

# Planetary-Scale Byzantine Fault Tolerance

Présentée le 12 janvier 2024

Faculté informatique et communications  
Laboratoire de calcul distribué  
Programme doctoral en informatique et communications

pour l'obtention du grade de Docteur ès Sciences

par

**Matteo MONTI**

Acceptée sur proposition du jury

Prof. K. Aberer, président du jury  
Prof. R. Guerraoui, directeur de thèse  
Prof. L. Alvisi, rapporteur  
Prof. R. Van Renesse, rapporteur  
Prof. B. Ford, rapporteur



*A Irene.*



# Preface

This thesis contains selected results of research performed during doctoral studies between September 2017 and August 2023, supervised by Prof. Rachid Guerraoui, at the Distributed Computing Laboratory (School of Computer and Communication Sciences) at EPFL in Lausanne, Switzerland. The results presented in this thesis appear in the following works:

- Rachid Guerraoui, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, Dragos-Adrian Seredinschi. "Scalable Byzantine Reliable Broadcast". *DISC 2019*. (Best Paper Award.)
- Martina Camaioni, Rachid Guerraoui, Matteo Monti, Manuel Vidigueira. "Oracular Byzantine Reliable Broadcast". *DISC 2022*.
- Martina Camaioni, Rachid Guerraoui, Matteo Monti, Pierre-Louis Roman, Manuel Vidigueira, Gauthier Voron. "Chop Chop: Byzantine Atomic Broadcast to the Network Limit". *Under Review-and-Resubmit for OSDI 2024*.

The following additional publications contain work that is not presented in this thesis, but are also a result of research performed during the same doctoral studies:

- Rachid Guerraoui, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, Dragos-Adrian Seredinschi. "The Consensus Number of a Cryptocurrency". *PODC 2019*.
- Daniel Collins, Rachid Guerraoui, Jovan Komatovic, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, Yvonne-Anne Pignolet, Dragos-Adrian Seredinschi, Andrei Tonkikh, Athanasios Xytkis. "Online Payments by Merely Broadcasting Messages". *DSN 2020*. (Runner-up for Best Paper Award.)
- Andrei Kucharavy, Matteo Monti, Rachid Guerraoui, Ljiljana Dolamic. "Byzantine-Resilient Learning Beyond Gradients: Distributing Evolutionary Search." *GECCO 2023*.
- Pierre Civit, Seth Gilbert, Rachid Guerraoui, Jovan Komatovic, Matteo Monti and Manuel Vidigueira. "Every Bit Counts in Consensus." *DISC 2023*. (Best Paper Award.)
- Martina Camaioni, Rachid Guerraoui, Jovan Komatovic, Matteo Monti, Pierre-Louis Roman, Manuel Vidigueira, Gauthier Voron. "Carbon: Scaling Trusted Payments with Untrusted Machines." *Under review for TDSC 2023*.

The following additional publications contain results from prior work in physics, but published during the same doctoral studies:

## Preface

---

- Armando Bazzani, Silvia Vitali, Carlo Emilio Montanari, Matteo Monti, Gastone Castellani. "Stochastic Properties of Colliding Particles in a Non-Equilibrium Thermal Bath." *Nonlocal and Fractional Operators 2019*.
- Silvia Vitali, Carlo Emilio Montanari, Matteo Monti, Gastone Castellani. "Event-Based Simulation of 2D Particles Gas in a Gradient of Temperature." *ICTT 2019*.

# Abstract

The scale and pervasiveness of the Internet make it a pillar of planetary communication, industry and economy, as well as a fundamental medium for public discourse and democratic engagement. In stark contrast with the Internet’s decentralized infrastructure, however, most hyper-scale online services adopt a centralized architecture, wherein millions or billions of clients entrust the same service provider with sensitive data to store, process and transmit, resulting in poor security and a dangerous asymmetry of power between users and system maintainer. Byzantine distributed computing has long held the promise to advance the security and transparency of our information infrastructure, replicating services across several, mutually mistrustful processes so as to uphold security and transparency even if a fraction of the processes deviate arbitrarily from the algorithm they are assigned. Despite extensive research, however, real-world replicated systems still fall short of global adoption. This is at least in part due to the limited scalability of Byzantine algorithms: the cost of replication is still too high, preventing the rise of Byzantine hyper-scalers.

This thesis attempts to overcome this limitation, contributing planetary-scale implementations for two classes of distributed abstractions: Reliable Broadcast, which has correct processes agree on messages issued by individual sources, and Consensus, which provides agreement on values contributed by multiple independent sources. The two classes of abstraction differ in power and requirements: Consensus is more powerful, enabling universal distributed computation, but Reliable Broadcast can be implemented without any assumption on network timeliness, a better fit for the public Internet’s large latency spikes.

First, we scale the number of servers that can take part in Reliable Broadcast. We present Contagion, the first probabilistic Reliable Broadcast protocol to achieve logarithmic per-process computation and communication complexity, enabling scalability to a practically unlimited number of servers. At the core of Contagion are samples, a probabilistic alternative to Byzantine quorums foregoing intersection guarantees for statistical representativeness. Contagion’s security evaluation is among the main technical contributions of this thesis, providing the first formal analysis of a probabilistic protocol in the Byzantine setting.

Second, we scale the number of clients that can concurrently submit messages to a Reliable Broadcast system. We present Draft, the first implementation of Reliable Broadcast to asymptotically amortize (whenever a set of commonplace conditions are met) all signature and communication overhead resulting from Byzantine fault tolerance, matching the complexity of a trusted, centralized solution. Draft’s efficiency is enabled by brokers, a novel

## Abstract

---

layer of untrusted processes we introduce between clients and servers to enhance system performance.

Finally, we scale the number of clients a Consensus-class system can handle. We do so with Chop Chop, the first system to introduce brokers to the Atomic Broadcast abstraction. When geo-deployed on 64 medium-sized servers, Chop Chop processes 43,600,000 messages per second, two orders of magnitude more than state-of-the-art alternatives. Even at maximum load, Chop Chop's performance is close to line-rate, putting 92% of server bandwidth towards transmitting application messages, negating nearly all cost resulting from Byzantine fault tolerance.

**Keywords.** Distributed systems, Byzantine fault tolerance, scalability, planetary scale, replication, reliable broadcast, consensus, randomized algorithms, broker-enhanced computation, oracularity, efficiency.



# Résumé

L'ampleur et l'omniprésence d'Internet en font un pilier de la communication, de l'industrie et de l'économie à l'échelle planétaire, ainsi qu'un support fondamental pour le discours public et l'engagement démocratique. Cependant, contrairement à l'infrastructure décentralisée d'Internet, la plupart des services en ligne à grande échelle adoptent une architecture centralisée, dans laquelle des millions ou des milliards de clients confient au même fournisseur de service des données sensibles à stocker, à traiter et à transmettre, ce qui entraîne une sécurité médiocre et une dangereuse asymétrie de pouvoir entre les utilisateurs et le responsable du système. Le calcul distribué Byzantin promet depuis longtemps de faire progresser la sécurité et la transparence de notre infrastructure d'information, en répliquant les services à travers plusieurs ordinateurs qui se méfient les uns des autres, de manière à maintenir la sécurité et la transparence même si une fraction des ordinateurs dévient arbitrairement de l'algorithme qui leur est assigné. Malgré des recherches approfondies, les systèmes répliqués n'ont pas encore été adoptés en pratique à l'échelle mondiale. Cela est dû, au moins en partie, à la capacité limitée pour des algorithmes Byzantins de passer à l'échelle : le coût de la réplication est encore trop élevé, ce qui empêche l'apparition d'*hyper-scalers* Byzantins.

Cette thèse tente de surmonter cette limitation par des implémentations, utilisables à l'échelle planétaire, pour deux classes d'abstractions distribuées : le *Reliable Broadcast*, qui permet aux processus corrects de se mettre d'accord sur des messages émis par des sources individuelles, et le *Consensus*, qui permet de se mettre d'accord sur des valeurs fournies par de multiples sources indépendantes. Les deux catégories d'abstractions diffèrent en termes de puissance et d'exigences : le *Consensus* est plus puissant et permet un calcul distribué universel, mais le *Reliable Broadcast* peut être mis en œuvre sans aucune hypothèse sur la latence du réseau, ce qui est mieux adapté aux importants pics de latence de l'Internet public.

Tout d'abord, nous augmentons le nombre de serveurs pouvant participer aux *Reliable Broadcast*. Nous présentons Contagion, le premier protocole probabiliste de *Reliable Broadcast* dont la complexité de calcul et de communication par processus est logarithmique, ce qui permet de l'étendre à un nombre pratiquement illimité de serveurs. Au cœur de Contagion se trouvent les échantillons, une alternative probabiliste aux quorums Byzantins qui renoncent aux garanties d'intersection pour la représentativité statistique. L'évaluation de la sécurité de Contagion est l'une des principales contributions techniques de cette thèse, car elle fournit la première analyse formelle d'un protocole probabiliste dans le cadre Byzantin.

Ensuite, nous augmentons le nombre de clients qui peuvent simultanément soumettre des

## Résumé

---

messages à un système de *Reliable Broadcast*. Nous présentons Draft, la première implémentation de *Reliable Broadcast* à amortir asymptotiquement, chaque fois qu'une série de conditions communes sont satisfaites, l'intergralité du surcoût des signatures et des communications résultant de la tolérance aux pannes Byzantines, et ainsi atteindre la complexité d'une solution centralisée et fiable. L'efficacité de Draft est rendue possible par les *brokers*, une nouvelle classe de processus ne nécessitant pas de garantie de fiabilité que nous introduisons entre les clients et les serveurs pour améliorer les performances du système.

Enfin, nous augmentons le nombre de clients qu'un système de type *Consensus* peut gérer. Nous le faisons avec Chop Chop, le premier système à introduire des *brokers* dans l'abstraction *Atomic Broadcast*. Lorsqu'il est géodéployé sur 64 serveurs de taille moyenne, Chop Chop traite 43 600 000 messages par seconde, soit deux ordres de grandeur de plus que l'état de l'art. Même à la charge maximale, les performances de Chop Chop sont proches du *line-rate*, consacrant 92% de la bande passante du serveur à la transmission des messages d'application, ce qui annule la quasi-totalité des coûts résultant de la tolérance aux pannes Byzantines.

# Contents

<b>Preface</b>	<b>i</b>
<b>Abstract (English / Français)</b>	<b>iii</b>
<b>Introduction</b>	<b>1</b>
<b>I Scalable Byzantine Reliable Broadcast</b>	<b>9</b>
<b>1 Overview</b>	<b>11</b>
1.1 Introduction . . . . .	11
1.1.1 Probabilistic Byzantine Reliable Broadcast . . . . .	12
1.1.2 Probability Analysis . . . . .	13
1.1.3 Security and Complexity Evaluation . . . . .	14
1.2 Related Work . . . . .	16
1.3 Model . . . . .	17
<b>2 Murmur</b>	<b>19</b>
2.1 Interface . . . . .	19
2.2 Algorithm . . . . .	20
2.3 No duplication, integrity and validity . . . . .	20
2.4 Totality . . . . .	22
<b>3 Sieve</b>	<b>27</b>
3.1 Interface . . . . .	27
3.2 Algorithm . . . . .	28
3.3 No duplication and integrity . . . . .	30
3.4 Total validity . . . . .	31
3.5 Preliminary lemmas . . . . .	32
3.6 Simplified Sieve . . . . .	38
3.6.1 Consistency-only broadcast . . . . .	38
3.6.2 Byzantine oracle . . . . .	39
3.6.3 Algorithm . . . . .	39
3.7 Adversarial execution . . . . .	44
3.7.1 Model (Sieve) . . . . .	44

## Contents

---

3.7.2	Model (Simplified Sieve)	45
3.7.3	Network scheduling	45
3.7.4	Interfaces	46
3.8	Simplified adversarial power	48
3.8.1	Preliminary definitions	48
3.8.2	Consistency of Simplified Sieve	50
3.9	Two-phase adversaries	58
3.10	Consistency	66
3.10.1	Two-phase adversaries	67
3.10.2	Random variables	68
3.10.3	Byzantine population, correct echoes, delivery	70
3.10.4	Second phase	72
3.10.5	First phase	78
3.11	Decorators	82
3.11.1	Auto-echo adversary	82
3.11.2	Process-sequential adversary	88
3.11.3	Sequential adversary	97
3.11.4	Non-redundant adversary	108
3.11.5	Sample-blind adversary	114
3.11.6	Byzantine-counting adversary	127
3.11.7	Single-response adversary	133
3.11.8	Two-phase adversary	139
<b>4</b>	<b>Contagion</b>	<b>147</b>
4.1	Interface	147
4.2	Algorithm	148
4.3	No duplication and integrity	150
4.4	Validity	151
4.5	Adversarial execution	152
4.5.1	Model	152
4.6	Epidemic processes	153
4.7	Threshold contagion (overview)	157
4.8	Preliminary lemmas	159
4.9	Consistency	161
4.10	Totality	163
4.10.1	Minimal operations	163
4.10.2	Delivery probability	164
4.10.3	C-step Threshold Contagion	165
4.11	Threshold Contagion	167
4.11.1	Epidemic processes	168
4.11.2	Threshold Contagion	169
4.11.3	Rules	169

4.11.4 Random variables . . . . .	172
4.11.5 Goal . . . . .	172
4.11.6 Sample space . . . . .	173
4.11.7 Random variables as sample functions . . . . .	178
4.11.8 Contagion step . . . . .	182
4.11.9 Final infection size . . . . .	190
<b>II Oracular Byzantine Reliable Broadcast</b>	<b>193</b>
<b>5 Overview</b>	<b>195</b>
5.1 Introduction . . . . .	195
5.2 Related Work . . . . .	201
5.3 Model & background . . . . .	203
5.3.1 Model . . . . .	203
5.3.2 Background . . . . .	204
<b>6 Draft</b>	<b>205</b>
6.1 Interface . . . . .	205
6.2 Algorithm . . . . .	206
6.2.1 Protocol & correctness overview . . . . .	206
6.2.2 Complexity overview . . . . .	212
6.2.3 Pseudocode (Client) . . . . .	214
6.2.4 Pseudocode (Broker) . . . . .	216
6.2.5 Pseudocode (Server) . . . . .	221
6.3 Correctness . . . . .	227
6.3.1 No duplication . . . . .	227
6.3.2 Consistency . . . . .	227
6.3.3 Totality . . . . .	232
6.3.4 Integrity . . . . .	234
6.3.5 Validity . . . . .	236
6.4 Complexity . . . . .	252
6.4.1 Auxiliary results . . . . .	252
6.4.2 Batching limit . . . . .	260
6.4.3 Protocol analysis . . . . .	261
<b>7 Dibs</b>	<b>267</b>
7.1 Interface . . . . .	267
7.2 Algorithm . . . . .	269
7.2.1 Pseudocode (Client) . . . . .	269
7.2.2 Pseudocode (Server) . . . . .	271
7.3 Correctness . . . . .	272
7.3.1 Correctness . . . . .	272
7.3.2 Signup integrity . . . . .	275

## Contents

---

7.3.3	Signup validity . . . . .	275
7.3.4	Self-knowledge . . . . .	277
7.3.5	Transferability . . . . .	278
7.3.6	Density . . . . .	279
<b>III</b>	<b>Byzantine Atomic Broadcast to the Network Limit</b>	<b>281</b>
<b>8</b>	<b>Chop Chop</b>	<b>283</b>
8.1	Overview . . . . .	283
8.2	Atomic Broadcast . . . . .	286
8.2.1	Cost of Atomic Broadcast . . . . .	286
8.2.2	Existing Mitigations . . . . .	287
8.3	Distilled Batches . . . . .	288
8.3.1	Distillation at a Glance . . . . .	289
8.3.2	Distillation Microbenchmark . . . . .	290
8.4	Chop Chop . . . . .	291
8.4.1	Architecture and Model . . . . .	292
8.4.2	Distillation Phase . . . . .	293
8.4.3	Submission Phase . . . . .	296
8.4.4	Correctness . . . . .	297
8.5	Implementation Details . . . . .	299
8.5.1	Broker . . . . .	299
8.5.2	Server . . . . .	300
8.6	Evaluation . . . . .	301
8.6.1	Baselines . . . . .	301
8.6.2	Setup . . . . .	302
8.6.3	RQ1 – Load Handling . . . . .	305
8.6.4	RQ2 – Distillation Benefits . . . . .	306
8.6.5	RQ3 – Number of Servers . . . . .	308
8.6.6	RQ4 – Overall Efficiency . . . . .	309
8.6.7	RQ5 – Chop Chop Under Failures . . . . .	309
8.6.8	RQ6 – Application Use Cases . . . . .	310
8.7	Related Work . . . . .	310
	<b>Conclusions</b>	<b>313</b>
	<b>Bibliography</b>	<b>315</b>

# Introduction

## A word of motivation

The Internet was designed with decentralization in mind, as a way to maintain reliable communications in the face of catastrophic infrastructural damage. Half a century after ARPANET's first deployment, the scale and pervasiveness of the Internet make it a pillar of planetary communication, industry and economy [int21], as well as a fundamental medium for public discourse and democratic engagement [ESWV22]. In stark contrast with the infrastructure on which they rely, however, most Internet-scale services adopt a centralized architecture, wherein millions or billions of clients entrust the same service provider with sensitive data to store, process and transmit [Fit22, Run20].

Centralized computer systems are fast, scalable and easy to deploy, but come at a severe cost in terms of security and transparency. On the one hand, single points of failure are notoriously bad for security, as even the best-maintained systems are prone to being hacked or compromised by random accidents. In 2023, cyber-attacks are projected to account for as much as 8 trillion US dollars in global economic damage [Far22] - over 8% of the world's economy, lost to improperly secured computer systems. In the last year alone, real-world attacks have been successfully dealt [cyb23] on energy systems, financial infrastructure, healthcare providers, governments and military organizations alike, both in war and peacetime. However well maintained, a centralized system is prone to being hacked into, physically tampered with, or compromised by random accidents.

On the other hand, centralized solutions force users to trust that the system's maintainer (often a self-interested, profit-driven organization) will behave in a way aligned with the users' best interests, e.g., by preserving data confidentiality or upholding accurate bookkeeping even when failing to do so would be advantageous and difficult to detect. Over the last decade, this power asymmetry resulted in a long sequence of abuses, with Internet-scale service providers violating user privacy [Reu20, Cer19, Pri19], mismanaging hate speech [Moz18], enabling democratic tampering [Con18] and harming the health of their most vulnerable users [Gay21]. Lacking ad-hoc transparency measures, centralized systems are opaque. Once cleartext data is handed off, it is not easy to determine how that data will be handled.

## Introduction

---

Byzantine distributed computing, on which this thesis focuses, has long<sup>1</sup> held the promise to advance the security and transparency of humanity’s information infrastructure. By replicating a service across several, reciprocally mistrusting *processes*, a Byzantine-resilient system upholds well-defined liveness and safety properties even if a fraction of the processes fail maliciously, deviating arbitrarily from the algorithm they are assigned. In theory, any service can be replicated, faithfully and transparently executed on a common system whose correctness and availability would be very hard to thwart. As such, Byzantine algorithms hold the promise of profound societal impact, securing critical infrastructure and empowering citizens with services whose operation is auditable and democratically established, beyond the influence of any individual party.

Despite extensive research and interest from industry and governments, however, real-world Byzantine systems still fall short of global adoption: from search engines to social networks, from messaging platforms to software distribution infrastructures, most planetary-scale services still adopt, at the time of writing, a centralized architecture. This lack of adoption, we argue, is at least in part<sup>2</sup> due to the limited scalability of Byzantine algorithms: the cost of replication is still too high, preventing the rise of Byzantine hyper-scalers.

The Internet’s foundational vision was that of a highly-available, decentralized, secure, universal computer shared by all. This thesis hopes to make a small step towards that vision, introducing practical distributed algorithms that billions of stakeholders can partake in as servers, or use as clients.

### Two levels of power, two ways to scale

This thesis focuses on two fundamental classes of distributed abstraction: Reliable Broadcast and Consensus. Abstractions within the same class can be implemented from each other (for example, Atomic Broadcast can be built atop Consensus, and vice versa). While Consensus can be used to implement Reliable Broadcast, however, the converse is provably impossible [68]. As such, abstractions in the Consensus class are strictly more powerful than those in the Reliable Broadcast class. Consensus-powered algorithms, however, are generally more complex, and require more stringent assumptions about the system on which they run.

**Reliable Broadcast.** Abstractions in the Reliable Broadcast class have correct processes agree on messages issued by individual sources. Reliable Broadcast [31], the class’ namesake, has one designated sender broadcast a single message. If the sender is correct, every correct process delivers the message it broadcasts. Even if the sender is Byzantine, all correct processes

---

<sup>1</sup>The seminal “*Byzantine Generals Problem*” [104] was published in 1982, over forty years prior to the writing of this thesis.

<sup>2</sup>Other factors might include the high latency inherent to WAN-based coordination, the lack of well-established, practical solutions for deploying real-world distributed systems, and the incomplete legislative frameworks regulating decentralised networks. This thesis aims to contribute technical results to the public discourse, but the road ahead is likely long and interdisciplinary.



deliver either the same message, or no message at all. Multi-shot variants [34] of Reliable Broadcast include Source-Order Broadcast, wherein every process maintains an independent, agreed-upon, append-only log, and Causal Broadcast, wherein no message is delivered before any of its causal antecedents. Reliable Broadcast has well-known applications [37, 61] as a building block in practical fault-tolerant systems. Most notably, the Asset-Transfer problem, which cryptocurrency-based systems traditionally solve using Consensus, was recently reduced [80] to Causal Broadcast, showing Reliable Broadcast’s relevance as a stand-alone solution. Abstractions in the Reliable Broadcast class are generally considered simpler to implement [50] than their Consensus-based counterparts. Maybe more importantly, Reliable Broadcast can be solved in the asynchronous setting [31]: all abstractions in the Reliable Broadcast class uphold liveness and safety even when no assumption can be made on how rapidly processes can exchange messages. This makes Reliable Broadcast solutions fast and secure in environments with large latency spikes, such as the public Internet. Simplicity and asynchrony, however, come at the cost of generality, as only a relatively narrow class of problems<sup>3</sup> can be solved by abstractions in the Reliable Broadcast class - Consensus is still required in the general case. We study abstractions in the Reliable Broadcast class in Parts I and II of this thesis.

**Consensus.** Abstractions in the Consensus class have correct processes agree on values contributed by multiple independent sources. Consensus [104], the class’ namesake, has each correct process first propose, then decide a value. All correct processes eventually decide the same value; if all correct processes propose the same value<sup>4</sup>, that value must be decided. The canonical multi-shot variant of Consensus is Atomic Broadcast, wherein all sources contribute messages to the same, totally-ordered, agreed-upon, append-only log. The most relevant application of Atomic Broadcast is arguably State Machine Replication, the fundamental distributed abstraction enabling universal Byzantine-resilient computation. In brief, State Machine Replication uses Atomic Broadcast to order requests, usually contributed by an external set of clients. Every correct process initializes a local copy of the same deterministic state machine to the same initial value, then serves the same requests in the same order. Because all correct processes deterministically transition through the same sequence of states, any response backed by sufficiently many processes is guaranteed to be correct - this technique can be employed to make any deterministic service safe and available despite a (suitably small) fraction of processes misbehaving arbitrarily. Abstractions in the Consensus class usually call for more complex implementations than their Reliable Broadcast counterparts. Additionally, Consensus is provably impossible in asynchrony [68]<sup>5</sup>, forcing Consensus algorithms to make assumptions about the timeliness of inter-process communication. We study abstractions in

---

<sup>3</sup>Reliable Broadcast can only solve problems whose *consensus number* is 1 Herlihy’s hierarchy [86].

<sup>4</sup>Several other flavors of *validity* are defined in literature [48], relating decisions to proposals. Weaker forms of validity, for example, only apply if all processes are correct.

<sup>5</sup>The impossibility of asynchronous Consensus applies only to the deterministic variant of Consensus. Probabilistic variants of Consensus, which deterministically uphold safety but only terminate with probability 1, do allow for asynchronous solutions.

## Introduction

---

the Consensus class in Part III.

The goal of this thesis is to contribute planetary-scale implementations for both classes of abstraction. For each class, we attempt in particular to scale the number of processes that can partake in the abstraction as *servers*, as well as the number of processes that can concurrently use the abstraction as *clients*.

**Scaling servers.** Designing a distributed algorithm that can scale to billions of servers presents some fundamental challenges. At the core of even the simplest distributed algorithms [31] is the concept of *quorums* [149], sets of servers that are large enough to intersect in at least one correct server. Classically, deterministic algorithms [22, 31, 43, 156] have at least one process exchange messages with a quorum of servers, resulting in an unworkable communication complexity: at planetary scale, no individual process can feasibly communicate with (a non-negligible fraction of) the whole system. We tackle this fundamental problem by shifting to the probabilistic setting, replacing quorums with **samples**. Unlike quorums, samples do not provide any intersection guarantees - a sample is just large enough to be (with high probability) statistically representative of the system as a whole. Studying the security of a probabilistic system while accounting for Byzantine behavior is one of the main challenges of this thesis. We focus on server scalability in Part I of this thesis.

**Scaling clients.** In order to reliably service billions of clients, a system will realistically need to sustain millions of requests per second. Designing a distributed algorithm for such a high throughput requires minimizing the amount of computation and communication required server-side to process each request. Techniques such as *batching* [43, 55, 70, 140] are available in literature to partially amortize the cost of coordination among servers. Even when batching is employed, however, the communication and computation cost incurred by a server when processing a high rate of requests is often dominated by operations (such as message authentication) that, we notice, are expensive to *perform*, but cheap to *verify*. To leverage this observation, we extend the classical client-server model with **brokers**, a permissionless, scalable set of processes whose only purpose is to alleviate server complexity. Unlike servers, a large fraction of which we must assume to be correct, all brokers but one can be faulty. Brokers act as an intermediary between clients and servers, taking upon themselves to perform expensive, verifiable operations. We focus on client scalability in Parts II and III of this thesis.

## Contributions of this thesis

This thesis (Figure 1) presents solutions to three of the four scalability challenges we outlined above: we scale Reliable Broadcast servers in Part I, Reliable Broadcast clients in Part II, and Consensus clients in Part III.

	Reliable Broadcast		Consensus	
	<i>Theory</i>	<i>Systems</i>	<i>Theory</i>	<i>Systems</i>
<b>Server Scalability</b>	Part I, [81]	<i>(unpublished)</i>		
<b>Client Scalability</b>	Part II, [41]	[39] <i>(submitted)</i>	<i>(work in progress)</i>	Part III, [40] <i>(minor revisions)</i>

**Figure 1: Contributions of this thesis.**

**Part I: Scalable Byzantine Reliable Broadcast.** We generalize single-shot Reliable Broadcast to the probabilistic setting, allowing each of its properties to be violated with a fixed, arbitrarily small probability. We leverage these relaxed guarantees in a protocol replacing quorums with stochastic samples. Unlike quorums, samples do not provide any intersection guarantee: a sample is just large enough to be statistically representative of the overall system (at least with high probability). We implement Reliable Broadcast with three algorithms, building on top of each other to progressively obtain stronger abstractions: Murmur for Probabilistic Broadcast, Sieve for Probabilistic Consistent Broadcast, and Contagion for Probabilistic Reliable Broadcast. To the best of our knowledge, Contagion is the first Reliable Broadcast protocol to achieve logarithmic per-process computation and communication complexity, enabling scalability to a practically unlimited number of servers. A major technical contribution of Part I is a fully formal analysis of our protocol, deriving closed-form bounds on the probability of each Reliable Broadcast property being compromised. The complexity of such an analysis is exacerbated in the Byzantine setting, as each bound must hold for every possible adversarial strategy. We tackle this problem using decorators, a novel suite of formal techniques designed to progressively narrow down the set of strategies that provably includes the optimal adversary.

The algorithms and techniques discussed in Part I were presented [81] at the International Symposium on Distributed Computing (DISC 2019). The paper won the Best Paper Award for that edition of the conference. In addition to the theoretical results presented in this thesis, a C++ implementation of Contagion was developed and tested on a global deployment of Amazon AWS machines. Preliminary results showed a throughput of upwards of a thousand messages per second even when Contagion was deployed on 2048 Micro instances (larger-scale testing was unfeasible due to budget constraints). Regrettably, the experimental results on Contagion are, at the time of writing, unpublished.

**Part II: Oracular Byzantine Reliable Broadcast.** We study Client-Server Byzantine Reliable Broadcast, a multi-shot variant of Reliable Broadcast whose interface is split between broadcasting clients and delivering servers. We present Draft, an optimally resilient implementation of Client-Server Byzantine Reliable Broadcast. Like most algorithms in the Reliable Broadcast class, Draft guarantees both liveness and safety in the asynchronous setting. Under good conditions, however, Draft stretches towards information-theoretical efficiency. In a moment

## Introduction

---

of synchrony, free from Byzantine misbehavior, and at the limit of infinitely many broadcasting clients, a Draft server delivers a  $b$ -bits payload at an asymptotic amortized cost of 0 signature verifications, and  $\log_2(c) + b$  bits exchanged, where  $c$  is the number of clients in the system. This is the minimum number of bits required to convey the payload ( $b$  bits, assuming it is incompressible) along with an identifier for its sender ( $\log_2(c)$  bits, necessary to enumerate any set of  $c$  elements, and optimal if broadcasting frequencies are uniform or unknown). Noting that this is the same complexity a server would incur by receiving plain, unauthenticated messages from a trusted oracle, we say that Draft can achieve *oracularity*, effectively achieving zero-cost Byzantine fault tolerance. We achieve oracularity by noting that a synchronous, interactive protocol can, in the absence of Byzantine behavior, reduce the metadata of a batch of messages down to a constant. Remarkably, this reduction protocol is expensive to perform, but publicly verifiable: ill-reduced batches are visibly malformed, and can be discarded without additional coordination. This observation allows us to offload reduction to brokers, a novel set of processes we add to the distributed computing model between clients and servers. Brokers play no safety role in the system (even if they all fail, only<sup>6</sup> liveness is lost), and can be trustlessly spun up to meet client demand, interacting with clients to produce reduced batches for the servers to oracularly process.

The results discussed in Part II were presented [41] at the International Symposium on Distributed Computing (DISC 2022). In addition to the theoretical results presented in this thesis, Draft was generalized into Carbon, a reconfigurable, garbage-collected, asynchronous asset-transfer system. When tested on a global deployment of 64 AWS machines, a Rust implementation of Carbon sustained upwards of one million messages per second. Our experimental results are currently under review for publication [39] in the IEEE Transactions on Dependable and Secure Computing (TDSC).

**Part III: Byzantine Atomic Broadcast to the Network Limit.** We present Chop Chop, a Byzantine Atomic Broadcast system that amortizes the cost of ordering, authenticating and deduplicating messages, achieving *line rate* (i.e., closely matching the complexity of a protocol that does not ensure any ordering, authentication or Byzantine resilience) even when processing messages as small as 8 bytes. Chop Chop attains this performance by means of *distillation*, a generalization of the reduction techniques we introduced in Part II for the Reliable Broadcast class. A distilled batch is a set of messages that are fast to authenticate and deduplicate, as well as order. Batches are distilled using a novel, interactive protocol based on brokers. In a geo-distributed deployment of 64 medium-sized servers, with clients situated cross-cloud, a Rust implementation of Chop Chop processes 43,600,000 messages per second with an average latency of 3.6 seconds. Under the same conditions, state-of-the-art alternatives offer two orders of magnitude less throughput for the same latency. We test three simple applications,

---

<sup>6</sup>The liveness of a critical system might arguably be as important as its safety. The trustless nature of brokers, however, helps protect liveness too. Should well-established brokers experience catastrophic failure, concerned stakeholders could, e.g., spin up their own brokers, bypassing lengthy security audits to quickly bring the system back online.

running on a replicated state machine powered by Chop Chop: a Payment system, an Auction house and an instance of the “Pixel War” game, respectively achieving 32, 2.3 and 35 million operations per second.

The results presented in Part III are under minor revisions for presentation [40] at the USENIX Symposium on Operating Systems Design and Implementation (OSDI 2024). In addition to the experimental results presented in this thesis, a paper is under development to generalize oracularity to the Atomic Broadcast abstraction, further establish brokers as a well-defined role in the distributed computing model, and prove Chop Chop’s correctness and oracularity to the fullest extent of formal detail.

**State of the art.** To the best of our knowledge, all three results presented in this thesis are state of the art at the time of writing. An implementation of Reliable Broadcast with sub-logarithmic complexity would likely require going past stochastic sampling: even assuming that only the sender is Byzantine, having every correct process query a sub-logarithmic sample would result in a non-negligible probability of non-representativeness for at least one sample, effectively poisoning a process’ view of the overall system. Concerning our oracular and line-rate implementations, improving efficiency would necessarily involve some compression scheme: assuming ids and messages are incompressible, both solutions asymptotically reach, to the bit, the information-theoretical lower bound.

**Towards Scalable Byzantine Atomic Broadcast.** A great deal of effort was put towards the development of a Byzantine Atomic Broadcast algorithm that (similarly to Contagion) would achieve logarithmic per-process computation and communication complexity. Doing so proved unfeasible within the time allowed by PhD research. While an algorithm, which we conjecture to be correct, was developed, a fully formal proof of its safety and liveness proved intractable within such a limited time frame. As we discuss in Part I, much of the complexity of Contagion’s analysis stems from the non-trivial interaction between the probabilistic behavior of correct processes and the arbitrary behavior of the Byzantine adversary. Adding to this the subtlety of time analysis (our Atomic Broadcast algorithm works in the partially synchronous model) further complicates our proofs. On the way to proving the algorithm’s correctness, a fully formal theory of decorators was developed, framing Byzantine behavior in the general scope of probabilistic games. Due to its length (comparable to the entire extent of this thesis) and dubious standing as stand-alone contribution, our theoretical work on decorators was omitted from this thesis, and is currently unpublished.



# Scalable Byzantine Reliable Broadcast **Part I**





# 1 Overview

## 1.1 Introduction

Broadcast is a popular abstraction in the distributed systems toolbox, allowing a process to transmit messages to a set of processes. The literature defines many flavors of broadcast, with different safety and liveness guarantees [34, 72, 82, 110, 125]. In this Part we focus on Byzantine reliable broadcast, as introduced by Bracha [31]. This abstraction is a central building block in practical Byzantine fault-tolerant (BFT) systems [37, 61, 80]. We tackle the problem of its scalability, namely reducing the complexity of Byzantine reliable broadcast, and seeking good performance despite a large number of participating processes.

In Byzantine reliable broadcast, a designated sender broadcasts a single message. Intuitively, the broadcast abstraction ensures that no two correct processes deliver different messages (*consistency*), either all correct processes deliver a message or none does (*totality*), and that, if the sender is correct, all correct processes eventually deliver the broadcast message (*validity*). This must hold despite a certain fraction of Byzantine processes, potentially including the sender. We denote by  $N$  the number of processes in the system, and  $f$  the fraction of processes that are Byzantine. Existing algorithms for Byzantine reliable broadcast scale poorly as they typically have  $O(N)$  per-process communication complexity [32, 110, 111, 146]. The root cause for the poor scalability of these algorithms is their use of quorums [112, 149], i.e., sets of processes that are large enough to always intersect in at least one correct process. The size of a quorum grows linearly with the size of the system [34].

To overcome the scalability limitation of quorum-based broadcast, Malkhi *et al.* [113] generalized quorums to the probabilistic setting. In this setting, two random quorums intersect with a fixed, arbitrarily high probability, allowing the size of each quorum to be reduced to  $O(\sqrt{N})$ . We are not aware of any Byzantine reliable broadcast algorithm building on probabilistic quorums; nevertheless, such an algorithm could have a per-process communication complexity reduced from  $O(N)$  to  $O(\sqrt{N})$ . The active<sub>t</sub> protocol of Malkhi *et al.* [110] uses a form of samples for an optimistic path, but relies on synchrony and has a linear worst-case complexity (that is arguably very likely to occur with only moderate amounts of faulty processes).

## Chapter 1. Overview

---

**Samples.** In this Part, we present a probabilistic gossip-based Byzantine reliable broadcast algorithm having  $O(\log N)$  per-process communication and computation complexity, at the expense of  $O(\log N / \log \log N)$  latency. Essentially, we propose *samples* as a replacement for quorums. Like a probabilistic quorum, a sample is a randomly selected set of processes. Unlike quorums, samples do not need to intersect. Any two quorums must have a large enough size to ensure their intersection in a correct node (at least with high probability). Samples, instead, can be significantly smaller than quorums, as each sample must be large enough only to be *representative* of the system with high probability.

As with quorums, a process can use its sample to gather information about the global state of the system. An old Italian saying provides an intuitive understanding of this shift of paradigm: *“To know if the sea is salty, one needs not drink all of it!”* Intuitively, we leverage the law of large numbers, trading performance for a fixed, arbitrarily small probability of non-representativeness. To get an intuition of the difference between quorums and samples, consider the emulation of a shared memory in message passing [9]. One writes in a quorum and reads from a quorum to fetch the last value written. Our algorithms are rather in the vein of “write all, read any”. Here we would “write” using a gossip primitive and “sample” the system to seek the last value.

Throughout this Part, we extensively use samples to estimate the number of processes satisfying a set of *yes-or-no* predicates, e.g., the number of processes that are ready to deliver a message  $m$ . Consider the case where a correct process  $\pi$  queries  $K$  randomly selected processes (a sample) for a predicate  $P$ . Assume a fraction  $p$  of correct processes from the whole system satisfy predicate  $P$ . Let  $x$  be the fraction of positive responses (out of  $K$ ) that  $\pi$  collects. By the Chernoff bound, the probability of  $|x - p| \geq f + \epsilon$  is smaller or equal to  $\exp(-\lambda(\epsilon)K)$ , where  $\lambda$  quickly increases with  $\epsilon$ . For sufficient  $K$ , the probability of  $x$  differing from  $p$  by more than  $f + \epsilon$  can be made exponentially small.

To obtain a sample, our algorithms use a *sampling oracle* that returns the identity of a process from the system picked with uniform probability. In a permissioned system (i.e., one where the set of participating processes is known) sampling reduces to picking with uniform probability an element from the set of processes. In a permissionless system subject to Byzantine failures and slow churn, a (nearly) uniform sampling mechanism is still achievable using gossip [28].

### 1.1.1 Probabilistic Byzantine Reliable Broadcast

Our probabilistic algorithm, Contagion, allows each property of Byzantine reliable broadcast to be violated with an arbitrarily small probability  $\epsilon$ . We show that  $\epsilon$  scales sub-quadratically with  $N$ , and decays exponentially in the size of the samples. As a result, for a fixed value of  $\epsilon$ , the per-node communication complexity of Contagion is logarithmic.

We build Contagion incrementally, relying on two sub-protocols, as we describe next.

First, **Murmur** is a probabilistic broadcast algorithm that uses simple message dissemination

to establish *validity* and *totality*. In this algorithm, each correct process relays the sender's message to a randomly picked *gossip sample* of other processes. For the sample size  $\Omega(\log N)$ , the resulting gossip network is a connected graph with  $O(\log N / \log \log N)$  diameter, with high probability [47, 65]. In case of a Byzantine sender, however, Murmur does not guarantee consistency.

Second, **Sieve** is a probabilistic consistent broadcast algorithm that guarantees *consistency*, i.e., no two correct processes deliver different messages. To do so, each correct process uses a randomly selected *echo sample*. Intuitively, if enough processes from any echo sample confirm a message  $m$ , then with high probability no correct processes in the system delivers a different message  $m'$ . Sieve, however, does not ensure totality. If a Byzantine sender broadcasts multiple conflicting messages, a correct process might be unable to gather sufficient confirmations for either of them from its echo sample, and consequently would not deliver any message, even if some correct process delivers a message.

Finally, **Contagion** is a probabilistic reliable broadcast algorithm that guarantees validity, consistency, and totality. The sender uses Sieve to disseminate a consistent message to a subset of the correct processes. In order to achieve totality, Contagion mimics the spreading of a contagious disease in a population. A process samples the system and if it observes enough other "infected" processes in its sample, it becomes infected itself. If a critical fraction of processes is initially infected by having received a message from the underlying Sieve layer, the message spreads to all correct processes with high probability. If a process observes enough other infected processes, it delivers. As in the original deterministic implementation by Bracha [31], the crucial point here is that "enough" for becoming infected is less than "enough" for delivering. This way, with high probability, either all correct processes deliver a message or none does—Contagion satisfies totality. The other two important properties (validity and consistency) are inherited from the underlying (Murmur and Sieve) layers.

### 1.1.2 Probability Analysis

A major technical contribution of this Part is a complete, formal analysis of the properties of our three algorithms. To the best of our knowledge, this is the first analysis of a probabilistic broadcast algorithm in the Byzantine fault model, and this turned out to be very challenging. Intuitively, providing a bound on the probability of a property being violated reduces to studying a joint distribution between the inherent randomness of the system and the behavior of the Byzantine adversary. Since the behavior of the adversary is arbitrary, the marginal distribution of the Byzantine's behavior is unknown.

We develop two novel strategies to bound the probability of a property being violated, which we use in the analysis of Sieve and Contagion respectively.

**Identifying the optimal adversarial strategy.** When evaluating the consistency of Sieve, we show that a bound holds for every *possibly optimal* adversarial strategy. Essentially, we identify a subset of adversarial strategies that we prove to include the optimal one, i.e., the one that has the highest probability of compromising Sieve’s consistency. We then prove that every possibly optimal adversarial strategy has a probability of compromising the consistency of Sieve smaller than some  $\epsilon$ .

Identifying a set of possibly optimal adversarial strategies that is narrow enough to be amenable to analytic bounds is the most technically involved challenge of this Part. We do so by means of a novel technique that involves modeling the adversary as an algorithm interacting with the system through a well-defined interface. We start from the set of all possible adversarial algorithms (trivially, this set includes an optimal adversary). We then reduce the set of possibly optimal adversaries in steps. At every step, we transform each adversary by means of a *decorator*. Intuitively, a decorator tweaks an adversary in a way that simplifies its behavior without reducing its effectiveness in compromising the system. With each decorator, the set of possibly optimal adversaries becomes narrower and more tractable, finally yielding adversaries whose behavior is so simple that the probability of consistency being compromised can be computed analytically.

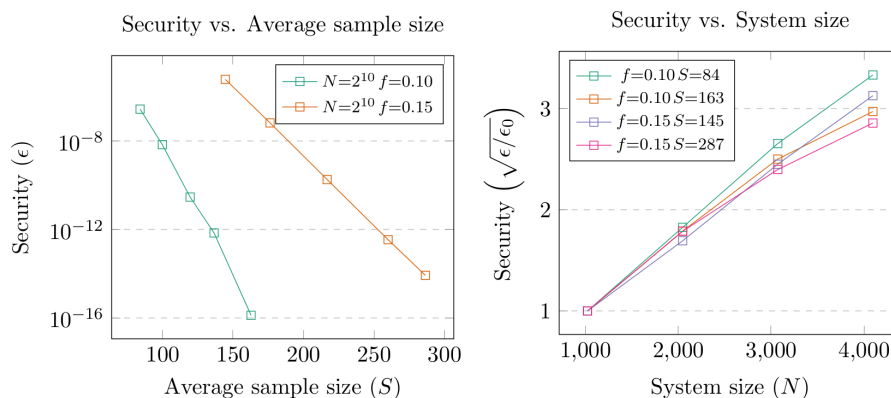
**Making adversarial strategy irrelevant.** When evaluating the totality of Contagion, we show that the adversarial strategy does not affect the outcome of the execution. Here, we show that any adversarial strategy reduces to a well-defined sequence of choices. We then prove that, due to the limited knowledge of the Byzantine adversary, every choice is equivalent to a random one.

### 1.1.3 Security and Complexity Evaluation

In Chapters 2 to 4 we present and study the security of three algorithms: Murmur, Sieve and Contagion, building upon each other to implement probabilistic Byzantine reliable broadcast. Each algorithm is configured by one or more security parameters (sample sizes and thresholds - we discuss each in detail in the relevant chapter). Sample sizes in particular are what determine the per-process communication complexity of each algorithm: in Murmur, a correct process gossips (on average) with  $G$  other processes to disseminate the sender’s message; in Sieve, a correct process samples  $E$  other processes to ensure consistency; in Contagion, a correct process decides whether or not to amplify / deliver a message based on two independent samples, respectively of size  $R$  and  $D$ .

The goal of Contagion’s security analysis is to derive closed-form expressions, relating sample sizes and thresholds to the probability  $\epsilon$  that Contagion will fail to uphold any of the properties of probabilistic Byzantine reliable broadcast.

We employ numerical techniques to minimize  $\epsilon$ , optimizing all security parameters under the



**Figure 1.1:** **Left** –  $\epsilon$ -security of Contagion, as a function of the average sample size  $S = \langle G, E, R, D \rangle$ . We use a system size of 1024 processes and fractions of tolerated Byzantine processes  $f = 0.1$  and  $f = 0.15$ . **Right** – Square root of the normalized  $\epsilon$ -security of Contagion, as a function of the system size  $N$ , for various fractions of Byzantine processes ( $f$ ) and average sample sizes ( $S$ ). We normalize the values in each series by the first element of that series. All lines appearing to grow sub-linearly with a square-rooted y-axis demonstrates that the normalized  $\epsilon$  security grows sub-quadratically.

constraint that the sum of all the sample sizes be constant ( $G + E + R + D = \text{const}$ ). Because a process communicates a constant number of times with each of the processes in its samples, this corresponds to a fixed communication complexity.

For a given system size  $N$  and fraction of Byzantine processes  $f$ , we relate this per-process communication complexity to the  $\epsilon$ -security of Contagion. As Figure 1.1 (left) shows, the probability  $\epsilon$  of compromising the security of Contagion decays exponentially in the average sample size  $S$ . We also study how the  $\epsilon$ -security of Contagion changes as a function of the system size  $N$ , for a fixed set of parameters ( $G, E, R, D$ ). Figure 1.1 (right) shows that the  $\epsilon$ -security is bounded by a quadratic function in  $N$ .

In summary, for a fixed security  $\epsilon$ , the average sample size (and consequently, the communication complexity of our algorithm) grows logarithmically with the system size  $N$ . For a practical choice of parameters, the probability of violating the properties of our algorithm can be brought down to  $10^{-16}$  for systems with thousands of processes.

The latency in terms of message delays between broadcasting and delivery of a message is  $O(\log N / \log \log N)$ . Specifically, the latency converges to  $O(\log N / \log \log N)$  message delays for gossip-based dissemination with Murmur (we prove this in Section 2.4, Theorem 5), and 2 message delays in total for Echo (Sieve) and Ready (Contagion) messages.

**Roadmap.** We discuss related work in Section 1.2. We state our model in Section 1.3. We discuss our probabilistic broadcast implementation Murmur in Chapter 2, our probabilistic consistent broadcast implementation Sieve in Chapter 3, and our probabilistic reliable

## Chapter 1. Overview

---

broadcast implementation Contagion in Chapter 4. Each chapter introduces the relevant abstraction, presents the relevant pseudocode, then formally studies the algorithm’s abidance to the abstraction’s properties. For deterministic properties, we provide a proof of correctness. For probabilistic properties, we analitically bound the probability that the property will be violated.

### 1.2 Related Work

At its base, our broadcast algorithm relies on gossip. There is a great body of literature studying various aspects of gossip, proposing flavors of gossip protocols for different environments and analyzing their complexities [6, 13, 18, 19, 20, 63, 67, 73, 74, 75, 76, 77, 84, 139, 148, 157]. However, to the best of our knowledge, we propose the first highly scalable gossip-based reliable broadcast protocol resilient to Byzantine faults with a thorough probabilistic analysis.

The communication pattern in the implementation of both our Sieve and Contagion algorithms can be traced back to the Asynchronous Byzantine Agreement (ABA) primitive of Bracha and Toueg [32] and the subsequent line of work [31, 37, 110, 132]. Indeed, our echo-based mechanism in Sieve resembles algorithms from classic quorum-based systems for Byzantine consistent broadcast [131, 146]. The ready-based mechanism in Contagion is inspired by a two-phase protocol appearing in several practical (quorum-based) systems [37, 61, 111]. Compared to classic work on this topic, the key feature of Contagion and Sieve is that they replace the building block of quorum systems with stochastic samples, thus enabling better scalability for the price of abandoning deterministic guarantees.

There is significant prior work on using epidemic algorithms to implement scalable *reliable* broadcast [23, 66, 89, 106]. Under benign failures or constant churn, these algorithms ensure, with high probability, that every broadcast message reaches all or none, and that all messages from correct senders are delivered. Our goal is to additionally provide *consistency* for broadcast messages, and tolerate *Byzantine* environments [32, 111, 146]. To the best of our knowledge, we are the first to apply the epidemic sample-based methodology in this context. Our main algorithm Contagion scales well to dynamic systems of thousands of nodes, some of which may be Byzantine. This makes it a suitable choice for *permissionless* settings that are gaining popularity with the advent of blockchains [122].

Distributed clustering techniques seek to group the processes of a system into clusters, sometimes called shards or quorums, of size  $O(\log N)$  [14, 79, 94, 95, 134]. This line of work has various goals (e.g., leader election, “almost everywhere” agreement, building an overlay network) and they also aim for scalable solutions. The overarching principle in clustering techniques is similar to our use of samples: build each cluster in a provably random manner so that the adversary cannot dominate any single cluster. Samples in our solution are private and individual on a per-process basis, in contrast to clusters which are typically public and global for the whole system.

The idea of *communication locality* appears in the context of secure multi-party computation (MPC) protocols [29, 46, 71]. This property captures the intuition that, in order to obtain scalable distributed protocols and permit a large number of participants, it is desirable to limit the number of participants each process must communicate with. All of our three algorithms have this communication locality property, since each process coordinates only with logarithmically-sized samples. In contrast to secure MPC protocols, our algorithms have different goals, system model, or assumptions (e.g., we do not assume a client-server model [71], nor do we seek to address privacy issues). Our algorithms can be used as building blocks towards helping tackle scalability in MPC protocols, and we consider this an interesting avenue for future work.

### 1.3 Model

In this section we discuss our system model and assumptions, and introduce some useful notation. We assume an asynchronous message-passing system where the set  $\Pi$  of  $N = |\Pi|$  processes partaking in an algorithm is fixed. Any two processes can communicate via a reliable authenticated point-to-point link.

We assume that each correct process has access to a local, unbiased, independent source of randomness. We assume that every correct process has direct access to an oracle  $\Omega$  that, provided with an integer  $n \leq N$ , yields the identities of  $n$  distinct processes, chosen uniformly at random from  $\Pi$ . Implementing  $\Omega$  is beyond the scope of this Part, but it is straightforward in practice. In a system where the set of participating processes is known, sampling reduces to picking with uniform probability an element from the set of processes. In a system without a global membership view, a (nearly) uniform sampling mechanism, *Brahms*, is available in literature due to Bortnikov *et al.* [28]. Remarkably, *Brahms* tolerates (slow) churn, a good fit for the usually permissionless nature of large-scale distributed systems. While we conjecture this amenability to reconfiguration would gracefully extend to our algorithms, a thorough analysis of churn is beyond the scope of this Part.

At most a fraction  $f$  of the processes are Byzantine, i.e., subject to arbitrary failures [104]. Byzantine processes may collude and coordinate their actions. Unless stated otherwise, we denote by  $\Pi_C \subseteq \Pi$  the set of correct processes and by  $C = |\Pi_C| = (1 - f)N$  the number of correct processes. We assume a static Byzantine adversary controlling the faulty processes, i.e., the set of processes controlled by the adversary is fixed at the beginning and does not change throughout the execution of the protocols.

We make standard cryptographic assumptions regarding the power of the adversary, namely that it cannot subvert cryptographic primitives, e.g., forge a signature. We finally assume that Byzantine processes are not aware of (1) the output of the local source of randomness of any correct process; and (2) which correct processes are communicating with each other. Both assumptions are crucial to the safety of sampling. Should the membership of a process' sample be revealed, an adversary could bias its view of the system, singling out the queried

## Chapter 1. Overview

---

processes with non-representative information. Anonymous communication is admittedly our system's most significant practical constraint. Even against ISP-grade adversaries, however, anonymous communication was achieved in practice (at the cost of additional latency) using techniques such as onion routing [58] or private messaging [147].



## 2 Murmur

In this chapter, we present the **probabilistic broadcast** abstraction and discuss its properties. We then present Murmur, an algorithm that implements probabilistic broadcast, and evaluate its **security** and **complexity** as a function of its **parameters**.

The probabilistic broadcast abstraction serves the purpose of reliably broadcasting a single message from a designated correct sender to all correct processes (**validity**, **totality**).

We use probabilistic broadcast in the implementation of Sieve (see Chapter 3) to initially distribute the message from the designated sender to all correct processes.

### 2.1 Interface

The **probabilistic broadcast** interface (instance  $pb$ , sender  $\sigma$ ) exports the following **events**:

- **Request**:  $\langle pb.Broadcast \mid m \rangle$ : Broadcasts a message  $m$  to all processes. This is only used by  $\sigma$ .
- **Indication**  $\langle pb.Deliver \mid m \rangle$ : Delivers a message  $m$  broadcast by process  $\sigma$ .

For any  $\epsilon \in [0, 1]$ , we say that probabilistic broadcast is  $\epsilon$ -secure if:

1. **No duplication**: No correct process delivers more than one message.
2. **Integrity**: If a correct process delivers a message  $m$ , and  $\sigma$  is correct, then  $m$  was previously broadcast by  $\sigma$ .
3.  $\epsilon$ -**Validity**: If  $\sigma$  is correct, and  $\sigma$  broadcasts a message  $m$ , then  $\sigma$  eventually delivers  $m$  with probability at least  $(1 - \epsilon)$ .
4.  $\epsilon$ -**Totality**: If a correct process delivers a message, then every correct process eventually delivers a message with probability at least  $(1 - \epsilon)$ .

### 2.2 Algorithm

Murmur (Algorithm 1) distributes a single message across the system by means of **gossip**: upon reception, a correct process relays the message to a set of randomly selected neighbors. The algorithm depends on one integer parameter,  $G$  (*expected gossip sample size*), whose value we discuss in Section 2.4.

**Initialization.** Upon initialization, (line 11) every correct process randomly samples a value  $\bar{G}$  from a *Poisson* distribution with expected value  $G$ , and uses the sampling oracle  $\Omega$  to select  $\bar{G}$  distinct processes that it will use to initialize its **gossip sample**  $\mathcal{G}$ .

**Link reciprocation.** Once its gossip sample is initialized, a correct process sends a `GossipSubscribe` message to all the processes in  $\mathcal{G}$  (line 13). Upon receiving a `GossipSubscribe` message from a process  $\pi$  (line 17), a correct process adds  $\pi$  to its own gossip sample (line 22), and sends back the gossiped message if it has already received it (line 20).

**Gossip.** When broadcasting the message (line 34), a correct designated sender  $\sigma$  signs the message and sends it to every process in its gossip sample  $\mathcal{G}$  (line 28). Upon receiving a correctly signed message from  $\sigma$  (line 37) for the first time (this is enforced by updating the value of *delivered*, line 25), a correct process delivers it (line 30) and forwards it to every process in its gossip sample (line 28).

### 2.3 No duplication, integrity and validity

We start by verifying that Murmur satisfies **no duplication**, **integrity** and **0-validity**, independently of  $G$ .

**Theorem 1.** *Murmur satisfies no duplication.*

*Proof.* Procedure *dispatch* explicitly checks (line 25) if the variable *delivered* is equal to  $\perp$  before delivering any message. Before a message is delivered (line 30), *delivered* is updated to a value different from  $\perp$  (line 26). Therefore a correct process only delivers one message.  $\square$

**Theorem 2.** *Murmur satisfies integrity.*

*Proof.* Upon receiving a `Gossip` message, a correct process checks its signature against the public key of the designated sender  $\sigma$  (line 37). Moreover, if  $\sigma$  is correct, it only signs *message* when broadcasting (line 34). Since we assume that cryptographic signatures cannot be forged, this implies that the message was previously broadcast by  $\sigma$ .  $\square$

---

**Algorithm 1** Murmur

---

```

1: Implements:
2:   ProbabilisticBroadcast, instance pb
3:
4: Uses:
5:   AuthenticatedPointToPointLinks, instance al
6:
7: Parameters:
8:    $G$ : expected gossip sample size
9:
10: upon event  $\langle pb.Init \rangle$  do
11:    $\mathcal{G} = \Omega(\text{Poisson}[G]);$ 
12:   for all  $\pi \in \mathcal{G}$  do
13:     trigger  $\langle al.Send \mid \pi, [\text{GossipSubscribe}] \rangle;$ 
14:   end for
15:    $delivered = \perp;$ 
16:
17: upon event  $\langle al.Deliver \mid \pi, [\text{GossipSubscribe}] \rangle$  do
18:   if  $delivered \neq \perp$  then
19:      $(message, signature) = delivered;$ 
20:     trigger  $\langle al.Send \mid \pi, [\text{Gossip}, message, signature] \rangle;$ 
21:   end if
22:    $\mathcal{G} \leftarrow \mathcal{G} \cup \{\pi\};$ 
23:
24: procedure  $dispatch(message, signature)$  is
25:   if  $delivered = \perp$  then
26:      $delivered \leftarrow (message, signature);$ 
27:     for all  $\pi \in \mathcal{G}$  do
28:       trigger  $\langle al.Send \mid \pi, [\text{Gossip}, message, signature] \rangle;$ 
29:     end for
30:     trigger  $\langle pb.Deliver \mid message \rangle$ 
31:   end if
32:
33: upon event  $\langle pb.Broadcast \mid message \rangle$  do ▷ only process  $\sigma$ 
34:    $dispatch(message, sign(message));$ 
35:
36: upon event  $\langle al.Deliver \mid \pi, [\text{Gossip}, message, signature] \rangle$  do
37:   if  $verify(\sigma, message, signature)$  then
38:      $dispatch(message, signature);$ 
39:   end if
40:

```

---

**Theorem 3.** *Murmur satisfies 0-validity.*

*Proof.* Upon broadcasting a message  $m$ , a correct sender calls the procedure  $dispatch(m, sign(m))$  (line 34). Since  $delivered$  is initialized to  $\perp$ , this immediately results in the delivery of  $m$  (line 30).

Since the validity property is satisfied deterministically, Murmur satisfies  $\epsilon$ -validity for  $\epsilon = 0$ .  $\square$

### 2.4 Totality

We now compute, given the parameter  $G$ , the  $\epsilon$ -**totality** of Murmur. To this end, we first prove some preliminary lemmas.

**Lemma 1.** *Let  $\rho$  and  $\pi$  be two correct processes, let  $\rho$  be in  $\pi$ 's gossip sample. Then  $\pi$  is eventually in  $\rho$ 's gossip sample.*

*Proof.* A gossip sample is updated only upon initialization (line 11) or when a GossipSubscribe message is received (line 22).

If  $\pi$  selected  $\rho$  upon initialization, then it also sent it a GossipSubscribe message (line 13). Since Byzantine network scheduling can only finitely delay the messages between correct processes,  $\rho$  eventually receives  $\pi$ 's message (line 17) and adds  $\pi$  to its gossip sample.

If  $\pi$  received a GossipSubscribe message from  $\rho$ , then (line 13)  $\rho$  selected  $\pi$  upon initialization, which means that  $\pi$  is already in  $\rho$ 's gossip sample.  $\square$

**Definition 1** (Correct gossip network). Let  $\pi, \rho$  be two correct processes, let  $\pi \leftrightarrow \rho$  denote the condition  $\rho$  is eventually in  $\pi$ 's gossip sample. Lemma 1 proves that

$$(\pi \leftrightarrow \rho) \Leftrightarrow (\rho \leftrightarrow \pi)$$

We define **correct gossip network** to be the undirected graph

$$\mathbb{G} = (\Pi_C, \{(\pi, \rho) \in \Pi_C^2 \mid \pi \leftrightarrow \rho\}) \quad (2.1)$$

**Lemma 2.** *If the correct gossip network is connected, then Murmur satisfies totality.*

*Proof.* We start by noting that a correct process eventually delivers a message (line 30) if and only if it eventually sets  $delivered$  to a value different from  $\perp$  (line 26).

Let  $\pi$  be a correct process for which eventually  $delivered \neq \perp$ . Upon setting  $delivered \leftarrow (m \neq \perp)$ ,  $\pi$  sends  $m$  to all the processes in its gossip sample (line 28). Moreover, upon receiving a GossipSubscribe message *after* setting  $delivered \leftarrow m$ ,  $\pi$  replies with  $m$  (line 20).

Therefore, every correct process that is eventually in  $\pi$ 's gossip sample eventually satisfies  $delivered \neq \perp$ . If  $\mathbb{G}$  is connected, then a path exists in  $\mathbb{G}$  between  $\pi$  and every other correct process, and they all eventually satisfy  $delivered \neq \perp$ , i.e., they deliver a message.  $\square$

From Lemma 2 it follows that Murmur satisfies  $\epsilon$ -totality if the probability of  $\mathbb{G}$  being disconnected is at most  $\epsilon$ .

**Notation 1** (Binomial distribution). We use  $\text{Bin}[N, p]$  to denote the **binomial distribution** with  $N$  trials and  $p$  probability of success.

**Notation 2** (Poisson distribution). We use  $\text{Poisson}[\lambda]$  to denote the **Poisson distribution** with expected value  $\lambda$ .

**Notation 3** (Probability). Let  $E, F$  be events. We use  $\mathcal{P}[E]$  to denote the probability of  $E$ . We use  $\mathcal{P}[E | F]$  to denote the probability of  $E$ , conditioned on the occurrence of  $F$ .

Let  $X, Y, Z$  be random variables. For example, we use the following expressions interchangeably:

$$\mathcal{P}[\bar{X}] \longleftrightarrow \mathcal{P}[X = \bar{X}]$$

Note how  $X$  is a random variable, while  $\bar{X}$  is an element in the codomain of  $X$ . Stand-ins can be combined. For example, we use the following expressions interchangeably:

$$\begin{aligned} \mathcal{P}[\bar{X}, \bar{Y}] &\longleftrightarrow \mathcal{P}[X = \bar{X}, Y = \bar{Y}] \\ \mathcal{P}[\bar{X} | \bar{Y}] &\longleftrightarrow \mathcal{P}[X = \bar{X} | Y = \bar{Y}] \\ \mathcal{P}[\bar{X}, \bar{Y} | \bar{Z}] &\longleftrightarrow \mathcal{P}[X = \bar{X}, Y = \bar{Y} | Z = \bar{Z}] \end{aligned}$$

Stand-ins are only used to express exact values. Whenever non-trivial expressions are needed, we use their explicit form. Explicit notation and stand-ins can be combined. For example, we use the following expressions interchangeably:

$$\begin{aligned} \mathcal{P}[\bar{X} | Y < K] &\longleftrightarrow \mathcal{P}[X = \bar{X} | Y < K] \\ \mathcal{P}[\bar{X} | X < K] &\longleftrightarrow \mathcal{P}[X = \bar{X} | X < K] \end{aligned}$$

**Lemma 3.** *In the limit  $N \rightarrow \infty$ ,  $\mathbb{G}$  is a  $G(C, p)$  Erdős–Rényi graph, with*

$$p = 1 - \left(1 - \frac{C}{N}\right)^2$$

*Proof.* It is a known result that, for large samples and small probabilities, a binomial distribution converges to a Poisson distribution:

$$\lim_{\substack{N \rightarrow \infty \\ Np = \text{const}}} \left[ \text{Bin}[N, p](n) = \binom{N}{n} p^n (1-p)^{N-n} \right] = \left[ \frac{(Np)^n}{n!} e^{-Np} = \text{Poisson}[Np](n) \right]$$

## Chapter 2. Murmur

---

therefore, in the limit  $N \rightarrow \infty$ ,

$$\text{Poisson}[G](n) \simeq \text{Bin}\left[N, \frac{G}{N}\right](n) \quad (2.2)$$

As we discussed in Section 2.2, a gossip sample  $\mathcal{G}$  is initialized upon initialization (line 11) by first sampling a value  $\bar{G}$  from a  $\text{Poisson}[G]$  distribution, then selecting  $\bar{G}$  distinct processes from  $\Pi$  with uniform probability.

Let  $\pi \in \Pi_C, \rho \in \Pi$ , let  $\mathcal{G}_\pi^{in}$  be  $\pi$ 's initial gossip sample, let  $q = G/N$ . By the law of total probability, and using Equation (2.2), we have for large  $N$

$$\begin{aligned} \mathcal{P}\left[\rho \in \mathcal{G}_\pi^{in}\right] &= \sum_{\bar{G}=0}^N \left( \mathcal{P}\left[\rho \in \mathcal{G}_\pi^{in} \mid \bar{G}\right] \mathcal{P}[\bar{G}] \right) \\ &= \sum_{\bar{G}=0}^N \left( \frac{\bar{G}}{N} \text{Poisson}[G](\bar{G}) \right) \simeq \sum_{\bar{G}=0}^N \left( \frac{\bar{G}}{N} \text{Bin}[N, q](\bar{G}) \right) \\ &= \sum_{\bar{G}=0}^N \left( \frac{\bar{G}}{N} \binom{N}{\bar{G}} q^{\bar{G}} (1-q)^{N-\bar{G}} \right) \\ &= \sum_{\bar{G}=0}^N \left( \frac{\bar{G}}{N} \frac{N!}{\bar{G}!(N-\bar{G})!} q^{\bar{G}} (1-q)^{N-\bar{G}} \right) \\ &= \sum_{\bar{G}=1}^N \left( \frac{(N-1)!}{(\bar{G}-1)!(N-\bar{G})!} q q^{\bar{G}-1} (1-q)^{N-\bar{G}} \right) \\ &= q \sum_{\bar{G}'=0}^{N-1} \left( \frac{(N-1)!}{\bar{G}'!(N-1-\bar{G}')!} q^{\bar{G}'} (1-q)^{N-1-\bar{G}'} \right) \\ &= q \sum_{\bar{G}'=0}^{N-1} \text{Bin}[N-1, q](\bar{G}') = q \end{aligned}$$

Let  $\rho_1, \dots, \rho_R$  be distinct processes, with  $R \leq N$ . Similar calculations yield

$$\mathcal{P}\left[\rho_1 \in \mathcal{G}_\pi^{in}, \dots, \rho_R \in \mathcal{G}_\pi^{in}\right] = q^R \quad (2.3)$$

Equation (2.3) proves that every process  $\rho \in \Pi$  has an independent probability  $q$  of being in  $\mathcal{G}_\pi^{in}$ . Since for any two  $\pi, \xi \in \Pi_C$  we have

$$(\pi \leftrightarrow \xi) \Leftrightarrow \left( \pi \in \mathcal{G}_\xi^{in} \vee \xi \in \mathcal{G}_\pi^{in} \right)$$

we can derive the probability  $p$  of any two correct processes being connected:

$$p = 1 - (1-q)^2 = 1 - \left(1 - \frac{G}{N}\right)^2 \quad (2.4)$$

Therefore, following Equations (2.3) and (2.4),  $\mathbb{G} = G(C, p)$  is an Erdős – Rényi graph with  $H$  nodes and  $p$  probability of connection between any two nodes.  $\square$

Lemma 3 allows us to bound the  $\epsilon$ -totality of Murmur, given  $G$ .

**Theorem 4.** *Murmur satisfies  $\epsilon_t$ -totality, with  $\epsilon_t$  bound by*

$$\epsilon_t \leq \sum_{k=1}^{C/2} \binom{C}{k} (1-p)^{k(C-k)} \quad (2.5)$$

*Proof.* It follows immediately from Lemma 3 and a known result ([2]) on the connectivity of Erdős–Rényi graphs.  $\square$

We prove an additional result on the latency of Murmur.

**Theorem 5.** *The latency of Murmur is asymptotically sub-logarithmic. More formally, the diameter  $D(C, G)$  of the correct gossip network limits to*

$$\lim_{C \rightarrow \infty} D(C, G) = \frac{\log(C)}{\log(2-2f) + \log(G)}$$

*Proof.* It is a known result ([47]) that the diameter of an Erdős–Rényi graph  $G(C, p)$  converges, for  $Cp \rightarrow \infty$ , to  $\log(C)/\log(Cp)$ .

Noting that

$$\lim_{C \rightarrow \infty} 1 - \left(1 - \frac{G}{N}\right)^2 = \frac{2G}{N}$$

we get

$$\begin{aligned} \lim_{C \rightarrow \infty} D(C, G) &= \frac{\log(C)}{\log(C) + \log(p)} \\ &= \frac{\log(C)}{\log(C) + \log(2) + \log(G) - \log(N)} \\ &= \frac{\log(C)}{\log\left(\frac{2C}{N}\right) + \log(G)} \\ &= \frac{\log(C)}{\log(2-2f) + \log(G)} \end{aligned}$$

which proves the lemma. For a fixed security  $\epsilon$ , we showed in Theorem 4 that  $G$  must scale logarithmically with the size of the system. As a result, for a fixed security  $\epsilon$ , the latency scales as  $O(\log(N)/\log(\log(N)))$   $\square$





## 3 Sieve

In this chapter, we present the **probabilistic consistent broadcast** abstraction and discuss its properties. We then present Sieve, an algorithm that implements probabilistic consistent broadcast, and evaluate its **security** and **complexity** as a function of its **parameters**.

The probabilistic consistent broadcast abstraction allows a subset of the correct processes to agree on a single message from a potentially Byzantine designated sender. Intuitively, probabilistic consistent broadcast offers a tradeoff with respect to probabilistic broadcast. Probabilistic broadcast guarantees (**totality**) that if any correct process delivers a message, every correct process delivers a message. Probabilistic consistent broadcast, instead, guarantees (**consistency**) that, even if the sender is Byzantine, no two correct processes deliver different messages. However, if the sender is Byzantine, it may happen with a non-negligible probability that only an intermediate fraction of the correct processes deliver the message.

We use probabilistic consistent broadcast in the implementation of Contagion (see Chapter 4) as a way to consistently broadcast messages.

### 3.1 Interface

The **probabilistic consistent broadcast** interface (instance  $pcb$ , sender  $\sigma$ ) exposes the following two **events**:

- **Request:**  $\langle pcb.Broadcast \mid m \rangle$ : Broadcasts a message  $m$  to all processes. This is only used by  $\sigma$ .
- **Indication:**  $\langle pcb.Deliver \mid m \rangle$ : Delivers a message  $m$  broadcast by process  $\sigma$ .

For any  $\epsilon \in [0, 1]$ , we say that probabilistic consistent broadcast is  $\epsilon$ -secure if:

1. **No duplication:** No correct process delivers more than one message.

## Chapter 3. Sieve

---

2. **Integrity:** If a correct process delivers a message  $m$ , and  $\sigma$  is correct, then  $m$  was previously broadcast by  $\sigma$ .
3.  **$\epsilon$ -Total validity:** If  $\sigma$  is correct, and  $\sigma$  broadcasts a message  $m$ , every correct process eventually delivers  $m$  with probability at least  $(1 - \epsilon)$ .
4.  **$\epsilon$ -Consistency:** Every correct process that delivers a message delivers the same message with probability at least  $(1 - \epsilon)$ .

### 3.2 Algorithm

---

**Algorithm 2 Procedure *sample*.**

---

```
1: procedure sample(message, size) is
2:    $\psi = \emptyset$ ;
3:   for size times do
4:      $\psi \leftarrow \psi \cup \Omega(1)$ ;
5:   end for
6:   for all  $\pi \in \psi$  do
7:     trigger  $\langle \text{al.Send} \mid \pi, [\text{message}] \rangle$ ;
8:   end for
9:   return  $\psi$ ;
10:
```

---

Algorithm 2 implements a *sample* procedure that we use both in the implementation of Sieve and Contagion. Procedure *sample*(*message*, *size*) uses  $\Omega$  to pick *size* processes with replacement, and sends them *message*.

Algorithm 3 implements Sieve. Sieve consistently distributes a single message across the system as follows:

- Initially, probabilistic broadcast distributes potentially conflicting copies of the message to every correct process.
- Upon receiving a message  $m$  from probabilistic broadcast, a correct process issues an Echo message for  $m$ .
- Upon receiving enough Echo messages for the message  $m$  it Echoed, a correct process delivers  $m$ .

A correct process collects Echo messages from a randomly selected *echo sample* of size  $E$ , and delivers the message it Echoed upon receiving  $\hat{E}$  Echoes for it. We discuss the values of the two parameters of Sieve in Section 3.10.

**Algorithm 3 Sieve**


---

```

1: Implements:
2:   ProbabilisticConsistentBroadcast, instance pcb
3:
4: Uses:
5:   AuthenticatedPointToPointLinks, instance al
6:   ProbabilisticBroadcast, instance pb
7:
8: Parameters:
9:    $E$ : echo sample size
10:   $\hat{E}$ : delivery threshold
11:
12: upon event  $\langle pcb.Init \rangle$  do
13:    $echo = \perp$ ;  $delivered = \text{False}$ ;  $\tilde{\mathcal{E}} = \emptyset$ ;
14:
15:    $\mathcal{E} = \text{sample}(\text{EchoSubscribe}, E)$ ;
16:    $replies = \{\perp\}^E$ ;
17:
18: upon event  $\langle al.Deliver \mid \pi, [\text{EchoSubscribe}] \rangle$  do
19:   if  $echo \neq \perp$  then
20:      $(message, signature) = echo$ ;
21:     trigger  $\langle al.Send \mid \pi, [\text{Echo}, message, signature] \rangle$ ;
22:   end if
23:    $\tilde{\mathcal{E}} \leftarrow \tilde{\mathcal{E}} \cup \{\pi\}$ ;
24:
25: upon event  $\langle pcb.Broadcast \mid message \rangle$  do ▷ only process  $\sigma$ 
26:   trigger  $\langle pb.Broadcast \mid [\text{Send}, message, sign(message)] \rangle$ ;
27:
28: upon event  $\langle pb.Deliver \mid [\text{Send}, message, signature] \rangle$  do
29:   if  $verify(\sigma, message, signature)$  then
30:      $echo \leftarrow (message, signature)$ ;
31:     for all  $\rho \in \tilde{\mathcal{E}}$  do
32:       trigger  $\langle al.Send \mid \rho, [\text{Echo}, message, signature] \rangle$ ;
33:     end for
34:   end if
35:
36: upon event  $\langle al.Deliver \mid \pi, [\text{Echo}, message, signature] \rangle$  do
37:   if  $\pi \in \tilde{\mathcal{E}}$  and  $replies[\pi] = \perp$  and  $verify(\sigma, message, signature)$  then
38:      $replies[\pi] \leftarrow (message, signature)$ ;
39:   end if
40:
41: upon  $|\{\rho \in \tilde{\mathcal{E}} \mid replies[\rho] = echo\}| \geq \hat{E}$  and  $delivered = \text{False}$  do
42:    $delivered \leftarrow \text{True}$ ;
43:   trigger  $\langle pcb.Deliver \mid message \rangle$ ;
44:

```

---

## Chapter 3. Sieve

---

**Sampling.** Upon initialization (line 12), a correct process randomly selects an **echo sample**  $\mathcal{E}$  of size  $E$ . Samples are selected with replacement by repeatedly calling  $\Omega$  (Algorithm 2, line 4). A correct process sends an `EchoSubscribe` message to all the processes in its echo sample (Algorithm 2, line 7).

**Publish-subscribe.** Unlike in the deterministic version of Authenticated Echo Broadcast, where a correct process broadcasts its Echo messages to the whole system, here each process only listens for messages coming from its echo sample (line 37).

A correct process maintains an **echo subscription set**  $\tilde{\mathcal{E}}$ . Upon receiving an `EchoSubscribe` message from a process  $\pi$ , a correct process adds  $\pi$  to  $\tilde{\mathcal{E}}$  (line 23). If a correct process receives an `EchoSubscribe` message *after* publishing its Echo message, it also sends back the previously published message (line 21).

A correct process only sends its Echo messages (line 32) to its echo subscription set.

**Echo.** The designated sender  $\sigma$  initially broadcasts its message using probabilistic broadcast (line 26). Upon `pb.Delivery` of a message  $m$  (correctly signed by  $\sigma$ ) (line 28), a correct process sends an Echo message for  $m$  to all the nodes in its echo subscription set (line 32).

**Delivery.** A correct process  $\pi$  that Echoed a message  $m$  delivers  $m$  (line 43) upon collecting at least  $\hat{E}$  Echo messages for  $m$  (line 41) from the processes in its echo sample.

### 3.3 No duplication and integrity

We start by verifying that Sieve satisfies both **no duplication** and **integrity**.

**Theorem 6.** *Sieve satisfies no duplication.*

*Proof.* A message is delivered (line 43) only if the variable *delivered* is equal to `False` (line 41). Before any message is delivered, *delivered* is set to `True`. Therefore no more than one message is ever delivered.  $\square$

**Theorem 7.** *Sieve satisfies integrity.*

*Proof.* Upon receiving an Echo message, a correct process checks its signature against the public of the designated sender  $\sigma$  (line 37), and the  $(message, signature)$  pair is added to the *replies* variable only if this check succeeds. Moreover, a message is delivered only if it is represented at least  $\hat{E} > 0$  times in *replies* (line 41).

If  $\sigma$  is correct, it only signs *message* when broadcasting (line 26). Since we assume that cryptographic signatures cannot be forged, this implies that the message was previously

broadcast by  $\sigma$ . □

### 3.4 Total validity

We now compute, given  $E$  and  $\hat{E}$ , the  $\epsilon$ -**total validity** of Sieve. To this end, we prove some preliminary lemmas.

**Lemma 4.** *In an execution of Sieve, if pb does not satisfy totality, then pcb does not satisfy total validity.*

*Proof.* A correct process delivers a message (line 43) only if the *echo* variable is different from  $\perp$ . Moreover, the *echo* variable is set to a value different from  $\perp$  (line 30) only upon pb.Delivery of a message (line 28).

Let  $m$  be the message broadcast by the correct sender  $\sigma$ . If pb does not satisfy totality, then at least one correct process never sets *echo* to  $m$ . Therefore, at least one correct process does not deliver the  $m$ , and the total validity of pcb is compromised. □

**Lemma 5.** *In an execution of Sieve, if pb satisfies totality and no correct process has more than  $E - \hat{E}$  Byzantine processes in its echo sample, then pcb satisfies total validity.*

*Proof.* Let  $m$  be the message broadcast by the correct sender  $\sigma$ . Since pb satisfies totality (it always satisfies validity), every correct process eventually issues an Echo( $m$ ) message (i.e., an Echo message for  $m$ ) (line 32).

Let  $\pi$  be a correct process that has no more than  $E - \hat{E}$  Byzantine processes in its echo sample. Obviously,  $\pi$  has at least  $\hat{E}$  correct processes in its echo sample. Therefore,  $\pi$  eventually receives at least  $\hat{E}$  Echo( $m$ ) messages (line 36), and delivers  $m$  (line 41). □

Lemmas 4 and 5 allow us to bound the  $\epsilon$ -total validity of Sieve, given  $E$  and  $\hat{E}$ .

**Theorem 8.** *Sieve satisfies  $\epsilon_v$ -total validity, with*

$$\begin{aligned} \epsilon_v &\leq \epsilon_t^{pb} + (1 - \epsilon_t^{pb})(1 - (1 - \epsilon_o)^C) \\ \epsilon_o &= \sum_{\bar{F}=E-\hat{E}+1}^E \text{Bin}[E, f](\bar{F}) \end{aligned} \tag{3.1}$$

*if the underlying abstraction of probabilistic broadcast satisfies  $\epsilon_t^{pb}$ -totality.*

*Proof.* Following from Lemmas 4 and 5, the total validity of pcb can be compromised only if the totality of pb is compromised as well, or if at least one correct process has more than  $E - \hat{E}$  Byzantine processes in its echo sample.

## Chapter 3. Sieve

---

Since procedure *sample* independently picks  $E$  processes with replacement, each element of a correct process' echo sample has an independent probability  $f$  of being Byzantine, i.e., the number of Byzantine processes in a correct echo sample is binomially distributed.

Therefore, a correct process has a probability  $\epsilon_o$  of having more than  $E - \hat{E}$  Byzantine processes in its echo sample. Since every correct process picks its echo sample independently, the probability of at least one correct process having more than  $E - \hat{E}$  Byzantine processes in its echo sample is  $1 - (1 - \epsilon_o)^C$ .  $\square$

### 3.5 Preliminary lemmas

In order to compute an upper bound for the probability of the consistency of Sieve being compromised, we will make use of some preliminary lemmas. The statements of these lemmas are independent from the context of Sieve. For the sake of readability, we therefore gather them in this section, and use them throughout the rest of this chapter.

**Lemma 6.** *Let  $A, B \in \mathbb{N}$ , let  $x, y \in \mathbb{N}$  such that  $x + y \leq B$ . Let  $X, Y$  be random variables defined by*

$$\begin{aligned}\mathcal{P}[\bar{X}] &= \text{Bin}\left[A, \frac{x}{B}\right](\bar{X}) \\ \mathcal{P}[\bar{Y} | \bar{X}] &= \text{Bin}\left[A - \bar{X}, \frac{y}{B - x}\right](\bar{Y})\end{aligned}$$

We have

$$\mathcal{P}[X + Y = K] = \text{Bin}\left[A, \frac{x + y}{B}\right](K)$$

*Proof.* Since  $X$  is binomially distributed, it can be expressed as a sum of independent Bernoulli random variables:

$$\begin{aligned}X &= X_1 + \dots + X_A \\ X_i &\sim \text{Bern}\left[\frac{x}{B}\right]\end{aligned}$$

Given the value of  $\bar{X}$ ,  $Y$  is also binomially distributed with probability  $y/(B - x)$  and  $E - \bar{X}$  trials. We can therefore express  $Y$  as the sum of  $E$  Bernoulli variables  $Y_1, \dots, Y_E$ :

$$\begin{aligned}Y &= Y_1 + \dots + Y_E \\ \mathcal{P}[Y_i = 1 | \bar{X}_i] &= \begin{cases} 0 & \text{iff } \bar{X}_i = 1 \\ \frac{y}{B - x} & \text{otherwise} \end{cases}\end{aligned}$$

We indeed note how, out of  $Y_1, \dots, Y_E$ :

- Only  $E - \bar{X}$  variables have a non-null probability of being equal to 1.

- Those variables that have a non-null probability of being equal to 1 have a probability  $y/(B-x)$  of being equal to 1.

We therefore have

$$X + Y = (X_1 + Y_1) + \dots + (X_A + Y_A)$$

and from the law of total probability we have

$$\begin{aligned} \mathcal{P}[X_i + Y_i = 1] &= \mathcal{P}[X_i = 1] + \mathcal{P}[Y_i = 1 \mid X_i = 0] \mathcal{P}[X_i = 0] \\ &= \frac{x}{B} + \left(1 - \frac{x}{B}\right) \left(\frac{y}{B-x}\right) \\ &= \frac{x}{B} + \left(\frac{(B-x)}{B} \frac{y}{(B-x)}\right) \\ &= \frac{x+y}{B} \end{aligned}$$

therefore

$$(X_i + Y_i) \sim \text{Bern}\left[\frac{x+y}{B}\right]$$

which proves the lemma. □

**Lemma 7.** Let  $A, B \in \mathbb{N}$  such that  $A \geq B$ , let  $p \in [0, 1]$ . Let  $X_1, \dots, X_B$  be random variables defined by

$$X_i \sim \text{Bin}[A - i, p]$$

We have that

$$\mathcal{P}[X_i \geq B - i]$$

is an increasing function of  $i$ .

*Proof.* We prove the lemma by induction by showing that, for any  $i < B$ ,

$$\mathcal{P}[X_i \geq B - i] \leq \mathcal{P}[X_{i+1} \geq B - (i + 1)]$$

In order to obtain the above, we expand

$$\begin{aligned} &\mathcal{P}[X_i \geq B - i] - \mathcal{P}[X_{i+1} \geq B - i - 1] \\ &= \sum_{n=B-i}^{A-i} \left( \frac{(A-i)!}{(A-i-n)!n!} p^n (1-p)^{A-i-n} \right) \\ &\quad - \sum_{n=B-i-1}^{A-i-1} \left( \frac{(A-i-1)!}{(A-i-1-n)!n!} p^n (1-p)^{A-i-1-n} \right) \\ &= (\star_1) \end{aligned}$$

### Chapter 3. Sieve

---

By shifting the index in the second sum we get

$$\begin{aligned}
 (\star_1) &= \sum_{n=B-i}^{A-i} \left( \frac{(A-i)!}{(A-i-n)!n!} p^n (1-p)^{A-i-n} \right) \\
 &\quad - \sum_{n=B-i}^{A-i} \left( \frac{(A-i-1)!}{(A-i-1-(n-1))!(n-1)!} \right. \\
 &\quad \quad \left. p^{(n-1)} (1-p)^{A-i-1-(n-1)} \right) \\
 &= \sum_{n=B-i}^{A-i} \left( \frac{(A-i)!}{(A-i-n)!n!} p^n (1-p)^{A-i-n} \right. \\
 &\quad \quad \left. - \frac{(A-i)!n}{(A-i)(A-i-n)!n!} \frac{p^n}{p} (1-p)^{A-i-n} \right) \\
 &= \sum_{n=B-i}^{A-i} \left( \left( \frac{(A-i)!}{(A-i-n)!n!} p^n (1-p)^{A-i-n} \right) \left( 1 - \frac{n}{(A-i)p} \right) \right) \\
 &= (\star_2)
 \end{aligned}$$

and by letting  $N = A - i$ ,  $M = B - i$  we get

$$(\star_2) = \sum_{n=M}^N \left( \text{Bin}[N, p](n) \left( 1 - \frac{n}{Np} \right) \right) = (\star_3)$$

Noticing that  $(1 - n/Np)$  is positive for  $n < Np$ , we have

$$\begin{aligned}
 (\star_3) &\leq \sum_{n=0}^N (\text{Bin}[N, p](n)) - \frac{1}{Np} \sum_{n=0}^N (n \text{Bin}[N, p](n)) \\
 &= 1 - \frac{Np}{Np} = 0
 \end{aligned}$$

which proves the lemma. □

**Notation 4** (Ranges). Let  $a, b \in \mathbb{N}$ , with  $b \geq a$ . We use  $a..b$  to denote the range of integers  $\{a, \dots, b\}$ .

**Lemma 8.** Let  $f, g : 0..K \rightarrow \mathbb{R}$ , with  $f$  increasing,  $g$  positive and

$$\sum_{x=0}^K g(x) = 1$$

we have

$$\sum_{x=0}^K (f(x)g(x)) \geq \sum_{x=0}^{K-1} \left( \frac{f(x)g(x)}{1-g(K)} \right)$$



*Proof.* We have

$$\begin{aligned}
 & \sum_{x=0}^K (f(x)g(x)) - \sum_{x=0}^{K-1} \left( \frac{f(x)g(x)}{1-g(K)} \right) \\
 &= f(K)g(K) + \sum_{x=0}^{K-1} (f(x)g(x)) \left( 1 - \frac{1}{1-g(K)} \right) \\
 &= f(K)g(K) - \sum_{x=0}^{K-1} (f(x)g(x)) \left( \frac{1-(1-g(K))}{1-g(K)} \right) \\
 &= g(K) \left( f(K) - \frac{1}{1-g(K)} \sum_{x=0}^{K-1} (f(x)g(x)) \right) \\
 &= \frac{g(K)}{1-g(K)} \left( f(K) - f(K)g(K) - \sum_{x=0}^{K-1} (f(x)g(x)) \right) \\
 &= \frac{g(K)}{1-g(K)} \left( f(K) \sum_{x=0}^K (f(x)g(x)) \right) \\
 &= (\star_1)
 \end{aligned}$$

and noting that  $g(K) \geq 0$ ,  $1-g(K) \geq 0$ , and  $f$  is increasing, we have

$$(\star_1) \geq \frac{g(K)}{1-g(K)} \left( f(K) - f(K) \sum_{x=0}^K g(x) \right) = (\star_2)$$

and since  $\sum g(x) = 1$  we get

$$\frac{g(K)}{1-g(K)} (f(K) - f(K)) = 0$$

□

**Corollary 1.** Let  $f, g : 0..K \rightarrow \mathbb{R}$ , with  $f$  increasing,  $g$  positive and

$$\sum_{x=0}^K g(x) = 1$$

for any  $l \in 0..(K-1)$ , we have

$$\sum_{x=0}^K (f(x)g(x)) \geq \frac{\sum_{x=0}^{K-l} f(x)g(x)}{\sum_{x=0}^{K-l} g(x)}$$

*Proof.* It follows immediately from applying Lemma 8  $l$  times. □

**Lemma 9.** Let  $f : -1..C \rightarrow \mathbb{R}$ , let  $g, h : -1..C \rightarrow [0, 1]$ , with:

- $f$  decreasing.
- $g, h$  increasing.
- $g(x) \leq h(x)$  for all  $x$ .

### Chapter 3. Sieve

---

- $g(-1) = h(-1) = 0$ .
- $g(C) = h(C) = 1$ .

We have

$$\sum_{x=0}^C (f(x)(g(x) - g(x-1))) \leq \sum_{x=0}^C (f(x)(h(x) - h(x-1)))$$

*Proof.* We have

$$\begin{aligned} & \sum_{x=0}^C f(x)(g(x) - g(x-1)) - \sum_{x=0}^C f(x)(h(x) - h(x-1)) \\ &= \sum_{x=0}^C f(x)((g(x) - h(x)) - (g(x-1) - h(x-1))) \\ &= \sum_{x=0}^C f(x)(g(x) - h(x)) - \sum_{x=0}^C f(x)(g(x-1) - h(x-1)) \\ &= (\star_1) \end{aligned}$$

By shifting the index in then second sum we get

$$\begin{aligned} (\star_1) &= \sum_{x=0}^C f(x)(g(x) - h(x)) - \sum_{x=-1}^{C-1} f(x+1)(g(x) - h(x)) \\ &= \sum_{x=0}^{C-1} (f(x) - f(x+1))(g(x) - h(x)) \\ &\quad + f(C)(g(C) - h(C)) - f(-1)(g(-1) - h(-1)) \\ &= (\star_2) \end{aligned}$$

and by noting that:

- Since  $f$  is decreasing,  $f(x) - f(x+1) \geq 0$ .
- By hypothesis,  $g(x) - h(x) \leq 0$ .
- By hypothesis,  $g(C) - h(C) = 1 - 1 = 0$ .
- By hypothesis,  $g(-1) - h(-1) = 0 - 0 = 0$ .

Consequently, all the terms of the sum in  $(\star_2)$  are negative, and the two terms out of the sum are null. Therefore,  $(\star_2) \leq 0$ .  $\square$

**Lemma 10.** Let  $N \in \mathbb{N}$ , let  $K < N$ , let  $h, p_1, \dots, p_T \in [0, 1]$  such that

$$h = \sum_i p_i \leq \frac{K - \sqrt{K}}{N}$$

let  $X_1, \dots, X_T$  be independent random variables defined by

$$\mathcal{P}[\bar{X}_i] = \text{Bin}[N, p_i](\bar{X}_i)$$

We have

$$\mathcal{P}\left[\bigvee_i (X_i > K)\right] \leq \left(\frac{eNh}{K}\right)^K e^{-Nh}$$

*Proof.* Let  $p \in [0, 1]$ , let  $X \sim \text{Bin}[N, p]$ . From the multiplicative form of the Chernoff bound we have

$$\begin{aligned} \mathcal{P}[X > K] &= \mathcal{P}[X > (1 + \delta)\mu] < \left(\frac{e^\delta}{(1 + \delta)^{(1 + \delta)}}\right)^\mu \\ \delta &= \left(\frac{K}{Np} - 1\right) \\ \mu &= Np \end{aligned}$$

From the above follows

$$\begin{aligned} \mathcal{P}[X > K] &< \left(\frac{\exp\left(\frac{K}{Np} - 1\right)}{\left(\frac{K}{Np}\right)^{\frac{K}{Np}}}\right)^\mu \\ &= \exp\left(Np\left(\frac{K}{Np} - 1 - \frac{K}{Np} \log\left(\frac{K}{Np}\right)\right)\right) \\ &= \exp\left(K - Np - K \log\left(\frac{K}{Np}\right)\right) \\ &= \underbrace{\exp(K - K \log K + K \log N)}_{(\star_a)} \underbrace{\exp(K \log p - Np)}_{(\star_b)} \end{aligned}$$

We now study the domain where  $(\star_b)$  is convex:

$$\frac{\partial^2}{\partial^2 p} \exp(K \log p - Np) = \underbrace{p^{K-2} e^{-Np}}_{\geq 0} (K^2 - K(2Np + 1) + N^2 p^2)$$

Therefore we require

$$N^2 p^2 - (2KN)p + (K^2 - K) \geq 0$$

Which reduces to

$$\begin{aligned} p &\leq \frac{2KN - \sqrt{4K^2N^2 - 4(N^2K^2 - N^2K)}}{2N^2} \\ &= \frac{K - \sqrt{K}}{N} \end{aligned}$$

From Boole's inequality we have

$$\mathcal{P}\left[\bigvee_i (X_i > K)\right] \leq \sum_i \mathcal{P}[X_i > K] = (\star_1)$$

which we can expand into

$$(\star_1) = \underbrace{\exp(K - K \log K + K \log N)}_{(\star_a)} \sum_i \underbrace{\exp(K \log p_i - N p_i)}_{(\star_b)} = (\star_2)$$

As we established,  $(\star_b)$  is convex on the range  $[0, \sum_i p_i]$ . Consequently,

$$\begin{aligned} (\star_2) &\leq \exp(K - K \log K + K \log N) \exp(K \log h - Nh) \\ &= \left(\frac{eNh}{K}\right)^K e^{-Nh} \end{aligned}$$

which proves the lemma. □

## 3.6 Simplified Sieve

In this section, we introduce Simplified Sieve, a modified version of Sieve.

Simplified Sieve is a strawman both from a performance and a safety point of view. Indeed, on the one hand Simplified Sieve has  $O(N^2)$  per-process communication complexity, which makes it unfit for any real-world, scalable deployment. On the other, we prove that it is strictly easier for any Byzantine adversary to compromise the consistency of Simplified Sieve than that of Sieve.

Unlike Sieve, however, Simplified Sieve allows for an analytic probabilistic analysis. A critical goal of this chapter is to compute a bound  $\epsilon_c$  on the probability of compromising the consistency of Simplified Sieve. Since the consistency of Simplified Sieve is weaker than that of Sieve,  $\epsilon_c$  is also a bound on the probability of compromising the consistency of Sieve.

### 3.6.1 Consistency-only broadcast

Simplified Sieve implements consistency-only broadcast, a minimal version of the probabilistic consistent broadcast abstraction, designed to only provide  $\epsilon$ -consistency. In particular, we

drop the no duplication property, i.e., we allow a correct process to deliver more than one message.

The **consistency-only broadcast** interface (instance  $cob$ , sender  $\sigma$ ) exposes the following two events:

- **Request**  $\langle cob.Broadcast \mid m \rangle$ : Broadcasts a message  $m$  to all processes. This is only used by  $\sigma$ .
- **Indication**  $\langle cob.Deliver \mid m \rangle$ : Delivers a message  $m$  broadcast by process  $\sigma$ .

For any  $\epsilon \in [0, 1]$ , we say that consistency-only broadcast is  $\epsilon$ -secure if:

1.  **$\epsilon$ -Consistency**: With probability at least  $(1 - \epsilon)$ , at most one message  $m$  exists, such that  $m$  is delivered by any correct process.

We note how the above definition of  $\epsilon$ -consistency is equivalent to the one we provided in Section 3.1, but adapted for a context where no duplication is not guaranteed. In consistency-only broadcast, consistency is compromised even if a single correct process delivers two or more different messages.

### 3.6.2 Byzantine oracle

In order to implement Simplified Sieve, we make an additional assumption about the system:

- **(Byzantine oracle)** Every correct process has direct access to an oracle  $\Psi$  that, provided with a process  $\pi$ , returns `True` if  $\pi$  is Byzantine, and `False` if  $\pi$  is correct.

This assumption is obviously unsatisfiable in any realistic distributed system. Indeed, a system where every faulty process is publicly flagged can hardly be considered being subject to arbitrary failures. It is therefore critical to underline that Assumption 3.6.2 is **not** a requirement for the implementation of Sieve. Indeed, no correct process invokes  $\Psi$  throughout any execution of Algorithm 3. Assumption 3.6.2 is purely a theoretical artifice to aid in our proof of correctness.

### 3.6.3 Algorithm

Before introducing the design principles behind Simplified Sieve, we prove a simple preliminary result.

**Lemma 11.** *No execution of probabilistic broadcast results in more than  $C$  different messages being delivered.*

## Chapter 3. Sieve

---

*Proof.* Following from Theorem 1, probabilistic broadcast satisfies no duplication, i.e., no correct process delivers more than one message. As we discussed in Section 1.3, the system is composed of  $C$  correct processes.  $\square$

Since the set of messages that are pb.Delivered by at least one correct process has no more than  $C$  elements, and noting that a correct process pcb.Delivers a message  $m$  only if it pb.Delivered  $m$ , it is not restrictive to introduce the following definition.

**Definition 2** (Message). A **message** is an element of the set

$$\mathcal{M} = 1..C$$

---

**Algorithm 4 Procedure** *mimic*.

---

```
1: procedure correct() is
2:   do
3:      $\rho = \Omega(1)$ 
4:   until  $\Psi(\rho) = \text{False}$ 
5:   return  $\rho$ ;
6:
7: procedure mimic(reference) is
8:    $\psi = \emptyset$ ;
9:   for all  $\rho \in \text{reference}$  do
10:    if  $\Psi(\rho) = \text{True}$  then
11:       $\psi \leftarrow \psi \cup \{\rho\}$ ;
12:    else
13:       $\psi \leftarrow \psi \cup \text{correct}()$ ;
14:    end if
15:  end for
16:  return  $\psi$ ;
17:
```

---

Algorithm 5 implements Simplified Sieve. Simplified Sieve bears multiple differences to Sieve:

- A correct process can deliver more than one message. No correct process, however, delivers the same message more than once.
- In order to cob.Deliver a message, a correct process does not need to pb.Deliver any message.
- A correct process maintains  $C$  **echo samples**  $\mathcal{E}[1..C]$ . The Echo messages collected from the processes in the  $i$ -th echo sample  $\mathcal{E}[i]$  determine whether or not message  $i \in \mathcal{M}$  is delivered.
- Echo messages have two fields: *sample* and *message*. Intuitively, an Echo( $s, m$ ) message (i.e., an Echo message with fields  $s$  and  $m$ ) represents the following statement: “within the context of message  $s$ , consider my Echo to be for message  $m$ ”.

**Algorithm 5 Simplified Sieve**


---

```

1: Implements:
2:   ConsistencyOnlyBroadcast, instance cob
3:
4: Uses:
5:   AuthenticatedPointToPointLinks, instance al
6:   ProbabilisticBroadcast, instance pb
7:
8: Parameters:
9:    $E$ : echo sample size
10:   $\hat{E}$ : delivery threshold
11:
12: upon event  $\langle cob.Init \rangle$  do
13:    $delivered = \{\text{False}\}^C$ ;  $reveal = \{\text{False}\}^C$ ;
14:    $revealed = \{\text{False}\}^C$ ;
15:    $replies = \{\perp\}^{C \times E}$ ; ▷  $C \times E$  table filled with  $\perp$ .
16:    $\mathcal{E} = \{\emptyset\}^C$ ;
17:    $\mathcal{E}[1] \leftarrow \text{sample}(\text{EchoSubscribe}, E)$ ;
18:
19:   for  $j \in 2..C$  do
20:      $\mathcal{E}[j] \leftarrow \text{mimic}(\mathcal{E}[1])$ ;
21:   end for
22:
23: upon event  $\langle cob.Broadcast \mid message \rangle$  do ▷ only process  $\sigma$ 
24:   trigger  $\langle pb.Broadcast \mid [\text{Send}, message] \rangle$ ;
25:
26: upon event  $\langle pb.Deliver \mid [\text{Send}, message] \rangle$  do
27:   for all  $\rho \in \Pi$  do
28:     for all  $sample \in \mathcal{M}$  do
29:       trigger  $\langle al.Send \mid \rho, [\text{Echo}, sample, message] \rangle$ ;
30:     end for
31:   end for
32:
33: upon event  $\langle al.Deliver \mid \rho, [\text{Echo}, sample, message] \rangle$  do
34:   if  $\rho \in \mathcal{E}[sample]$  and  $replies[sample][\rho] = \perp$  then
35:      $replies[sample][\rho] \leftarrow message$ ;
36:      $revealed[sample] \leftarrow \text{False}$ ;
37:   end if
38:
39: upon exists  $message$  such that  $|\{\rho \in \mathcal{E}[message] \mid replies[message][\rho] = message\}| \geq$ 
    $\hat{E}$  and  $delivered[message] = \text{False}$  do
40:    $delivered[message] \leftarrow \text{True}$ ;
41:    $reveal[message] \leftarrow \text{True}$ ;
42:   trigger  $\langle cob.Deliver \mid message \rangle$ ;
43:

```

---

### Chapter 3. Sieve

---

```

44: upon exists message such that reveal[message] = True and revealed[message] =
    False do
45:   revealed[message] ← True;
46:   sample = {ρ ∈ ℰ[message] | replies[message][ρ] ≠ ⊥};
47:   for all  $\pi \in \Pi$  do
48:     trigger  $\langle \text{al.Send} \mid \pi, [\text{Reveal}, \text{message}, \text{sample}] \rangle$ ;
49:   end for
50:
51: upon event  $\langle \text{al.Deliver} \mid \rho, [\text{Reveal}, \text{message}, \text{sample}] \rangle$  do
52:   if  $\Psi(\rho) = \text{False}$  then
53:     reveal[message] ← True;
54:   end if
55:

```

---

Upon `pb.Delivering` a message  $m$ , a correct process sends  $C$  Echo messages to each other process, one `Echo( $s, m$ )` message for every  $s \in \mathcal{M}$ . In other words, the correct behavior is to echo  $m$  across all contexts  $s \in \mathcal{M}$ . A Byzantine process, however, can in principle send to the same process a set of Echo messages echoing different messages in different contexts (e.g, `Echo( $s, m$ )` and `Echo( $s', m' \neq m$ )`).

- When a correct process  $\pi$  collects at least  $\hat{E}$  `Echo( $m, m$ )` messages from the processes in  $\mathcal{E}[m]$ ,  $\pi$  delivers  $m$ .

**Mimic.** Algorithm 4 presents two utility procedures for manipulating samples with respect to their Byzantine component:

- (*correct*, line 1) Procedure *correct* returns a correct process, picked with uniform probability. It does so by invoking  $\Omega$  to select a process  $\rho$  with uniform probability (line 3), then using  $\Psi$  to pick again if  $\rho$  is Byzantine (line 4).
- (*mimic*, line 7) Provided with a sample *reference*, procedure *mimic* returns a sample  $\psi$  that shares with *reference* all Byzantine processes. It does so by looping over each process  $\rho$  in *reference*. If  $\rho$  is Byzantine (line 11),  $\rho$  is added to  $\psi$ . If  $\rho$  is correct (line 13), procedure *correct*() is used to add a random correct process to  $\psi$ .

**Samples.** Upon initialization (line 12), a correct process initializes  $C$  echo samples  $\mathcal{E}[1..C]$  that share the same set of Byzantine processes. It does so by using procedure *sample*(...) to randomly pick  $\mathcal{E}[1]$  (line 17), then using *mimic*( $\mathcal{E}[1]$ ) to pick samples 2 to  $C$  (line 20).

We underline how  $\mathcal{E}[1]$  is selected using the *sample* procedure we defined in Algorithm 2. As a result, upon initialization, a correct process sends an `EchoSubscribe` message to each process in  $\mathcal{E}[1]$ . However, a correct process does not handle the `al.Delivery` of an `EchoSubscribe` message. This is done on purpose. The only goal of those `EchoSubscribe` messages is to let



the Byzantine adversary know which Byzantine processes are in  $\mathcal{E}[1]$  (and, consequently, in every other sample).

**Broadcast.** Upon `cob.Broadcast`ing a message *message* (line 23), the correct designated sender uses `pb.Broadcast` to distribute *message*.

**Echo.** When a correct process `pb.Delivers` a message *message* (line 26), it sends to each process  $\rho$  an `Echo(sample, message)` message, for every *sample* in  $\mathcal{M}$  (line 29). In other words, the correct behavior of a correct process that `pb.Delivered message` is to echo *message* across all samples.

We note how Simplified Sieve does not make use of echo subscription sets. A correct process sends its Echo messages to every process in the system. The goal of Simplified Sieve, indeed, is not performance, but probabilistic tractability.

**Delivery.** A correct process maintains a table *replies* to keep track of the Echo messages received by each node in its echo samples. Upon receiving an `Echo(sample, message)` message from a process  $\rho$  for the first time (line 33), if  $\rho$  is in  $\mathcal{E}[sample]$ , a correct process sets *replies[sample][ $\rho$ ]* to *message* (line 35).

Upon receiving at least  $\hat{E}$  `Echo(message, message)` messages from the processes in  $\mathcal{E}[message]$  (line 39) (this is checked using the *replies* table), a correct process `cob.Delivers message` (line 42).

**Reveal.** A correct process maintains a *reveal* array to keep track of which echo samples it should *reveal*. When, for some *message*, *reveal[message]* = True (line 44), a correct process sends to every process a `Reveal` message, containing the set of processes in  $\mathcal{E}[message]$  that issued an `Echo(message, message')` message for some *message'*  $\in \mathcal{M}$  (line 48). In other words, whenever *reveal[message]* = True, a correct process reveals the set of processes in its echo sample for *message* that issued an Echo message for that sample.

If, after revealing its sample for *message*, a correct process receives additional Echo messages from the processes in  $\mathcal{E}[message]$ , the reveal procedure is performed again. This is enforced by setting a *revealed* flag back to `False` (line 36) every time a new Echo message is received.

A correct process sets *reveal[message]* to True under two circumstances: when it `cob.Delivers message` (line 41) and when it receives a `Reveal` message for *message* from a correct process (line 53). As a result, whenever any correct process delivers *message*, every correct process reveals its sample for *message*, regardless of whether or not it delivered *message*.

Like `EchoSubscribe`, the `Reveal` message serves the only purpose to provide information to

the Byzantine adversary.

### 3.7 Adversarial execution

In this section, we define the model underlying an adversarial execution of Sieve and Simplified Sieve, and identify the set of Byzantine adversaries for either algorithm. Here, a Byzantine adversary is an agent that acts upon a system with the goal of compromising its consistency. Throughout the rest of this chapter, we use the term **pcb adversary** to denote a Byzantine adversary for Sieve, and the term **cob adversary** (or just **adversary**) to denote a Byzantine adversary for Simplified Sieve.

The main goal of this section is to formalize the information available both to the pcb and the cob adversary, and the set of actions that they can perform on the system throughout an adversarial execution of either algorithm.

Throughout the rest of this chapter, we bound the probability of compromising the consistency of Sieve by assuming that, if the totality of pb is compromised, then the consistency of pcb is compromised as well. In what follows, therefore, we assume that pb satisfies totality.

#### 3.7.1 Model (Sieve)

Let  $\pi$  be any correct process. We make the following assumptions about an adversarial execution of Sieve:

- As we established in Section 1.3, the pcb adversary does not know which correct processes are in  $\pi$ 's echo sample. The pcb adversary knows, however, which Byzantine processes are in  $\pi$ 's echo sample.
- At any time, the pcb adversary knows if  $\pi$  delivered a message. If  $\pi$  delivered a message, then the pcb adversary knows which message did  $\pi$  deliver.
- The pcb adversary can cause  $\pi$  to pb.Deliver any message. As we established with Theorem 1,  $\pi$  will, however, pb.Deliver only one message throughout an execution of Sieve.

Throughout an adversarial execution of Sieve, an adversary performs a sequence of minimal operations on the system. Each operation consists of either of the following:

- Selecting a correct process that did not pb.Deliver any message, and causing it to pb.Deliver a message.
- Selecting a Byzantine process and causing it to send an Echo message to a correct process.

As a result of each operation, zero or more correct processes deliver a message. The pcb adversary is successful if, at the end of the adversarial execution, at least two different messages are delivered by at least one correct process.

#### 3.7.2 Model (Simplified Sieve)

Let  $\pi$  be any correct process. We make the following assumptions about an adversarial execution of Simplified Sieve:

- As we established in Section 1.3, the cob adversary does not know which correct processes are in  $\pi$ 's echo samples. The cob adversary knows, however, which Byzantine processes are in  $\pi$ 's echo samples.
- At any time, the cob adversary knows if  $\pi$  delivered a message. If  $\pi$  delivered a message, then the cob adversary knows which message did  $\pi$  deliver. Moreover, if  $\pi$  delivered a message  $m$ , then at any time the cob adversary also knows the processes in  $\pi$ 's echo sample for  $m$  that sent an  $\text{Echo}(m, m')$  message to  $\pi$ , for some message  $m'$ .
- The cob adversary can cause  $\pi$  to pb.Deliver any message. As we established with Theorem 1,  $\pi$  will, however, pb.Deliver only one message throughout an execution of Simplified Sieve.

Throughout an adversarial execution of Simplified Sieve, an adversary performs a sequence of minimal operations on the system. Each operation consists of either of the following:

- Selecting a correct process that did not pb.Deliver any message, and causing it to pb.Deliver a message.
- Selecting a Byzantine process and causing it to send an Echo message to a correct process.

As a result of each operation, zero or more correct processes deliver a message. The cob adversary is successful if, at the end of the adversarial execution, at least two different messages are delivered by at least one correct process.

#### 3.7.3 Network scheduling

In this section, we discuss the behavior of the adversary in relation to network scheduling. As we discussed in Section 1.3, the system is asynchronous, i.e., every message is eventually delivered but can be delayed by an arbitrary, finite amount of time.

## Chapter 3. Sieve

---

**Gossip messages.** As we stated in Section 3.7, throughout this section we assume that the pb instance used by Sieve and Simplified Sieve satisfies totality. While this means that the adversary cannot prevent any correct process from eventually pb.Delivering a message, the adversary can indeed arbitrarily choose which correct process pb.Delivers which message.

This can be achieved by delaying the delivery of the Gossip messages issued by correct processes. Noting that a correct process will accept a Gossip message from any source, the adversary can then cause any of the processes it controls to quickly send a Gossip message with arbitrary content to any correct process, effectively causing it to pb.Deliver an arbitrary message.

**Echo messages.** As we stated in Sections 3.7.1 and 3.7.2, the two minimal operations a (pcb or cob) adversary can perform essentially reduce to causing a Byzantine process to either send a Gossip or an Echo message to a correct process. We can see that those operations are indeed minimal: a correct process atomically al.Delivers a message (i.e., a message is the minimal amount of information that can be meaningfully transferred on the network), and a correct process will ignore any message that is not a Gossip or an Echo message.

Upon pb.Delivering a message, a correct process will issue zero or more Echo messages. As we discussed in Section 1.3, the adversary can arbitrarily delay those messages, but they will eventually be delivered. As a result, the outcome of an adversarial execution is solely determined by the sequence of operations performed by the adversary, and is not affected by network scheduling.

While the adversary could delay the delivery of Echo messages issued by correct processes, the only effect this would have is to prevent the adversary from knowing the effect of an operation on the system before performing the next one. An optimal adversary, therefore, performs an operation, then waits until all the Echo messages issued by correct processes are delivered before performing the next operation.

### 3.7.4 Interfaces

In Sections 3.7.1 and 3.7.2, we defined the model underlying an adversarial execution of Sieve and Simplified Sieve respectively. In Section 3.7.3, we discussed the behavior of the Byzantine adversary in relation to network scheduling. Throughout the rest of this chapter, we concretely model a (pcb or cob) adversary as an *algorithm* that interacts with a *system*.

As we discussed, a (pcb or cob) adversary works in steps: at every step, the adversary either performs one operation on the system, or queries the system for information about its state. In this section, we model this interaction by defining four interfaces, respectively implemented by the (pcb or cob) adversary and the (pcb or cob) system.

Both the **pcb adversary** and the **cob adversary** interfaces (instance *adv*) expose the following **procedures**:

- *Init()*: It is called once, at the beginning of the adversarial execution, before any operation is performed on the system. Here the (pcb or cob) adversary setups its internal state.
- *Step()*: It is called repeatedly, until the adversarial execution is completed. Here the (pcb or cob) adversary performs one operation on the system. The execution fails (e.g., an exception is raised) if a call to *adv.Step()* does not result in one, and only one, call to *sys.Deliver(...)*, *sys.Echo(...)* or *sys.End()* (as we define them below).

The **pcb system** interface (instance *sys*) exposes the following **procedures**:

- *Byzantine(process ∈ Π<sub>C</sub>)*: Returns a list of all the Byzantine processes in *process*' echo sample. The pcb adversary can invoke this procedure an unlimited number of times both from the *Init()* and the *Step()* procedure.
- *State()*: Returns a list of pairs ( $\pi \in \Pi_C, m \in \mathcal{M}$ ), representing which correct process currently delivered which message. The pcb adversary can invoke this procedure an unlimited number of times from the *Step()* procedure.
- *Deliver(process ∈ Π<sub>C</sub>, message ∈ ℳ)*: Causes *process* to pb.Deliver *message*. The execution fails if *Deliver* is provided with the same *process* argument more than once: a correct process does not pb.Deliver more than one message. The procedure does not return any value.
- *Echo(process ∈ Π<sub>C</sub>, source ∈ Π \ Π<sub>C</sub>, message ∈ ℳ)*: Causes *source* to send an *Echo(message)* message to *process*. The execution fails if *Echo* is provided with the same *process* and *source* arguments more than once: a correct process does not consider more than one *Echo* message from the same source. The procedure does not return any value.
- *End()*: Causes the execution to gracefully terminate. The execution fails if *End()* is called before *Deliver(...)* is invoked exactly *C* times: under the assumption that pb satisfies totality, every correct process eventually pb.Delivers a message. The procedure does not return any value.

The **cob system** interface (instance *sys*) exposes the following **procedures**:

- *Byzantine(process ∈ Π<sub>C</sub>)*: Returns a list of all the Byzantine processes in the first echo sample of *process*. The cob adversary can invoke this procedure an unlimited number of times both from the *Init()* and the *Step()* procedure.

- *State()*: Returns a list of pairs  $(\pi \in \Pi_C, m \in \mathcal{M})$ , representing which correct process currently delivered which message. The cob adversary can invoke this procedure an unlimited number of times from the *Step()* procedure.
- *Sample*( $process \in \Pi_C, message \in \mathcal{M}$ ): Returns the processes that are in the echo sample for message *message* of process *process* and that sent an *Echo*(*message*, *message'*) to *process*, for some message *message'*. The cob adversary can invoke this procedure an unlimited number of times from the *Step()* procedure. The execution fails if no correct process has *cob.Delivered message*: a correct process does not *reveal* its echo sample for *message* before *message* is delivered by at least one correct process.
- *Deliver*( $process \in \Pi_C, message \in \mathcal{M}$ ): Causes *process* to *pb.Deliver message*. The execution fails if *Deliver* is provided with the same *process* argument more than once: a correct process does not *pb.Deliver* more than one message. The procedure does not return any value.
- *Echo*( $process \in \Pi_C, sample \in \mathcal{M}, source \in \Pi \setminus \Pi_C, message \in \mathcal{M}$ ): Causes *source* to send an *Echo*(*sample*, *message*) message to *process*. The execution fails if, throughout an execution, *Echo* is provided with the same *process*, *sample* and *source* arguments more than once: a correct process does not consider more than one *Echo* message for the same sample from the same source. The procedure does not return any value.
- *End()*: Causes the execution to gracefully terminate. The execution fails if *End()* is called before *Deliver*(...) is invoked exactly  $C$  times: under the assumption that *pb* satisfies totality, every correct process eventually *pb.Delivers* a message. The procedure does not return any value.

### 3.8 Simplified adversarial power

In this section, we prove that an optimal consistency-only broadcast adversary is more powerful than an optimal probabilistic consistent broadcast adversary. This result is intuitive: a correct process in Simplified Sieve can deliver more than one message, and in general more information is available to the cob adversary than to the pcb adversary.

#### 3.8.1 Preliminary definitions

Before proving that an optimal cob adversary is more powerful than an optimal pcb adversary, we provide some definitions on pcb and cob systems and adversaries.

**Definition 3** (Pcb system). A **pcb system**  $\sigma$  is an element of the set

$$\begin{aligned} \mathcal{S}_{pcb} &= \mathcal{E}_{pcb}^C \\ \mathcal{E}_{pcb} &= \Pi^E \end{aligned}$$

Intuitively, a system  $\sigma \in \mathcal{S}_{pcb}$  is defined by the echo sample of each of its  $C$  correct processes. The echo sample of a correct process is a vector of  $E$  processes.

Let  $\sigma \in \mathcal{S}_{pcb}$ . We use  $\sigma[\pi \in \Pi_C][i \in 1..E]$  to denote the  $i$ -th process in  $\pi$ 's echo sample.

**Definition 4** (Cob system). A **cob system**  $\sigma$  is an element of the set

$$\begin{aligned} \mathcal{S}_{cob} &= \mathcal{E}_{cob}^C \\ \mathcal{E}_{cob} &= \{(e_1, \dots, e_C \in \Pi^E) \mid \mathcal{M}(e_i, e_j) \ \forall i, j \in 1..C\} \\ \mathcal{M}(e, e') &: \forall k, (e_k \in \Pi \setminus \Pi_C) \implies (e'_k = e_k) \end{aligned}$$

Intuitively, a system  $\sigma \in \mathcal{S}_{cob}$  is defined by the echo samples of each of its  $C$  correct processes. Each correct process has  $C$  echo samples  $e_1, \dots, e_C$  (one per message), each represented by a vector of  $E$  processes. Any two echo samples  $e_i, e_j$  of a given process satisfy  $\mathcal{M}(e_i, e_j)$ , i.e., they share the same set of Byzantine processes.

We also use just  $\mathcal{S}$  to denote the set of cob systems  $\mathcal{S}_{cob}$ . Let  $\sigma \in \mathcal{S}_{cob}$ , we use  $\sigma[\pi \in \Pi_C][m \in \mathcal{M}][i \in 1..E]$  to denote the  $i$ -th process in  $\pi$ 's echo sample for  $m$ .

**Definition 5** (Adversary). A **pcb adversary (cob adversary)** is a terminating algorithm that exposes the pcb adversary (cob adversary) interface and does not cause the adversarial execution to fail (see Section 3.7.4) when coupled with any system  $\sigma \in \mathcal{S}_{pcb}$  ( $\sigma \in \mathcal{S}_{cob}$ ).

Let  $\alpha, \alpha'$  be two pcb (cob) adversaries such that, for every  $\sigma \in \mathcal{S}_{pcb}$  ( $\sigma \in \mathcal{S}_{cob}$ ), the execution of  $\alpha$  coupled with  $\sigma$  is identical to the execution of  $\alpha'$  coupled with  $\sigma$ . We consider  $\alpha$  and  $\alpha'$  to be functionally the same adversary.

We use  $\mathcal{A}_{pcb}$  to denote the set of pcb adversaries. We use  $\mathcal{A}_{cob}$  (or just  $\mathcal{A}$ ) to denote the set of cob adversaries.

**Definition 6** (Adversarial power). Let  $\alpha$  be a pcb (cob) adversary. The **adversarial power** of  $\alpha$  is the probability of  $\alpha$  compromising the consistency of a pcb (cob) system, picked with uniform probability from  $\mathcal{S}_{pcb}$  ( $\mathcal{S}_{cob}$ ).

**Definition 7** (Optimal adversary). Let  $\alpha$  be a pcb (cob) adversary. We say that  $\alpha$  is an **optimal adversary** if its adversarial power is greater or equal to that of any other pcb (cob) adversary.

We note that Definition 7 is well defined: indeed, both  $\mathcal{A}_{pcb}$  and  $\mathcal{A}_{cob}$  are finite sets, and therefore admit a maximum for the adversarial power.

**Definition 8** (Optimal set of adversaries). Let  $\mathcal{A}'$  be a set of pcb (cob) adversaries. We say that  $\mathcal{A}'$  is an **optimal set of adversaries** if  $\mathcal{A}'$  includes an optimal pcb (cob) adversary.

**Definition 9** (Pcb invocation/response pair). The pair  $(i, r)$  is a **pcb invocation/response pair**

if

$$\begin{aligned}
 i = (\text{Byzantine}, \pi \in \Pi_C) & & r = (\xi_1, \dots, \xi_k \in \Pi \setminus \Pi_C) \\
 i = (\text{State}) & & r = ((\pi_1, m_1), \dots, \\
 & & (\pi_k \in \Pi_C, m_k \in \mathcal{M})) \\
 i = (\text{Deliver}, \pi \in \Pi_C, m \in \mathcal{M}) & & r = \perp \\
 i = (\text{Echo}, \pi \in \Pi_C, \xi \in \Pi \setminus \Pi_C, m \in \mathcal{M}) & & r = \perp
 \end{aligned}$$

**Definition 10** (Cob invocation/response pair). The pair  $(i, r)$  is a **cob invocation/response pair** if

$$\begin{aligned}
 i = (\text{Byzantine}, \pi \in \Pi_C) & & r = (\xi_1, \dots, \xi_k \in \Pi \setminus \Pi_C) \\
 i = (\text{Sample}, \pi \in \Pi_C, m \in \mathcal{M}) & & r = (\rho_1, \dots, \rho_k \in \Pi) \\
 i = (\text{State}) & & r = ((\pi_1, m_1), \dots, \\
 & & (\pi_k \in \Pi_C, m_k \in \mathcal{M})) \\
 i = (\text{Deliver}, \pi \in \Pi_C, m \in \mathcal{M}) & & r = \perp \\
 i = (\text{Echo}, \pi \in \Pi_C, s \in \mathcal{M}, \xi \in \Pi \setminus \Pi_C, m \in \mathcal{M}) & & r = \perp
 \end{aligned}$$

**Definition 11** (Trace). A **pcb trace (cob trace)** is a finite sequence of pcb (cob) invocation/response pairs. Let  $\alpha$  be a (pcb or cob) adversary, let  $\sigma$  be a (pcb or cob, correspondingly) system. We use  $\tau(\alpha, \sigma)$  to denote the trace produced by  $\alpha$  coupled with  $\sigma$ . We use  $\mathcal{T}$  to denote the set of traces.

**Notation 5** (Power set). Let  $X$  be a set. We use  $\mathbb{P}(X)$  to denote the **power set** of  $X$ . We use  $\mathbb{P}^K(X) = \{x \in \mathbb{P}(X) \mid |x| = K\}$  to denote the elements in  $\mathbb{P}(X)$  that have  $K$  elements. We use  $\mathbb{P}^{K+}(X) = \{x \in \mathbb{P}(X) \mid |x| \geq K\}$  to denote the elements in  $\mathbb{P}(X)$  that have at least  $K$  elements.

### 3.8.2 Consistency of Simplified Sieve

We can now prove that the  $\epsilon$ -consistency of Simplified Sieve is strictly weaker than that of Sieve.

**Lemma 12.** *An optimal cob adversary is more powerful than an optimal pcb adversary.*

*Proof.* Let  $\alpha^*$  be an optimal pcb adversary. In order to prove that an optimal cob adversary is more powerful than  $\alpha^*$ , we just need to find a cob adversary  $\alpha^+$  that is more powerful than  $\alpha^*$ . We achieve this using a **pcb-to-cob decorator**, i.e., an algorithm that acts as an interface between a pcb adversary and cob system. A pcb adversary coupled with a pcb-to-cob decorator effectively implements a cob adversary. Here we show that a pcb-to-cob decorator  $\Delta_{cob}$  exists such that, for every  $\alpha \in \mathcal{A}_{pcb}$ , the cob adversary  $\alpha' = \Delta_{cob}(\alpha)$  is more powerful than  $\alpha$ . If this is true, the lemma is proved: indeed,  $\alpha^+ = \Delta_{cob}(\alpha^*)$  is more powerful than  $\alpha^*$ .



---

**Algorithm 6** Cob decorator.

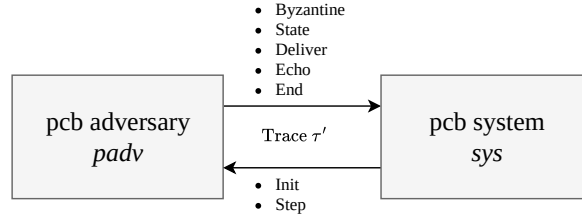
---

```

1: Implements:
2:   CobAdversary + PcbSystem, instance cadv
3:
4: Uses:
5:   PcbAdversary, instance padv, system cadv
6:   CobSystem, instance sys
7:
8: procedure cadv.Init() is
9:   deliveries =  $\{\perp\}^C$ ;
10:  padv.Init();
11:
12: procedure cadv.Step() is
13:  padv.Step();
14:
15: procedure cadv.Byzantine(process) is
16:  return sys.Byzantine(process);
17:
18: procedure cadv.State() is
19:  state =  $\emptyset$ ;
20:
21:  for all  $(\pi, m) \in \text{sys.State}()$  do
22:    if deliveries[ $\pi$ ] = m then
23:      state  $\leftarrow$  state  $\cup$   $\{(\pi, m)\}$ ;
24:    end if
25:  end for
26:
27:  return state;
28:
29: procedure cadv.Deliver(process, message) is
30:  deliveries[process]  $\leftarrow$  message;
31:  sys.Deliver(process, message);
32:
33: procedure cadv.Echo(process, source, message) is
34:  sys.Echo(process, message, source, message);
35:
36: procedure cadv.End() is
37:  sys.End();
38:

```

---



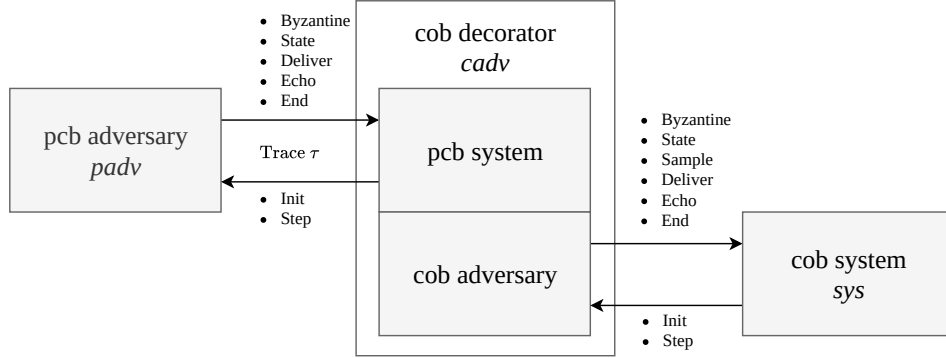
**Figure 3.1: An execution without decorator.**

**Decorator.** Algorithm 6 implements `Cob_decorator`, a `pcb-to-cob` decorator. Provided with a `pcb adversary` *padv*, `Cob_decorator` acts as an interface between *padv* and a `cob system` *sys*, effectively implementing a `cob adversary` *cadv*. `Cob_decorator` exposes both the `cob adversary` and the `pcb system` interfaces: the underlying `pcb adversary` *padv* uses *cadv* as its system.

`Cob_decorator` works as follows:

- Procedure *cadv.Init()* initializes a *deliveries* array that is used to keep track of the message `pb.Delivered` by each correct process, and a *gap* set that it uses to keep track of the messages `cob.Delivered` by each correct process in *sys*.
- Procedure *cadv.Step()* simply forwards the call to *padv.Step()*.
- Procedure *cadv.State()* returns a list of pairs  $(\pi \in \Pi_C, m \in \mathcal{M})$  such that  $\pi$  both `pb.Delivered` and delivered *m* in *sys*. This is achieved by querying *sys.State()*, then looping over each element  $(\pi, m)$  of the response and checking if *deliveries*[ $\pi$ ] = *m*.
- Procedure *cadv.Byzantine(process)* simply forwards the call to *sys.Byzantine(process)*.
- Procedure *cadv.Deliver(process, message)* sets *deliveries*[*process*] to *message* (to signify that *process* `pb.Delivered message`). It then forwards the call to *sys.Deliver(process, message)*, causing *process* to `pb.Deliver message`.
- Procedure *cadv.Echo(process, source, message)* forwards the call to *sys.Echo(process, message, source, message)*, causing *source* to send an *Echo(message, message)* message to *process*.
- Procedure *cadv.End()* simply forwards the call to *sys.End()*.

Let  $\Delta_{cob} : \mathcal{A}_{pcb} \rightarrow \mathcal{A}_{cob}$  denote the function that `Cob_decorator` implements, mapping `pcb adversaries` into `cob adversaries`. We want to prove that, for every  $\alpha \in \mathcal{A}_{pcb}$ , the adversarial power of  $\alpha' = \Delta_{cob}(\alpha)$  is greater than that of  $\alpha$ .



**Figure 3.2: A decorator exposing both its system interface to the pcb adversary and its adversary interface to the cob system.**

**System translation.** Let  $\alpha$  be a pcb adversary. We start by noting that, since  $\alpha$  is correct,  $\alpha$  always causes every correct process to pb.Deliver a message. We can therefore define a function

$$\mu : \mathcal{A}_{pcb} \times \mathcal{S}_{pcb} \times \Pi_C \rightarrow \mathcal{M}$$

such that  $\mu(\alpha, \sigma, \pi) = m$  if and only if  $\alpha$  eventually causes  $\pi$  to pb.Deliver  $m$ , when  $\alpha$  is coupled with  $\sigma$ .

We then define a **system translation** function  $\Psi[\alpha] : \mathcal{S}_{pcb} \rightarrow \mathbb{P}(\mathcal{S}_{cob})$  that maps a pcb system into a set of cob systems:

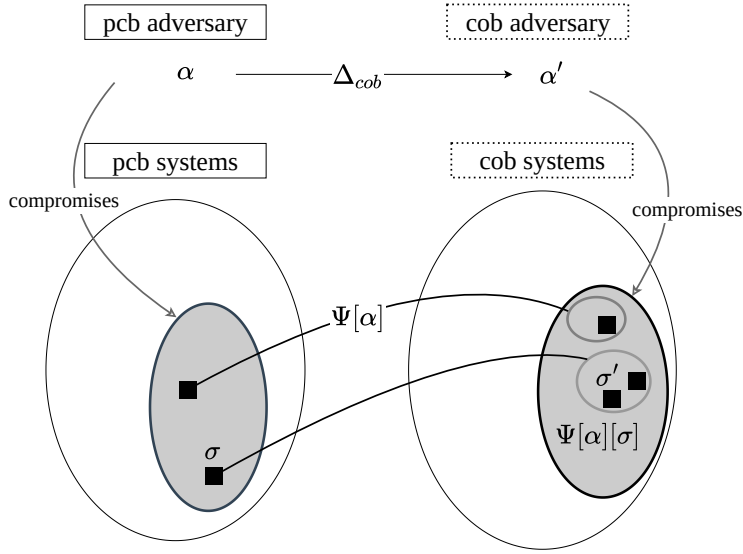
$$(\sigma' \in \Psi[\alpha](\sigma)) \iff (\forall \pi \in \Pi_C, \sigma[\pi] = \sigma'[\pi][\mu(\alpha, \sigma, \pi)])$$

Let  $\sigma$  be a pcb system, let  $\sigma'$  be a cob system, let  $\pi$  be any correct process, let  $m$  be the message that  $\alpha$  eventually causes  $\pi$  to pb.Deliver, when  $\alpha$  is coupled with  $\sigma$ . Intuitively,  $\sigma'$  is in  $\Psi[\alpha](\sigma)$  if  $\pi$ 's echo sample for  $m$  in  $\sigma'$  is identical to  $\pi$ 's echo sample in  $\sigma$ .

**Roadmap.** Let  $\alpha \in \mathcal{A}_{pcb}$ ,  $\alpha' = \Delta_{cob}(\alpha)$ . Let  $\sigma \in \mathcal{S}_{pcb}$  such that  $\alpha$  compromises the consistency of  $\sigma$ . In order to prove that  $\alpha'$  is more powerful than  $\alpha$ , we prove that:

- For every  $\sigma' \in \Psi[\alpha](\sigma)$ ,  $\alpha'$  compromises the consistency of  $\sigma'$ .
- The probability of  $\Psi[\alpha](\sigma)$  is equal to the probability of  $\sigma$ .
- For every  $\hat{\sigma} \in \mathcal{S}_{pcb}$  such that  $\hat{\sigma} \neq \sigma$ , the sets  $\Psi[\alpha](\sigma)$  and  $\Psi[\alpha](\hat{\sigma})$  are disjoint.

Indeed, if all of the above are true, then the probability of  $\alpha'$  compromising the consistency of a random cob system  $\sigma'$  is greater or equal to the probability of  $\alpha$  compromising the consistency of a random system  $\sigma$ , and the lemma is proved.



**Figure 3.3:** An illustration of the steps needed to prove that the adversarial power of  $\alpha$  is greater than that of  $\alpha'$ .

**Trace.** We start by noting that, if we couple Cob decorator with  $\sigma'$ , we effectively obtain a pcb system interface  $\delta$  with which  $\alpha$  directly exchanges invocations and responses. Here we show that the trace  $\tau(\alpha, \sigma)$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\sigma$ , or it has been coupled with  $\sigma'$ , with Cob decorator acting as an interface. We prove this by induction.

Let us assume

$$\begin{aligned} \tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \delta) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . By hypothesis, at least one of the echo samples of  $\pi$  in  $\sigma'$  is identical to the echo sample of  $\pi$  in  $\sigma$ . Moreover, all  $\pi$ 's echo samples in  $\sigma'$  share the same set of Byzantine processes. Therefore, the first of  $\pi$ 's echo samples in  $\sigma'$  contains the same Byzantine processes as  $\pi$ 's echo sample in  $\sigma$ . Finally, the decorator simply forwards the call to  $cadv.\text{Byzantine}(\pi)$  to  $sys.\text{Byzantine}(\pi)$ . Consequently,  $r_{n+1} = r'_{n+1}$ .

Before considering the case  $i_{n+1} = (\text{State})$ , we prove some auxiliary results. Let  $\pi$  be a correct process, let  $\rho$  be a process, let  $\xi$  be a Byzantine process, let  $m$  be a message. For every  $j \leq n+1$ , as we established, we have  $i_j = i'_j$ . Therefore, after the  $(n+1)$ -th invocation, the following hold true:

- $\pi$  pb.Delivered  $m$  in  $\sigma'$  if and only if  $\pi$  pb.Delivered  $m$  in  $\sigma$ . Indeed,  $\text{cadv.Deliver}(\pi, m)$  was invoked if and only if  $\text{sys.Deliver}(\pi, m)$  was invoked as well.
- $\pi$  pb.Delivered  $m$  in  $\sigma'$  if and only if  $\text{deliveries}[\pi] = m$ . Indeed,  $\text{cadv.Deliver}(\pi, m)$  was invoked if and only if  $\text{deliveries}[\pi]$  was set to  $m$ .
- $\xi$  sent an  $\text{Echo}(m)$  to  $\pi$  in  $\sigma'$  if and only if  $\xi$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma$ . Indeed,  $\text{cadv.Echo}(\pi, \xi, m)$  was invoked if and only if  $\text{sys.Echo}(\pi, m, \xi, m)$  was invoked as well.
- If  $\pi$  pb.Delivered  $m$  in  $\sigma$ , then  $\pi$ 's echo sample for  $m$  in  $\sigma'$  is identical to  $\pi$ 's echo sample in  $\sigma$ . This follows from the definition of  $\Psi$  (we recall that  $\sigma' \in \Psi[\alpha](\sigma)$ ).
- If  $\pi$  delivered  $m$  in  $\sigma$ , it also delivered  $m$  in  $\sigma'$ . Indeed, since  $\pi$  pb.Delivered  $m$  in  $\sigma$ ,  $\pi$ 's echo sample for  $m$  in  $\sigma'$  is identical to  $\pi$ 's echo sample in  $\sigma$ . Moreover, if  $\pi$  received an  $\text{Echo}(m)$  message from  $\rho$  in  $\sigma$ , then it also received an  $\text{Echo}(m, m)$  message from  $\rho$  in  $\sigma'$ .
- If  $\pi$  both pb.Delivered and delivered  $m$  in  $\sigma'$ , it also delivered  $m$  in  $\sigma$ . Indeed, since  $\pi$  pb.Delivered  $m$  in  $\sigma'$ , then it also pb.Delivered  $m$  in  $\sigma$ , and  $\pi$ 's echo sample in  $\sigma$  is identical to  $\pi$ 's echo sample for  $m$  in  $\sigma'$ . Moreover, if  $\pi$  received an  $\text{Echo}(m)$  message from  $\rho$  in  $\sigma'$ , then it also received an  $\text{Echo}(m, m)$  message from  $\rho$  in  $\sigma$ .

Let us assume  $i_{n+1} = (\text{State})$ . We start by noting that  $\text{cadv.State}()$  returns all the pairs  $(\pi', m')$  in  $\text{sys.State}()$  that satisfy  $\text{deliveries}[\pi'] = m'$ . If  $(\pi, m) \in r_{n+1}$ , then  $\pi$  both pb.Delivered and delivered  $m$  both in  $\sigma$ . Therefore,  $\pi$  both pb.Delivered and delivered  $m$  in  $\sigma'$ , and  $\text{deliveries}[\pi] = m$ . Consequently,  $(\pi, m) \in r'_{n+1}$ . If  $(\pi, m) \in r'_{n+1}$ , then  $(\pi, m)$  was returned from  $\text{sys.State}()$ , and  $\text{deliveries}[\pi] = m$ . Therefore,  $\pi$  both pb.Delivered and delivered  $m$  in  $\sigma'$ . Consequently,  $\pi$  delivered  $m$  in  $\sigma$ , and  $(\pi, m) \in r_{n+1}$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ .

**Consistency of  $\sigma'$ .** We proved that  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ . Moreover, we proved that if a correct process  $\pi$  eventually pcb.Delivers a message  $m$  in  $\sigma$ , then  $\pi$  also cob.Delivers  $m$  in  $\sigma'$ .

### Chapter 3. Sieve

---

Since  $\alpha$  compromises the consistency of  $\sigma$ , two correct processes  $\pi, \pi'$  and two distinct messages  $m, m' \neq m$  exist such that, in  $\sigma$ ,  $\pi$  pcb.Delivered  $m$  and  $\pi'$  pcb.Delivered  $m'$ . Therefore, in  $\sigma'$ ,  $\pi$  cob.Delivered  $m$  and  $\pi'$  cob.Delivered  $m'$ . Therefore  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Translation probabilities.** We now prove that, for every  $\sigma \in \mathcal{S}_{pcb}$ , the probability of  $\Psi[\alpha](\sigma)$  is equal to the probability of  $\sigma$ .

The probability of  $\sigma$  is

$$\mathcal{P}[\sigma] = \mathcal{P}[\sigma[\pi_1][1] = \pi_{1,1}, \dots, \sigma[\pi_C][E] = \pi_{C,E}] = N^{-EC}$$

and the probability of  $\Psi[\alpha](\sigma)$  is

$$\begin{aligned} \mathcal{P}[\Psi[\alpha](\sigma)] &= \\ \mathcal{P}[\sigma[\pi_1][\mu(\alpha, \sigma, \pi_1)][1] = \pi_{1,1}, \dots, \sigma[\pi_C][\mu(\alpha, \sigma, \pi_C)[E] = \pi_{C,E}] &= N^{-EC} \end{aligned}$$

which proves the result.

**Translation disjunction.** We now prove that, for any two  $\sigma_a, \sigma_b \neq \sigma_a$ , we have  $\Psi[\alpha](\sigma_a) \cap \Psi[\alpha](\sigma_b) = \emptyset$ . We prove this by contradiction. Suppose a system  $\sigma'$  exists such that  $\sigma' \in \Psi[\alpha](\sigma_a)$  and  $\sigma' \in \Psi[\alpha](\sigma_b)$ . We want to prove that  $\sigma_a = \sigma_b$ .

We start by noting that, if  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ , then  $\sigma_a = \sigma_b$ . Indeed, we have

$$\begin{aligned} \tau(\alpha, \sigma_a) &= \tau(\alpha, \sigma_b) \\ \Rightarrow \mu(\alpha, \sigma_a, \pi) &= \mu(\alpha, \sigma_b, \pi) \quad \forall \pi \in \Pi_C \\ \Rightarrow \sigma_a[\pi] &= \sigma'[\pi][\mu(\alpha, \sigma_a, \pi)] \\ &= \sigma'[\pi][\mu(\alpha, \sigma_b, \pi)] \\ &= \sigma_b[\pi] \quad \forall \pi \\ \Rightarrow \sigma_a &= \sigma_b \end{aligned}$$

We prove that  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$  by induction. Let us assume

$$\begin{aligned} \tau(\alpha, \sigma_a) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \sigma_b) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j & \quad \forall j \leq n \end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is

a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . By hypothesis, among the echo samples of  $\pi$  in  $\sigma'_a$ , at least one is identical to the echo sample of  $\pi$  in  $\sigma_a$ , and at least one is identical to the echo sample of  $\pi$  in  $\sigma_b$ . Noting that  $\pi$ 's echo samples share the same set of Byzantine processes, we immediately have that the Byzantine processes in  $\sigma_a[\pi]$  are the same as in  $\sigma_b[\pi]$ , and  $r_{n+1} = r'_{n+1}$ .

Before considering the case  $i_{n+1} = (\text{State})$ , we prove some auxiliary results. Let  $\pi$  be a correct process, let  $\rho$  be a process, let  $\xi$  be a Byzantine process, let  $m$  be a message. For every  $j \leq n+1$ , as we established, we have  $i_j = i'_j$ . Therefore, after the  $(n+1)$ -th invocation, the following hold true:

- $\pi$  pb.Delivered  $m$  in  $\sigma_a$  if and only if  $\pi$  pb.Delivered  $m$  in  $\sigma_b$ .
- $\xi$  sent an Echo( $m$ ) message to  $\pi$  in  $\sigma_a$  if and only if  $\xi$  sent an Echo( $m$ ) message to  $\pi$  in  $\sigma_b$ .
- If  $\pi$  pb.Delivered  $m$  (both in  $\sigma_a$  and  $\sigma_b$ ), then  $\sigma_a[\pi] = \sigma_b[\pi]$ . Indeed,

$$\begin{aligned} \sigma_a[\pi] &= \sigma'[\pi][\mu(\alpha, \sigma_a, \pi)] \\ &= \sigma'[\pi][\mu(\alpha, \sigma_b, \pi)] \\ &= \sigma_b[\pi] \end{aligned}$$

- $\pi$  delivered  $m$  in  $\sigma_a$  if and only if  $\pi$  delivered  $m$  in  $\sigma_b$ . Indeed, if  $\pi$  delivered  $m$  in  $\sigma_a$ , then it also pb.Delivered  $m$  in  $\sigma_a$  and, consequently,  $\sigma_b$ . Therefore,  $\pi$ 's echo sample in  $\sigma_a$  is identical to  $\pi$ 's echo sample in  $\sigma_b$ . Since  $\pi$  received the same Echo messages in  $\sigma_a$  and  $\sigma_b$  then  $\pi$  delivered  $m$  in  $\sigma_b$ . The argument can be trivially reversed to prove that, if  $\pi$  delivered  $m$  in  $\sigma_b$ , then  $\pi$  also delivered  $m$  in  $\sigma_a$ .

Let us consider the case  $i_{n+1} = (\text{State})$ . From the above follows  $r_{n+1} = r'_{n+1}$ .

Noting that procedures *Deliver*(...) and *Echo*(...) never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ .

Therefore,  $\sigma_a = \sigma_b$ , which contradicts the hypothesis and thus proves that the sets  $\Psi[\alpha](\sigma_a)$  and  $\Psi[\alpha](\sigma_b)$  are disjoint.

□

### 3.9 Two-phase adversaries

In Section 3.8 we proved the important result that it is easier to compromise the consistency of Simplified Sieve than that of Sieve. Throughout the rest of this chapter, we compute a bound on the  $\epsilon$ -security of Simplified Sieve.

It is easy to see that the  $\epsilon$ -security of Simplified Sieve is equal to the adversarial power of an optimal adversary. Therefore,  $\epsilon_c$  is a bound on the  $\epsilon$ -security of Simplified Sieve if  $\epsilon_c$  bounds the adversarial power of every adversary in an optimal set of adversaries.

In this section, we derive a set  $\mathcal{A}_{tp} \subseteq \mathcal{A}$  of *two-phase* adversaries that we prove to be optimal. Unlike  $\mathcal{A}$ ,  $\mathcal{A}_{tp}$  is small enough to be probabilistically tractable. In the next sections, we compute a bound on the adversarial power of every  $a \in \mathcal{A}_{tp}$ .

In a similar way to Lemma 12, the proofs of optimality of most of the sets of adversaries presented in this section make extensive use of decorators, and are in general lengthy and non-trivial. For the sake of readability, in this section we only state our results, and defer each explicit proof to Section 3.11.

#### Auto-echo adversary

As we introduced in Section 3.6.3, an Echo message in Simplified Sieve has two fields: a sample  $s$  and a message  $m$ . Intuitively, an  $\text{Echo}(s, m)$  message represents the following statement: “within the context of message  $s$ , consider my Echo to be for message  $m$ ”.

Upon pb.Delivering a message  $m$ , a correct process sends to every other process an  $\text{Echo}(s, m)$  for every  $s$ . In other words, a correct process supports the message it pb.Delivers across all samples. A Byzantine process, however, is not constrained to do this.

A correct process cob.Delivers a message  $m$  upon collecting enough  $\text{Echo}(m, m)$  messages from its echo sample for  $m$ . It is easy to see, therefore, that the probability of a correct process  $\pi$  cob.Delivering  $m$  increases if all the Byzantine processes send an  $\text{Echo}(m, m)$  message to  $\pi$ .

**Definition 12** (Auto-echo adversary). An adversary  $a \in \mathcal{A}$  is an **auto-echo adversary** if, at the beginning of its execution, it causes  $\xi$  to send an  $\text{Echo}(m, m)$  message to  $\pi$ , for every  $\pi \in \Pi_c$ ,  $\xi \in \Pi \setminus \Pi_c$  and  $m \in \mathcal{M}$ . We use  $\mathcal{A}_{ae}$  to denote the set of auto-echo adversaries.

In Section 3.11.1, we formally prove this intuition, i.e., we prove that the set of auto-echo adversaries  $\mathcal{A}_{ae}$  is optimal.

#### Process-sequential adversary

As we discussed in Section 3.6.3, a correct process reveals its sample for a message  $m$  only after delivering  $m$ . At the beginning of the execution, the adversary only knows which Byzantine



processes are in each correct process' echo samples. In Section 3.9, however, we proved that this does not affect the optimal adversary's strategy: the set of Byzantine processes in a correct process' echo samples don't play any role in an optimal adversarial execution.

Intuitively, therefore, an optimal adversary has effectively no meaningful way to distinguish any two correct processes based on the outcome that their actions will have on the system.

**Definition 13** (Correct process enumeration). We define a bijection

$$\zeta : 1..C \leftrightarrow \Pi_C$$

that uniquely maps an integer identifier  $i \in 1..C$  to a correct process.

**Definition 14** (Process-sequential adversary). An auto-echo adversary  $\alpha \in \mathcal{A}_{ae}$  is a **process-sequential adversary** if it never causes  $\zeta(i)$  to pb.Deliver a message before any  $\zeta(j < i)$ . We use  $\mathcal{A}_{ps}$  to denote the set of process-sequential adversaries.

In Section 3.11.2, we formally prove this intuition, i.e., we prove that the set of process-sequential adversaries  $\mathcal{A}_{ps}$  is optimal.

#### Sequential adversary

As we introduced in Section 3.6.3, in Simplified Sieve a correct process independently selects  $C$  echo samples, one for every message in  $\mathcal{M}$ . Moreover, every echo sample shares the same set of Byzantine processes. Finally, let  $\pi$  be a correct process, let  $m$  be a message, no correct process in  $\pi$ 's echo sample for  $m$  is known to the adversary before  $\pi$  delivers  $m$ .

Intuitively, therefore, an adversary has effectively no meaningful way of distinguishing two messages, based on the outcome that their pb.Delivery will have on the system.

**Definition 15** (Poisoned process). Let  $\sigma$  be a system, let  $\pi$  be a correct process. We say that  $\pi$  is **poisoned in**  $\sigma$  if and only if at least  $\hat{E}$  processes in  $\pi$ 's first echo sample in  $\sigma$  are Byzantine.

**Definition 16** (Sequential adversary). A process-sequential adversary  $\alpha \in \mathcal{A}_{ps}$  is a **sequential adversary** if it never causes a correct process to pb.Deliver  $m \in \mathcal{M}$  before causing every  $l < m \in \mathcal{M}$  to be pb.Delivered by at least one correct process. We use  $\mathcal{A}_{sq}$  to denote the set of sequential adversaries.

In Section 3.11.3, we formally prove this intuition, i.e., we prove that the set of sequential adversaries  $\mathcal{A}_{sq}$  is optimal.

#### Non-redundant adversary

As we established in Section 3.6.1, the consistency of consistency-only broadcast is compromised if and only if at least two messages are delivered by at least one correct process.

## Chapter 3. Sieve

---

It is easy to see, therefore, that an adversary that has already caused at least one correct process to deliver a message  $m$  gains no advantage from causing more correct processes to `pb.Deliver`  $m$ . Indeed, doing so would not increase the probability of at least one correct process delivering  $m$  (that condition is verified with probability 1): an optimal adversary should focus its remaining `pb.Deliveries` on achieving the goal to cause at least one other message to be delivered by at least one correct process.

**Definition 17** (Non-redundant adversary). A sequential adversary  $\alpha \in \mathcal{A}_{sq}$  is a **non-redundant adversary** if, whenever exactly one message  $m$  has been delivered, it never causes any additional correct process to `pb.Deliver`  $m$ . We use  $\mathcal{A}_{nr}$  to denote the set of non-redundant adversaries.

In Section 3.11.4, we formally prove this intuition, i.e., we prove that the set of non-redundant adversaries  $\mathcal{A}_{nr}$  is optimal.

### Sample-blind adversary

In Section 3.6.3, we discussed how, in Simplified Sieve, a correct process reveals its echo sample for a message after at least one correct process delivered that message. Throughout Section 3.10.1, we extensively used `Reveal` messages (through the `Sample(...)` system interface) to build a sequence of decorators that improved the power of any adversary in their domain.

In this section, we prove the counter-intuitive result that the information contained in a `Reveal` message is actually useless to an optimal adversary. Indeed, the decorators we developed leveraged `Reveal` messages to *correct* the sub-optimal behavior of a generic adversary. However, for every decorator that we developed, we argue that we could develop an adversary in the codomain of that decorator that never uses the information provided by `Reveal` messages.

An intuitive insight on `Reveal` messages can be provided by the observation that the information they provide is disclosed in the moment it ceases to actually be useful. Indeed, a correct process reveals the content of its echo sample for a message  $m$  only after at least one correct process delivered  $m$ . As we proved in Section 3.9, causing additional processes to deliver  $m$  gives no advantage to the adversary. Moreover, since the correct processes in each echo sample are picked independently from each other, the knowledge of a correct process  $\pi$ 's echo sample for  $m$  does not grant any advantage in causing  $\pi$  to deliver  $m' \neq m$ .

**Notation 6** (Undefined minima and maxima). Let  $X \subset \mathbb{N}$ , with  $X$  finite, let  $S : X \rightarrow \{\text{True}, \text{False}\}$  be a predicate on  $X$ . We use

$$\begin{aligned} (\min n \in X \mid S(n)) &= \perp \\ (\max n \in X \mid S(n)) &= \perp \end{aligned}$$

to denote that

$$\nexists n \in X \mid S(n)$$

**Definition 18** (Trace compatibility). Let  $\tau$  be a trace, let  $\sigma$  be a system. We say that  $\tau$  is **compatible with**  $\sigma$ , or  $\tau \sim \sigma$ , if the sequence of invocations in  $\tau$ , applied in order to  $\sigma$ , produces the corresponding sequence of responses in  $\tau$ .

**Notation 7** (Consistency compromission). Let  $\alpha$  be an adversary, let  $\sigma$  be a system, let  $\tau$  be a trace. We use  $\alpha \searrow \sigma$  to signify that  $\alpha$  compromises the consistency of  $\sigma$ . We use  $\tau \searrow \sigma$  to signify that the sequence of invocations in  $\tau$  compromises the consistency of  $\sigma$ .

**Definition 19** (Sample-blind adversary). A non-redundant adversary  $\alpha \in \mathcal{A}_{nr}$  is a **sample-blind adversary** if it never invokes *Sample(...)*. We use  $\mathcal{A}_{sb}$  to denote the set of sample-blind adversaries.

In Section 3.11.5, we formally prove this intuition, i.e., prove that the set of sample-blind adversaries  $\mathcal{A}_{sb}$  is optimal.

### Byzantine-counting adversary

In Section 3.7.2 we discussed how an adversary for Simplified Sieve knows which Byzantine processes are in the first echo sample of any correct process. In Section 3.9, however, we proved that the optimal adversarial behavior with respect to Echo messages is always to cause every Byzantine process to send an *Echo(m, m)* message to every correct process, for every message  $m \in \mathcal{M}$ .

Intuitively, therefore, a correct process gains no advantage from knowing specifically which Byzantine processes are in the first echo sample of any correct process.

**Definition 20** (Byzantine-counting adversary). A sample-blind adversary  $\alpha \in \mathcal{A}_{sb}$  is a **Byzantine-counting adversary** if, whenever it invokes *Byzantine( $\pi \in \Pi_C$ )*, it invokes  $|Byzantine(\pi)|$ . In other words, the behavior of a Byzantine-counting adversary does not depend on the specific set of Byzantine processes in the first echo sample of any correct process. We use  $\mathcal{A}_{bc}$  to denote the set of Byzantine-counting adversaries.

In Section 3.11.6, we formally prove this intuition, i.e., we prove that the set of Byzantine-counting adversaries  $\mathcal{A}_{bc}$  is optimal.

### Single-response adversary

As we introduced in Section 3.7.2, the goal of a cob adversary is to compromise the consistency of a cob system by causing two distinct messages to be delivered by at least one correct process each. In order to achieve this, it acts upon the system in steps, causing correct processes to pb.Deliver a sequence of messages, until the consistency is compromised.

We distinguish two phases of an adversarial execution.

**Definition 21** (Trace phases). Let  $\alpha$  be an adversary, let  $\sigma$  be a system. We call **first phase of**  $\tau(\alpha, \sigma)$  the sequence  $\tau(\alpha, \sigma)_1, \dots, \tau(\alpha, \sigma)_n$  with  $n$  given by

$$n = \begin{cases} \min_j |S(j)| & \text{iff } \exists j |S(j)| \\ |\tau(\alpha, \sigma)| & \text{otherwise} \end{cases}$$

$$S(j) = (\tau(\alpha, \sigma)_j = (\text{State}, r_h), r_h \neq \emptyset)$$

We call  $\tau(\alpha, \sigma)_{n+1}, \dots, \tau(\alpha, \sigma)_{|\tau(\alpha, \sigma)|}$  the **second phase of**  $\tau(\alpha, \sigma)$ . We call  $\eta(\alpha, \sigma)$  the first phase of  $\tau(\alpha, \sigma)$ . We call  $\theta(\alpha, \sigma)$  the second phase of  $\tau(\alpha, \sigma)$ .

The first phase of a trace ends when, for the first time, a call to  $State()$  returns a non-empty set. Intuitively, the first phase ends when the adversary becomes aware that at least one correct process delivered a message.

Let us focus on the second phase of an adversarial execution carried out by a Byzantine-counting adversary. We know that, at the beginning of the second phase, at least one message has been delivered by at least one correct process. If more than one message has been delivered, the adversary already compromised the consistency of the system, and the invocations in the second phase are irrelevant to its success.

If exactly one message has been delivered, an optimal adversary will issue a sequence of invocations that, given the information available on the system, maximizes the probability of at least one more message being delivered by at least one correct process. Since the adversary is non-redundant, the response provided by any invocation to  $State()$  will not change until the consistency is compromised. Intuitively, therefore, the information available to the adversary throughout the second phase does not change until consistency is compromised. Since any invocation issued by the adversary after consistency is compromised is irrelevant to its success, an optimal adversary does not need to invoke  $State()$  throughout the second phase of any adversarial execution.

**Definition 22** (Single-response adversary). A Byzantine-counting adversary  $\alpha \in \mathcal{A}_{bc}$  is a **single-response adversary** if it never invokes  $State()$  throughout the second phase of any adversarial execution. We use  $\mathcal{A}_{sr}$  to denote the set of single-response adversaries.

In Section 3.11.7, we formally prove this intuition, i.e., we prove that the set of single-response adversaries  $\mathcal{A}_{sr}$  is optimal.

### State-polling adversary

In Section 3.9, we proved that an optimal adversary does not need to invoke  $State()$  in the second phase of an adversarial execution, i.e., after at least one message has been delivered by

at least one correct process.

It is easy to see, however, that, throughout the first phase, the information provided by  $State()$  is useful to the adversary. Intuitively, the sooner a single-response adversary becomes aware that at least one correct process delivered a message, the sooner it can focus its strategy to cause the delivery of a second, distinct message.

In this section, we prove this intuition, i.e., we formally prove that the set of *state-polling adversaries* is optimal.

**Definition 23** (State-polling adversary). A single-response adversary  $\alpha \in \mathcal{A}_{sr}$  is a **state-polling adversary** if it invokes  $State()$  before the first invocation of  $Deliver(\dots)$  and after each invocation of  $Deliver(\dots)$ , until  $State()$  returns a non-empty set. We use  $\mathcal{A}_{sp}$  to denote the set of state-polling adversaries.

**Lemma 13.** *The set of state-polling adversaries  $\mathcal{A}_{sp}$  is optimal.*

*Proof.* It follows immediately from the observation that, for any adversary, not invoking  $State()$  is equivalent to invoking  $State()$  and ignoring its response.  $\square$

#### Two-phase adversary

In Section 3.9, we proved that: throughout the first phase, an optimal adversary invokes  $State()$  before the first invocation of  $Deliver(\dots)$  and after each invocation of  $Deliver(\dots)$ ; throughout the second phase, an optimal adversary never needs to invoke  $State()$ .

As we discussed, if the first phase is concluded with more than one message being delivered, the adversary already compromised the consistency of the system, and the invocations in the second phase are irrelevant to its success.

Let us consider the case where, at the beginning of the second phase, exactly one message  $m^*$  has been delivered. In Section 3.6.3, we discussed how a correct process selects the correct component of each echo sample independently. Intuitively, therefore, the knowledge of which processes delivered  $m^*$  is useless to the adversary, as it provides no information about the correct component of any echo sample for a message  $m \neq m^*$ . In other words, an optimal adversary only needs to know *when* the first phase of the execution is concluded, but not *how*.

**Definition 24** (Two-phase adversary). A state-polling adversary  $\alpha \in \mathcal{A}_{sp}$  is a **two-phase adversary** if, whenever it invokes  $State()$ , it only invokes  $(State() \neq \emptyset)$ . In other words, the behavior of a two-phase adversary does not depend on the content of  $State()$ , but only on whether or not  $State()$  is empty.

**Lemma 14.** *Let  $\alpha$  be a state-polling adversary, let  $\sigma, \sigma'$  be systems such that*

$$|\eta(\alpha, \sigma)| = |\eta(\alpha, \sigma')|$$

### Chapter 3. Sieve

---

and, for all  $\pi \in \Pi_C$ ,  $m \in \mathcal{M}$ ,

$$|\{n \in 1..E \mid \sigma[\pi][m][n] \in \Pi_C\}| = |\{n \in 1..E \mid \sigma'[\pi][m][n] \in \Pi_C\}|$$

We have

$$\forall n < |\eta(\alpha, \sigma)|, \tau(\alpha, \sigma)_n = \tau(\alpha, \sigma')_n$$

*Proof.* The lemma is proved by induction. Let us assume

$$\begin{aligned} \tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \sigma') &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

We start by noting that, since  $\alpha$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . By hypothesis, the number of Byzantine processes in  $\pi$ 's first echo sample is identical in  $\sigma$  and  $\sigma'$ : with a minor abuse of notation we effectively have  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{State})$ . By hypothesis,  $n+1 < |\eta(\alpha, \sigma)| = |\eta(\alpha, \sigma')|$ , and we immediately get  $r_{n+1} = r'_{n+1} = \emptyset$ .

Since  $\alpha$  is a sample-blind adversary, we have  $i_{n+1} \neq (\text{Sample}, \pi, m)$ .

Noting that procedures *Deliver*(...) and *Echo*(...) never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have that, for every  $n < |\eta(\alpha, \sigma)|$ ,  $\tau(\alpha, \sigma)_n = \tau(\alpha, \delta)_n$ .  $\square$

**Lemma 15.** *Let  $\alpha$  be a state-polling adversary, let  $\sigma, \sigma'$  be systems such that*

$$\eta(\alpha, \sigma) = \eta(\alpha, \sigma')$$

and, for all  $\pi \in \Pi_C$ ,  $m \in \mathcal{M}$ ,

$$|\{n \in 1..E \mid \sigma[\pi][m][n] \in \Pi_C\}| = |\{n \in 1..E \mid \sigma'[\pi][m][n] \in \Pi_C\}|$$

We have

$$\tau(\alpha, \sigma) = \tau(\alpha, \sigma')$$

*Proof.* The proof is similar to the proof of Lemma 14, and we omit it for the sake of brevity. The lemma is proved by induction and noting that, since  $\alpha$  is a single-response adversary, it never invokes  $State()$  throughout the second phase of an adversarial execution.  $\square$

In Section 3.11.8, we formally prove that the set of two-phase adversaries  $\mathcal{A}_{tp}$  is optimal.

Before moving on to computing a bound on the adversarial power of  $\mathcal{A}_{tp}$ , we prove two additional lemmas on the behavior of two-phase adversaries.

**Lemma 16.** *Let  $\alpha$  be a two-phase adversary. Let  $\eta^{(i)}(\alpha, \sigma)$  denote the sequence of invocations in  $\eta(\alpha, \sigma)$ . Let  $\sigma, \sigma'$  be systems such that, for all  $\pi \in \Pi_C$ ,*

$$|\sigma.\text{Byzantine}(\pi)| = |\sigma'.\text{Byzantine}(\pi)|$$

*We have*

$$\forall n \leq \min(|\eta(\alpha, \sigma)|, |\eta(\alpha, \sigma')|), \eta^{(i)}(\alpha, \sigma)_n = \eta^{(i)}(\alpha, \sigma')_n$$

*Proof.* The proof is similar to the proof of Lemma 14, and we omit it for the sake of brevity. The lemma is proved by induction and noting that, except for the last one, every response to  $(\text{State})$  in  $\eta(\alpha, \sigma), \eta(\alpha, \sigma')$  is, by definition,  $\emptyset$ .  $\square$

**Lemma 17.** *Let  $\alpha$  be a two-phase adversary. Let  $\sigma, \sigma'$  be systems such that  $|\eta(\alpha, \sigma)| = |\eta(\alpha, \sigma')|$ . Let  $\theta^{(i)}(\alpha, \sigma)$  denote the sequence of invocations in  $\theta(\alpha, \sigma)$  and, for all  $\pi \in \Pi_C$ ,*

$$|\sigma.\text{Byzantine}(\pi)| = |\sigma'.\text{Byzantine}(\pi)|$$

*We have*

$$\theta^{(i)}(\alpha, \sigma) = \theta^{(i)}(\alpha, \sigma')$$

*Proof.* The proof is again similar to the proof of Lemma 14, and we omit it for the sake of brevity. The lemma is proved by induction and noting that:

- Since  $\alpha$  is two-phase, it only invokes  $State() \neq \emptyset$ , the content of the  $|\eta(\alpha, \sigma)|$ -th response does not affect its behavior.
- Since  $\alpha$  is single-response, it never invokes  $State()$  throughout the second phase.

$\square$

### 3.10 Consistency

In this section, we finally achieve the main goal of this chapter, i.e., to compute a bound on the  $\epsilon$ -consistency of Sieve. In order to achieve this, in Section 3.6, we introduced Simplified Sieve, a strawman algorithm designed to be analytically tractable.

In Section 3.8, we proved that the consistency of Simplified Sieve is weaker than the consistency of Sieve. More precisely, we proved that an optimal adversary has a greater probability of compromising the consistency of Simplified Sieve than that of Sieve.

In doing so, we reduced the problem of bounding the  $\epsilon$ -consistency of Sieve to that of bounding the adversarial power of a set of adversaries for Simplified Sieve that provably includes an optimal adversary.

Throughout Section 3.9, we employed a sequence of decorators to iteratively reduce the size of the set that provably includes an optimal adversary. Specifically, we proved that the set  $\mathcal{A}_{1p}$  of two-phase adversaries is optimal. Intuitively, we proved that the behavior of an optimal adversary reduces to:

- (Echo phase): Causing every Byzantine process to send an  $\text{Echo}(m, m)$  message to every correct process, for every message  $m$ .
- (First phase): In sequence, causing correct processes to deliver a predefined sequence of messages until at least one correct process delivers a message.
- (Second phase): In sequence, causing the remaining set of correct process to deliver a predefined sequence of messages, determined only by the number of correct processes that pb.Delivered a message throughout the first phase.

In particular, the only information that we did not prove to be unnecessary to the Byzantine adversary is:

- The number of Byzantine processes in the first echo sample of each correct process  $\pi$ . This information is available to the adversary from the beginning of the adversarial execution, and does not change throughout the execution. We conjecture this information to still be of no use to the adversary, but we don't rely on this conjecture in proving what follows.
- The number of correct processes that pb.Deliver a message throughout the first phase of the adversarial execution, i.e., before at least one correct process delivers a message.

In this section, we redefine a two-phase adversary as a table of messages. In doing so, we provide a sound structure to a set of adversaries that provably includes an optimal one. We



then use this structure to analitically bound the probability of any two-phase adversary compromising the consistency of a random Simplified Sieve system.

First, we focus on the second phase of an adversarial execution, and study the probability of any two-phase adversary compromising the consistency of Simplified Sieve, given the number of correct processes that pb.Delivered, throughout the first phase, the message that was delivered by at least one correct process at the end of the first phase.

We then focus on the first phase of an adversarial execution, and study the probability of any two-phase adversary concluding the first phase of an adversarial execution having caused less than  $n$  correct processes to pb.Deliver  $m$ ,  $m$  being the message that at least one correct process delivers at the end of the first phase.

We finally join the two above results to compute a bound  $\epsilon_c$  on the probability of a two-phase adversary compromising the consistency of Simplified Sieve. Since at least one two-phase adversary is provably optimal, Simplified Sieve satisfies  $\epsilon_c$ -consistency. Since the  $\epsilon$ -consistency of Sieve is provably bound by the  $\epsilon$ -consistency of Simplified Sieve, Sieve satisfies  $\epsilon_c$ -consistency.

#### 3.10.1 Two-phase adversaries

In Section 3.9, we proved that the set  $\mathcal{A}_{tp}$  is optimal. In this section, we use Lemmas 16 and 17 to re-define the set of two-phase adversaries as a set of *triangular message tables*.

**Definition 25** (Byzantine population). A **Byzantine population** is a vector in the set

$$\mathcal{F} = (0..E)^{\Pi_C}$$

Let  $\sigma$  be a system. We define the **Byzantine population** of  $\sigma$  by

$$\forall \pi, F(\sigma)_\pi = |\sigma.Byzantine(\pi)|$$

**Definition 26** (Two-phase adversary). A **two-phase adversary**  $\alpha \in \mathcal{A}_{tp}$  is a *triangular table* defined by:

$$\begin{aligned} \alpha[F]_i &\in \mathcal{M} & F \in \mathcal{F}, i \in 1..C \\ \alpha[F]_i^n &\in \mathcal{M} & F \in \mathcal{F}, n \in 0..C, i \in 1..(C-n) \end{aligned}$$

Coupled with a system  $\sigma$ , a two-phase adversary  $\alpha$ :

- (Echo phase) Causes every Byzantine process to send an Echo( $m$ ,  $m$ ) message to every correct process in  $\sigma$ , for every message  $m$ .
- (First phase) Sequentially causes  $\zeta(1)$  to pb.Deliver  $\alpha[F(\sigma)]_1$ ,  $\zeta(2)$  to pb.Deliver  $\alpha[F(\sigma)]_2$ ,

... in  $\sigma$ , until, as a result of the  $n$ -th pb.Delivery, at least one correct process delivers a message in  $\sigma$ . We note that, if  $\sigma$  is poisoned, then at least one correct process delivers a message in  $\sigma$  as a result of the echo phase, and  $n = 0$ .

- (Second phase) Sequentially causes  $\zeta(n+1)$  to pb.Deliver  $\alpha[F(\sigma)]_1^n$ , ...,  $\zeta(C)$  to pb.Deliver  $\alpha[F(\sigma)]_{C-n}^n$  in  $\sigma$ .

### 3.10.2 Random variables

Let  $\alpha$  be a two-phase adversary. In the next sections, we compute a bound on the probability of  $\alpha$  compromising the consistency of a random, non-poisoned system. To this end, in this section we introduce a set of random variables.

**Notation 8** (Delivery indicator). Let  $\sigma$  be a system, let  $m, m_1, \dots, m_n$  be messages. We use

$$\delta_m[m_1, \dots, m_n](\sigma) \in \{\text{True}, \text{False}\}$$

to indicate whether or not at least one correct process delivers  $m$  in  $\sigma$ , if  $\zeta(1)$  pb.Delivers  $m_1, \dots, \zeta(n)$  pb.Delivers  $m_n$  in  $\sigma$ . We additionally define

$$\delta[m_1, \dots, m_n](\sigma) = \bigvee_{m \in \mathcal{M}} \delta_m[m_1, \dots, m_n](\sigma)$$

Let  $\sigma$  be a random, non-poisoned system. We define:

- **Byzantine population**  $F_{\pi \in \Pi_C}(\sigma)$ : represents the number of Byzantine processes in the first echo sample of  $\pi$  in  $\sigma$ .
- **First phase duration**  $\eta(\sigma)$ : represents the number of correct processes that pb.Deliver a message in the first phase, when  $\alpha$  is coupled with  $\sigma$ . More formally,

$$\eta(\sigma) = \min n \mid (\delta[\alpha[F(\sigma)]_1, \dots, \alpha[F(\sigma)]_n] = \text{True} \vee n = C)$$

- **First-phase deliveries**  $S_{m \in \mathcal{M}}(\sigma)$ : represents the number of correct processes that pb.Deliver message  $m$  throughout the first phase, when  $\alpha$  is coupled with  $\sigma$ . More formally,

$$S_m(\sigma) = \left| \left\{ n \in 1.. \eta(\sigma) \mid \alpha[F(\sigma)]_n = m \right\} \right|$$

- **Second-phase deliveries**  $T_{m \in \mathcal{M}}(\sigma)$ : represents the number of correct processes that pb.Deliver message  $m$  throughout the second phase, when  $\alpha$  is coupled with  $\sigma$ . More formally,

$$T_m(\sigma) = \left| \left\{ n \in 1..(C - \eta(\sigma)) \mid \alpha[F(\sigma)]_n^{\eta(\sigma)} = m \right\} \right|$$

- **Deliveries**  $N_{m \in \mathcal{M}}(\sigma)$ : represents the number of correct processes that pb.Deliver mes-

sage  $m$ , when  $\alpha$  is coupled with  $\sigma$ . More formally,

$$N_m(\sigma) = S_m(\sigma) + T_m(\sigma)$$

- **First delivered message**  $H(\sigma) \in \mathcal{M} \cup \{\perp\}$ : if, when  $\alpha$  is coupled with  $\sigma$ , at least one correct process delivers a message,  $H(\sigma)$  represents the first message to be delivered by at least one correct process in  $\sigma$ . Otherwise,  $H(\sigma) = \perp$ . More formally,

$$H(\sigma) = \begin{cases} \alpha[F(\sigma)]_{\eta(\sigma)} & \text{iff } \delta[\alpha[F(\sigma)]_1, \dots, \alpha[F(\sigma)]_C] = \text{True} \\ \perp & \text{otherwise} \end{cases}$$

- **Correct echoes**  $E_{m \in \mathcal{M}}^{k \in 0..C}[\pi](\sigma) \in 0..E \cup \{\perp\}$ : if  $k \leq N_i(\sigma)$ , then  $E_i^k[\pi](\sigma)$  represents the number of correct processes in  $\pi$ 's echo sample for  $m$  that sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma$ , when exactly  $k$  correct processes pb.Delivered  $m$  in  $\sigma$ . Otherwise,  $E_m^k[\pi](\sigma) = \perp$ .
- **Delivery**  $A_{m \in \mathcal{M}}^{k \in 0..C}[\pi \in \Pi_C](\sigma) \in \{\text{True}, \text{False}, \perp\}$ : if  $k \leq N_m(\sigma)$ ,  $A_m^k[\pi]$  represents, when  $\alpha$  is coupled with  $\sigma$ , whether or not  $\pi$  delivered  $m$  after  $k$  correct processes pb.Delivered  $m$ . More formally,

$$A_m^k[\pi](\sigma) = \begin{cases} E_m^k[\pi] \geq \hat{E} - F_\pi & \text{iff } k \leq N_m(\sigma) \\ \perp & \text{otherwise} \end{cases}$$

- **Global delivery**  $A_{m \in \mathcal{M}}^{k \in 0..C}(\sigma)$ : if  $k \leq N_m(\sigma)$ ,  $A_i^k$  represents, when  $\alpha$  is coupled with  $\sigma$ , whether or not at least one process delivered  $m$  after  $k$  correct processes pb.Delivered  $m$ . More formally,

$$A_m^k(\sigma) = \begin{cases} \bigvee_{\pi \in \Pi_C} A_m^k[\pi](\sigma) & \text{iff } k \leq N_m(\sigma) \\ \perp & \text{otherwise} \end{cases}$$

- **First phase plan**  $L_m(\sigma)$ : represents the number of times  $m$  appears in the sequence

$$\alpha[F(\sigma)]_1, \dots, \alpha[F(\sigma)]_C$$

Intuitively,  $L_m$  represents the number of correct processes that  $\alpha$  would eventually cause to pb.Deliver  $m$ , if no correct process ever delivered any message.

- **Adversarial success**  $W$ :  $W$  represents whether or not the adversary successfully compromises the consistency of the system.

We additionally define:

$$\begin{aligned}
 E_m[\pi](\sigma) &= E_m^{N_m(\sigma)}[\pi](\sigma) \\
 E_m^{(s)}[\pi](\sigma) &= E_m^{S_m(\sigma)}[\pi](\sigma) \\
 E_m^{(t)}[\pi](\sigma) &= E_m[\pi](\sigma) - E_m^{(s)}[\pi](\sigma) \\
 A_m[\pi](\sigma) &= A_m^{N_m(\sigma)}[\pi](\sigma) \\
 A_m(\sigma) &= A^{N_m(\sigma)}(\sigma)
 \end{aligned}$$

### 3.10.3 Byzantine population, correct echoes, delivery

In this section, we compute the probability distributions underlying Byzantine population. Given the Byzantine population, we then compute the number of correct echoes and the probability of delivery.

**Byzantine population.** As we discussed in Section 3.6.3, every correct process selects its first echo sample using the *Sample*(...) procedure, which, in turn, picks each element independently from the set of processes. Therefore, the number of correct processes in the first echo sample of each correct process is independently binomially distributed:

$$\mathcal{D}[\bar{F}_\pi] = \text{Bin}[E, f](\bar{F}_\pi)$$

**Correct echoes.** Let  $\pi$  be a correct process, let  $m$  be a message. If  $\pi$  has  $\bar{F}_\pi$  Byzantine processes in its first echo sample and exactly  $k$  correct processes pb.Delivered  $m$ , then each of the  $E - \bar{F}_\pi$  correct process in  $\pi$ 's echo sample for  $m$  has an independent probability  $k/C$  of having pb.Delivered  $m$ .

Consequently, we have

$$\mathcal{D}[\bar{E}_m^k[\pi] \mid \bar{F}_\pi] = \begin{cases} \text{Bin}\left[E - \bar{F}_\pi, \frac{k}{C}\right](\bar{E}_m^k[\pi]) \mathcal{D}[k \leq N_m \mid \bar{F}_\pi] & \text{iff } \bar{E}_m^k[\pi] \neq \perp \\ \mathcal{D}[k > N_m \mid \bar{F}_\pi] & \text{otherwise} \end{cases}$$

We underline that the above holds true only because the adversary  $\alpha$  is non-redundant. Indeed, since  $\alpha$  knows the first phase duration  $\eta$ , it also knows  $H$  (this immediately follows from  $H = \alpha[F]_\eta$ ). Therefore, if  $\alpha$  was not non-redundant, the value of  $E_m^k[\pi]$  would not necessarily be independent from the event  $k \leq N_m$ .

We can see this with an example. With a minor slip of notation, consider an adversary  $\alpha$  such

that

$$\begin{aligned}\alpha[\{0\}^{\Pi_C}]_1 &= 1 \\ \alpha[\{0\}^{\Pi_C}]_{i>1} &\neq 1 \\ \alpha[\{0\}^{\Pi_C}]_1^1 &= 1\end{aligned}$$

We can immediately see that  $\alpha$  is not non-redundant: if no correct process has any Byzantine process in its echo samples, and at least one correct process delivers 1 as an immediate result of  $\zeta(1)$  pb. Delivering 1,  $\alpha$  causes  $\zeta(2)$  to pb. Deliver 1 again. If  $\eta = 1$ , then  $l \geq 1$  correct process  $\pi_1^*, \dots, \pi_l^*$  exists such that  $\zeta(1)$  appears at least  $\hat{E}$  times in  $\pi_i^*$ 's echo sample for 1. Since a correct process  $\pi$  has a probability  $l/C$  of being among  $\pi_1^*, \dots, \pi_l^*$ , if  $N_1 \geq 2$  the distribution of  $\bar{E}_m^k[\pi]$  becomes

$$\begin{aligned}\mathcal{P}\left[\bar{E}_m^k[\pi] \mid F[\pi] = 0, N_1 \geq 2\right] \\ = \text{Bin}\left[E, \frac{k}{C}\right]\left(\bar{E}_m^k[\pi]\right)\left(\frac{C-l}{C} + \frac{l}{C} \frac{I(\bar{E}_m^k[\pi] \geq \hat{E})}{\sum_{e=\hat{E}}^E \text{Bin}\left[E, \frac{k}{C}\right](e)}\right)\end{aligned}$$

which is clearly not a binomial. Intuitively, if  $\alpha$  was not non-redundant, the value of  $N_m$  may come to depend on whether or not  $m$  was delivered by at least one correct process, which obviously correlates with the value of  $E_m^k[\pi]$ .

Since  $\alpha$  is non-redundant, however, and every correct process picks each echo sample independently, the value of  $E_m^k[\pi]$  is indeed independent from the event  $k \leq N_m$ .

**Delivery.** Noting that a correct process  $\pi$  delivers a message  $m$  if it collects at least  $\hat{E}$  Echo( $m$ ,  $m$ ) messages from its echo sample for  $m$ , we can use the distribution underlying the correct echoes to obtain

$$\mathcal{P}\left[A_m^k[\pi] \mid \bar{F}_\pi\right] = \sum_{\bar{E}_m^k[\pi] = \hat{E} - \bar{F}_\pi}^{E - \bar{F}_\pi} \mathcal{P}\left[\bar{E}_m^k[\pi] \mid \bar{F}_\pi\right]$$

and, using the law of total probability, we get

$$\mathcal{P}\left[A_m^k[\pi]\right] = \sum_{\bar{F}_\pi=0}^E \mathcal{P}\left[A_m^k[\pi] \mid \bar{F}_\pi\right] \mathcal{P}\left[\bar{F}_\pi\right]$$

Finally, since the above holds independently for every process  $\pi$ , we have

$$\mathcal{P}\left[A_m^k\right] = 1 - \prod_{\pi \in \Pi_C} \left(1 - \mathcal{P}\left[A_m^k[\pi]\right]\right) = 1 - \left(1 - \mathcal{P}\left[A_m^k[\zeta(1)]\right]\right)^C$$

### 3.10.4 Second phase

In the previous sections, we computed the probability of any correct process delivering a message  $m$ , given that  $k$  correct processes pb.Delivered  $m$ . In Section 3.10.1, we discussed how an optimal adversarial execution unfolds in two phases: the first takes place before any correct process delivers any message; throughout the second, the goal of the adversary is to cause at least one correct process to deliver one additional message.

In this section, we focus on the second phase. We assume that a message  $H$  has already been delivered by at least one correct process. Given the number of correct processes that pb.Delivered each message throughout the first phase, we compute (where possible) a bound on the probability of any message different from  $H$  being delivered before the end of the adversarial execution, i.e., the probability of the adversary successfully compromising the consistency of the system.

**Correct echoes for a non-delivered message.** Let  $\pi$  be a correct process that has  $\bar{F}$  Byzantine processes in its first echo sample. Let  $m$  be a message such that  $\pi$  does not deliver  $m$  after  $k$  correct processes pb.Delivered  $m$ . Here we use Bayes' theorem to compute the probability distribution underlying the number of correct echoes received by  $\pi$  for  $m$ .

**Notation 9** (Indicator function). We use  $I$  to denote the **indicator function**. Let  $c$  be a predicate, then

$$I(c) = \begin{cases} 1 & \text{iff } c \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

$$\begin{aligned} \mathcal{P}\left[\bar{E}_m^k[\pi] \mid \bar{A}_m^k[\pi], \bar{F}_\pi\right] &= \mathcal{P}\left[\bar{E}_m^k[\pi] \mid E_m^k[\pi] < \hat{E} - \bar{F}_\pi, \bar{F}_\pi\right] \\ &= \frac{\mathcal{P}\left[E_m^k[\pi] < \hat{E} - \bar{F}_\pi \mid \bar{E}_m^k[\pi], \bar{F}_\pi\right] \mathcal{P}\left[\bar{E}_m^k[\pi] \mid \bar{F}_\pi\right]}{\mathcal{P}\left[E_m^k[\pi] < \hat{E} - \bar{F}_\pi \mid \bar{F}_\pi\right]} \\ &= \frac{I(\bar{E}_m^k[\pi] < \hat{E} - \bar{F}_\pi) \mathcal{P}\left[\bar{E}_m^k[\pi] \mid \bar{F}_\pi\right]}{\sum_{e=0}^{\hat{E}-\bar{F}_\pi-1} \mathcal{P}\left[E_m^k[\pi] = e \mid \bar{F}_\pi\right]} \end{aligned}$$

where the numerator of the last term includes an indicator function because any condition  $A < B$ , given  $\bar{A}$  and  $\bar{B}$ , is always satisfied deterministically.

**Conditions.** Let  $\pi$  be a correct process, let  $m$  be a message. Throughout the rest of this section, we compute the probability of  $\pi$  eventually delivering  $m$  under the following conditions:

- $\bar{F}_\pi$  processes in  $\pi$ 's first echo sample are Byzantine.
- $m$  is not the message that is delivered at the end of the first phase, i.e.,  $H \neq m$ .

- $\bar{S}_m$  correct processes pb.Deliver  $m$  throughout the first phase.
- $\bar{T}_m$  correct processes pb.Deliver  $m$  throughout the second phase.

**First phase correct echoes.** Here we compute, under the above conditions, the probability distribution underlying  $E_m^{(s)}[\pi]$ , i.e., the number of correct echoes that  $\pi$  collects for  $m$  throughout the first phase.

Since  $H \neq m$ ,  $\pi$  does not deliver  $m$  throughout the first phase. In other words,  $\pi$  does not deliver  $m$  after  $\bar{S}_m$  correct processes pb.Delivered  $m$ , and we immediately have

$$\mathcal{P}[\bar{E}_m^{(s)}[\pi] | H \neq m, \bar{S}_m, \bar{F}_\pi] = \mathcal{P}\left[E_m^{\bar{S}_m}[\pi] = \bar{E}_m^{(s)}[\pi] | \bar{A}_m^{\bar{S}_m}[\pi], \bar{F}_\pi\right]$$

**Second phase correct echoes.** Here we compute, under the above conditions and given  $\bar{E}_m^{(s)}[\pi]$ , the probability distribution underlying  $\bar{E}_m^{(t)}[\pi]$ , i.e., the number of correct echoes that  $\pi$  collects for  $m$  throughout the second phase.

We start by noting that, out of the  $E$  elements in  $\pi$ 's echo sample for  $m$ :

- $\bar{F}_\pi$  are Byzantine.
- $\bar{E}_m^{(s)}[\pi]$  belong to the set of  $\bar{S}_m$  processes that pb.Delivered  $m$  throughout the first phase.
- $E - \bar{F}_\pi - \bar{E}_m^{(s)}[\pi]$  belong to the set of  $C - \bar{S}_m$  processes that did not pb.Deliver  $m$  throughout the first phase.

Moreover, out of the  $C - \bar{S}_m$  processes that did not pb.Deliver  $m$  throughout the first phase,  $\bar{T}_m$  pb.Delivered  $m$  throughout the second phase. Therefore, each of the processes in  $\pi$ 's echo sample for  $m$  that did not pb.Deliver  $m$  throughout the first phase has an independent probability  $\bar{T}_m / (C - \bar{S}_m)$  of pb.Delivering  $m$  throughout the second phase.

Consequently,  $\bar{E}_m^{(t)}$  is binomially distributed:

$$\mathcal{P}[\bar{E}_m^{(t)} | \bar{E}_m^{(s)}, \bar{S}_m, \bar{T}_m, \bar{F}_\pi] = \text{Bin}\left[E - \bar{F}_\pi - \bar{E}_m^{(s)}, \frac{\bar{T}_m}{C - \bar{S}_m}\right](\bar{E}_m^{(t)}[\pi])$$

**Delivery probability (given message).** We can finally compute, under the above conditions, the probability of  $\pi$  eventually delivering  $m$ .

We start by expanding the definition of  $A_m[\pi]$  to get

$$\begin{aligned} \mathcal{P}[A_m[\pi] | H \neq m, \bar{S}_m, \bar{T}_m, \bar{F}_\pi] \\ = \mathcal{P}[E_m[\pi] \geq \hat{E} - \bar{F}_\pi | H \neq m, \bar{S}_m, \bar{T}_m, \bar{F}_\pi] = (\star_1) \end{aligned}$$

### Chapter 3. Sieve

---

and then expand the definition of  $E_m[\pi]$  to get

$$(\star_1) = \mathcal{P}[E_m^{(t)}[\pi] \geq \hat{E} - \bar{F}_\pi - E_m^{(s)}[\pi] \mid H \neq m, \bar{S}_m, \bar{T}_m, \bar{F}_\pi] = (\star_2)$$

Finally, using the law of total probability on each possible value of  $E_m^{(s)}[\pi]$ , we get

$$(\star_2) = \sum_{\bar{E}_m^{(s)}[\pi]=0}^{E-\bar{F}_\pi} \left( \underbrace{\mathcal{P}[E_m^{(t)}[\pi] \geq \hat{E} - \bar{F}_\pi - \bar{E}_m^{(s)}[\pi] \mid \bar{E}_m^{(s)}[\pi], \bar{S}_m, \bar{T}_m, \bar{F}_\pi]}_{(\star_a)} \cdot \underbrace{\mathcal{P}[\bar{E}_m^{(s)}[\pi] \mid H \neq m, \bar{S}_m, \bar{F}_\pi]}_{(\star_b)} \right)$$

As we previously established,

$$\mathcal{P}[\bar{E}_m^{(t)}[\pi] \mid E_m^{(s)}[\pi] = i, \bar{S}_m, \bar{T}_m, \bar{F}_\pi] = \mathcal{P}[X_i = \bar{E}_m^{(t)}[\pi]]$$

with

$$\begin{aligned} X_i &\sim \text{Bin}[A - i, p] \\ A &= E - \bar{F}_\pi \\ p &= \frac{\bar{T}_m}{(C - \bar{S}_m)} \end{aligned}$$

Moreover,

$$(\star_a) = \mathcal{P}[X_i \geq B - i]$$

with

$$B = \hat{E} - \bar{F}_\pi \leq E - \bar{F}_\pi = A$$

Therefore, following from Lemma 7,  $(\star_a)$  is an increasing function of  $\bar{E}_m^{(s)}[\pi]$ . Moreover, as we previously established,

$$\begin{aligned} \mathcal{P}[\bar{E}_m^{(s)}[\pi] \mid H \neq m, \bar{S}_m, \bar{F}_\pi] &= \mathcal{P}\left[E_m^{\bar{S}_m}[\pi] = \bar{E}_m^{(s)}[\pi] \mid A_m^{\bar{S}_m}[\pi], \bar{F}_\pi\right] \\ &= \frac{I(\bar{E}_m^{\bar{S}_m}[\pi] < \hat{E} - \bar{F}_\pi) \mathcal{P}\left[\bar{E}_m^{\bar{S}_m}[\pi] \mid \bar{F}_\pi\right]}{\sum_{e=0}^{\hat{E}-\bar{F}_\pi-1} \mathcal{P}\left[E_m^{\bar{S}_m}[\pi] = e \mid \bar{F}_\pi\right]} \end{aligned}$$



and  $(\star_2)$  can be restated as

$$\begin{aligned}
 (\star_2) &= \frac{\sum_{x=0}^{K-l} f(x)g(x)}{\sum_{x=0}^{K-l} g(x)} \\
 K &= E - \bar{F}_\pi \\
 l &= E - \hat{E} + 1 \\
 f(x) &= \mathcal{P}\left[E_m^{(t)}[\pi] \geq \hat{E} - \bar{F}_\pi - x \mid E_m^{(s)}[\pi] = x, \bar{S}_m, \bar{T}_m, \bar{F}_\pi\right] \\
 g(x) &= \mathcal{P}\left[E_m^{\bar{S}_m}[\pi] = x \mid \bar{F}_\pi\right]
 \end{aligned}$$

with  $f(x)$  increasing and  $\sum_{x=0}^K g(x) = 1$ . Following from Corollary 1, we therefore have

$$\begin{aligned}
 &\mathcal{P}\left[A_m[\pi] \mid H \neq m, \bar{S}_m, \bar{T}_m, \bar{F}_\pi\right] \\
 &\leq \sum_{\bar{E}_m^{(s)}[\pi]=0}^{E-\bar{F}_\pi} \mathcal{P}\left[\bar{E}_m^{(s)}[\pi] + E_m^{(t)}[\pi] \geq \hat{E} - \bar{F}_\pi \mid \bar{E}_m^{(s)}[\pi], \bar{S}_m, \bar{T}_m, \bar{F}_\pi\right] \\
 &\quad \cdot \mathcal{P}\left[E_m^{\bar{S}_m}[\pi] = \bar{E}_m^{(s)}[\pi] \mid \bar{F}_\pi\right] \\
 &= (\star_3)
 \end{aligned}$$

which, as we previously established, can be restated as

$$\begin{aligned}
 (\star_3) &= \mathcal{P}[X + Y \geq H] = \sum_{K=H}^A \mathcal{P}[X + Y = K] \\
 \mathcal{P}[\bar{X}] &= \text{Bin}\left[A, \frac{x}{B}\right](\bar{X}) \\
 \mathcal{P}[\bar{Y} \mid \bar{X}] &= \text{Bin}\left[A - \bar{X}, \frac{y}{B-x}\right](\bar{Y}) \\
 H &= \hat{E} - \bar{F}_\pi \\
 A &= E - \bar{F}_\pi \\
 B &= C \\
 x &= \bar{S}_m \\
 y &= \bar{T}_m
 \end{aligned}$$

which, using Lemma 6, yields the bound

$$\mathcal{P}\left[A_m[\pi] \mid H \neq m, \bar{S}_m, \bar{T}_m, \bar{F}_\pi\right] \leq \sum_{e=\hat{E}-\bar{F}_\pi} \text{Bin}\left[E - \bar{F}_\pi, \frac{\bar{S}_m + \bar{T}_m}{C}\right](e) \quad (3.2)$$

**Delivery probability (any message).** We now move on to compute the probability that a correct process  $\pi$  will eventually deliver any message other than  $H$ , under the following assumptions:

### Chapter 3. Sieve

---

- The first phase of the adversarial execution is concluded.
- The number  $\bar{F}_\pi$  of Byzantine processes in the first echo sample of  $\pi$  is given.
- The number  $\bar{S}_m, \bar{T}_m$  of correct processes that pb.Delivered each message  $m$  throughout the first and second phase respectively is given.

Since every echo sample is picked independently, from Equation (3.2) follows

$$\begin{aligned} \mathcal{P} \left[ \bigvee_{m \neq \bar{H}} A_m[\pi] \mid \bar{S}_1, \dots, \bar{S}_C, \bar{T}_1, \dots, \bar{T}_C, \bar{F}_\pi \right] &\leq \mathcal{P} \left[ \bigvee_{i \neq m} (X_i \geq K) \right] \\ \mathcal{P}[\bar{X}_i] &= \text{Bin}[N, p_i](\bar{X}_i) \\ N &= E - \bar{F}_\pi \\ K &= \hat{E} - \bar{F}_\pi \\ p_i &= \frac{\bar{S}_i + \bar{T}_i}{C} \end{aligned}$$

and noting that

$$\sum_{m \neq \bar{H}} \frac{\bar{S}_m + \bar{T}_m}{C} = \sum_{n \neq \bar{H}} \bar{N}_m = C - \bar{N}_{\bar{H}}$$

we can use Lemma 10 to obtain the bound

$$\mathcal{P} \left[ \bigvee_{m \neq \bar{H}} A_m[\pi] \mid \bar{S}_1, \dots, \bar{S}_C, \bar{T}_1, \dots, \bar{T}_C, \bar{F}_\pi \right] \leq \phi(\bar{N}_{\bar{H}}, \bar{F}_\pi)$$

with

$$\phi(\bar{N}_{\bar{H}}, \bar{F}_\pi) = \begin{cases} \alpha(\bar{N}_{\bar{H}}, \bar{F}_\pi) \cdot \beta(\bar{N}_{\bar{H}}, \bar{F}_\pi) & \text{iff } \frac{C - \bar{N}_{\bar{H}}}{C} \leq \frac{(\hat{E} - \bar{F}_\pi) - \sqrt{\hat{E} - \bar{F}_\pi}}{E - \bar{F}_\pi} \\ 1 & \text{otherwise} \end{cases} \quad (3.3)$$

where

$$\begin{aligned} \alpha(\bar{N}_{\bar{H}}, \bar{F}_\pi) &= \left( \frac{e(E - \bar{F}_\pi) \frac{C - \bar{N}_{\bar{H}}}{C}}{\hat{E} - \bar{F}_\pi} \right)^{(\hat{E} - \bar{F}_\pi)} \\ \beta(\bar{N}_{\bar{H}}, \bar{F}_\pi) &= \exp \left( -(E - \bar{F}_\pi) \frac{C - \bar{N}_{\bar{H}}}{C} \right) \end{aligned}$$

At a first glance, the second branch of the bound above could seem unreasonably lax. We underline, however, that for a large enough  $(\hat{E} - \bar{F}_\pi)$ ,

$$\frac{(\hat{E} - \bar{F}_\pi) - \sqrt{\hat{E} - \bar{F}_\pi}}{E - \bar{F}_\pi} \simeq \frac{\hat{E} - \bar{F}_\pi}{E - \bar{F}_\pi}$$

and, since the median of  $\text{Bin}[N, p]$  is either  $\lceil Np \rceil$  or  $\lfloor Np \rfloor$ ,

$$\sum_{e=\hat{E}-\bar{F}_\pi}^{E-\bar{F}_\pi} \text{Bin}\left[E-\bar{F}_\pi, \frac{\hat{E}-\bar{F}_\pi}{E-\bar{F}_\pi}\right](e) \simeq \frac{1}{2}$$

Therefore, even in the second branch, the bound introduces a limited multiplicative error. Moreover, as we will see in the numerical analysis, the error introduced by the bound is non-negligible only for extremely unlikely values of  $\bar{N}_H$ .

**Adversarial success probability.** Throughout this section, we computed the probability that a correct process  $\pi$  will deliver a message different from the message that was delivered throughout the first phase.

We showed that such probability can be bound by a function that only depends on the number of Byzantine processes in the first echo sample of  $\pi$ , and the number of correct processes that pb.Delivered  $H$  throughout the first phase.

We therefore have

$$\mathcal{P}\left[\bigvee_{m \neq H} A_m[\pi] \mid \bar{N}_H, \bar{F}_\pi\right] \leq \phi(\bar{N}_H, \bar{F}_\pi)$$

By the law of total probability we have

$$\begin{aligned} \mathcal{P}\left[\bigvee_{m \neq H} A_m[\pi] \mid \bar{N}_H\right] &= \sum_{m \neq H} \mathcal{P}\left[\bigvee_{m \neq H} A_m[\pi] \mid \bar{N}_H, \bar{F}_\pi\right] \mathcal{P}[\bar{F}_\pi \mid \bar{N}_H] \\ &\leq \sum_{m \neq H} \phi(\bar{N}_H, \bar{F}_\pi) \mathcal{P}[\bar{F}_\pi \mid \bar{N}_H] \end{aligned}$$

Since  $H$  was delivered by at least one correct process at the end of the first phase, we know that:

- One correct process  $\pi^+$  delivered  $H$  immediately after  $\bar{N}_H$  correct processes pb.Delivered  $H$ .
- Every other correct process did not deliver  $H$  before  $\bar{N}_H$  correct processes pb.Delivered  $H$ .

We start by computing the probability distribution underlying  $\bar{F}_{\pi^+}$ . Using Bayes' theorem we

get

$$\begin{aligned} \mathcal{P}[\bar{F}_{\pi^+} | \bar{N}_H] &= \mathcal{P}\left[\bar{F}_{\pi^+} | A_H^{\bar{N}_H}[\pi^+], \cancel{A_H^{\bar{N}_H-1}[\pi^+]}\right] \\ &= \frac{\mathcal{P}\left[A_H^{\bar{N}_H}[\pi^+], \cancel{A_H^{\bar{N}_H-1}[\pi^+]} | \bar{F}_{\pi^+}\right] \mathcal{P}[\bar{F}_{\pi^+}]}{\mathcal{P}\left[A_H^{\bar{N}_H}[\pi^+], \cancel{A_H^{\bar{N}_H-1}[\pi^+]}\right]} \end{aligned}$$

and noting that  $A_H^{\bar{N}_H-1}[\pi^+] \implies A_H^{\bar{N}_H}[\pi^+]$ , we have

$$\begin{aligned} \mathcal{P}\left[A_H^{\bar{N}_H}[\pi^+], \cancel{A_H^{\bar{N}_H-1}[\pi^+]}\right] &= \mathcal{P}\left[A_H^{\bar{N}_H}[\pi^+]\right] - \mathcal{P}\left[A_H^{\bar{N}_H-1}[\pi^+]\right] \\ \mathcal{P}\left[A_H^{\bar{N}_H}[\pi^+], \cancel{A_H^{\bar{N}_H-1}[\pi^+]} | \bar{F}_{\pi^+}\right] &= \mathcal{P}\left[A_H^{\bar{N}_H}[\pi^+] | \bar{F}_{\pi^+}\right] \\ &\quad - \mathcal{P}\left[A_H^{\bar{N}_H-1}[\pi^+] | \bar{F}_{\pi^+}\right] \end{aligned}$$

Similarly, for  $\pi^- \neq \pi^+$ , we get

$$\begin{aligned} \mathcal{P}[\bar{F}_{\pi^-} | \bar{N}_H] &= \mathcal{P}\left[\bar{F}_{\pi^-} | \cancel{A_H^{\bar{N}_H-1}[\pi^-]}\right] \\ &= \frac{\mathcal{P}\left[\cancel{A_H^{\bar{N}_H-1}[\pi^-]} | \bar{F}_{\pi^-}\right] \mathcal{P}[\bar{F}_{\pi^-}]}{\mathcal{P}\left[\cancel{A_H^{\bar{N}_H-1}[\pi^-]}\right]} \end{aligned}$$

Since each correct process picks its echo sample independently, we have

$$\begin{aligned} \mathcal{P}[W | \bar{N}_H] &= 1 - \prod_{\pi \in \Pi_C} \left(1 - \mathcal{P}\left[\bigvee_{m \neq H} A_m[\pi] | \bar{N}_H\right]\right) \\ &\leq 1 - (1 - \phi^+(\bar{N}_H))(1 - \phi^-(\bar{N}_H))^{C-1} \end{aligned} \tag{3.4}$$

with

$$\begin{aligned} \phi^+(\bar{N}_H) &= \sum_{\bar{F}_{\pi^+}=0}^E \phi(\bar{N}_H, \bar{F}_{\pi^+}) \mathcal{P}[\bar{F}_{\pi^+} | \bar{N}_H] \\ \phi^-(\bar{N}_H) &= \sum_{\bar{F}_{\pi^-}=0}^E \phi(\bar{N}_H, \bar{F}_{\pi^-}) \mathcal{P}[\bar{F}_{\pi^-} | \bar{N}_H] \end{aligned}$$

### 3.10.5 First phase

In the previous section, we computed, given the number of correct processes that pb.Delivered the first delivered message, the probability of a two-phase adversary successfully compromising the consistency of a system.

In this section, we compute the probability distribution underlying the number of correct processes that pb.Deliver the first delivered message.

**Definition 27** (Deafened adversary). Let  $\alpha$  be a two-phase adversary. We define  $\Delta(\alpha)$  the **deafened version of  $\alpha$**  if:

- $\Delta(\alpha)$  is a process-sequential adversary.
- Coupled with a system  $\sigma$ ,  $\Delta(\alpha)$  sequentially causes the pb.Delivery of  $\alpha[F(\sigma)]_1, \dots, \alpha[F(\sigma)]_C$ .

Intuitively, the deafened version of a two-phase adversary  $\alpha$  is an adversary whose adversarial execution would be identical to  $\alpha$ 's, if no correct process ever delivered any message.

**Lemma 18.** *Let  $\alpha$  be a two-phase adversary, let  $\sigma$  be a system. We have*

$$\eta(\alpha, \sigma) = \eta(\Delta(\alpha), \sigma)$$

*Proof.* It follows immediately from Definition 27:  $\Delta(\alpha)$  causes the same processes to pb.Deliver the same messages as  $\alpha$  throughout the first phase.  $\square$

**Definition 28** (Delivery cost). Let  $\alpha$  be an auto-echo adversary, let  $\sigma$  be a non-poisoned system, let  $m$  be a message such that, when  $\alpha$  is coupled with  $\sigma$ , at least one correct process delivers  $m$ . We define the **delivery cost of  $m$**   $\lambda(\alpha, \sigma, m)$  as the minimum  $\lambda \in 1..C$  such that, when  $\alpha$  is coupled with  $\sigma$ , at least one correct process delivers  $m$  after  $\lambda$  correct processes pb.Delivered  $m$ .

**Lemma 19.** *Let  $\alpha$  be a two-phase adversary, let  $\sigma$  be a non-poisoned system such that, when coupled with  $\sigma$ ,  $\alpha$  causes at least one correct process to deliver one message. Let  $\bar{H}$  be the first message delivered by at least one correct process, when  $\alpha$  is coupled with  $\sigma$ .*

*We have that*

$$\lambda(\alpha, \sigma, \bar{H}) \geq \min_{m \in \mathcal{M}} \lambda(\Delta(\alpha), \sigma, m)$$

*Proof.* Following from Lemma 18  $\eta(\alpha, \sigma) = \eta(\Delta(\alpha), \sigma)$ . Therefore, at least one correct process delivers  $\bar{H}$  after  $\lambda(\alpha, \sigma, \bar{H})$  processes pb.Deliver  $\bar{H}$ , when  $\Delta(\alpha)$  is coupled with  $\sigma$ .  $\square$

In this section, we bound the cumulative probability  $\mathcal{P}[N_H \leq L]$  for an adversary  $\alpha$  by bounding the probability that the deafened adversary  $\Delta(\alpha)$  will cause the delivery of at least one message  $m$ , with a cost smaller or equal to  $L$ .

Let  $m$  be a message. We start by noting that, by definition,  $\Delta(\alpha)$  eventually causes  $L_m$  correct processes to pb.Deliver  $m$ . Let  $\pi$  be a correct process, let  $\bar{F}_\pi$  be the number of Byzantine processes in  $\pi$ 's first echo sample.

### Chapter 3. Sieve

---

We denote with  $\Lambda_m[\pi]$  the random variable representing the minimum number of correct processes that pb.Deliver  $m$ , before  $\pi$  delivers  $m$ . If  $\pi$  never delivers  $m$ , we set  $\Lambda_m[\pi] = \infty$ .

Let  $L \in 1..C$ . Using the tools we developed in the previous section, we immediately get

$$\mathcal{P}[\Lambda_m[\pi] \leq L \mid \bar{L}_m, \bar{F}_\pi] = \sum_{e=\hat{E}-\bar{F}_\pi}^{E-\bar{F}_\pi} \text{Bin}\left[E - \bar{F}_\pi, \frac{\min(\bar{L}_m, L)}{C}\right](e)$$

and using the independence of echo samples, we get

$$\begin{aligned} \mathcal{P}\left[\bigvee_{m \in \mathcal{M}} \Lambda_m[\pi] \leq L \mid \bar{L}_1, \dots, \bar{L}_C, \bar{F}_\pi\right] &= \mathcal{P}\left[\bigvee_{i \in \mathcal{M}} (X_i \geq K)\right] \\ \mathcal{P}[\bar{X}_i] &= \text{Bin}[N, p_i](\bar{X}_i) \\ N &= E - \bar{F}_\pi \\ K &= \hat{E} - \bar{F}_\pi \\ p_i &= \frac{\min(\bar{L}_i, L)}{C} \end{aligned}$$

and we can use Lemma 10 to obtain the bound

$$\begin{aligned} &\mathcal{P}\left[\bigvee_{m \in \mathcal{M}} \Lambda_m[\pi] \leq L \mid \bar{L}_1, \dots, \bar{L}_C, \bar{F}_\pi\right] \\ &\leq 1 - (1 - \psi(L, \bar{F}_\pi))^{\lfloor \frac{C}{L} \rfloor} (1 - \psi(C \bmod L, \bar{F}_\pi)) \end{aligned}$$

with

$$\psi(M, \bar{F}_\pi) = \begin{cases} \alpha(M, \bar{F}_\pi) \cdot \beta(M, \bar{F}_\pi) & \text{iff } \frac{M}{C} \leq \frac{(\hat{E} - \bar{F}_\pi) - \sqrt{\hat{E} - \bar{F}_\pi}}{E - \bar{F}_\pi} \\ 1 & \text{otherwise} \end{cases} \quad (3.5)$$

where

$$\begin{aligned} \alpha(M, \bar{F}_\pi) &= \left(\frac{e(E - \bar{F}_\pi) \frac{M}{C}}{\hat{E} - \bar{F}_\pi}\right)^{(\hat{E} - \bar{F}_\pi)} \\ \beta(M, \bar{F}_\pi) &= \exp\left(- (E - \bar{F}_\pi) \frac{M}{C}\right) \end{aligned}$$

Noting that the bound holds for any value of  $\bar{L}_1, \dots, \bar{L}_C$ , we can use again the law of total probability to obtain

$$\mathcal{P}\left[\bigvee_{m \in \mathcal{M}} \Lambda_m[\pi] \leq L\right] \leq \psi(L)$$

with

$$\psi(L) = \sum_{\bar{F}_\pi=0}^E 1 - (1 - \psi(L, \bar{F}_\pi))^{\lfloor \frac{C}{L} \rfloor} (1 - \psi(C \bmod L, \bar{F}_\pi)) \mathcal{P}[\bar{F}_\pi] \quad (3.6)$$

and using the independence of echo samples across correct processes we finally get

$$\mathcal{P} \left[ \bigvee_{\pi \in \Pi_C, m \in \mathcal{M}} \Lambda_m[\pi] \leq L \right] \leq 1 - (1 - \psi(L))^C \quad (3.7)$$

We now have all the elements to prove

**Theorem 9.** *Sieve satisfies  $\epsilon_c$ -consistency, with*

$$\begin{aligned} \epsilon_c &\leq \epsilon_p + \sum_{L=0}^C \tilde{\psi}(L) \tilde{\phi}(L) \\ \tilde{\psi}(L) &= \begin{cases} (1 - (1 - \psi(L))^C) - (1 - (1 - \psi(L-1))^C) & \text{iff } L \in 1..C \\ 0 & \text{iff } L \in \{-1, 0\} \\ 1 & \text{iff } L = C \end{cases} \\ \tilde{\phi}(L) &= (1 - (1 - \phi^+(L))(1 - \phi^-(L))^{C-1}) \\ \epsilon_p &= 1 - \left( 1 - \sum_{\bar{F}=\hat{E}}^E \text{Bin}[E, f](\bar{F}) \right)^C \end{aligned}$$

*Proof.* Following from Lemma 19, we have

$$\mathcal{P}[N_H \leq L] \leq \mathcal{P} \left[ \bigvee_{\pi \in \Pi_C, m \in \mathcal{M}} \Lambda_m[\pi] \leq L \right] \quad (3.8)$$

By the law of total probability, we have

$$\begin{aligned} \mathcal{P}[W] &= \sum_{x=0}^C (f(x)(g(x) - g(x-1))) \\ f(x) &= \mathcal{P}[W \mid \bar{N}_H = x] \\ g(x) &= \mathcal{P}[N_H \leq x] \end{aligned}$$

and from Lemma 9 we get

$$\begin{aligned} \mathcal{P}[W] &\leq \sum_{x=0}^C (f(x)(h(x) - h(x-1))) \\ h(x) &= \mathcal{P} \left[ \bigvee_{\pi \in \Pi_C, m \in \mathcal{M}} \Lambda_m[\pi] \leq x \right] \end{aligned}$$

The probability of compromising a non-poisoned system is obtained by applying the bounds in Equations (3.4), (3.7) and (3.8).

It is easy to see that  $\epsilon_p$  represents the probability of a random system being poisoned: indeed,

each correct process has an independent probability

$$\sum_{\bar{F}=\hat{E}}^E \text{Bin}[E, f](\bar{F})$$

of having more than  $\hat{E}$  Byzantine processes in its first echo sample, i.e., of being poisoned.

Therefore, the bound on  $\epsilon_c$  bounds the probability of any two-phase adversary compromising the consistency of a cob system. Due to Lemma 27, the set  $\mathcal{A}_{tp}$  of two-phase adversaries is optimal. Therefore, Simplified Sieve satisfies  $\epsilon_c$ -consistency.

Due to Lemma 12, the adversarial power of an optimal pcb adversary is bound by the adversarial power of an optimal cob adversary, and the theorem is proved.  $\square$

### 3.11 Decorators

In this section, we provide the proof that each of the sets of cob adversaries presented in Section 3.9 is optimal.

#### 3.11.1 Auto-echo adversary

**Lemma 20.** *The set of auto-echo adversaries  $\mathcal{A}_{ae}$  is optimal.*

*Proof.* We prove the result using a **decorator**, i.e., an algorithm that acts as an interface between an adversary and a system. An adversary coupled with a decorator effectively implements an adversary. Here we show that a decorator  $\Delta_{ae}$  exists such that, for every  $\alpha \in \mathcal{A}$ , the adversary  $\alpha' = \Delta_{ae}(\alpha)$  is an auto-echo adversary, and more powerful than  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the auto-echo adversary  $\alpha^+ = \Delta_{ae}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 7 implements **Auto-echo decorator**, a decorator that transforms an adversary into an auto-echo adversary. Provided with an adversary  $adv$ , Auto-echo decorator acts an interface between  $adv$  and a system  $sys$ , effectively implementing an auto-echo adversary  $aeadv$ . Auto-echo decorator exposes both the adversary and the system interfaces: the underlying adversary  $adv$  uses  $aeadv$  as its system.

Auto-echo decorator works as follows:

- Procedure  $aeadv.Init()$  initializes the following variables:
  - A *queue* list that contains every combination of  $(\pi, m, \xi)$ ,  $\pi$  being a correct process,  $m$  being a message and  $\xi$  being a Byzantine process: *queue* is used to initially



**Algorithm 7 Auto-echo decorator.**


---

```

1: Implements:
2:   AutoEchoAdversary + CobSystem, instance aeadv
3:
4: Uses:
5:   CobAdversary, instance adv, system aeadv
6:   CobSystem, instance sys
7:
8: procedure aeadv.Init() is
9:   queue =  $\emptyset$ ;
10:
11:   for all  $\pi \in \Pi_C$  do
12:     for all  $m \in \mathcal{M}$  do
13:       for all  $\xi \in \Pi \setminus \Pi_C$  do
14:         queue  $\leftarrow$  queue  $\cup$   $\{(\pi, m, \xi)\}$ ;
15:       end for
16:     end for
17:   end for
18:
19:   echoes =  $\{\perp\}^{C \times C \times N}$ ;  $\triangleright C \times C \times N$  table filled with  $\perp$ .
20:   executed = False;
21:   adv.Init();
22:
23: procedure aeadv.Step() is
24:   if queue  $\neq \emptyset$  then
25:      $(\pi, m, \xi) = \textit{queue}[1]$ ;
26:     sys.Echo( $\pi, m, \xi, m$ );
27:     queue  $\leftarrow$  queue  $\setminus \{(\xi, \pi, m)\}$ ;
28:   else
29:     executed  $\leftarrow$  False;
30:     while executed = False do
31:       adv.Step();
32:     end while
33:   end if
34:
35: procedure aeadv.Byzantine(process) is
36:   return sys.Byzantine(process);
37:

```

---

### Chapter 3. Sieve

---

```
38: procedure aeadv.State() is
39:   state =  $\emptyset$ ;
40:
41:   for all  $(\pi, m) \in \text{sys.State}()$  do
42:     n = 0;
43:
44:     for all  $\rho \in \text{sys.Sample}(\pi, m)$  do
45:       if echoes[ $\pi$ ][m][ $\rho$ ] = m then
46:         n  $\leftarrow$  n + 1;
47:       end if
48:     end for
49:
50:     if n  $\geq \hat{E}$  then
51:       state  $\leftarrow$  state  $\cup$   $\{(\pi, m)\}$ ;
52:     end if
53:   end for
54:
55:   return state;
56:
57: procedure aeadv.Sample(process, message) is
58:   sample =  $\emptyset$ ;
59:
60:   for all  $\rho \in \text{sys.Sample}(process, message)$  do
61:     if echoes[process][message][ $\rho$ ]  $\neq \perp$  then
62:       sample  $\leftarrow$  sample  $\cup$   $\{\rho\}$ ;
63:     end if
64:   end for
65:
66:   return sample;
67:
68: procedure aeadv.Deliver(process, message) is
69:   executed = True;
70:
71:   for all  $\pi \in \Pi_C$  do
72:     for all  $m \in \mathcal{M}$  do
73:       echoes[ $\pi$ ][m][process] = message;
74:     end for
75:   end for
76:
77:   sys.Deliver(process, message, flag);
78:
79: procedure aeadv.Echo(process, sample, source, message) is
80:   echoes[process][sample][source] = message;
81:
```

---

---

82: **procedure** *aeadv.End()* **is**

83:     *executed* = True;

84:     *sys.End()*;

85:

---

cause every Byzantine process  $\xi$  to send an  $\text{Echo}(m, m)$  message to every correct process  $\pi$ , for every message  $m$ .

- An *echoes* table, initialized with  $\perp$  values: *echoes* is used to keep track of all the Echo messages that would have been sent to each correct process in *sys*, if *adv* was playing instead of *aeadv*.

- Procedure *aeadv.Step()* checks if *queue* is not empty. If it is not empty, it pops (i.e., picks and removes) its first element  $(\pi, m, \xi)$ , with  $\xi \in \Pi \setminus \Pi_C$ ,  $\pi \in \Pi_C$  and  $m \in \mathcal{M}$ . It then causes  $\xi$  to send  $\pi$  an  $\text{Echo}(m, m)$  message.

If *queue* is empty instead, the procedure calls *adv.Step()* until either *sys.Deliver(...)* or *sys.End()* are called: this is achieved using the *executed* flag.

- Procedure *aeadv.Byzantine(process)* simply forwards the call to *sys.Byzantine(process)*.
- Procedure *aeadv.State()* returns a list of pairs  $(\pi \in \Pi_C, m \in \mathcal{M})$  such that  $\pi$  delivered  $m$  in *sys*, and  $\pi$  would have delivered  $m$  in *sys*, if *adv* was playing instead of *aeadv*.

This is achieved by querying *sys.State()*, then looping over each element  $(\pi, m)$  of the response. For each  $(\pi, m)$ , the procedure loops over every element  $\rho$  of *sys.Sample( $\pi, m$ )*, and computes the number  $n$  of  $\text{Echo}(m, m)$  messages that  $\pi$  would have received from its echo sample for  $m$  in *sys*, if *adv* was playing instead of *aeadv*. This is achieved using the *echoes* table. If  $n$  is greater or equal to  $\hat{E}$ ,  $(\pi, m)$  is included in the list returned by the procedure.

- Procedure *aeadv.Sample(process, message)* returns every process in *sys.Sample(process, message)* that would have sent an  $\text{Echo}(message, message')$  message for some message  $message'$  to *process* in *sys*, if *adv* was playing instead of *aeadv*. This is achieved using the *echoes* table.
- Procedure *aeadv.Deliver(process, message)* updates the *echo* table to reflect all the Echo messages that *process* will send, as a result of having pb.Delivered *message*. It then forwards the call to *sys.Deliver(process, message)*, causing *process* to pb.Deliver *message*.
- Procedure *aeadv.Echo(process, sample, source, message)* updates the *echo* table to include the  $\text{Echo}(sample, message)$  message that *process* would receive from *source*, if *adv* was playing instead of *aeadv*.
- Procedure *aeadv.End()* simply forwards the call to *sys.End()*.

## Chapter 3. Sieve

---

**Correctness.** We start by proving that no adversary, coupled with Auto-echo decorator, causes the execution to fail.

We start by establishing a preliminary result. Let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ . If  $(\pi, m)$  is returned from  $aeadv.State()$ , then  $\pi$  delivered  $m$  in  $sys$ . Indeed,  $(\pi, m)$  is returned from  $aeadv.State()$  only if  $(\pi, m)$  is returned from  $sys.State()$ .

Let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ . The following hold true:

- An invocation to  $aeadv.Step()$  results in one and only one call to  $sys.Deliver(...)$ ,  $sys.Echo(...)$  or  $sys.End()$ . Indeed, if  $queue$  is not empty, exactly one call to  $sys.Echo(...)$  is issued. Otherwise,  $adv.Step()$  is called until  $executed = \text{True}$ , and  $executed$  is set to  $\text{True}$  only after an invocation to  $sys.Deliver(...)$  or  $sys.End()$ .
- Procedure  $aeadv.State()$  never causes the execution to fail. Indeed,  $sys.Sample(\pi, m)$  is called only if  $(\pi, m)$  is returned from  $sys.State()$ . This means that  $sys.Sample(\pi, m)$  is called only if  $\pi$  delivered  $m$  in  $sys$ . Therefore,  $sys.Sample(\pi, m)$  is never invoked from  $aeadv.State()$  unless at least one correct process delivered  $m$  in  $sys$ .
- No invocation of  $aeadv.Sample(...)$  causes the execution to fail. Noting that  $adv$  is correct, it will never invoke  $aeadv.Sample(\pi, m)$  unless  $(\pi', m)$  was returned from a previous invocation of  $aeadv.State()$ , for some  $\pi' \in \Pi_C$ . As we previously established,  $(\pi', m)$  is returned from  $aeadv.State()$  only if  $\pi'$  delivered  $m$  in  $sys$ . Therefore,  $sys.Sample(\pi, m)$  is never invoked from  $aeadv.Sample(...)$  unless at least one correct process delivered  $m$  in  $sys$ .

**Auto-echo.** It is easy to prove that Auto-echo decorator always implements an auto-echo adversary. Indeed, every call to  $aeadv.Step()$  results in a call to  $sys.Echo(\pi, m, \xi, m)$ , causing the Byzantine process  $\xi$  to send an  $\text{Echo}(m, m)$  message to the correct process  $\pi$ , until  $queue$  is exhausted.

Therefore, only  $sys.Echo(...)$  is invoked until  $\xi$  sent an  $\text{Echo}(m, m)$  message to  $\pi$ , for every  $\pi \in \Pi_C$ , every  $m \in \mathcal{M}$ , and every  $\xi \in \Pi \setminus \Pi_C$ .

**Roadmap.** Let  $\alpha \in \mathcal{A}$ , let  $\alpha' = \Delta_{ae}(\alpha)$ . Let  $\sigma$  be a system such that  $\alpha$  compromises the consistency of  $\sigma$ . Let  $\sigma'$  be an identical copy of  $\sigma$ . In order to prove that  $\alpha'$  is more powerful than  $\alpha$ , we prove that  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Trace.** We start by noting that, if we couple Auto-echo decorator with  $\sigma'$ , we effectively obtain a system instance  $\delta$  with which  $\alpha$  directly exchanges invocations and responses. Here we show that the trace  $\tau(\alpha, \sigma)$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\sigma$ , or it has been coupled with  $\sigma'$ , with Auto-echo decorator acting as an interface. We prove this by induction.

Let us assume

$$\begin{aligned}\tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \delta) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n\end{aligned}$$

We start by noting that, since  $\alpha$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . Since  $\text{aeadv.Byzantine}(\pi)$  simply forwards the call to  $\text{sys.Byzantine}(\pi)$ , and  $\sigma'$  is an identical copy of  $\sigma$ , we immediately have  $r_{n+1} = r'_{n+1}$ .

Before considering the remaining possible values of  $i_{n+1}$ , we prove some auxiliary results. Let  $\pi$  be a correct process, let  $\xi$  be a Byzantine process, let  $\rho$  be a process, let  $s, m$  be messages. For every  $j \leq n + 1$ , as we established, we have  $i_j = i'_j$ . Therefore, after the  $(n + 1)$ -th invocation, the following hold true:

- $\rho$  sent an  $\text{Echo}(s, m)$  message to  $\pi$  in  $\sigma$  if and only if  $\text{echoes}[\pi][s][\rho] = m$ . Indeed, if  $\rho$  is correct,  $\rho$  sent an  $\text{Echo}(s, m)$  message to  $\pi$  in  $\sigma$  if and only if  $\text{aeadv.Deliver}(\rho, m)$  was invoked. In turn,  $\text{echo}[\pi][s][\rho]$  was set to  $m$  for every  $\pi' \in \Pi_C, s' \in \mathcal{M}$  if and only if  $\text{aeadv.Deliver}(\rho, m)$  was invoked. If  $\rho$  is Byzantine,  $\rho$  sent an  $\text{Echo}(s, m)$  message to  $\pi$  in  $\sigma$  if and only if  $\text{aeadv.Echo}(\pi, s, \rho, m)$  was invoked. In turn,  $\text{echo}[\pi][s][\rho]$  was set to  $m$  if and only if  $\text{aeadv.Echo}(\pi, s, \rho, m)$  was invoked.
- If  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma$  for some  $m' \in \mathcal{M}$ , then  $\rho$  sent an  $\text{Echo}(m, m'')$  message to  $\pi$  in  $\sigma'$  as well, for some  $m'' \in \mathcal{M}$ . Indeed, if  $\rho$  is correct, then  $\text{aeadv.Deliver}(\rho, m')$  was invoked. As a result,  $\text{sys.Deliver}(\rho, m')$  was called, and  $\rho$  sent an  $\text{Echo}(s', m')$  message to  $\pi'$  for every  $\pi' \in \Pi_C, s' \in \mathcal{M}$ . If  $\rho$  is Byzantine, then it sent an  $\text{Echo}(m''', m''')$  message to  $\pi'$ , for every  $\pi' \in \Pi_C, m''' \in \mathcal{M}$ .
- If  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma$ , then  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma'$  as well. Indeed, if  $\rho$  is correct, then  $\text{aeadv.Deliver}(\rho, m)$  was invoked. As a result,  $\text{sys.Deliver}(\rho, m)$  was called, and  $\rho$  sent an  $\text{Echo}(s', m)$  message to  $\pi'$  for every  $\pi' \in \Pi_C, s' \in \mathcal{M}$ . If  $\rho$  is Byzantine, then it sent an  $\text{Echo}(m', m')$  message to  $\pi'$ , for every  $\pi' \in \Pi_C, m' \in \mathcal{M}$ .
- If  $\pi$  delivered  $m$  in  $\sigma$ , then  $\pi$  delivered  $m$  in  $\sigma'$  as well. This follows from the above and the fact that  $\sigma'$  is an identical copy of  $\sigma$  (i.e.,  $\pi$ 's echo sample for  $m$  in  $\sigma$  is identical to  $\pi$ 's echo sample in  $\sigma'$ ).

## Chapter 3. Sieve

---

Let us assume that  $i_{n+1} = (\text{State})$ . Let  $\pi$  be a correct process, let  $m$  be message. We start by noting that  $\text{aeadv.State}()$  returns  $(\pi, m)$  if and only if  $\pi$  delivered  $m$  in  $\sigma'$ , and  $\pi$  delivered  $m$  in  $\sigma$ . Indeed,  $(\pi, m)$  is added to the return list of  $\text{aeadv.State}()$  if and only if  $(\pi, m)$  is returned from  $\text{sys.State}()$ , and at least  $\hat{E}$  processes sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma$ . If  $(\pi, m) \in r_{n+1}$ , then  $\pi$  delivered  $m$  in  $\sigma$ , and  $\pi$  delivered  $m$  in  $\sigma'$  as well. Therefore  $(\pi, m) \in r'_{n+1}$ . If  $(\pi, m) \in r'_{n+1}$ , then we immediately have that  $\pi$  delivered  $m$  in  $\sigma$ , and  $(\pi, m) \in r_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Sample}, \pi, m)$ . Let  $\rho$  be a process. We start by noting that  $\text{aeadv.Sample}(\pi, m)$  returns  $\rho$  if and only if  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma'$  for some  $m' \in \mathcal{M}$ , and  $\text{echoes}[\pi][m][\rho] \neq \perp$ . If  $\rho \in r_{n+1}$ , then  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma$ , for some  $m' \in \mathcal{M}$ . Therefore,  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma'$ , for some  $m' \in \mathcal{M}$ , and  $\text{echoes}[\pi][m][\rho] = m' \neq \perp$ . Consequently,  $\rho \in r'_{n+1}$ . If  $\rho \in r'_{n+1}$ , then  $\text{echoes}[\pi][m][\rho] = m' \neq \perp$  for some  $m' \in \mathcal{M}$ . Therefore,  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma$ , and  $\rho \in r_{n+1}$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ .

**Consistency of  $\sigma'$ .** We proved that  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ . Moreover, we proved that if a correct process  $\pi$  eventually delivers a message  $m$  in  $\sigma$ , then  $\pi$  also delivers  $m$  in  $\sigma'$ .

Since  $\alpha$  compromises the consistency of  $\sigma$ , two correct processes  $\pi, \pi'$  and two distinct messages  $m, m' \neq m$  exist such that, in  $\sigma$ ,  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$ . Therefore, in  $\sigma'$ ,  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$ . Therefore  $\alpha'$  compromises the consistency of  $\sigma'$ .

Consequently, the adversarial power of  $\alpha$  is smaller or equal to the adversarial power of  $\alpha' = \Delta_{ae}(a)$ , and the lemma is proved.  $\square$

### 3.11.2 Process-sequential adversary

**Lemma 21.** *The set of process-sequential adversaries  $\mathcal{A}_{ps}$  is optimal.*

*Proof.* We again prove the result using a decorator, i.e., an algorithm that acts as an interface between an adversary and a system. An adversary coupled with a decorator effectively implements an adversary. Here we show that a decorator  $\Delta_{ps}$  exists such that, for every  $\alpha \in \mathcal{A}_{ae}$ , the adversary  $\alpha' = \Delta_{ps}(\alpha)$  is a process-sequential adversary, and as powerful as  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the process-sequential  $\alpha^+ = \Delta_{ps}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 8 implements **Process-sequential decorator**, a decorator that transforms an auto-echo adversary into a process-sequential adversary. Provided with an auto-echo

**Algorithm 8 Process-sequential decorator.**


---

```

1: Implements:
2:   ProcessSequentialAdversary + CobSystem, instance psadv
3:
4: Uses:
5:   AutoEchoAdversary, instance aeadv, system psadv
6:   CobSystem, instance sys
7:
8: procedure psadv.Init() is
9:   perm =  $\{\perp\}^C$ ;   cursor = 1;
10:  aeadv.Init();
11:
12: procedure psadv.Step() is
13:  aeadv.Step();
14:
15: procedure psadv.Byzantine(process) is
16:  return sys.Byzantine(process);
17:
18: procedure psadv.State() is
19:  return sys.State();
20:
21: procedure psadv.Sample(process, message) is
22:  sample =  $\emptyset$ ;
23:
24:  for all  $\rho \in \text{sys.Sample}(process, message)$  do
25:    if  $\rho \in \Pi_C$  then
26:      sample  $\leftarrow$  sample  $\cup$   $\{\zeta(\text{perm}[\zeta^{-1}(\rho)])\}$ 
27:    else
28:      sample  $\leftarrow$  sample  $\cup$   $\{\rho\}$ ;
29:    end if
30:  end for
31:
32:  return sample;
33:
34: procedure psadv.Deliver(process, message) is
35:  perm[cursor] =  $\zeta^{-1}(process)$ ;
36:  sys.Deliver( $\zeta(cursor)$ , message);
37:  cursor  $\leftarrow$  cursor + 1;
38:
39: procedure psadv.Echo(process, sample, source, message) is
40:  sys.Echo(process, sample, source, message);
41:
42: procedure psadv.End() is
43:  sys.End();
44:

```

---

### Chapter 3. Sieve

---

adversary  $aeadv$ , Process-sequential decorator acts as an interface between  $aeadv$  and a system  $sys$ , effectively implementing a process-sequential adversary  $psadv$ . Process-sequential decorator exposes both the adversary and the system interfaces: the underlying adversary  $aeadv$  uses  $psadv$  as its system.

Process-sequential decorator works as follows:

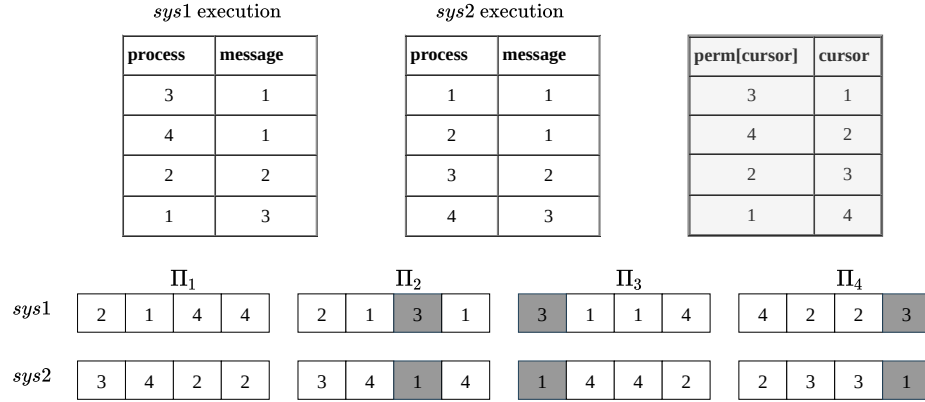
- Procedure  $psadv.Init()$  initializes the following variables:
  - A  $perm$  array of  $C$  elements:  $perm$  is used to consistently translate process identifiers between  $aeadv$  and  $sys$ .
  - A  $cursor$  variable, initially set to 1: at any time,  $cursor$  identifies the next process that will pb.Deliver a message in  $sys$ .
- Procedure  $psadv.Step()$  simply forwards the call to  $aeadv.Step()$ .
- Procedure  $psadv.Byzantine(process)$  simply forwards the call to  $sys.Byzantine(process)$ .
- Procedure  $psadv.State()$  simply forwards the call to  $sys.State()$ .
- Procedure  $psadv.Sample(process, message)$  returns the list of processes returned by  $sys.Sample(process, message)$ , translated through  $perm$ . More specifically, for every process  $\rho$  in  $sys.Sample(process, message)$ : if  $\rho$  is correct, it is translated to  $\zeta(perm[\zeta^{-1}(\rho)])$ ; if  $\rho$  is Byzantine, it is left unchanged.
- Procedure  $psadv.Deliver(process, message)$  sets  $perm[cursor]$  to  $\zeta^{-1}(process)$ , then forwards the call to  $sys.Deliver(\zeta(cursor), message)$ . Finally, it increments  $cursor$ . This serves the purpose to sequentially cause  $\zeta(1), \zeta(2), \dots$  to deliver a message, while storing the translation in  $perm$  in order for  $psadv.Sample(\dots)$  to provide a response consistent with any previous invocation of  $psadv.Deliver(\dots)$ .
- Procedure  $psadv.Echo(process, sample, source, message)$  simply forwards the call to  $sys.Echo(process, sample, source, message)$ .
- Procedure  $psadv.End()$  simply forwards the call to  $sys.End()$ .

**Correctness.** We start by proving that no adversary, coupled with Process-sequential decorator, causes the execution to fail.

The following hold true:

- No invocation of  $psadv.Sample(\dots)$  causes the execution to fail. Noting that  $aeadv$  is correct, it will never invoke  $psadv.Sample(\pi, m)$  unless  $(\pi', m)$  was returned from a





**Figure 3.4: Two systems with (one of) their respective echo samples.** The table on the right shows the permutation from *sys1* to *sys2*. Clearly both systems are equally likely. Moreover, the effect of process 3 delivering message 1 (grey) in *sys1*, is equal to process 1 delivering the same message in *sys2*. It can be seen that this holds for all further message deliveries. Intuitively this shows why we can restrict the adversary to always deliver in sequence.

previous invocation of  $psadv.State()$ , for some  $\pi' \in \Pi_C$ . Moreover, since  $psadv.State()$  simply forwards the call to  $sys.State()$ , if  $(\pi', m)$  was returned from  $psadv.State()$ , then  $\pi'$  delivered  $m$  in  $sys$ . Therefore,  $sys.Sample(\pi, m)$  is never invoked from  $psadv.Sample(\dots)$  unless at least one correct process delivered  $m$  in  $sys$ .

- Procedure  $sys.Sample(\dots)$  never calls  $\zeta(\perp)$ . We defer the proof of this result to a later section of this lemma.
- Procedure  $sys.Deliver(\dots)$  is never invoked twice on the same process. Indeed, by definition,  $\zeta$  is a bijection between  $1..C$  and  $\Pi_C$ , and  $cursor$  is incremented every time  $sys.Deliver(\dots)$  is called.

**Process-sequential.** It is easy to prove that Process-sequential decorator always implements a process-sequential adversary. Indeed,  $sys.Deliver(\dots)$  is invoked sequentially on  $\zeta(1), \zeta(2), \dots$  as  $cursor$  is incremented, regardless of the process originally provided to  $psadv.Deliver(\dots)$ .

**System translation.** Let  $\alpha$  be an adversary. We start by noting that, since  $\alpha$  is correct,  $\alpha$  always causes every correct process to  $pb.Deliver$  a message. We can therefore define a function

$$\mu: \mathcal{A} \times \mathcal{S} \times \Pi_C \rightarrow 1..C$$

such that  $\mu(\alpha, \sigma, \pi) = d$  if and only if  $\pi$  is the  $d$ -th process that  $\alpha$  causes to  $pb.Deliver$  a message, when  $\alpha$  is coupled with  $\sigma$ . We additionally define

$$(\mu^{-1}(\alpha, \sigma, d) = \pi) \stackrel{def}{\iff} (\mu(\alpha, \sigma, \pi) = d)$$

### Chapter 3. Sieve

---

We then define a **system translation function**  $\Psi[\alpha] : \mathcal{S} \rightarrow \mathcal{S}$  such that, for every system  $\sigma$ , every correct process  $\pi$ , every message  $m$ , and every  $e \in 1..E$ ,

$$\Psi[\alpha](\sigma)[\pi][m][e] = \begin{cases} \zeta(\mu(\alpha, \sigma, \sigma[\pi][m][e])) & \text{iff } \sigma[\pi][m][e] \in \Pi_C \\ \sigma[\pi][m][e] & \text{otherwise} \end{cases}$$

Let  $\sigma$  be a system, let  $\sigma' = \Psi[\alpha](\sigma)$ . Intuitively,  $\sigma'$  is obtained from  $\sigma$  simply by relabeling every correct process in every echo sample. Whenever a correct process  $\pi$  appears in an echo sample in  $\sigma$ , it is replaced with  $\zeta(d)$ ,  $d$  being the position of  $\pi$  in the ordered list of processes that  $\alpha$  causes to pb.Deliver a message, when coupled with  $\sigma$ . Byzantine processes are left unchanged.

**Roadmap.** Let  $\alpha \in \mathcal{A}_{ae}$ , let  $\alpha' = \Delta_{ps}(\alpha)$ . Let  $\sigma \in \mathcal{S}$  such that  $\alpha$  compromises the consistency of  $\sigma$ . In order to prove that  $\alpha'$  is as powerful as  $\alpha$ , we prove that:

- $\alpha'$  compromises the consistency of  $\sigma' = \Psi[\alpha](\sigma)$ .
- $\Psi[\alpha](\sigma)$  is a permutation on  $\mathcal{S}$ .

Indeed, if the above are true, then the probability of  $\alpha'$  compromising the consistency of a random system  $\sigma'$  is equal to the probability of  $\alpha$  compromising the consistency of a random system  $\sigma$ , and the lemma is proved.

**Trace.** We start by noting that, if we couple Process-sequential decorator with  $\sigma'$ , we effectively obtain a system interface  $\delta$  with which  $\alpha$  directly exchanges invocations and responses. Here we show that the trace  $\tau(\alpha, \sigma)$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\sigma$ , or it has been coupled with  $\sigma'$ , with Process-sequential decorator acting as an interface. We prove this by induction.

Let us assume

$$\begin{aligned} \tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \delta) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . Let  $\xi$  be a Byzantine process. If  $\xi \in r_{n+1}$  then, by definition,  $\xi \in \sigma[\pi][1]$ , i.e., for at least one  $e \in 1..E$ ,  $\sigma[\pi][1][e] = \xi$ . Therefore,  $\sigma'[\pi][1][e] = \xi$ , and  $\xi \in r'_{n+1}$ . If  $\xi \notin r_{n+1}$  then, for all  $e \in 1..E$ ,  $\sigma[\pi][1][e] \neq \xi$ . If  $\sigma[\pi][1][e] \in \Pi_C$ , then  $\sigma'[\pi][1][e] \in \Pi_C$  as well, so  $\sigma'[\pi][1][e] \neq \xi$ . If  $\sigma[\pi][1][e] \in \Pi \setminus \Pi_C$ , then  $\sigma'[\pi][1][e] = \sigma[\pi][1][e] \neq \xi$ . Therefore,  $\xi \notin r'_{n+1}$ .

Before considering the remaining possible values of  $i_{n+1}$ , we prove some auxiliary results. Let  $\pi$  be a correct process, let  $m$  be a message, let  $d \in 1..C$ , let  $e \in 1..E$ . For every  $j \leq n+1$ , as we established, we have  $i_j = i'_j$ . Therefore, after the  $(n+1)$ -th invocation, the following hold true:

- $\pi$  pb.Delivered  $m$  in  $\sigma$  if and only if  $\zeta(\mu(\alpha, \sigma, \pi))$  pb.Delivered  $m$  in  $\sigma'$ . Indeed:
  - If  $\pi$  pb.Delivered  $m$  in  $\sigma$ , then  $psadv.Deliver(\pi, m)$  was invoked. Moreover, by definition,  $psadv.Deliver(\pi, m)$  was the  $\mu(\alpha, \sigma, \pi)$ -th invocation of  $psadv.Deliver(\dots)$ . Noting that *cursor* is incremented at each invocation of  $psadv.Deliver(\dots)$ , when  $psadv.Deliver(\pi, m)$  was invoked we had  $cursor = \mu(\alpha, \sigma, \pi)$ . Finally,  $psadv.Deliver(\pi, m)$  forwards the call to  $sys.Deliver(\zeta(cursor), m)$ . Consequently,  $\zeta(\mu(\alpha, \sigma, \pi))$  pb.Delivered  $m$  in  $\sigma'$ .
  - If  $\zeta(\mu(\alpha, \sigma, \pi))$  pb.Delivered  $m$  in  $\sigma'$  then  $sys.Deliver(\zeta(cursor), m)$  was invoked, with  $cursor = \mu(\alpha, \sigma, \pi)$ . Noting that *cursor* is incremented after each invocation of  $sys.Deliver(\dots)$ , we have that  $psadv.Deliver(\dots)$  was invoked at least  $\mu(\alpha, \sigma, \pi)$  times. By definition, this means that  $psadv.Deliver(\pi, m)$  was invoked. Consequently,  $\pi$  pb.Delivered  $m$  in  $\sigma$ .
- If  $\pi$  pb.Delivered a message in  $\sigma'$ , then  $perm[\zeta^{-1}(\pi)] \neq \perp$ . Indeed, noting that *cursor* is incremented every time  $psadv.Deliver(\dots)$  is invoked, we have that  $\rho$  pb.Delivered a message in  $\sigma'$  as a result of the  $\zeta^{-1}(\pi)$ -th invocation of  $psadv.Deliver(\dots)$ . As a result,  $perm[\zeta^{-1}(\pi)]$  was set to a value other than  $\perp$ . From this follows that procedure  $sys.Sample(\dots)$  never calls  $\zeta(\perp)$ .
- If  $perm[d] \neq \perp$ , then  $perm[d] = \zeta^{-1}(\mu^{-1}(\alpha, \sigma, d))$ . Indeed, noting that *cursor* is incremented every time  $psadv.Deliver(\dots)$  is invoked,  $perm[d]$  was set to a value other than  $\perp$  upon the  $d$ -th invocation of  $psadv.Deliver(\dots)$ . By the definition of  $\mu$ , the  $d$ -th invocation of  $psadv.Deliver(\dots)$  is  $psadv.Deliver(\mu^{-1}(\alpha, \sigma, d), m')$ , for some  $m' \in \mathcal{M}$ .
- $\sigma[\pi][m][e]$  sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma$  if and only if  $\sigma'[\pi][m][e]$  sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma'$ . Indeed, if  $\sigma[\pi][m][e] \in \Pi_C$ , then  $\sigma'[\pi][m][e] = \zeta(\mu(\alpha, \sigma, \sigma[\pi][m][e]))$ . Therefore,  $\sigma[\pi][m][e]$  pb.Delivered  $m$  in  $\sigma$  if and only if  $\sigma'[\pi][m][e]$  pb.Delivered  $m$  in  $\sigma'$ . Noting that  $\alpha$  is an auto-echo adversary, if  $\sigma[\pi][m][e] \in \Pi \setminus \Pi_C$ , then  $\sigma'[\pi][m][e] = \sigma[\pi][m][e]$ , and both sent an Echo( $m, m$ ) message to  $\pi$  (in  $\sigma$  and  $\sigma'$ , respectively).

### Chapter 3. Sieve

---

- $\pi$  delivered  $m$  in  $\sigma$  if and only if  $\pi$  delivered  $m$  in  $\sigma'$ . This immediately follows from the above.

Let us assume  $i_{n+1} = (\text{State})$ . From the above immediately follows  $r_{n+1} = r'_{n+1}$ .

Let us assume  $i_{n+1} = (\text{Sample}, \pi, m)$ . Let  $\rho$  be a process. The following hold true:

- If  $\rho \in r_{n+1}$ , then  $\rho \in r'_{n+1}$ . Indeed, if  $\rho \in \Pi_C$ , then  $\rho \in \sigma[\pi][m]$  and  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma$ , for some  $m' \in \mathcal{M}$ . Therefore,  $\rho$  delivered  $m'$  in  $\sigma$ . By definition,  $\zeta(\mu(\alpha, \sigma, \rho)) \in \sigma'[\pi][m]$ . Moreover,  $\zeta(\mu(\alpha, \sigma, \rho))$  delivered  $m'$  in  $\sigma'$  and, as a result, it sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma'$ . Therefore,  $\zeta(\mu(\alpha, \sigma, \rho)) \in \text{sys.Sample}(\pi, m)$ . Finally,  $\text{perm}[\mu(\alpha, \sigma, \rho)] = \zeta^{-1}(\rho)$ . Consequently,  $\rho \in r'_{n+1}$ . If  $\rho \in \Pi \setminus \Pi_C$  then  $\rho \in \sigma[\pi][m]$  and  $\rho \in \sigma'[\pi][m]$ . Moreover,  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$ , for some  $m' \in \mathcal{M}$ , both in  $\sigma$  and  $\sigma'$ . Consequently,  $\rho \in r'_{n+1}$ .
- If  $\rho \in r'_{n+1}$ , then  $\rho \in r_{n+1}$ . Indeed, if  $\rho \in \Pi_C$ , then  $\zeta(\text{perm}^{-1}[\zeta^{-1}(\rho)])^1$  was returned from  $\text{sys.Sample}(\pi, m)$ , in other words  $\zeta(\text{perm}^{-1}[\zeta^{-1}(\rho)])$  pb.Delivered some message  $m' \in \mathcal{M}$  in  $\sigma'$ . Moreover, using our auxiliary result on  $\text{perm}$  we obtain

$$\zeta(\text{perm}^{-1}[\zeta^{-1}(\rho)]) = \zeta(\mu(\alpha, \sigma, \rho))$$

therefore  $\zeta(\mu(\alpha, \sigma, \rho))$  pb.Delivered  $m'$  in  $\sigma'$ , and  $\rho$  pb.Delivered  $m'$  in  $\sigma$ . Finally, since  $\zeta(\mu(\alpha, \sigma, \rho)) \in \sigma'[\pi][m]$ , then by definition  $\rho \in \sigma[\pi][m]$ . Consequently,  $\rho \in r_{n+1}$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ .

**Consistency of  $\sigma'$ .** We proved that  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ . Moreover, we proved that if a correct process  $\pi$  eventually delivers a message  $m$  in  $\sigma$ , then  $\zeta(\mu(\alpha, \sigma, \pi))$  also delivers  $m$  in  $\sigma'$ .

Since  $\alpha$  compromises the consistency of  $\sigma$ , two correct processes  $\pi, \pi'$  and two distinct messages  $m, m' \neq m$  exist such that, in  $\sigma$ ,  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$ . Therefore, in  $\sigma'$ ,  $\zeta(\mu(\alpha, \sigma, \pi))$  delivered  $m$  and  $\zeta(\mu(\alpha, \sigma, \pi'))$  delivered  $m'$ . Therefore  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Translation permutation.** We now prove that, for any two  $\sigma_a, \sigma_b \neq \sigma_a$ , we have  $\Psi[\alpha](\sigma_a) \neq \Psi[\alpha](\sigma_b)$ . We prove this by contradiction. Suppose a system  $\sigma'$  exists such that  $\sigma' = \Psi[\alpha](\sigma_a) = \Psi[\alpha](\sigma_b)$ . We want to prove that  $\sigma_a = \sigma_b$ .

---

<sup>1</sup>Noting that  $\text{perm}$  is injective, we define  $\text{perm}^{-1}[b] = a \iff \text{perm}[a] = b$ .

We start by noting that, if  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ , then  $\sigma_a = \sigma_b$ . Indeed, if  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ , then for every  $\pi \in \Pi_C$  and every  $d \in 1..C$  we have

$$\begin{aligned}\mu(\alpha, \sigma_a, \pi) &= \mu(\alpha, \sigma_b, \pi) \\ \mu^{-1}(\alpha, \sigma_a, d) &= \mu^{-1}(\alpha, \sigma_b, d)\end{aligned}$$

and since, by definition, for every  $\pi \in \Pi_C$ ,  $m \in \mathcal{M}$  and  $e \in 1..E$ , we have

$$\begin{aligned}\sigma_a[\pi][m][e] &= \begin{cases} \mu^{-1}(\alpha, \sigma_a, \zeta^{-1}(\sigma'[\pi][m][e])) & \text{iff } \sigma'[\pi][m][e] \in \Pi_C \\ \sigma'[\pi][m][e] & \text{otherwise} \end{cases} \\ \sigma_b[\pi][m][e] &= \begin{cases} \mu^{-1}(\alpha, \sigma_b, \zeta^{-1}(\sigma'[\pi][m][e])) & \text{iff } \sigma'[\pi][m][e] \in \Pi_C \\ \sigma'[\pi][m][e] & \text{otherwise} \end{cases}\end{aligned}$$

we get

$$\sigma_a[\pi][m][e] = \sigma_b[\pi][m][e]$$

and  $\sigma_a = \sigma_b$ .

We prove that  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$  by induction. Let us assume

$$\begin{aligned}\tau(\alpha, \sigma_a) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \sigma_b) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n\end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . Let  $\xi$  be a Byzantine process. if  $\xi \in r_{n+1}$ , then for at least one  $e \in 1..E$  we have  $\sigma_a[\pi][m][e] = \xi$ . Therefore,  $\sigma'[\pi][m][e] = \xi$ , and  $\sigma_b[\pi][m][e] = \xi$ . Consequently,  $\xi \in r'_{n+1}$ . The argument can be reversed to prove  $\xi \in r'_{n+1} \implies \xi \in r_{n+1}$ .

Before considering the remaining possible values of  $i_{n+1}$ , we prove some auxiliary result. Let  $\pi$  be a correct process, let  $m$  be a message, let  $e \in 1..E$ . For every  $j \leq n+1$ , as we established, we have  $i_j = i'_j$ . Therefore, after the  $(n+1)$ -th invocation, the following hold true:

- $\pi$  pb.Delivered  $m$  in  $\sigma_a$  if and only if  $\pi$  pb.Delivered  $m$  in  $\sigma_b$ . Indeed, if  $\pi$  pb.Delivered  $m$  in  $\sigma_a$ , then some  $j \leq (n+1)$  exists such that  $i_j = (\text{Deliver}, \pi, m)$ . Since  $i'_j = i_j$ ,

$\pi$  pb.Delivered  $m$  in  $\sigma_b$  as well. The argument can be reversed to prove that, if  $\pi$  pb.Delivered  $m$  in  $\sigma_b$ , then  $\pi$  pb.Delivered  $m$  in  $\sigma_a$  as well.

- If  $\pi$  pb.Delivered  $m$  in  $\sigma_a$  (or, equivalently,  $\sigma_b$ ), then  $\mu(\alpha, \sigma_a, \pi) = \mu(\alpha, \sigma_b, \pi)$ . Indeed, some  $j \leq (n + 1)$  exists such that  $i_j = i'_j = (\text{Deliver}, \pi, m)$ . Since, for all  $h < j$ , we also have  $i_h = i'_h$ , then

$$\begin{aligned} & |\{h \in 1..(j-1) \mid i_h = (\text{Deliver}, \pi' \in \Pi_C, m' \in \mathcal{M})\}| \\ &= |\{h \in 1..(j-1) \mid i'_h = (\text{Deliver}, \pi' \in \Pi_C, m' \in \mathcal{M})\}| \end{aligned}$$

- $\sigma_a[\pi][m][e]$  sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma_a$  if and only if  $\sigma_b[\pi][m][e]$  sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma_b$ . We prove this by cases:

- Let us assume that  $\sigma_a[\pi][m][e]$  is correct, and pb.Delivered  $m$  in  $\sigma_a$ . By definition, we have

$$\begin{aligned} \sigma'[\pi][m][e] &= \zeta(\mu(\alpha, \sigma_a, \sigma_a[\pi][m][e])) \\ \sigma'[\pi][m][e] &= \zeta(\mu(\alpha, \sigma_b, \sigma_b[\pi][m][e])) \end{aligned}$$

and from the above we have

$$\zeta(\mu(\alpha, \sigma_a, \sigma_a[\pi][m][e])) = \zeta(\mu(\alpha, \sigma_b, \sigma_b[\pi][m][e]))$$

Equating the two above we get

$$\zeta(\mu(\alpha, \sigma_b, \sigma_a[\pi][m][e])) = \zeta(\mu(\alpha, \sigma_b, \sigma_b[\pi][m][e]))$$

and noting that  $\mu$  is always injective, we have  $\sigma_a[\pi][m][e] = \sigma_b[\pi][m][e]$ . Therefore  $\sigma_b[\pi][m][e]$  pb.Delivered  $m$  in  $\sigma_b$ .

The argument can be inverted to prove that, if  $\sigma_b[\pi][m][e]$  is correct, and pb.Delivered  $m$  in  $\sigma_b$ , then  $\sigma_a[\pi][m][e]$  pb.Delivered  $m$  in  $\sigma_a$  as well.

- Let us assume that  $\sigma_a[\pi][m][e]$  is correct, but did not pb.Deliver  $m$ . From the definition of  $\Psi[\alpha]$ , we know that  $\sigma_b[\pi][m][e]$  is correct as well. By contradiction, following from the above, we have that if  $\sigma_b[\pi][m][e]$  pb.Delivered  $m$  in  $\sigma_b$ ,  $\sigma_a[\pi][m][e]$  would have pb.Delivered  $m$  in  $\sigma_a$  as well.

The argument can be inverted to prove that, if  $\sigma_b[\pi][m][e]$  is correct, but did not pb.Deliver  $m$  in  $\sigma_b$ , then  $\sigma_a[\pi][m][e]$  did not pb.Deliver  $m$  in  $\sigma_a$  either.

- Let us assume that  $\sigma_a[\pi][m][e]$  is Byzantine. Then, from the definition of  $\Psi[\alpha]$ , we immediately have  $\sigma_b[\pi][m][e] = \sigma_a[\pi][m][e]$  and, since  $\alpha$  is an auto-echo adversary, both sent an Echo( $m, m$ ) message to  $\pi$  (in their respective systems).
- $\pi$  delivered  $m$  in  $\sigma_a$  if and only if  $\pi$  delivered  $m$  in  $\sigma_b$  as well. This follows immediately from the above.

Let us assume  $i_{n+1} = (\text{State})$ . From the above immediately follows  $r_{n+1} = r'_{n+1}$ .

Let us assume  $i_{n+1} = (\text{Sample}, \pi, m)$ . Let  $\rho$  be a process. If  $\rho \in r_{n+1}$ , then for some  $e \in 1..E$ ,  $\sigma_a[\pi][m][e] = \rho$ , and  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma_a$ , for some  $m' \in \mathcal{M}$ . Following from the above, we have  $\sigma_b[\pi][m][e] = \rho$  as well, and  $\rho$  sent an  $\text{Echo}(m, m')$  message to  $\pi$  in  $\sigma_b$  as well. Therefore,  $\rho \in r'_{n+1}$ . The argument can be inverted to prove that, if  $\rho \in r'_{n+1}$ , then  $\rho \in r_{n+1}$  as well.

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ .

Therefore,  $\sigma_a = \sigma_b$ , which contradicts the hypothesis.  $\square$

### 3.11.3 Sequential adversary

**Lemma 22.** *The set of sequential adversaries  $\mathcal{A}_{sq}$  is optimal.*

*Proof.* We again prove the result using a decorator. Here we show that a decorator  $\Delta_{sq}$  exists such that, for every  $\alpha \in \mathcal{A}_{ps}$ , the adversary  $\alpha' = \Delta_{sq}(\alpha)$  is a sequential adversary, and as powerful as  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the sequential  $\alpha^+ = \Delta_{sq}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 9 implements **Sequential decorator**, a decorator that transforms a process-sequential adversary into a sequential adversary. Provided with a process-sequential adversary  $psadv$ , Sequential decorator acts as an interface between  $psadv$  and a system  $sys$ , effectively implementing a sequential adversary  $sqadv$ . Sequential decorator exposes both the adversary and the system interfaces: the underlying adversary  $psadv$  uses  $sqadv$  as its system.

Sequential decorator works as follows:

- Procedure  $sqadv.Init()$  initializes the following variables:
  - A  $perm$  array of  $C$  elements:  $perm$  is used to consistently translate messages between  $psadv$  and  $sys$ .
  - A  $cursor$  variable, initially set to 1: at any time,  $cursor$  identifies the next message that will be pb.Delivered in  $sys$ , if  $psadv$  will invoke the delivery of a process whose delivery  $psadv$  never invoked before.
  - A  $poisoned$  variable:  $poisoned$  is set to True if and only if at least one correct process in  $sys$  is poisoned. This condition is verified by looping over  $sys.Byzantine(\pi)$  for every correct process  $\pi$ .

## Chapter 3. Sieve

---

### Algorithm 9 Sequential decorator.

---

```
1: Implements:
2:   SequentialAdversary + CobSystem, instance sqadv
3:
4: Uses:
5:   ProcessSequentialAdversary, instance psadv, system sqadv
6:   CobSystem, instance sys
7:
8: procedure sqadv.Init() is
9:    $perm = \{\perp\}^C$ ;    $cursor = 1$ ;    $step = 0$ ;
10:
11:    $poisoned = \text{False}$ ;
12:   for all  $\pi \in \Pi_C$  do
13:     if  $|sys.Byzantine(\pi)| \geq \hat{E}$  then
14:        $poisoned \leftarrow \text{True}$ ;
15:     end if
16:   end for
17:
18:   psadv.Init();
19:
20: procedure sqadv.Step() is
21:    $step \leftarrow step + 1$ ;
22:
23:   if  $poisoned = \text{False}$  or  $step \leq (N - C)C^2$  then
24:     psadv.Step();
25:   else if  $step \leq (N - C)C^2 + C$  then
26:      $sys.Deliver(\zeta(step - (N - C)C^2), 1)$ ;
27:   else
28:      $sys.End()$ ;
29:   end if
30:
31: procedure sqadv.Byzantine(process) is
32:   return  $sys.Byzantine(process)$ ;
33:
34: procedure sqadv.State() is
35:    $state = \emptyset$ ;
36:
37:   for all  $(\pi, m) \in sys.State()$  do
38:      $state \leftarrow state \cup \{(\pi, perm[m])\}$ ;
39:   end for
40:
41:   return  $state$ ;
42:
```

---



---

```

43: procedure sqadv.Sample(process, message) is
44:   return sys.Sample(process, perm-1[message]);
45:
46: procedure sqadv.Deliver(process, message) is
47:   if message ∈ perm then
48:     sys.Deliver(process, perm-1[message]);
49:   else
50:     perm[cursor] = message;
51:     sys.Deliver(process, cursor);
52:     cursor ← cursor + 1;
53:   end if
54:
55: procedure sqadv.Echo(process, sample, source, message) is
56:   sys.Echo(process, sample, source, message);
57:
58: procedure sqadv.End() is
59:   sys.End();
60:

```

---

- A *step* variable, initially set to 0: at any time, *step* counts how many times *sqadv.Step*() has been invoked.
- Procedure *sqadv.Step*() increments *step*, then implements two different behaviors depending on the value of *poisoned*:
  - If *poisoned* = True, it forwards the call to *psadv.Step*() for the first  $(N - C)C^2$  times. For the next  $C$  steps, it sequentially invokes *sys.Deliver*( $\zeta(1), 1$ ), ..., *sys.Deliver*( $\zeta(C), 1$ ). Finally, it calls *sys.End*().
  - If *poisoned* = False, it forwards the call to *psadv.Step*().
- Procedure *sqadv.Byzantine*(*process*) simply forwards the call to *sys.Byzantine*(*process*).
- Procedure *sqadv.State*() returns the list of process / message pairs returned by *sys.State*(), with each message translated through *perm*. More specifically, *sqadv.State*() returns  $(\pi, perm[m])$  for every  $(\pi, m)$  in *sys.State*().
- Procedure *sqadv.Sample*(*process, message*) simply forwards the call to *sys.Sample*(*process, perm*<sup>-1</sup>[*message*]).
- Procedure *sqadv.Deliver*(*process, message*) checks if *psadv* has already invoked the delivery of *message* (this is achieved by checking if *message* is in *perm*). If so, it forwards the call to *sys.Deliver*(*process, perm*<sup>-1</sup>[*message*]). Otherwise, it sets *perm*[*cursor*] to *message*, then forwards the call to *sys.Deliver*(*process, cursor*). Finally, it increments *cursor*. This mechanism serves two purposes:

### Chapter 3. Sieve

---

- To consistently translate a  $sqadv.Deliver(\dots)$  invocation to a  $sys.Deliver(\dots)$  invocation. More specifically, the set of invocations  $psadv.Deliver(\pi_1, m), \dots, psadv.Deliver(\pi_k, m)$  is always translated to  $sys.Deliver(\pi_1, m'), \dots, sys.Deliver(\pi_k, m')$ .
- To never cause the pb.Delivery of a message  $b$  in  $sys$  before every message  $a < b$  has been pb.Delivered in  $sys$  at least once.
- Procedure  $sqadv.Echo(process, sample, source, message)$  simply forwards the call to  $sys.Echo(process, sample, source, message)$ .
- Procedure  $sqadv.End()$  simply forwards the call to  $sys.End()$ .

**Correctness.** We start by proving that no adversary, coupled with Sequential decorator, causes the execution to fail. We distinguish two cases, based on the value of  $poisoned$ .

Let us assume  $poisoned = \text{True}$ . When  $sqadv.Step()$  is invoked, the call is forwarded to  $psadv.Step()$  only for the first  $(N - C)C^2$  times. Noting that  $psadv$  is an auto-echo adversary, every call to  $psadv.Step()$  results in a call to  $sqadv.Echo(\dots)$ . For the next  $C$  steps,  $sqadv.Step()$  sequentially causes  $\zeta(1), \zeta(2), \dots$  to pb.Deliver message 1. Finally,  $sqadv.Step()$  invokes  $sys.End()$ . Therefore,  $sqadv$  never causes the execution to fail, and it implements a process-sequential adversary.

Let us assume  $poisoned = \text{False}$ . Let  $\pi$  be a correct process, let  $m$  be a message. The following hold true:

- Procedure  $sqadv.State()$  never returns a  $(\pi, \perp)$  pair. Indeed, if  $(\pi, m) \in sys.State()$ , then  $\pi$  pb.Delivered  $m$  in  $sys$ . Since  $\pi$  is not poisoned,  $\pi$  received at least one  $Echo(m, m)$  message from a correct process. Consequently, if  $(\pi, m)$  is returned from  $sys.State()$ , then at least one correct process pb.Delivered  $m$  in  $sys$ , i.e.,  $sys.Deliver(\pi', m)$  was invoked for some  $\pi' \in \Pi_C$ . The statement is proved by noting that, whenever  $sys.Deliver(\pi', m)$  is invoked for some  $\pi' \in \Pi_C$ , we have  $perm[m] \neq \perp$ : indeed, either  $sys.Sample(process, perm^{-1}[message])$  is invoked, and  $message \in perm$ , or  $sys.Sample(process, cursor)$  is invoked, and  $perm[cursor] = message \neq \perp$ .
- No invocation of  $sqadv.Sample(\dots)$  causes the execution to fail. Noting that  $psadv$  is correct, it will never invoke  $sqadv.Sample(\pi, m)$  unless  $(\pi', m)$  was returned from a previous invocation of  $sqadv.State()$ , for some  $\pi' \in \Pi_C$ . Since  $\pi'$  is not poisoned,  $(\pi', perm^{-1}[m])$  was returned from  $sys.State()$ , therefore  $\pi'$  delivered  $perm^{-1}[m]$  in  $sys$ . Therefore  $sys.Sample(\pi, m)$  is never invoked from  $sqadv.Sample(\dots)$  unless at least one correct process delivered  $m$  in  $sys$ .

**Sequential.** It is easy to prove that Sequential decorator always implements a sequential adversary. Indeed, if  $poisoned = \text{True}$ ,  $sqadv$  simply causes every correct process to pb.Deliver

message 1 (which trivially implements a sequential adversary). If  $poisoned = \text{False}$ , then whenever  $sys.Deliver(\pi, m)$  is invoked, either of the following holds true:

- $m = perm^{-1}[message]$  for some  $message \in perm$ . In this case  $sys.Deliver(\dots)$  was previously invoked on  $m$  (i.e., some process  $\pi'$  exists such that  $sys.Deliver(\pi', m)$  was previously invoked).
- $m = cursor$ . Then  $sys.Deliver(\dots)$  was never invoked on  $m$ . Noting that, whenever  $sys.Deliver(\dots)$  is invoked on a new message,  $cursor$  is incremented, we have that every message  $l < m$  was previously pb.Delivered by at least one correct process in  $sys$ .

**System translation.** Let  $\alpha$  be an adversary. We can define a function

$$\mu : \mathcal{A} \times \mathcal{S} \times \mathcal{M} \rightarrow 1..C \cup \{\perp\}$$

such that:

- $\mu(\alpha, \sigma, m) = (d \in 1..C)$  if and only if  $m$  is the  $d$ -th distinct message that  $\alpha$  causes at least one correct process to pb.Deliver, when  $\alpha$  is coupled with  $\sigma$ .
- $\mu(\alpha, \sigma, m) = \perp$  if and only if  $\alpha$  never causes any correct process to pb.Deliver  $m$ , when  $\alpha$  is coupled with  $\sigma$ .

We additionally define  $v : \mathcal{A} \times \mathcal{S} \rightarrow 1..C$  by

$$v(\alpha, \sigma) = \max_{m \in \mathcal{M}} \mu(\alpha, \sigma, m)$$

and

$$(\mu^{-1}(\alpha, \sigma, d) = m) \stackrel{def}{\iff} (\mu(\alpha, \sigma, m) = d)$$

for all  $d \leq v(\alpha, \sigma)$ . Here  $v(\alpha, \sigma)$  counts the number of distinct messages that  $\alpha$  causes at least one correct process to pb.Deliver, when coupled with  $\sigma$ . It is immediate to see that  $\mu(\alpha, \sigma, d) = \perp$  for all  $d > v(\alpha, \sigma)$ .

We then define a **message permutation function**  $\chi : \mathcal{A} \times \mathcal{S} \times \mathcal{M} \rightarrow \mathcal{M}$  as follows:

$$\chi(\alpha, \sigma, d) = \begin{cases} \mu^{-1}(\alpha, \sigma, d) & \text{iff } d \leq v(\alpha, \sigma) \\ \max m \in \mathcal{M} \mid |\{l \leq m : \mu(\alpha, \sigma, l) = \perp\}| = d - v(\alpha, \sigma) & \text{otherwise} \end{cases}$$

For a given  $\alpha$  and  $\sigma$ , the permutation  $\chi$  maps  $d$  to the  $d$ -th distinct message that is pb.Delivered when  $\alpha$  is coupled with  $\sigma$ , if such a message exists. If such a message does not exist,  $\chi$  simply enumerates sequentially the messages that are never pb.Delivered when  $\alpha$  is

### Chapter 3. Sieve

---

coupled with  $\sigma$ .

For example, let us consider the case where  $C = 10$  and  $\alpha$  coupled with  $\sigma$  causes the pb.Delivery of messages 3, 7, 1, 4 (in this order of first appearance). Then  $\chi$  will assume the following values for  $d \in 1..C$ : 3, 7, 1, 4, 2, 5, 6, 8, 9, 10.

Finally, we define a **system translation function**  $\Psi[\alpha]: \mathcal{S} \rightarrow \mathcal{S}$  such that, for system  $\sigma$ , every correct process  $\pi$  and every message  $m$ ,

$$\Psi[\alpha](\sigma)[\pi][m] = \begin{cases} \sigma[\pi][m] & \text{iff } \exists \pi' \in \Pi_C \mid \pi' \text{ is poisoned in } \sigma \\ \sigma[\pi][\chi(\alpha, \sigma, m)] & \text{otherwise} \end{cases}$$

Let  $\sigma$  be a system, let  $\sigma' = \Psi[\alpha](\sigma)$ . Intuitively, if at least one correct process is poisoned in  $\sigma$ , then  $\sigma' = \sigma$ . Otherwise,  $\sigma'$  is obtained from  $\sigma$  by permuting the echo samples of each correct process in  $\sigma$  using  $\chi$ .

**Roadmap.** Let  $\alpha \in \mathcal{A}_{ps}$ , let  $\alpha' = \Delta_{sq}(\alpha)$ . Let  $\sigma \in \mathcal{S}$  such that  $\alpha$  compromises the consistency of  $\sigma$ . In order to prove that  $\alpha'$  is as powerful as  $\alpha$ , we prove that:

- $\alpha'$  compromises the consistency of  $\sigma' = \Psi[\alpha](\sigma)$ .
- $\Psi[\alpha](\sigma)$  is a permutation on  $\mathcal{S}$ .

Indeed, if the above are true, then the probability of  $\alpha'$  compromising the consistency of a random system  $\sigma'$  is equal to the probability of  $\alpha$  compromising the consistency of a random system  $\sigma$ , and the lemma is proved.

**Poisoned case.** We start by considering the case where *poisoned* = True. Let  $\pi$  be a correct process that is poisoned in  $\sigma$ . Noting that *psadv* is an auto-echo adversary,  $\pi$  eventually delivers every message. Indeed, every Byzantine process eventually sends to  $\pi$  an Echo( $m$ ,  $m$ ) message, for every  $m \in \mathcal{M}$ . Since all of  $\pi$ 's echo samples share the same set of at least  $\hat{E}$  Byzantine processes,  $\pi$  eventually delivers every message.

As a result, if at least one correct process in  $\sigma$  is poisoned, the consistency of  $\sigma$  is compromised by any auto-echo adversary. Noting that  $\sigma' = \Psi[\alpha](\sigma) = \sigma$ , and  $\Delta_{sq}(\alpha)$  is an auto-echo adversary, we immediately have that  $\alpha'$  compromises the consistency of  $\sigma'$  as well.

In the next sections of this proof, we consider the case *poisoned* = False.

**Trace.** We start by noting that, if we couple Process-sequential decorator with  $\sigma'$ , we effectively obtain a system interface  $\delta$  with which  $\alpha$  directly exchanges invocations and responses. Here we show that, if  $poisoned = \text{False}$ , the trace  $\tau(\alpha, \sigma)$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that, if  $poisoned = \text{False}$ ,  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\sigma$ , or it has been coupled with  $\sigma'$ , with Process-sequential decorator acting as an interface. We prove this by induction.

Let us assume  $poisoned = \text{False}$ , and

$$\begin{aligned} \tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \delta) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . We can note that  $sqadv.\text{Byzantine}(\text{process})$  simply forwards the call to  $sys.\text{Byzantine}(\text{process})$ , and  $\chi$  defines a permutation over  $1..C$ . Therefore, a message  $m$  exists such that  $\pi$ 's first echo sample for in  $\sigma'$  is identical to  $\pi$ 's echo sample for  $m$  in  $\sigma$ . Moreover, all of  $\pi$ 's echo samples in  $\sigma$  share the same set of Byzantine processes. Consequently,  $r_{n+1} = r'_{n+1}$ .

Before considering the remaining possible values of  $i_{n+1}$ , we prove some auxiliary results. We start by noting the following:

- Let  $d \in 1..C$ . At any time, if  $perm[d] \neq \perp$ , then  $perm[d] = \mu^{-1}(\alpha, \sigma, d)$ . Indeed, at any time, a message  $m$  is in  $perm$  if and only if  $sqadv.\text{Deliver}(\dots)$  was previously invoked on  $m$ . Moreover, whenever  $sqadv.\text{Deliver}(\dots)$  is invoked on a message  $m$  that is not in  $perm$ ,  $m$  is added to  $perm$  and  $cursor$  is incremented. Therefore  $perm[cursor]$  is set to  $m$  if and only if  $sqadv.\text{Deliver}(\dots)$  was never invoked on  $m$ , and  $sqadv.\text{Deliver}(\dots)$  was previously invoked on exactly  $cursor - 1$  distinct messages. Moreover, by definition, when  $sqadv.\text{Deliver}(\dots)$  is invoked on  $m$  for the first time,  $sqadv.\text{Deliver}(\dots)$  was previously invoked on exactly  $\mu(\alpha, \sigma, m) - 1$  distinct messages. Consequently,  $cursor = \mu(\alpha, \sigma, m)$ , and  $m = \mu^{-1}(\alpha, \sigma, cursor)$ .
- No two values of  $perm$  are equal to each other. Indeed, a message  $m$  is added to  $perm$  only if  $m \notin perm$ .
- Let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ . If  $\pi$  delivered  $m$ , then at least one correct process pb.Delivered  $m$ .

This separately holds true both in  $\sigma$  and  $\sigma'$ . Indeed, if  $\pi$  delivered  $m$ , then it received at least  $\hat{E}$   $\text{Echo}(m, m)$  messages from its echo sample for  $m$  and, since no correct process is poisoned in neither  $\sigma$  nor  $\sigma'$ , at least one of them must have come from a correct process.

Let  $\pi$  be a correct process, let  $\rho$  be a process, let  $m, s$  be messages. For every  $j \leq n + 1$ , as we established, we have  $i_j = i'_j$ . By hypothesis,  $\alpha$  is an auto-echo adversary, so  $i_j = i'_j, r_j = r'_j = \perp$  for every  $j \leq (N - C)C^2$ . Let us consider the non-trivial case  $n \geq (N - C)C^2$ . After the  $(n + 1)$ -th invocation, the following hold true:

- $\pi$  pb.Delivered  $m$  in  $\sigma$  if and only if  $\pi$  pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Indeed:
  - If  $\pi$  pb.Delivered  $m$  in  $\sigma$ , then  $\text{sqadv.Deliver}(\pi, m)$  was invoked. If  $\text{sqadv.Deliver}(\pi, m)$  was the first invocation of  $\text{sqadv.Deliver}(\dots)$  on  $m$ , then  $m$  was not in  $\text{perm}$ ,  $\text{perm}[\text{cursor}]$  was set to  $m$ , and  $\text{sys.Deliver}(\pi, \text{cursor})$  was invoked. As we previously proved, however, we have  $\text{perm}[\text{cursor}] = \mu^{-1}(\alpha, \sigma, m)$ , so  $\text{cursor} = \mu(\alpha, \sigma, m)$ . Consequently,  $\text{sys.Deliver}(\pi, \mu(\alpha, \sigma, m))$  was invoked, and  $\pi$  pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . If  $\text{sqadv.Deliver}(\pi, m)$  was not the first invocation of  $\text{sqadv.Deliver}(\dots)$  on  $m$ , then  $m$  was in  $\text{perm}$ , and  $\text{sys.Deliver}(\pi, \text{perm}^{-1}(m))$  was invoked. Due to the above, we have again  $\text{perm}^{-1}[m] = \mu(\alpha, \sigma, m)$ . Consequently,  $\text{sys.Deliver}(\pi, \mu(\alpha, \sigma, m))$  was invoked, and  $\pi$  pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ .
  - If  $\pi$  pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ , then  $\text{sys.Deliver}(\pi, \mu(\alpha, \sigma, m))$  was invoked. If  $\text{sys.Deliver}(\pi, \text{cursor})$  was invoked, we have that  $\text{cursor} = \mu(\alpha, \sigma, m)$ , and  $\text{sqadv.Deliver}(\pi, m')$  was invoked for some  $m' \notin \text{perm}$ . As a result,  $\text{perm}[\text{cursor}]$  was set to  $m'$ . As we previously established, however,

$$m' = \mu^{-1}(\alpha, \sigma, \text{cursor}) = \mu^{-1}(\alpha, \sigma, \mu(\alpha, \sigma, m)) = m$$

and  $\text{sqadv.Deliver}(\pi, m)$  was invoked. As a result,  $\pi$  pb.Delivered  $m$  in  $\sigma$ . If  $\text{sys.Deliver}(\pi, \text{perm}^{-1}(m'))$  was invoked for some  $m' \in \text{perm}$ , we have  $\text{perm}^{-1}[m'] = \mu(\alpha, \sigma, m)$ , and again  $m' = m$ . Consequently,  $\text{sqadv.Deliver}(\pi, m)$  was invoked, and  $\pi$  pb.Delivered  $m$  in  $\sigma$ .

- $\pi$  received an  $\text{Echo}(m, m)$  message from  $\rho$  in  $\sigma$  if and only if  $\pi$  received an  $\text{Echo}(\mu(\alpha, \sigma, m), \mu(\alpha, \sigma, m))$  message from  $\rho$  in  $\sigma'$ . Indeed:
  - If  $\rho$  is a correct process, from the above we have that  $\pi$  pb.Delivered  $m$  in  $\sigma$  if and only if  $\pi$  pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Therefore,  $\rho$  sent to  $\pi$  an  $\text{Echo}(m, m)$  message if and only if  $\rho$  sent to  $\pi$  an  $\text{Echo}(\mu(\alpha, \sigma, m), \mu(\alpha, \sigma, m))$  message.
  - If  $\rho$  is a Byzantine process then, noting that  $\alpha$  is an auto-echo adversary,  $\rho$  sent to  $\pi$  an  $\text{Echo}(m, m)$  both in  $\sigma$  and  $\sigma'$ .

- $\pi$  received an  $\text{Echo}(s, m')$  message for some  $m' \in \mathcal{M}$  from  $\rho$  in  $\sigma$  if and only if  $\pi$  received an  $\text{Echo}(\mu(\alpha, \sigma, s), m'')$  message for some  $m'' \in \mathcal{M}$  from  $\rho$  in  $\sigma'$ . Indeed:
  - If  $\rho$  is correct, it sent an  $\text{Echo}(s', m')$  message for every  $s' \in \mathcal{S}$  and some  $m' \in \mathcal{S}$  to  $\pi$  in  $\sigma$  if and only if  $\rho$  pb.Delivered a message in  $\sigma$ . Moreover,  $\rho$  pb.Delivered a message in  $\sigma$  if and only if  $\rho$  pb.Delivered a message in  $\sigma'$ . Finally,  $\rho$  pb.Delivered a message in  $\sigma'$  if and only if  $\rho$  sent an  $\text{Echo}(s'', m'')$  message for every  $s'' \in \mathcal{S}$  and some  $m'' \in \mathcal{S}$  to  $\pi$  in  $\sigma'$ .
  - If  $\rho$  is Byzantine, it sent an  $\text{Echo}(m', m')$  message for every  $m' \in \mathcal{S}$  both in  $\sigma$  and  $\sigma'$ .
- $\pi$  delivered  $m$  in  $\sigma$  if and only if  $\pi$  delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Indeed, if  $\pi$  delivered  $m$  in  $\sigma$ , then at least one correct process pb.Delivered  $m$  in  $\sigma$ , and at least one correct process pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ ; if  $\pi$  delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ , then at least one correct process pb.Delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ , and at least one correct process pb.Delivered  $m$  in  $\sigma$ . Following from the definition of  $\chi$ ,  $\pi$ 's echo sample for  $m$  in  $\sigma$  is identical to  $\pi$ 's echo sample for  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Moreover,  $\pi$  received an  $\text{Echo}(m, m)$  message from  $\rho$  in  $\sigma$  if and only if  $\pi$  received an  $\text{Echo}(\mu(\alpha, \sigma, m), \mu(\alpha, \sigma, m))$  message from  $\rho$  in  $\sigma'$ . Therefore  $\pi$  delivered  $m$  in  $\sigma$  if and only if  $\pi$  delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ .

Let us assume  $i_{n+1} = (\text{State})$ . Let  $\pi$  be a correct process, let  $m$  be a message. The following hold true:

- If  $(\pi, m) \in r_{n+1}$ , then  $(\pi, m) \in r'_{n+1}$ . Indeed,  $\pi$  delivered  $m$  in  $\sigma$ , therefore  $\pi$  delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Moreover,  $\text{sqadv.Deliver}(\dots)$  was invoked at least once on  $m$ , and  $\text{perm}[\mu(\alpha, \sigma, m)] = m$ . Finally,  $\text{perm}[\mu(\alpha, \sigma, m)]$  was returned from  $\text{sqadv.State}()$ , i.e.,  $(\pi, m) \in r'_{n+1}$ .
- If  $(\pi, m) \in r'_{n+1}$ , then  $\pi$  delivered  $\text{perm}^{-1}[m]$  in  $\sigma'$ . Since  $\text{perm}[m] = \mu^{-1}(\alpha, \sigma, m)$ ,  $\pi$  delivered  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Therefore  $\pi$  delivered  $m$  in  $\sigma$ , and  $(\pi, m) \in r_{n+1}$ .

Let us assume  $i_{n+1} = (\text{Sample}, \pi, m)$ . At least one correct process delivered  $m$  in  $\sigma$ . Since no correct process is poisoned, at least one correct process pb.Delivered  $m$  in  $\sigma$ , and  $\text{sqadv.Deliver}(\dots)$  was invoked at least once on  $m$ . Therefore,  $\text{perm}[m] = \mu^{-1}(\alpha, \sigma, m)$ . Moreover, from the definition of  $\chi$ , we have that  $\pi$ 's echo sample for  $m$  in  $\sigma$  is identical to  $\pi$ 's echo sample for  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ . Finally, every process that sent an  $\text{Echo}(s, m')$  for some  $m' \in \mathcal{M}$  to  $\pi$  in  $\sigma$  sent an  $\text{Echo}(\mu(\alpha, \sigma, s), m'')$  for some  $m'' \in \mathcal{M}$  to  $\pi$  in  $\sigma'$ . Since  $\text{sqadv.Sample}(\pi, m)$  forwards the call to  $\text{sys.Sample}(\pi, \text{perm}^{-1}(m))$ , we again have  $r_{n+1} = r'_{n+1}$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ .

### Chapter 3. Sieve

---

**Consistency of  $\sigma'$ .** We proved that, if *poisoned* = False, then  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ . Moreover, we proved that if a correct process  $\pi$  eventually delivers a message  $m$  in  $\sigma$ , then  $\pi$  delivers  $\mu(\alpha, \sigma, m)$  in  $\sigma'$ .

Since  $\alpha$  compromises the consistency of  $\sigma$ , two correct processes  $\pi, \pi'$  and two distinct messages  $m, m' \neq m$  exist such that, in  $\sigma$ ,  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$ . Therefore, in  $\sigma'$ ,  $\pi$  delivered  $\mu(\alpha, \sigma, m)$  and  $\pi'$  delivered  $\mu(\alpha, \sigma, m') \neq \mu(\alpha, \sigma, m)$  (since  $\mu$  is a permutation). Therefore  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Translation permutation.** We now prove that, for any two  $\sigma_a, \sigma_b \neq \sigma_a$ , we have  $\Psi[\alpha](\sigma_a) \neq \Psi[\alpha](\sigma_b)$ . We prove this by contradiction. Suppose a system  $\sigma'$  exists such that  $\sigma' = \Psi[\alpha](\sigma_a) = \Psi[\alpha](\sigma_b)$ . We want to prove that  $\sigma_a = \sigma_b$ .

Following from the definition of  $\Psi[\alpha]$ , if at least one correct process in  $\sigma'$  is poisoned, then we immediately have  $\sigma_a = \sigma' = \sigma_b$ . Consequently, no correct process in  $\sigma'$  is poisoned.

We start by noting that, if  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ , then  $\sigma_a = \sigma_b$ . Indeed, if  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ , then for every  $\pi \in \Pi_C$  and every  $m \in \mathcal{M}$  we have

$$\mu(\alpha, \sigma_a, m) = \mu(\alpha, \sigma_b, m)$$

from which immediately follows

$$\chi(\alpha, \sigma_a, m) = \chi(\alpha, \sigma_b, m)$$

and, since no correct process is poisoned, for every  $\pi \in \Pi_C$  and every  $m \in \mathcal{M}$  we have

$$\begin{aligned} \sigma_a[\pi][m] &= \sigma'[\pi][\chi^{-1}(\alpha, \sigma_a, m)] \\ &= \sigma'[\pi][\chi^{-1}(\alpha, \sigma_b, m)] \\ &= \sigma_b[\pi][m] \end{aligned}$$

therefore  $\sigma_a = \sigma_b$ .

We prove that  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$  by induction. Let us assume

$$\begin{aligned} \tau(\alpha, \sigma_a) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \sigma_b) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$



and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . As we previously established, the Byzantine processes in  $\pi$ 's echo samples in  $\sigma'$  are identical to the Byzantine processes in  $\pi$ 's echo samples in  $\sigma_a$  and  $\sigma_b$ . Therefore  $r_{n+1} = r'_{n+1}$ .

Before considering the remaining possible values of  $i_{n+1}$ , we prove some auxiliary result. Let  $\pi$  be a correct process, let  $\rho$  be a process, let  $m, s$  be messages. For every  $j \leq n+1$ , as we established, we have  $i_j = i'_j$ . By hypothesis,  $\alpha$  is an auto-echo adversary, so  $i_j = i'_j, r_j = r'_j = \perp$  for every  $j \leq (N-C)C^2$ . Let us consider the non-trivial case  $n \geq (N-C)C^2$ . After the  $(n+1)$ -th invocation, the following hold true:

- $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma_a$  if and only if  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma_b$ . Indeed, if  $\rho$  is a correct process, and  $\rho$  pb.Delivered  $m$  in  $\sigma_a$ , then some  $j \leq (n+1)$  exists such that  $i_j = (\text{Deliver}, \rho, m)$ . Since  $i'_j = i_j$ ,  $\rho$  pb.Delivered  $m$  in  $\sigma_b$  as well. If  $\rho$  is Byzantine, and  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma_a$ , then some  $j \leq (n+1)$  exists such that  $i_j = (\text{Echo}, \pi, m, \rho, m)$ . Since  $i'_j = i_j$ ,  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma_b$  as well. Both arguments can be reversed to prove that, if  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma_b$ , then  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma_a$  as well.
- $\rho$  sent an  $\text{Echo}(s, m')$  for some  $m' \in \mathcal{M}$  to  $\pi$  in  $\sigma_a$  if and only if  $\rho$  sent an  $\text{Echo}(s, m')$  message for some  $m'' \in \mathcal{M}$  to  $\pi$  in  $\sigma_b$ . Indeed:
  - If  $\rho$  is correct, and it sent an  $\text{Echo}(s, m')$  message to  $\pi$  in  $\sigma_a$ , then it pb.Delivered  $m'$  in both  $\sigma_a$  and  $\sigma_b$ . Consequently,  $\rho$  sent an  $\text{Echo}(s, m')$  message to  $\pi$  in  $\sigma_b$  as well. The argument can be inverted to prove that, if  $\rho$  is correct and it sent an  $\text{Echo}(s, m'')$  message for some  $m'' \in \mathcal{M}$  to  $\pi$  in  $\sigma_b$ , then  $\rho$  sent an  $\text{Echo}(s, m')$  message for some  $m' \in \mathcal{M}$  in  $\sigma_a$ .
  - If  $\rho$  is Byzantine, then it sent an  $\text{Echo}(m', m')$  message for every  $m' \in \mathcal{M}$ , both in  $\sigma$  and  $\sigma'$ .
- If at least one correct process pb.Delivered  $m$  in  $\sigma_a$  (or, equivalently,  $\sigma_b$ ), then  $\mu(\alpha, \sigma_a, m) = \mu(\alpha, \sigma_b, m)$ . Indeed, let  $j$  be the minimum index such that  $i_j = i'_j = (\text{Deliver}, \pi', m)$  for some  $\pi' \in \Pi_C$ . By definition, we have

$$\begin{aligned}
 \mu(\alpha, \sigma_a, m) &= |\{m \in \mathcal{M} \mid \exists k \leq j, \pi' \in \Pi_C : i_k = (\text{Deliver}, \pi', m)\}| \\
 &= |\{m \in \mathcal{M} \mid \exists k \leq j, \pi' \in \Pi_C : i'_k = (\text{Deliver}, \pi', m)\}| \\
 &= \mu(\alpha, \sigma_b, m)
 \end{aligned}$$

- $\pi$  delivered  $m$  in  $\sigma_a$  if and only if  $\pi$  delivered  $m$  in  $\sigma_b$ . Indeed, if  $\pi$  delivered  $m$  in  $\sigma_a$ , then at least one correct process pb.Delivered  $m$  both in  $\sigma_a$  and  $\sigma_b$ , and  $\mu(\alpha, \sigma_a, m) = \mu(\alpha, \sigma_b, m)$ . From the definition of  $\chi$ , we immediately get  $\chi(\alpha, \sigma_a, m) = \chi(\alpha, \sigma_b, m)$  and,

## Chapter 3. Sieve

---

as we previously established,  $\pi$ 's echo sample for  $m$  in  $\sigma_a$  is identical to  $\pi$ 's echo sample for  $m$  in  $\sigma_b$ . Since  $\pi$  received the same  $\text{Echo}(m, m)$  messages in  $\sigma_a$  and  $\sigma_b$ ,  $\pi$  delivered  $m$  in  $\sigma_b$  as well. The argument can be reversed to prove that, if  $\pi$  delivered  $m$  in  $\sigma_b$ , then  $\pi$  delivered  $m$  in  $\sigma_a$  as well.

Let us assume  $i_{n+1} = (\text{State})$ . From the above it immediately follows  $r_{n+1} = r'_{n+1}$ .

Let us assume  $i_{n+1} = (\text{Sample}, \pi, m)$ . As we established,  $\pi$  receives an  $\text{Echo}(m, m')$  message for some  $m' \in \mathcal{M}$  from the same set of processes in  $\sigma_a$  and  $\sigma_b$ . Moreover, since at least one correct process pb.Delivered  $m$  in both  $\sigma_a$  and  $\sigma_b$ ,  $\pi$ 's echo sample for  $m$  in  $\sigma_a$  is identical to  $\pi$ 's echo sample for  $m$  in  $\sigma_b$ . Therefore,  $r_{n+1} = r'_{n+1}$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma_a) = \tau(\alpha, \sigma_b)$ .

Therefore,  $\sigma_a = \sigma_b$ , which contradicts the hypothesis.

□

### 3.11.4 Non-redundant adversary

**Lemma 23.** *The set of non-redundant adversaries  $\mathcal{A}_{nr}$  is optimal.*

*Proof.* We again prove the result using a decorator. Here we show that a decorator  $\Delta_{nr}$  exists such that, for every  $\alpha \in \mathcal{A}_{sq}$ , the adversary  $\alpha' = \Delta_{sq}(\alpha)$  is a non-redundant adversary, and more powerful than  $\alpha$ . If this is true, then indeed the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the sequential  $\alpha^+ = \Delta_{nr}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 10 implements **Non-redundant decorator**, a decorator that transforms a sequential adversary into a non-redundant adversary. Provided with a sequential adversary  $sqadv$ , Non-redundant decorator acts as an interface between  $sqadv$  and a system  $sys$ , effectively implementing a non-redundant adversary  $nradv$ . Non-redundant decorator exposes both the adversary and the system interfaces: the underlying adversary  $sqadv$  uses  $nradv$  as its system.

Non-redundant decorator works as follows:

- Procedure  $nradv.Init()$  initializes a *deliveries* array that is used to keep track of the message each correct process would have delivered, if  $sqadv$  was playing instead of  $nradv$ .
- Procedure  $nradv.Step()$  simply forwards the call to  $sqadv.Step()$ ;

**Algorithm 10 Non-redundant decorator.**


---

```

1: Implements:
2:   NonRedundantAdversary + CobSystem, instance nradv
3:
4: Uses:
5:   SequentialAdversary, instance sqadv, system nradv
6:   CobSystem, instance sys
7:
8: procedure nradv.Init() is
9:   deliveries =  $\{\perp\}^C$ ;
10:  sqadv.Init();
11:
12: procedure nradv.Step() is
13:  sqadv.Step();
14:
15: procedure nradv.Byzantine(process) is
16:  return sys.Byzantine(process);
17:
18: procedure nradv.State() is
19:  state =  $\emptyset$ ;
20:
21:  for all  $(\cdot, m) \in \text{sys.State}()$  do
22:    for all  $\pi \in \Pi_C$  do
23:       $n = 0$ ;
24:
25:      for all  $\rho \in \text{sys.Sample}(\pi, m)$  do
26:        if  $\rho \in \Pi \setminus \Pi_C$  or  $\text{deliveries}[\rho] = m$  then
27:           $n \leftarrow n + 1$ ;
28:        end if
29:      end for
30:
31:      if  $n \geq \hat{E}$  then
32:         $\text{state} \leftarrow \text{state} \cup \{(\pi, m)\}$ ;
33:      end if
34:    end for
35:  end for
36:
37:  return state;
38:
39: procedure nradv.Sample(process, message) is
40:  return sys.Sample(process, message);
41:

```

---

### Chapter 3. Sieve

---

```
42: procedure nradv.Deliver(process, message) is
43:   state =  $\emptyset$ ;
44:
45:   for all  $(\cdot, m) \in \text{sys.State}()$  do
46:     state  $\leftarrow$  state  $\cup$   $\{m\}$ ;
47:   end for
48:
49:   if state =  $\{message\}$  then
50:     sys.Deliver(process, message + 1);
51:   else
52:     sys.Deliver(process, message);
53:   end if
54:
55:   deliveries[process] = message;
56:
57: procedure nradv.Echo(process, sample, source, message) is
58:   sys.Echo(process, sample, source, message);
59:
60: procedure nradv.End() is
61:   sys.End();
62:
```

---

- Procedure *nradv.Byzantine*(*process*) simply forwards the call to *sqadv.Byzantine*(*process*).
- Procedure *nradv.State*() returns a list of pairs  $(\pi \in \Pi_C, m \in \mathcal{M})$  such at least one correct process delivered *m* in *sys*, and  $\pi$  would have delivered *m* in *sys*, if *sqadv* was playing instead of *nradv*.

This is achieved by querying *sys.State*(), then looping over each message *m* in the response. For every  $\pi \in \Pi_C$ , the procedure loops over every element  $\rho$  of *sys.Sample*( $\pi, m$ ), and computes the number *n* of *Echo*(*m, m*) messages that  $\pi$  would have received from its echo sample for *m* in *sys*, if *sqadv* was playing instead of *nradv*. This is achieved using the *deliveries* table, and the hypothesis that *sqadv* is an auto-echo adversary. If *n* is greater or equal to  $\hat{E}$ ,  $(\pi, m)$  is included in the list returned by the procedure.

- Procedure *nradv.Sample*(*process*, *message*) simply forwards the call to *sys.Sample*(*process*, *message*).
- Procedure *nradv.Deliver*(*process*, *message*) uses *sys.State*() to determine which messages have been delivered by at least one correct process in *sys*. If *message* is the only message that was delivered, the procedure forwards the call to *sys.Deliver*(*process*, *message* + 1). Otherwise, it forwards the call to *sys.Deliver*(*process*, *message*). Finally, it updates the *deliveries* array to reflect

the fact that *process* would have `pb.Delivered message` in *sys*, if *sqadv* was playing instead of *nradv*.

- Procedure *nradv.Echo(process, sample, source, message)* simply forwards the call to *sys.Echo(process, sample, source, message)*.
- Procedure *nradv.End()* simply forwards the call to *sys.End()*.

**Correctness.** We start by proving that no adversary, coupled with Non-redundant decorator, causes the execution to fail.

Let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ . The following hold true:

- Procedure *nradv.State()* never causes the execution to fail. Indeed, *sys.Sample( $\pi, m$ )* is called only if  $(\pi', m)$  was returned from *sys.State()*, for some  $\pi' \in \Pi_C$ . This means that *sys.Sample( $\pi, m$ )* is called only if at least one correct process delivered  $m$  in *sys*.
- No invocation of *nradv.Sample(...)* causes the execution to fail. Noting that *sqadv* is correct, it will never invoke *nradv.Sample( $\pi, m$ )* unless  $(\pi', m)$  was returned from a previous invocation of *nradv.State()*, for some  $\pi' \in \Pi_C$ . Moreover,  $(\pi', m)$  is returned from *nradv.State()* is and only if, for some  $\pi'' \in \Pi_C$ ,  $(\pi'', m)$  is returned from *sys.State()*. Therefore, *sys.Sample( $\pi, m$ )* is never invoked unless at least one correct process delivered  $m$  in *sys*.
- Procedure *nradv.Deliver(...)* never calls *sys.Deliver(...)* on a message greater than  $C$ . Let  $m \in \mathcal{M}$ . If  $m$  is the only message that was delivered in *sys*, then no correct process is poisoned in *sys*: indeed, as we proved, if a correct message was poisoned in *sys*, it would have delivered every message. Therefore, at least one correct process `pb.Delivered  $m$` . Moreover, since *sqadv* is a sequential adversary, it invokes *nradv.Deliver(...)* for the  $n$ -th time only on a message  $m \leq n$ . Since *nradv.Deliver(...)* is invoked at most  $C$  times, we have  $n \leq C$ , and, since  $m$  was delivered as a result of a previous invocation of *nradv.Deliver(...)*, we have  $m \leq n - 1$ . Consequently,  $m + 1 \leq C$ .

We further prove that *nradv* is a sequential adversary. Let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ . Since *sqadv* is sequential, it invokes *nradv.Deliver( $\pi, m$ )* only if it previously invoked *nradv.Deliver(...)* on every message  $l < m \in \mathcal{M}$ . Therefore, *nradv* can be a non-sequential adversary only as a result of a call to *sys.Deliver( $\pi, m + 1$ )*. If  $m$  is the only message that was delivered by at least one correct process in *sys*, then no correct process is poisoned in *sys*. Therefore, if *sys.Deliver( $\pi, m + 1$ )* is invoked, then, as we established, at least one correct process `pb.Delivered  $m$`  in *sys*. Noting that the set of messages that are delivered by at least one correct process in *sys* is non-decreasing, if no correct process is poisoned in *sys* then *sqadv* invoked *nradv.Deliver(...)* on  $m$  at least once when no correct process had delivered  $m$ . Consequently, for every  $l < m$ , *nradv.Deliver(...)*, and as a result *sys.Deliver(...)*, was invoked on  $l$ .

### Chapter 3. Sieve

---

**Non-redundant.** It is easy to prove that Non-redundant decorator always implements a non-redundant adversary. Indeed, let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ ,  $sys.Deliver(\pi, m)$  is never invoked if  $m$  is the only message that was delivered.

**Roadmap.** Let  $\alpha \in \mathcal{A}_{sq}$ , let  $\alpha' = \Delta_{nr}(\alpha)$ . Let  $\sigma$  be a system such that  $\alpha$  compromises the consistency of  $\sigma$ . Let  $\sigma'$  be an identical copy of  $\sigma$ . In order to prove that  $\alpha'$  is more powerful than  $\alpha$ , we prove that  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Trace.** We start by noting that, if we couple Non-redundant decorator with  $\sigma'$ , we effectively obtain a system instance  $\delta$  with which  $\alpha$  directly exchanges invocations and responses. Here we show that the trace  $\tau(\alpha, \sigma)$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\sigma$ , or it has been coupled with  $\sigma'$ , with Non-redundant decorator acting as an interface. We prove this by induction.

Let us assume

$$\begin{aligned} \tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \delta) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

We start by noting that, since  $\alpha$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = \text{Byzantine}(\pi)$ . Since  $nradv.Byzantine(\pi)$  simply forwards the call to  $sys.Byzantine(\pi)$ , and  $\sigma'$  is an identical copy of  $\sigma$ , we immediately have  $r_{n+1} = r'_{n+1}$ .

Before considering the remaining possible values of  $i_{n+1}$ , we prove some auxiliary results. Let  $\pi$  be a correct process, let  $\rho$  be a process, let  $m$  be a message. For every  $j \leq n+1$ , as we established, we have  $i_j = i'_j$ . Therefore, after the  $(n+1)$ -th invocation, the following hold true:

- $\pi$  pb.Delivered  $m$  in  $\sigma$  if and only if  $delivers[\pi] = m$ . This follows immediately from the fact that, whenever  $nradv.Deliver(\pi, m)$  is invoked,  $delivers[\pi]$  is set to  $m$ .
- $\pi$  pb.Delivered a message in  $\sigma$  if and only if  $\pi$  pb.Delivered a message in  $\sigma'$ . This follows immediately from the fact that every  $nradv.Deliver(\pi, m)$  is always either forwarded to  $sys.Deliver(\pi, m)$  or  $sys.Deliver(\pi, m+1)$ .
- If  $m$  was delivered by at least one correct process in  $\sigma$ , then  $m$  was delivered by at least

one correct process in  $\sigma'$  as well. Indeed:

- If at least one correct process is poisoned, then it delivered every message both in  $\sigma$  and  $\sigma'$ .
- If no correct process is poisoned then, for some  $j^* \leq n+1$ , after the  $j$ -th invocation, exactly one message  $m^*$  was delivered by at least one correct process in  $\sigma$ . This follows from the fact that a non-poisoned process delivers  $m$  only as a result of receiving an  $\text{Echo}(m, m)$  message, and no two messages  $\text{Echo}(m, m)$ ,  $\text{Echo}(m' \neq m, m')$  are ever issued as a result of a single invocation.
- If no correct process is poisoned, and  $m = m^*$ , then some correct process  $\pi^*$  delivered  $m$  in  $\sigma$  as a result of the  $j^*$ -th invocation. It is easy to see that, up to the  $j^*$ -th invocation, every call to  $\text{nradv.Deliver}(\pi, m)$  was simply forwarded to  $\text{sys.Deliver}(\pi, m)$ . Therefore, noting that  $\pi^*$ 's echo sample for  $m$  is identical in  $\sigma$  and  $\sigma'$ ,  $\pi^*$  delivered  $m$  in  $\text{sys}$  as well.
- If no correct process is poisoned, and  $m \neq m^*$ , then no invocation of  $\text{nradv.Deliver}(\dots)$  sees  $m$  as the only message delivered by at least one correct process in  $\text{sys}$ . Therefore, all calls to  $\text{nradv.Deliver}(\pi, m)$  are simply forwarded to  $\text{sys.Deliver}(\pi, m)$ . Consequently, noting that  $\pi$ 's echo sample for  $m$  is identical in  $\sigma$  and  $\sigma'$ , if  $\pi$  delivered  $m$  in  $\sigma$ , then  $\pi$  delivered  $m$  in  $\sigma'$  as well.

Let us assume that  $i_{n+1} = (\text{State})$ . Let  $\pi$  be a correct process, let  $m$  be message. The following hold true:

- If  $(\pi, m) \in r_{n+1}$ , then  $\pi$  delivered  $m$  in  $\sigma$ . Therefore, at least one correct process delivered  $m$  in  $\sigma'$ . Let  $\pi'$  be a correct process in  $\pi$ 's echo sample for  $m$  that pb.Delivered  $m$  in  $\sigma$ : as we established, we have  $\text{deliveries}[\pi'] = m$  and  $\pi'$  pb.Delivered a message in  $\sigma'$ . Therefore,  $\pi' \in \text{sys.Sample}(\pi, m)$ . Since  $\text{nradv.State}(\dots)$  counts the processes in  $\text{sys.Sample}(\pi, m)$  that are either Byzantine or have their  $\text{deliveries}$  value set to  $m$ , we have  $(\pi, m) \in r'_{n+1}$ .
- If  $(\pi, m) \notin r_{n+1}$ , then less than  $\hat{E}$  processes in  $\pi$ 's echo sample for  $m$  are either Byzantine or have pb.Delivered  $m$ . Therefore, less than  $\hat{E}$  processes in  $\pi$ 's echo sample are either Byzantine or have their  $\text{deliveries}$  value set to  $m$ . Since  $\text{sys.Sample}(\pi, m)$  is a subset of  $\pi$ 's echo sample for  $m$ , and since  $\text{nradv.State}(\dots)$  counts the processes in  $\text{sys.Sample}(\pi, m)$  that are either Byzantine or have their  $\text{deliveries}$  value set to  $m$ ,  $(\pi, m) \notin r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Sample}, \pi, n)$ . By hypothesis,  $\pi$ 's echo sample for  $m$  is identical in  $\sigma$  and  $\sigma'$ . Moreover, the set of processes that pb.Delivered a message is identical in  $\sigma$  and  $\sigma'$ . Noting that  $\text{nradv.Sample}(\pi, m)$  simply forwards the call to  $\text{sys.Sample}(\pi, m)$ , we immediately get  $r_{n+1} = r'_{n+1}$ .

Noting that procedures  $Deliver(\dots)$  and  $Echo(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (Deliver, \pi, m)$  or  $i_{n+1} = (Echo, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ .

**Consistency of  $\sigma'$ .** We proved that  $\tau(\alpha, \sigma) = \tau(\alpha, \delta)$ . Moreover, we proved that if a message  $m$  is eventually delivered by at least a correct process in  $\sigma$ , then  $m$  is eventually delivered by at least a correct process in  $\sigma'$  as well.

Since  $\alpha$  compromises the consistency of  $\sigma$ , two distinct messages  $m, m' \neq m$  exist such that, in  $\sigma$ , both  $m$  and  $m'$  are delivered by at least one correct process. Therefore, in  $\sigma'$ , both  $m$  and  $m'$  are delivered by at least one correct process as well. Therefore,  $\alpha'$  compromises the consistency of  $\sigma'$ .

Consequently, the adversarial power of  $\alpha$  is smaller or equal to the adversarial power of  $\alpha' = \Delta_{nr}(\alpha)$ , and the lemma is proved.  $\square$

### 3.11.5 Sample-blind adversary

**Lemma 24.** *The set of sample-blind adversaries  $\mathcal{A}_{sb}$  is optimal.*

*Proof.* We again prove the result using a decorator. Here we show that a decorator  $\Delta_{sm}$  exists such that, for every  $\alpha \in \mathcal{A}_{nr}$ , the adversary  $\alpha' = \Delta_{sm}^{C^2}(\alpha)$  is a sample-blind adversary, and more powerful than  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the sample-blind  $\alpha^+ = \Delta_{sm}^{C^2}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 11 implements **Sample-masking decorator**, a decorator that masks every invocation of  $Sample(\pi, m)$  issued by a non-redundant adversary, if  $Sample(\pi, m)$  is the first invocation of  $Sample(\dots)$  issued by that adversary.

Provided with a non-redundant adversary  $nradv$ , Sample-masking decorator acts as an interface between  $nradv$  and a system  $sys$ . Sample-masking decorator is only guaranteed to mask any invocation to  $sys.Sample(\pi, m)$ , for one process  $\pi$  and one message  $m$ . Noting that  $|\Pi_C| = C$  and  $|\mathcal{M}| = C$ , we have that, for every  $\alpha \in \mathcal{A}_{nr}$ ,  $\alpha' = \Delta_{sm}^{C^2}(\alpha)$  is a sample-blind adversary: indeed, all of  $\alpha'$ 's  $C^2$  possible calls to  $Sample(\dots)$  are necessarily masked.

Sample-masking decorator exposes both the adversary and the system interfaces: the underlying adversary  $nradv$  uses  $smadv$  as its system. Sample-masking decorator works as follows:

- Procedure  $smadv.Init()$  initializes the following variables:
  - An *index* and a *cache* variable, both initially set to  $\perp$ : *index* is used to store the pair  $(\pi \in \Pi_C, m \in \mathcal{M})$  that was provided as argument to the first invocation of



**Algorithm 11 Sample-masking decorator.**


---

```

1: Implements:
2:   SampleMaskedAdversary + CobSystem, instance smadv
3:
4: Uses:
5:   NonRedundantAdversary, instance nradv, system smadv
6:   CobSystem, instance sys
7:
8: procedure smadv.Init() is
9:   index =  $\perp$ ;      cache =  $\perp$ ;
10:  trace = [];
11:  deliveries =  $\{\perp\}^C$ ;
12:  nradv.Init();
13:
14: procedure optimize(process, message) is
15:  smadv.Byzantine(process);
16:  best.sample =  $\perp$ ;  best.probability = 0;
17:
18:  for all sample  $\in \Pi^E$  do
19:    systems = 0;
20:    compromissions = 0;
21:    for all  $\sigma \in \mathcal{S}$  do
22:      if trace  $\sim \sigma$  and  $\sigma[\textit{process}][\textit{message}] = \textit{sample}$  then
23:        systems  $\leftarrow \textit{systems} + 1$ ;
24:        if NonRedundantAdversary  $\setminus \sigma$  then
25:          compromissions  $\leftarrow \textit{compromissions} + 1$ ;
26:        end if
27:      end if
28:    end for
29:
30:    if systems > 0 and compromissions/systems > best.probability then
31:      best.sample = sample;
32:      best.probability = compromissions/systems;
33:    end if
34:  end for
35:
36:  return best.sample;
37:
38: procedure smadv.Step() is
39:  nradv.Step();
40:
41: procedure smadv.Byzantine(process) is
42:  trace  $\leftarrow \textit{trace} + [(\textit{Byzantine}, \textit{process}, \textit{sys.Byzantine}(\textit{process}))]$ ;
43:  return sys.Byzantine(process);
44:

```

---

### Chapter 3. Sieve

---

```
45: procedure smadv.State() is
46:   trace  $\leftarrow$  trace + [(State, sys.State())];
47:   state = sys.State() \ {index};
48:
49:   if index  $\neq$   $\perp$  then
50:     ( $\pi$ , m) = index;
51:
52:     n = 0;
53:     for all  $\rho \in$  cache do
54:       if  $\rho \in \Pi \setminus \Pi_C$  or deliveries[ $\rho$ ] = m then
55:         n  $\leftarrow$  n + 1;
56:       end if
57:     end for
58:
59:     if n  $\geq$   $\hat{E}$  then
60:       state  $\leftarrow$  state  $\cup$  {( $\pi$ , m)};
61:     end if
62:   end if
63:
64:   return state;
65:
66: procedure smadv.Sample(process, message) is
67:   if index =  $\perp$  then
68:     index  $\leftarrow$  (process, message);
69:     cache  $\leftarrow$  optimize(process, message);
70:   end if
71:
72:   if (process, message) = index then
73:     sample = {};
74:     for all  $\pi \in$  cache do
75:       if deliveries[ $\pi$ ]  $\neq$   $\perp$  then
76:         sample  $\leftarrow$  sample  $\cup$  { $\pi$ };
77:       end if
78:     end for
79:     return sample;
80:   else
81:     return sys.Sample(process, message);
82:   end if
83:
84: procedure smadv.Deliver(process, message) is
85:   trace  $\leftarrow$  trace + [(Deliver, (process, message))];
86:   deliveries[process] = message;
87:   sys.Deliver(process, message);
88:
```

---

---

```

89: procedure smadv.Echo(process, sample, source, message) is
90:   trace ← trace + [(Echo, (process, sample, source, message))];
91:   sys.Echo(process, sample, source, message);
92:
93: procedure smadv.End() is
94:   sys.End();
95:

```

---

*smadv.Sample*(...); *cache* is used to store the content of the echo sample *smadv* generates for  $(\pi, m)$  when *smadv.Sample*(...) is invoked for the first time. This guarantees that subsequent invocations of *smadv.Sample*( $\pi, m$ ) are provided with consistent responses throughout the entire adversarial execution.

- A *trace* array: *trace* is used to store the sequence of invocations and responses exchanged between *nradv* and *sys*.
- A *deliveries* array of  $C$  elements: *deliveries* is used to track the message pb.Delivered in *sys* by each correct process.
- Procedure *optimize*(*process, message*) returns the sample *sample* for (*process, message*) that maximizes the probability of *nradv* winning against a random system  $\sigma$  that is compatible with *trace*, and satisfies  $\sigma[\textit{process}][\textit{message}] = \textit{sample}$ . This is achieved as follows:
  - The procedure calls *smadv.Byzantine*(*process*), causing an invocation to *sys.Byzantine*(*process*) to be appended to *trace* along with its response. This is necessary because, if *smadv.Byzantine*(*process*) was never invoked before, the set of Byzantine processes in the generated *sample* might differ from the Byzantine processes in *process*' echo sample for *message* in *sys*. Noting that all of *process*' echo samples in *sys* share the same set of Byzantine processes, a subsequent call to *smadv.Byzantine*(*process*) could return a set of Byzantine processes that is inconsistent with the *sample*, causing undefined behavior on *nradv*.
  - The procedure loops over every possible value of *sample*. For each value of *sample*, it counts the number *systems* of systems  $\sigma$  that are compatible with *trace*, and satisfy  $\sigma[\textit{process}][\textit{message}] = \textit{sample}$ . Among the systems that satisfy those two constraints, the procedure counts the number *compromissions* of systems whose consistency the adversary would compromise.
  - The procedure returns the value of *sample* that satisfies  $\textit{systems} > 0$ , and maximizes  $\textit{compromissions}/\textit{systems}$ . In other words, the procedure returns a sample *sample* that is compatible with at least with one of the systems that are compatible with *trace*, and maximizes the probability that the adversary would compromise the consistency of a randomly selected system compatible with *trace*, picked among those where *process*' echo sample for *message* is *sample*.

- Procedure *smadv.Byzantine(process)* appends to *trace* the invocation of *sys.Byzantine(process)* along with its response. It then forwards the call to *sys.Byzantine(process)*.
- Procedure *smadv.State()* appends to *trace* the invocation of *sys.State()* along with its response. It then returns the response of *sys.State()*, modified to be compatible with any previous masked invocation of *sys.Sample(...)*. More specifically, if  $index = (\pi, m) \neq \perp$  (i.e., *nradv*'s first invocation of *smadv.Sample(...)* was *smadv.Sample( $\pi, m$ )*), then  $(\pi, m)$  is included in the set of pairs returned by *smadv.State()* only if  $\pi$  would have delivered *m* in *sys*, if  $\pi$ 's echo sample for *m* was *cache*. This is achieved by looping over every process in *cache*, and counting the number *n* of those processes that are either Byzantine, or pb.Delivered *m* (this is achieved using the *deliveries* array).
- Procedure *smadv.Sample(process, message)* determines whether *smadv.Sample(...)* has ever been invoked before by checking the value of *index*. If it has not, it sets *index* to  $(process, message)$ , and generates a sample for  $(process, message)$  by setting *cache* to the value returned by *optimize(process, sample)*.

If  $(process, message)$  is equal to *index*, the procedure returns the set of processes in *cache* that pb.Delivered a message in *sys*. This is achieved by looping over every process  $\rho$  in *cache*, and adding  $\rho$  to the response if  $\rho$  is either Byzantine, or satisfy  $deliveries[\rho] \neq \perp$ .

If  $(process, message)$  is not equal to *index*, the call is forwarded to *sys.Sample(process, message)*.

- Procedure *smadv.Deliver(process, message)* appends to *trace* the invocation of *sys.Deliver(process, message)*. To reflect the fact that *process* pb.Delivered *m* in *sys*, it then updates the *deliveries* array. Finally, it forwards the call to *sys.Deliver(process, message)*.
- Procedure *smadv.Echo(process, sample, source, message)* appends to *trace* the invocation of *sys.Echo(process, sample, source, message)*. It then forwards the call to *sys.Echo(process, sample, source, message)*.
- Procedure *smadv.End()* simply forwards the call to *sys.End()*.

**Correctness.** We start by proving that no adversary has undefined behavior when coupled with *Sample-masked* decorator. An adversary has undefined behavior if, at any point, the sequence of invocations and responses it exchanges with *smadv* is incompatible with every system.

Let  $\pi \in \Pi_C$ , let  $m \in \mathcal{M}$ , let us assume that the first invocation to *smadv.Sample(...)* is *smadv.Sample( $\pi, m$ )*. We start by noting that every invocation in *smadv* is forwarded to

the corresponding invocation in *sys* except for *smadv.State()* and *smadv.Sample( $\pi, m$ )*. Moreover, before the first invocation of *smadv.Sample( $\pi, m$ )*, *index* is set to  $\perp$  and, as a result, *smadv.State()* effectively forwards to *sys.State()*. Therefore, the trace exchanged between *nradv* and *smadv* is trivially compatible with *sys* before the first invocation of *smadv.Sample( $\pi, m$ )*.

When *smadv.Sample( $\pi, m$ )* is invoked for the first time, *cache* is set to *optimize( $\pi, m$ )*. When *optimize( $\pi, m$ )* is called, it calls *smadv.Byzantine( $\pi$ )*, which appends the invocation and the corresponding response to *trace*. After that, the set of systems that are compatible with *trace* is non empty, as it trivially includes *sys*. The procedure *optimize( $\pi, m$ )* returns a sample *sample* only if at least one system  $\sigma$  is compatible with *trace*, and satisfies  $\sigma[\pi][m] = \text{sample}$ . Since *sys.Byzantine( $\pi$ )* is in *trace*, the Byzantine component of *sample* is identical to *sys.Byzantine( $\pi$ )*: indeed, any system  $\sigma$  where the Byzantine component of  $\sigma[\pi][m]$  is different from *sys.Byzantine( $\pi$ )* is incompatible with  $\sigma$ .

Therefore, the system obtained by replacing  $\pi$ 's echo sample for *m* in *sys* with *cache* is a valid system, and it is compatible with *trace* up to the first invocation of *smadv.Sample( $\pi, m$ )*. Moreover, *trace* will always be compatible with such system. Indeed:

- Every subsequent call to *smadv.Sample( $\pi, m$ )* uses the *deliveries* table to determine which processes in *cache* pb.Delivered a message in *sys*, thus returning a response that is consistent with  $\pi$ 's echo sample for *m* being *cache*.
- Every subsequent call to *smadv.State()* includes  $(\pi, m)$  in its response only if at least  $\hat{E}$  processes in *cache* are either Byzantine or pb.Delivered *m* in *sys* (this is verified using the *deliveries* table).

This proves that that no adversary, coupled with Sample-masked decorator, has undefined behavior.

**Sample-blind.** It is easy to see that Sample-masking decorator masks the first invocation to *Sample(...)* issued by the decorated adversary. Indeed, if *smadv.Sample( $\pi, m$ )* is the first invocation of *smadv.Sample(...)* issued by *nradv*, then *index* is set to  $(\pi, m)$ , and *sys.Sample( $\pi, m$ )* is never be invoked.

Let  $\alpha$  be a non-redundant adversary, we have that  $\Delta_{sb}(\alpha)$  issues calls to *Sample(...)* for at most  $C^2 - 1$  pairs  $(\pi' \in \Pi_C, m' \in \mathcal{M})$ . The same argument can be applied again to see that, by composing Sample-masking decorator with itself  $C^2$  times, all possible calls to *Sample(...)* are masked. Therefore,  $\alpha' = \Delta_{sb}^{C^2}(\alpha)$  is a sample-blind adversary.

### Chapter 3. Sieve

---

**Sample replacement.** Let  $\alpha$  be an adversary, let  $\sigma$  be a system. We define a function  $v : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{N} \cup \{\perp\}$  by

$$v(\alpha, \sigma) = \min n \mid (\tau(\alpha, \sigma)_n = ((\text{Sample}, \pi \in \Pi_C, m \in \mathcal{M}), \perp))$$

Intuitively,  $v(\alpha, \sigma)$  returns the index of the first invocation of  $\text{Sample}(\dots)$  in  $\tau(\alpha, \sigma)$  if such invocation exists, and  $\perp$  otherwise. We additionally define  $\pi(\alpha, \sigma)$  and  $m(\alpha, \sigma)$  by

$$\tau(\alpha, \sigma)_{v(\alpha, \sigma)} = ((\text{Sample}, \pi(\alpha, \sigma), m(\alpha, \sigma)), \perp)$$

if  $v(\alpha, \sigma) \neq \perp$ , and by

$$\pi(\alpha, \sigma) = m(\alpha, \sigma) = \perp$$

if  $v(\alpha, \sigma) = \perp$ . Whenever at least an invocation to  $\text{Sample}(\dots)$  is issued when  $\alpha$  is coupled with  $\sigma$ ,  $\pi(\alpha, \sigma)$  and  $m(\alpha, \sigma)$  are the arguments to that invocation.

We then define  $\sigma^- : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{N} \cup \{\perp\}$ ,  $\sigma^+ : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{N} \cup \{\perp\}$ . If  $v(\alpha, \sigma) \neq \perp$

$$\begin{aligned} \sigma^-(\alpha, \sigma) &= \max n < v(\alpha, \sigma) \mid \tau(\alpha, \sigma)_n = ((\text{State}), r_n), \Psi(r_n, \alpha, \sigma) \\ \sigma^+(\alpha, \sigma) &= \min n < v(\alpha, \sigma) \mid \tau(\alpha, \sigma)_n = ((\text{State}), r_n), \Psi(r_n, \alpha, \sigma) \end{aligned}$$

where  $\Psi$  is a predicate defined as

$$\Psi(r_n, \alpha, \sigma) = (\pi(\alpha, \sigma), m(\alpha, \sigma)) \in r_n$$

Otherwise, i.e. if  $v(\alpha, \sigma) = \perp$

$$\sigma^-(\alpha, \sigma) = \sigma^+(\alpha, \sigma) = \perp$$

otherwise. Intuitively, when  $v(\alpha, \sigma) \neq \perp$ :  $\sigma^-(\alpha, \sigma)$  returns the index of the last invocation of  $\text{State}()$  prior to  $v(\alpha, \sigma)$  that did not include  $(\pi(\alpha, \sigma), m(\alpha, \sigma))$  in its response;  $\sigma^+(\alpha, \sigma)$  returns the index of the first invocation of  $\text{State}()$  prior to  $v(\alpha, \sigma)$  that included  $(\pi(\alpha, \sigma), m(\alpha, \sigma))$  in its response.

We additionally define  $\delta : \mathcal{A} \times \mathcal{S} \times \mathcal{M} \times \mathbb{N} \rightarrow \mathbb{P}(\Pi_C)$  by

$$\pi \in \delta(\alpha, \sigma, m, n) \stackrel{\text{def}}{\iff} \exists j < n \mid \tau(\alpha, \sigma)_j = ((\text{Deliver}, \pi, m), \perp) \quad (3.9)$$

Intuitively,  $\pi$  is in  $\delta(\alpha, \sigma, m, n)$  if  $\alpha$  invokes  $\text{Deliver}(\pi, m)$  before the  $n$ -th invocation it issues, when coupled with  $\sigma$ . In other words,  $\delta(\alpha, \sigma, m, n)$  represents the set of correct processes that pb.Deliver  $m$  before the  $n$ -th invocation, when  $\alpha$  is coupled with  $\sigma$ .

Finally, we define  $\delta^- : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{P}(\Pi_C)$ ,  $\delta^+ : \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{P}(\Pi_C)$  by

$$\delta^-(\alpha, \sigma) = \begin{cases} \delta(\alpha, \sigma, m(\alpha, \sigma), \sigma^-(\alpha, \sigma)) & \text{iff } \sigma^-(\alpha, \sigma) \neq \perp \\ \emptyset & \text{otherwise} \end{cases}$$

$$\delta^+(\alpha, \sigma) = \begin{cases} \delta(\alpha, \sigma, m(\alpha, \sigma), \sigma^+(\alpha, \sigma)) & \text{iff } \sigma^+(\alpha, \sigma) \neq \perp \\ \Pi_C & \text{otherwise} \end{cases}$$

Intuitively:

- When  $\sigma^-(\alpha, \sigma) \neq \perp$ ,  $\delta^-(\alpha, \sigma)$  represents the set of processes that pb.Delivered  $m(\alpha, \sigma)$  before  $\sigma^-(\alpha, \sigma)$ . Intuitively,  $\delta^-$  is designed to guarantee that less than  $\hat{E}$  elements of  $\sigma[\pi(\alpha, \sigma)][m(\alpha, \sigma)]$  are either Byzantine or included in  $\delta^-(\alpha, \sigma)$ . If this was not the case, the  $\sigma^-(\alpha, \sigma)$ -th invocation of  $State()$  would have included  $(\pi(\alpha, \sigma), m(\alpha, \sigma))$  in its response.
- When  $\sigma^+(\alpha, \sigma) \neq \perp$ ,  $\delta^+(\alpha, \sigma)$  represents the set of processes that pb.Delivered  $m(\alpha, \sigma)$  before  $\sigma^+(\alpha, \sigma)$ . Intuitively,  $\delta^+$  is designed to guarantee that at least  $\hat{E}$  elements of  $\sigma[\pi(\alpha, \sigma)][m(\alpha, \sigma)]$  are either Byzantine or included in  $\delta^+(\alpha, \sigma)$ . If this was not then case, the  $\sigma^+(\alpha, \sigma)$ -th invocation of  $State()$  would not have included  $(\pi(\alpha, \sigma), m(\alpha, \sigma))$  in its response.

All the above definitions allow us to define a **sample replacement function**  $\mathcal{E}[\alpha] : \mathcal{S} \rightarrow \mathbb{P}(\Pi^E)$  by

$$\mathcal{E}[\alpha](\sigma) = \emptyset$$

if  $v(\alpha, \sigma) = \perp$  and

$$\bar{E} \in \mathcal{E}[\alpha](\sigma) \stackrel{\text{def}}{\iff} \begin{cases} \sigma[\pi(\alpha, \sigma)][m(\alpha, \sigma)][n] \in \Pi \setminus \Pi_C \implies \\ (\bar{E}[n] = \sigma[\pi(\alpha, \sigma)][m(\alpha, \sigma)][n]) \\ |\{n \in 1..E \mid \bar{E}[n] \in \delta^-(\alpha, \sigma) \cup (\Pi \setminus \Pi_C)\}| < \hat{E} \\ |\{n \in 1..E \mid \bar{E}[n] \in \delta^+(\alpha, \sigma) \cup (\Pi \setminus \Pi_C)\}| \geq \hat{E} \end{cases}$$

otherwise. Intuitively,  $\mathcal{E}[\alpha]$  is designed so that, if  $\alpha$  is non-redundant, when  $v(\alpha, \sigma) \neq \perp$ , a sample  $E$  is in  $\mathcal{E}[\alpha](\sigma)$  if, by replacing  $\pi(\alpha, \sigma)$ 's echo sample for  $m(\alpha, \sigma)$  in  $\sigma$  with  $E$ , we obtain a system  $\sigma'$  that is *interchangeable* with  $\sigma$ , i.e., a system that cannot be distinguished from  $\sigma$  up to the  $v(\alpha, \sigma)$ -th invocation, and whose consistency is compromised by the same set of traces. We prove these two properties in the next section of this proof.

More specifically, a sample  $\bar{E}$  is in  $\mathcal{E}[\alpha](\sigma)$  if it satisfies the following conditions:

- $\bar{E}$  shares the set of Byzantine processes in  $\sigma[\pi(\alpha, \sigma)][m(\alpha, \sigma)]$ .

### Chapter 3. Sieve

---

- Less than  $\hat{E}$  processes in  $\bar{E}$  pb.Deliver  $m(\alpha, \sigma)$  before the last invocation of  $State(\dots)$  in  $\tau(\alpha, \sigma)$  (before  $v(\alpha, \sigma)$ ) that does not include  $(\pi(\alpha, \sigma), m(\alpha, \sigma))$  in its response.
- At least  $\hat{E}$  processes in  $\bar{E}$  pb.Deliver  $m(\alpha, \sigma)$  before the first invocation of  $State(\dots)$  in  $\tau(\alpha, \sigma)$  (before  $v(\alpha, \sigma)$ ) that includes  $(\pi(\alpha, \sigma), m(\alpha, \sigma))$  in its response.

**Sample interchangeability.** Let  $\alpha$  be a non-redundant adversary, let  $\sigma$  be a system such that  $v(\alpha, \sigma) \neq \perp$ . Let  $\pi^* = \pi(\alpha, \sigma)$ , let  $m^* = m(\alpha, \sigma)$ . Let  $\sigma'$  be a system such that, for every pair  $(\pi, m) \neq (\pi^*, m^*)$  (i.e.,  $\pi \neq \pi^*$  or  $m \neq m^*$ ), the two following statements hold:

$$\begin{aligned}\sigma'[\pi^*][m^*] &\in \mathcal{E}[\alpha](\sigma) \\ \sigma'[\pi][m] &= \sigma[\pi][m]\end{aligned}$$

In this section, we prove the following:

$$\begin{aligned}\forall n < v(\alpha, \sigma), \tau(\alpha, \sigma)_n &= \tau(\alpha, \sigma')_n \\ (\alpha \setminus \sigma) &\implies (\tau(\alpha, \sigma) \setminus \sigma')\end{aligned}$$

We establish the first result by induction. Let us assume

$$\begin{aligned}\tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \sigma') &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n\end{aligned}$$

with  $n \geq 0$  (here  $n = 0$  means that this is  $\alpha$ 's first invocation). We start by noting that, since  $a$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us consider the case  $i_{n+1} = (\text{Byzantine}, \pi, m)$ . Following from the definition of  $\mathcal{E}[\alpha](\sigma)$ ,  $\pi^*$ 's echo sample for  $m^*$  in  $\sigma'$  includes the same set of Byzantine processes as  $\pi^*$ 's echo sample for  $m^*$  in  $\sigma$ . Since all other echo samples are trivially identical in  $\sigma$  and  $\sigma'$ , we have  $r_{n+1} = r'_{n+1}$ .

Let us consider the case  $i_{n+1} = (\text{State})$ . Let  $\pi \in \Pi_C$ , let  $\rho \in \Pi$ , let  $m \in \mathcal{M}$ . Noting that  $i_j = i'_j \forall j \leq n+1$ , we trivially have that  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma$  if and only if  $\rho$  sent an  $\text{Echo}(m, m)$  message to  $\pi$  in  $\sigma$ . Noting that all echo samples but  $\pi^*$ 's echo sample for  $m^*$  are identical in  $\sigma$ , we immediately get that the symmetric difference between  $r_{n+1}$  and  $r'_{n+1}$  can only include  $(\pi^*, m^*)$ . The following hold true:



- If  $(\pi^*, m^*) \in r_{n+1}$ , then  $(\pi^*, m^*) \in r'_{n+1}$ . Indeed, if  $(\pi^*, m^*) \in r_{n+1}$ , then by definition  $\sigma^+(\alpha, \sigma) \leq n + 1$ . Therefore, by definition, every correct process in  $\delta^+(\alpha, \sigma)$  pb.Delivered  $m^*$  (both in  $\sigma$  and  $\sigma'$ ). Noting that  $\alpha$  is an auto-echo adversary, every process in  $\delta^+(\alpha, \sigma) \cup (\Pi \setminus \Pi_C)$  sent an  $\text{Echo}(m^*, m^*)$  message to  $\pi^*$ , both in  $\sigma$  and  $\sigma'$ . Finally, by definition,  $\mathcal{E}[\alpha](\sigma)$  includes at least  $\hat{E}$  processes in  $\delta^+(\alpha, \sigma) \cup (\Pi \setminus \Pi_C)$ . Therefore  $\pi^*$  delivered  $m^*$  in  $\sigma'$ , and  $(\pi^*, m^*) \in r_{n+1}$ .
- If  $(\pi^*, m^*) \notin r_{n+1}$ , then  $(\pi^*, m^*) \notin r'_{n+1}$ . Indeed, if  $(\pi^*, m^*) \in r_{n+1}$ , then by definition  $\sigma^-(\alpha, \sigma) \geq n + 1$ . Therefore, by definition, every correct process that pb.Delivered  $m^*$  (both in  $\sigma$  and  $\sigma'$ ) is included in  $\delta^-(\alpha, \sigma)$ . Finally, by definition,  $\mathcal{E}[\alpha](\sigma)$  includes less than  $\hat{E}$  processes in  $\delta^-(\alpha, \sigma) \cup (\Pi \setminus \Pi_C)$ . Therefore  $\pi^*$  did not deliver  $m^*$  in  $\sigma'$ , and  $(\pi^*, m^*) \notin r_{n+1}$ .

which proves  $r_{n+1} = r'_{n+1}$ .

Noting that, by definition,  $n < \nu(\alpha, \sigma)$ ,  $i_{n+1}$  cannot be  $(\text{Sample}, \pi, m)$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have

$$\forall n < \nu(\alpha, \sigma), \tau(\alpha, \sigma) = \tau(\alpha, \sigma')$$

Let us assume that  $\alpha$  compromises the consistency of  $\sigma$ . We want to prove that  $\tau(\alpha, \sigma)$  compromises the consistency of  $\sigma'$ .

We start by noting that, since by definition  $\alpha$ 's  $\nu(\alpha, \sigma)$ -th invocation in  $\tau(\alpha, \sigma)$  is  $(\text{Sample}, \pi^*, m^*)$  then, since  $\alpha$  is correct, for some  $j < \nu(\alpha, \sigma)$ , the  $j$ -th invocation in  $\tau(\alpha, \sigma)$  is  $(\text{State})$ , and its response includes  $(\pi, m^*)$  for some  $\pi \in \Pi_C$ . Therefore, before the  $\nu(\alpha, \sigma)$ -th invocation, at least one correct process in  $\sigma$  delivered  $m^*$ .

We previously proved, however, that since  $j < \nu(\alpha, \sigma)$ , we have  $\tau(\alpha, \sigma)_j = \tau(\alpha, \sigma')_j$ . Therefore, at least one correct process delivered  $m^*$  in  $\sigma'$  as well.

Since  $\alpha$  compromises the consistency of  $\sigma$ , at least one correct process  $\pi'$  eventually delivers a message  $m' \neq m^*$  in  $\sigma$ . Noting that  $\pi'$ 's echo sample for  $m'$  is identical in  $\sigma$  and  $\sigma'$ , we immediately have that  $\pi'$  delivers  $m'$  in  $\sigma'$  as well.

**System optimization.** Let  $\alpha$  be a non-redundant adversary, let  $\sigma$  be a system. In the previous section of this proof, we proved that, if we replace  $\pi(\alpha, \sigma)$ 's echo sample for  $m(\alpha, \sigma)$  in  $\sigma$  with any sample in  $\mathcal{E}[\alpha](\sigma)$ , we obtain a system  $\sigma'$  such that  $\tau(\alpha, \sigma)_n = \tau(\alpha, \sigma')_n$  for all  $n < \nu(\alpha, \sigma)$ .

We start by defining a function  $\mathcal{N} : \mathcal{A} \rightarrow \mathbb{P}(\mathcal{S})$  by

$$\mathcal{N}(\alpha) = \{\sigma \in \mathcal{S} \mid \nu(\alpha, \sigma) \neq \perp\}$$

### Chapter 3. Sieve

---

Provided with an adversary  $\alpha$ ,  $\mathcal{N}$  returns the set of systems coupled with which  $\alpha$  issues at least one invocation to  $Sample(\dots)$ .

We then define a function  $\mathcal{S}[\alpha] : \mathcal{N}(\alpha) \rightarrow \mathbb{P}(\mathcal{N}(\alpha))$  by

$$\mathcal{S}[\alpha](\sigma) = \{\sigma' \in \mathcal{S} \mid \tau(\alpha, \sigma)_{1..(v(\alpha, \sigma)-1)} \sim \sigma'\}$$

Intuitively, when  $v(\alpha, \sigma) \neq \perp$ ,  $\mathcal{S}[\alpha](\sigma)$  returns the set of systems that  $\alpha$  cannot distinguish from  $\sigma$ , before the first invocation of  $Sample(\dots)$ .

Let  $\sigma$  be a system such that  $v(\alpha, \sigma) \neq \perp$ , let  $\sigma' \in \mathcal{S}[\alpha](\sigma)$ . Noting that  $\alpha$  is a deterministic adversary, we immediately get

$$\tau(\alpha, \sigma')_n = \tau(\alpha, \sigma)_n \quad \forall n < v(\alpha, \sigma)$$

and

$$v(\alpha, \sigma') = v(\alpha, \sigma)$$

from which immediately follows

$$\mathcal{S}[\alpha](\sigma') = \mathcal{S}[\alpha](\sigma)$$

Let  $\alpha$  be a non-redundant adversary, let  $\sigma, \sigma'$  be systems in  $\mathcal{N}(\alpha)$ . Let  $\sigma \mathcal{S}[\alpha] \sigma'$  denote the relationship

$$\sigma' \in \mathcal{S}[\alpha](\sigma)$$

Since  $\tau(\alpha, \sigma) \sim \sigma$ , we immediately have that  $\mathcal{S}[\alpha]$  is reflexive. Since we established  $\mathcal{S}[\alpha](\sigma') = \mathcal{S}[\alpha](\sigma)$ ,  $\mathcal{S}[\alpha]$  is also symmetric and transitive. Therefore,  $\mathcal{S}[\alpha]$  is an equivalence relation on  $\mathcal{N}(\alpha)$ .

Let

$$\mathcal{S}[\alpha]_1, \dots, \mathcal{S}[\alpha]_h = \frac{\mathcal{N}(\alpha)}{\mathcal{S}[\alpha]}$$

intuitively, each  $\mathcal{S}[\alpha]_i$  is a distinct set of systems that are indistinguishable to  $\alpha$ , before the first invocation of  $Sample(\dots)$ .

Let  $i \in 1..h$ . Let  $\sigma \in \mathcal{S}[\alpha]_i$ , let  $E \in \mathcal{E}[\alpha](\sigma)$ , let  $\sigma'$  be identical to  $\sigma$ , with the exception of  $\pi(\alpha, \sigma)$ 's echo sample for  $m(\alpha, \sigma)$ , which is replaced with  $E$ . As we previously proved,  $\tau(\alpha, \sigma)_{1..(v(\alpha, \sigma)-1)} \sim \sigma'$ , therefore have  $\sigma' \in \mathcal{S}[\alpha]_i$ . Moreover, we proved that for every  $\sigma$  in  $\mathcal{S}[\alpha]_i$ ,  $\mathcal{E}[\alpha](\sigma)$  yields the same set of samples.

Let  $\sigma, \sigma'$  be systems in  $\mathcal{S}[\alpha]_i$ , let  $\pi^* = \pi(\alpha, \sigma) = \pi(\alpha, \sigma')$ , let  $m^* = m(\alpha, \sigma) = m(\alpha, \sigma')$ . Let  $\sigma \mathcal{E}[\alpha] \sigma'$  denote the relationship

$$\sigma[\pi^*][m^*] = \sigma'[\pi^*][m^*] \in (\mathcal{E}[\alpha](\sigma) = \mathcal{E}[\alpha](\sigma'))$$

from its definition we can immediately see that  $\mathcal{E}[\alpha]$  is an equivalence relation, and we can partition

$$\mathcal{E}[\alpha]_1^i, \dots, \mathcal{E}[\alpha]_l^i = \frac{\mathcal{S}[\alpha]_i}{\mathcal{E}[\alpha]}$$

with

$$|\mathcal{E}[\alpha]_1^i| = \dots = |\mathcal{E}[\alpha]_l^i|$$

Let  $\mathcal{C}[\alpha]_1^i, \dots, \mathcal{C}[\alpha]_l^i$  denote the probability of  $\alpha$  compromising a random element of  $\mathcal{E}[\alpha]_1^i, \dots, \mathcal{E}[\alpha]_l^i$ :

$$\mathcal{C}[\alpha]_j^i = \frac{|\{\sigma \in \mathcal{E}[\alpha]_j^i \mid \alpha \searrow \sigma\}|}{|\mathcal{E}[\alpha]_j^i|}$$

we can determine the subset whose consistency  $\alpha$  has the highest probability of compromising by

$$\mathcal{C}[\alpha]_*^i = \operatorname{argmax}_j \mathcal{C}[\alpha]_j^i$$

Finally, we define an **optimization function**  $\mathcal{O}[\alpha] : \mathcal{N}(\alpha) \rightarrow \mathcal{N}(\alpha)$ . Let  $\sigma \in \mathcal{S}[\alpha]_i$ , we define  $\mathcal{O}[\alpha]$  by

$$\mathcal{O}[\alpha](\sigma)[\pi][m] = \begin{cases} \mathcal{E}[\alpha](\sigma)_{\mathcal{C}[\alpha]_*^i} & \text{iff } \pi = \pi(\alpha, \sigma), m = m(\alpha, \sigma) \\ \sigma[\pi][m] & \text{otherwise} \end{cases}$$

As we previously proved, every  $\mathcal{E}[\alpha]_j^i$  has the same number of elements. Moreover,  $\mathcal{O}[\alpha]$  maps a system  $\sigma$  in  $\mathcal{E}[\alpha]_j^i$  to the corresponding system  $\sigma'$  in  $\mathcal{E}[\alpha]_{\mathcal{C}[\alpha]_*^i}^i$  that is identical to  $\sigma$ , except for  $\pi(\alpha, \sigma)$ 's echo sample for  $m(\alpha, \sigma)$ , which is replaced with  $\mathcal{E}[\alpha](\sigma)_{\mathcal{C}[\alpha]_*^i}$ .

Therefore, for every  $\sigma, \sigma' \in \mathcal{E}[\alpha]_{\mathcal{C}[\alpha]_*^i}^i$ ,

$$|\mathcal{O}[\alpha]^{-1}(\sigma)| = |\mathcal{O}[\alpha]^{-1}(\sigma')| = \frac{|\mathcal{S}[\alpha]_i|}{|\mathcal{E}[\alpha]_1^i|}$$

**System masking.** Let  $\alpha$  be a non-redundant adversary, let  $\alpha' = \Delta_{sb}(\alpha)$ , let  $\sigma$  be a system.

We start by noting that, if  $\nu(\alpha, \sigma) = \perp$ , then  $\tau(\alpha, \sigma) = \tau(\alpha', \sigma)$ . Indeed, if  $\alpha$  never invokes *Sample(...)* when coupled with  $\sigma$ , all calls to *smadv* are simply forwarded to the corresponding calls in *sys*. Therefore, if  $\alpha$  compromises the consistency of  $\sigma$ , then trivially  $\alpha'$  compromises the consistency of  $\sigma$  as well.

Let us assume that  $\nu(\alpha, \sigma) \neq \perp$ . Let  $\sigma'$  be an identical copy of  $\sigma$ . We start by noting that, if we couple Sample-masking decorator with  $\sigma'$ , we effectively obtain a system instance  $\delta$

### Chapter 3. Sieve

---

with which  $\alpha$  directly exchanges invocations and responses. Here we show that the trace  $\tau(\alpha, \mathcal{O}[\alpha](\sigma))$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\mathcal{O}[\alpha](\sigma)$ , or it has been coupled with  $\sigma'$ , with Non-redundant decorator acting as an interface.

We previously proved that the trace exchanged between *nradv* and *smadv* is identical to the trace that *nradv* would exchange with *sys*, if  $\pi(\alpha, \sigma)$ 's echo sample for  $m(\alpha, \sigma)$  in *sys* was replaced with *cache*.

Let  $i \in \mathbb{N}$  such that  $\sigma \in \mathcal{S}[\alpha]_i$ . Procedure *optimize* explicitly loops over all possible values of *sample*  $\in \Pi^E$ . For every value of *sample*, it loops over all the systems  $\bar{\sigma}$  that are compatible with *trace*, and satisfy  $\bar{\sigma}[\pi(\alpha, \sigma)][m(\alpha, \sigma)] = \textit{sample}$ . If, at the end of the loop, *systems*  $\neq 0$ , then *compromissions* effectively represents, for some  $j$ , the number of systems in  $\mathcal{E}[\alpha]_j^i$  that  $\alpha$  compromises. Since *optimize* selects the value of *sample* that maximizes *compromissions/system*, the value that is eventually assigned to *cache* is effectively  $\mathcal{E}[\alpha](\sigma)_{\mathcal{E}[\alpha]_*^i}$ , which proves the statement.

We previously proved that, if  $\alpha$  compromises the consistency of  $\mathcal{O}[\alpha](\sigma)$ , then  $\tau(\alpha, \mathcal{O}[\alpha](\sigma))$  compromises the consistency of  $\sigma$  as well. Noting that every invocation to *smadv.Deliver(...)* or *smadv.Echo(...)* is respectively forwarded to *sys.Deliver(...)* or *sys.Echo(...)*, we finally obtain that if  $\alpha$  compromises the consistency of  $\mathcal{O}[\alpha](\sigma)$ , then  $\alpha'$  compromises the consistency of  $\sigma$  as well.

**Adversarial power.** We can finally show that the adversarial power of  $\alpha'$  is greater than the adversarial power of  $\alpha$ . Let  $\sigma$  be a system.

As we previously established, if  $\sigma \notin \mathcal{N}(\alpha)$ , then the probability of  $\alpha$  compromising  $\sigma$  is identical to the probability of  $\alpha$  compromising  $\sigma'$ .

Let us assume that  $\sigma \in \mathcal{N}(\alpha)$ . Let  $i, j \in \mathbb{N}$  such that  $\sigma \in \mathcal{E}[\alpha]_j^i$ . The probability of  $\alpha$  compromising the consistency of  $\sigma$  is

$$\mathcal{P}[\alpha \searrow \sigma] = \mathcal{C}[\alpha]_j^i$$

and, since  $\alpha'$  compromises the consistency of  $\sigma$  if  $\alpha$  compromises the consistency of  $\mathcal{O}[\alpha](\sigma)$ , the probability of  $\alpha'$  compromising the consistency of  $\sigma$  is

$$\mathcal{P}[\alpha' \searrow \sigma] = \mathcal{P}[\alpha \searrow \mathcal{O}[\alpha](\sigma)] = \mathcal{C}[\alpha]_{\mathcal{E}[\alpha]_*^i}^i \geq \mathcal{C}[\alpha]_j^i = \mathcal{P}[\alpha \searrow \sigma]$$

Which proves that the adversarial power of  $\alpha'$  is greater or equal to the adversarial power of  $\alpha$ .  $\square$

### 3.11.6 Byzantine-counting adversary

**Lemma 25.** *The set of Byzantine-counting adversaries  $\mathcal{A}_{bc}$  is optimal.*

*Proof.* We again prove the result using a decorator. Here we show that a decorator  $\Delta_{bc}$  exists such that, for every  $\alpha \in \mathcal{A}_{sb}$ , the adversary  $\alpha' = \Delta_{bc}(\alpha)$  is a Byzantine-counting adversary, and more powerful than  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the Byzantine-counting  $\alpha^+ = \Delta_{bc}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 12 implements **Byzantine-counting decorator**, a decorator that transforms a sample-blind adversary into a Byzantine-counting adversary. Provided with a sample-blind adversary  $sbadv$ , Byzantine-counting decorator acts as an interface between  $sbadv$  and a system  $sys$ , effectively implementing a Byzantine-counting adversary  $bcadv$ . Byzantine-counting decorator exposes both the adversary and the system interface: the underlying adversary  $sbadv$  uses  $bcadv$  as its system.

Byzantine-counting decorator works as follows:

- Procedure  $bcadv.Init()$  generates  $best.byzantine$ , an array of  $C$  pre-computed responses that  $bcadv$  will provide to any subsequent invocation of  $bcadv.Byzantine(\dots)$ , optimized to maximize the probability of compromising  $sys$ . This is achieved as follows:
  - The procedure loops over every correct process  $\pi$ , and queries  $|sys.Byzantine(\pi)|$  to determine how many Byzantine processes there are in the first echo sample of  $\pi$ . For each  $\pi$ , the procedure sets variable  $space[\pi]$  to the set of all possible responses to  $bcadv.Byzantine(\pi)$  that satisfy the condition  $|bcadv.Byzantine(\pi)| = |sys.Byzantine(\pi)|$ .
  - The procedure loops over every possible array  $byzantine$  of  $C$  responses that, for every  $\pi \in \Pi_C$ , satisfies  $byzantine[\pi] \in space[\pi]$ . It then counts the number of systems  $\sigma$  that are compatible with  $byzantine$  (i.e., that satisfy, for every  $\pi \in \Pi_C$ ,  $\sigma.Byzantine(\pi) = byzantine[\pi]$ ) and whose consistency is compromised by the underlying adversary  $SampleBlindAdversary$ .
  - The procedure sets  $best.byzantine$  to the array  $byzantine$  that maximizes the number of systems compatible with  $byzantine$  whose consistency is compromised by  $SampleBlindAdversary$ .
- Procedure  $bcadv.Byzantine(process)$  simply returns  $best.byzantine[process]$ .
- Procedure  $bcadv.State()$  simply forwards the call to  $sys.State()$ .
- Procedure  $bcadv.Sample(process, message)$  is never called. This is due to the fact that  $sbadv$  is sample-blind.

**Algorithm 12 Byzantine-counting decorator.**

---

```

1: Implements:
2:   ByzantineCountingAdversary + CobSystem, instance bcadv
3:
4: Uses:
5:   SampleBlindAdversary, instance sbadv, system bcadv
6:   CobSystem, instance sys
7:
8: procedure bcadv.Init() is
9:   best.byzantine =  $\perp$ ;   best.compromissions = 0;
10:  space =  $\{\perp\}^C$ ;
11:
12:  for all  $\pi \in \Pi_C$  do
13:    count =  $|sys.Byzantine(\pi)|$ ;
14:    space $[\pi]$  =  $(\Pi \setminus \Pi_C)^{count}$ ;
15:  end for
16:
17:  for all byzantine  $\in space[\pi_1] \times \dots \times space[\pi_C]$  do
18:    compromissions = 0;
19:    for all  $\sigma \in \mathcal{S}$  do
20:      match = True;
21:      for all  $\pi \in \Pi_C$  do
22:        if  $\sigma.Byzantine(\pi) \neq byzantine[\pi]$  then
23:          match  $\leftarrow$  False;
24:        end if
25:      end for
26:
27:      if match and SampleBlindAdversary  $\setminus \sigma$  then
28:        compromissions  $\leftarrow$  compromissions + 1;
29:      end if
30:    end for
31:
32:    if compromissions > best.compromissions then
33:      best.byzantine  $\leftarrow$  byzantine;
34:      best.compromissions = compromissions;
35:    end if
36:  end for
37:  sbadv.Init();
38:
39: procedure bcadv.Step() is
40:  sbadv.Step();
41:
42: procedure bcadv.Byzantine(process) is
43:  return best.byzantine[process];
44:

```

---

---

```

45: procedure bcadv.State() is
46:   return sys.State();
47:
48: procedure bcadv.Sample(process, message) is
49:   raise error;
50:
51: procedure bcadv.Deliver(process, message) is
52:   sys.Deliver(process, message);
53:
54: procedure bcadv.Echo(process, sample, source, message) is
55:   sys.Echo(process, sample, source, message);
56:
57: procedure bcadv.End() is
58:   sys.End();
59:

```

---

- Procedure *bcadv.Deliver(process, message)* simply forwards the call to *sys.Deliver(process, message)*.
- Procedure *bcadv.Echo(process, sample, source, message)* simply forwards the call to *sys.Echo(process, sample, source, message)*.
- Procedure *bcadv.End()* simply forwards the call to *sys.End()*.

**Correctness.** We start by proving that no adversary has undefined behavior when coupled with `Byzantine-counting` decorator. An adversary has undefined behavior if, at any point, the sequence of invocations and responses it exchanges with *bcadv* is incompatible with every system.

Upon initialization, *bcadv* generates an array *best.byzantine* of  $C$  responses, one for every call to *bcadv.Byzantine*( $\pi \in \Pi_C$ ). For every correct process  $\pi$ , *best.byzantine*[ $\pi$ ] contains only Byzantine processes and satisfies  $|\text{best.byzantine}[\pi]| = |\text{sys.Byzantine}(\pi)|$ . Let *sys'* be the system obtained by replacing the Byzantine component of each correct process  $\pi$ 's echo samples in *sys* with *best.byzantine*[ $\pi$ ]. The trace exchanged between *sbadv* and *bcadv* is always compatible with *sys'*. Indeed:

- Every call to *bcadv.Byzantine*( $\pi$ ) returns *best.byzantine*[ $\pi$ ] which is equal, by definition, to *sys.Byzantine*( $\pi$ ).
- Every call to *bcadv.State*() is simply forwarded to *sys.State*(). Let  $\pi$  be a correct process, let  $m$  be a message. Since that *bcadv* is an auto-echo adversary, when *bcadv.State*() is invoked, every Byzantine process in  $\pi$ 's echo sample for  $m$  has sent an `Echo( $m, m$ )` message both in *sys* and *sys'*. Moreover, the number of Byzantine processes in  $\pi$ 's echo sample for  $m$  is identical in *sys* and *sys'*. Finally, set of cor-

### Chapter 3. Sieve

---

rect processes in  $\pi$ 's echo sample for  $m$  is identical in  $sys$  and  $sys'$ . Consequently,  $bcadv.State() = sys.State() = sys'.State()$ .

**Byzantine-counting.** It is immediate to see that Byzantine-counting decorator always implements a Byzantine-counting adversary. Indeed, for any  $\pi \in \Pi_C$ ,  $sys.Byzantine(\pi)$  is only invoked from  $|sys.Byzantine(\pi)|$ .

**Byzantine interchangeability.** Let  $\alpha$  be a sample-blind system. Let  $\sigma$  be a system, let  $\sigma'$  be a system such that, for every correct process  $\pi$ , every message  $m$ , and every  $n \in 1..E$ ,

$$\begin{aligned} (\sigma[\pi][m][n] \in \Pi_C) &\implies (\sigma'[\pi][m][n] = \sigma[\pi][m][n]) \\ (\sigma[\pi][m][n] \notin \Pi_C) &\implies (\sigma'[\pi][m][n] \notin \Pi_C) \end{aligned}$$

In other words, for every  $\pi \in \Pi_C$  and every  $m \in \mathcal{M}$ , the set of correct processes in  $\pi$ 's echo sample for  $m$  is identical in  $\sigma$  and  $\sigma'$ .

Here we prove that, if  $\alpha$  compromises  $\sigma$ , then  $\tau(\alpha, \sigma)$  compromises  $\sigma'$ . In order to do this, we first establish some auxiliary results.

Let us consider the case where  $\alpha$  is run against  $\sigma$  and  $\tau(\alpha, \sigma)$  is applied to  $\sigma'$ . Let  $\pi$  be a correct process, let  $\rho$  be a process, let  $m$  be a message. At the end of both adversarial executions, the following hold true:

- If  $\pi$  pb.Delivered  $m$  in  $\sigma$ , then  $\pi$  pb.Delivered  $m$  in  $\sigma'$  as well. This follows immediately from the fact that  $\tau(\alpha, \sigma)$  is applied to  $\sigma'$ , and  $((\text{Deliver}, \pi, m), \perp) \in \tau(\alpha, \sigma)$ .
- If  $\rho$  sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma$ , then  $\rho$  sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma'$ . Indeed, if  $\rho$  is a correct process, and it sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma$ , then it pb.Delivered  $m$  both in  $\sigma$  and  $\sigma'$ . Therefore, it sent an Echo( $m, m$ ) message to  $\pi$  in  $\sigma'$  as well. If  $\rho$  is a Byzantine process then, since  $\alpha$  is an auto-echo adversary,  $\rho$  sent an Echo( $m, m$ ) message to  $\pi$  both in  $\sigma$  and  $\sigma'$ .
- If  $\pi$  delivered  $m$  in  $\sigma$ , then  $\pi$  also delivered  $m$  in  $\sigma'$ . This follows from the above, and the fact that the correct processes in  $\pi$ 's echo sample for  $m$  are identical in  $\sigma$  and  $\sigma'$ .

If  $\alpha$  compromises the consistency of  $\sigma$ , then two correct processes  $\pi, \pi'$  and two distinct messages  $m, m' \neq m$  exist such that  $\pi$  delivered  $m$ , and  $\pi'$  delivered  $m'$  in  $\sigma$ . From the above, however,  $\pi$  delivered  $m$ , and  $\pi'$  delivered  $m'$ , in  $\sigma'$  as well. Consequently,  $\tau(\alpha, \sigma)$  compromises the consistency of  $\sigma'$ .



**System optimization.** Let  $\sigma, \sigma'$  be systems. We define the relationship  $\overset{[F]}{\sim}$  by

$$\left(\sigma \overset{[F]}{\sim} \sigma'\right) \stackrel{def}{\iff} (\forall \pi \in \Pi_C, \forall n \in 1..E, \sigma[\pi][1][n] \in \Pi_C \Leftrightarrow \sigma'[\pi][1][n] \in \Pi_C)$$

In other words,  $\sigma \overset{[F]}{\sim} \sigma'$  if, for every  $\pi$  and for every  $n \in 1..E$ , the  $n$ -th element of the first of  $\pi$ 's echo samples is either correct both in  $\sigma$  and  $\sigma'$ , or Byzantine both in  $\sigma$  and  $\sigma'$ .

It is immediate to see that  $\overset{[F]}{\sim}$  is an equivalence relation. We can therefore partition  $\mathcal{S}$  with  $\overset{[F]}{\sim}$  to obtain

$$\mathcal{S}_1, \dots, \mathcal{S}_h = \frac{\mathcal{S}}{\overset{[F]}{\sim}}$$

Let  $\sigma, \sigma'$  be systems. We define the relationship  $\overset{F}{\sim}$  by

$$\left(\sigma \overset{F}{\sim} \sigma'\right) \stackrel{def}{\iff} (\forall \pi \in \Pi_C, \forall n \in 1..E, \sigma[\pi][1][n] \notin \Pi_C \Leftrightarrow \sigma'[\pi][1][n] = \sigma[\pi][1][n])$$

Intuitively,  $\sigma \overset{F}{\sim} \sigma'$  if the Byzantine processes in each echo sample are identical in  $\sigma$  and  $\sigma'$ . Again,  $\overset{F}{\sim}$  is an equivalence relation that we can use to partition  $\mathcal{S}_i$ :

$$\mathcal{S}_1^i, \dots, \mathcal{S}_l^i = \frac{\mathcal{S}_i}{\overset{F}{\sim}}$$

and noting that, in Simplified Sieve, every correct process selects independently the correct processes in its echo samples, we have

$$|\mathcal{S}_1^i| = \dots = |\mathcal{S}_l^i|$$

Let  $\alpha$  be a sample-blind adversary. We define  $\mathcal{C}[\alpha]_j^i$  as the fraction of systems in  $\mathcal{S}_j^i$  whose consistency is compromised by  $\alpha$ :

$$\mathcal{C}[\alpha]_j^i = \frac{|\{\sigma \in \mathcal{S}_j^i \mid \alpha \searrow \sigma\}|}{|\mathcal{S}_j^i|}$$

From  $\mathcal{C}[\alpha]_j^i$  we can define

$$\mathcal{C}[\alpha]_*^i = \operatorname{argmax}_j \mathcal{C}[\alpha]_j^i$$

Intuitively,  $\mathcal{C}[\alpha]_*^i$  identifies the partition of  $\mathcal{S}_i$  that  $\alpha$  has the highest probability of compromising consistency.

### Chapter 3. Sieve

---

Finally, we define an **optimization function**  $\mathcal{O}[\alpha] : \mathcal{S} \rightarrow \mathcal{S}$ . Let  $\sigma \in \mathcal{S}_i$ , we define  $\mathcal{O}[\alpha]$  by

$$\begin{aligned} \mathcal{O}[\alpha](\sigma) &\in \mathcal{S}_{\mathcal{C}[\alpha]_*^i}^i \\ \sigma[\pi][m][n] \in \Pi_C &\implies \mathcal{O}[\alpha](\sigma)[\pi][m][n] = \sigma[\pi][m][n] \end{aligned}$$

As we previously proved, every  $\mathcal{S}_j^i$  has the same number of elements. Moreover,  $\mathcal{O}[\alpha]$  maps a system  $\sigma$  in  $\mathcal{S}_j^i$  to the corresponding system  $\sigma'$  in  $\mathcal{S}_{\mathcal{C}[\alpha]_*^i}^i$  such that every correct process in an echo sample in  $\sigma$  is identical to the corresponding process in  $\sigma'$ .

Therefore, for every  $\sigma, \sigma' \in \mathcal{S}_{\mathcal{C}[\alpha]_*^i}^i$ ,

$$|\mathcal{O}[\alpha]^{-1}(\sigma)| = |\mathcal{O}[\alpha]^{-1}(\sigma')| = \frac{|\mathcal{S}_i|}{|\mathcal{S}_1^i|}$$

**System masking.** Let  $\alpha$  be a sample-blind adversary, let  $\alpha' = \Delta_{bc}(\alpha)$ , let  $\sigma$  be a system, let  $\sigma'$  be an identical copy of  $\sigma$ . We start by noting that, if we couple Byzantine-counting decorator with  $\sigma'$ , we effectively obtain a system instance  $\delta$  with which  $\alpha$  directly exchanges invocations and responses. Here we show that the trace  $\tau(\alpha, \mathcal{O}[\alpha](\sigma))$  is identical to the trace  $\tau(\alpha, \delta)$ . Intuitively, this means that  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\mathcal{O}[\alpha](\sigma)$ , or it has been coupled with  $\sigma'$ , with Byzantine-counting decorator acting as an interface.

We previously proved that the trace exchanged between *sbadv* and *bcadv* is identical to the trace that *sbadv* would exchange with the system *sys'* that is obtained by replacing the Byzantine component of each correct process  $\pi$ 's echo samples in *sys* with *best.byzantine*[\mathcal{C}[\alpha]\_\*^i].

Let  $i \in \mathbb{N}$  such that  $\sigma \in \mathcal{S}_i$ . Procedure *bcadv.Init()* explicitly loops over all the possible values of *byzantine* that satisfy the condition  $|\text{byzantine}[\pi]| = |\text{sys.Byzantine}(\pi)|$  for all  $\pi \in \Pi_C$ . It then loops over every system  $\sigma$  that satisfies  $\sigma.\text{Byzantine}(\pi) = \text{byzantine}[\pi]$ , and counts the number of systems that  $\alpha$  compromises. It finally selects the value of *byzantine* that maximizes the number of compromissions. In doing so, *bcadv.Init()* is effectively looping over every  $\mathcal{S}_j^i$ , and selecting the  $j$  that maximizes the probability of  $\alpha$  compromising a random element of  $\mathcal{S}_j^i$ . Since *bcadv.Init()* is effectively masking  $\sigma$  with the element of  $\mathcal{S}_{\mathcal{C}[\alpha]_*^i}^i$  with which  $\sigma$  shares the correct component of every sample, the trace  $\tau(\alpha, \mathcal{O}[\alpha](\sigma))$  is identical to the trace  $\tau(\alpha, \delta)$ .

We previously proved that, if  $\alpha$  compromises the consistency of  $\mathcal{O}[\alpha](\sigma)$ , then  $\tau(\alpha, \mathcal{O}[\alpha](\sigma))$  compromises the consistency of  $\sigma$  as well. Noting that every invocation to *bcadv.Deliver(...)* or *bcadv.Echo(...)* is respectively forwarded to *sys.Deliver(...)* or *sys.Echo(...)*, we finally obtain that if  $\alpha$  compromises the consistency of  $\mathcal{O}[\alpha](\sigma)$ , then  $\alpha'$  compromises the consistency of  $\sigma$  as well.

**Adversarial power.** We can finally show that the adversarial power of  $\alpha'$  is greater than the adversarial power of  $\alpha$ . Let  $\sigma$  be a system.

Let  $i, j \in \mathbb{N}$  such that  $\sigma \in \mathcal{S}_j^i$ . The probability of  $\alpha$  compromising the consistency of  $\sigma$  is

$$\mathcal{P}[\alpha \searrow \sigma] = \mathcal{C}[\alpha]_j^i$$

and, since  $\alpha'$  compromises the consistency of  $\sigma$  if  $\alpha$  compromises the consistency of  $\mathcal{O}[\alpha](\sigma)$ , the probability of  $\alpha'$  compromising the consistency of  $\sigma$  is

$$\mathcal{P}[\alpha' \searrow \sigma] = \mathcal{P}[\alpha \searrow \mathcal{O}[\alpha](\sigma)] = \mathcal{C}[\alpha]_{\mathcal{C}[\alpha]_j^i}^i \geq \mathcal{C}[\alpha]_j^i = \mathcal{P}[\alpha \searrow \sigma]$$

Which proves that the adversarial power of  $\alpha'$  is greater or equal to the adversarial power of  $\alpha$ .  $\square$

### 3.11.7 Single-response adversary

**Lemma 26.** *The set of single-response adversaries  $\mathcal{A}_{sr}$  is optimal.*

*Proof.* We again prove the result using a decorator. Here we show that a decorator  $\Delta_{sr}$  exists such that, for every  $\alpha \in \mathcal{A}_{bc}$ , the adversary  $\alpha' = \Delta_{sr}(\alpha)$  is a single-response adversary, and as powerful as  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the sequential  $\alpha^+ = \Delta_{sr}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 13 implements **Single-response decorator**, a decorator that transforms a Byzantine-counting adversary into a single-response adversary. Provided with a Byzantine-counting adversary  $bcadv$ , Single-response decorator acts as an interface between  $bcadv$  and a system  $sys$ , effectively implementing a single-response adversary  $sradv$ . Single-response decorator exposes both the adversary and the system interfaces: the underlying adversary  $bcadv$  uses  $sradv$  as its system.

Single-response decorator works as follows:

- Procedure  $sradv.Init()$  initializes the following variables:
  - A *cache* set, initially empty: *cache* is used to store the first non-empty set returned from  $sys.State()$ .
  - A *poisoned* variable: *poisoned* is set to True if and only if at least one correct process in  $sys$  is poisoned. This condition is verified by looping over  $sys.Byzantine(\pi)$  for every correct process  $\pi$ .

## Chapter 3. Sieve

---

### Algorithm 13 Single-response decorator.

---

```
1: Implements:
2:   SingleResponseAdversary + CobSystem, instance sradv
3:
4: Uses:
5:   ByzantineCountingAdversary, instance bcadv, system sradv
6:   CobSystem, instance sys
7:
8: procedure sradv.Init() is
9:   cache =  $\emptyset$ ;   poisoned = False;   step = 0;
10:
11:  for all  $\pi \in \Pi_C$  do
12:    if  $|\text{sys.Byzantine}(\pi)| \geq \hat{E}$  then
13:      poisoned  $\leftarrow$  True;
14:    end if
15:  end for
16:
17:  bcadv.Init();
18:
19: procedure sradv.Step() is
20:  step  $\leftarrow$  step + 1;
21:
22:  if poisoned = False or step  $\leq (N - C)C^2$  then
23:    bcadv.Step();
24:  else if step  $\leq (N - C)C^2 + C$  then
25:    sys.Deliver( $\zeta(\text{step} - (N - C)C^2), 1$ );
26:  else
27:    sys.End();
28:  end if
29:
30: procedure sradv.Byzantine( $\pi$ ) is
31:  count =  $|\text{sys.Byzantine}(\pi)|$ ;
32:  return  $\{\perp\}^{\text{count}}$ ;
33:
34: procedure sradv.State() is
35:  return cache;
36:
37: procedure sradv.Sample(process, message) is
38:  raise error;
39:
```

---

---

```

40: procedure sradv.Deliver(process, message) is
41:   sys.Deliver(process, message);
42:
43:   if cache =  $\emptyset$  then
44:     cache  $\leftarrow$  sys.State();
45:   end if
46:
47: procedure sradv.Echo(process, sample, source, message) is
48:   sys.Echo(process, sample, source, message);
49:
50: procedure sradv.End() is
51:   sys.End();
52:

```

---

– A *step* variable, initially set to 0: at any time, *step* counts how many times *sradv.Step*() has been invoked.

- Procedure *sradv.Step*() increments *step*, then implements two different behaviors depending on the value of *poisoned*:
  - If *poisoned* = True, it forwards the call to *bcadv.Step*() for the first  $(N - C)C^2$  times. For the next  $C$  steps, it sequentially invokes *sys.Deliver*( $\zeta(1), 1$ ), ..., *sys.Deliver*( $\zeta(C), 1$ ). Finally, it calls *sys.End*().
  - If *poisoned* = False, it forwards the call to *bcadv.Step*().
- Procedure *sradv.Byzantine*(*process*) returns an array of *count* elements, *count* being the number of elements returned by *sys.Byzantine*(*process*). The array is filled with  $\perp$  values: since *bcadv* is Byzantine-counting, the content of the array is irrelevant.
- Procedure *sradv.State*() simply returns *cache*.
- Procedure *sradv.Sample*(*process, message*) is never called. This is due to the fact that *bcadv* is sample-blind.
- Procedure *sradv.Deliver*(*process, message*) forwards the call to *sys.Deliver*(*process, message*). Then, if *cache* is empty, it updates *cache* with *sys.State*().
- Procedure *sradv.Echo*(*process, sample, source, message*) simply forwards the call to *sys.Echo*(*process, sample, source, message*).
- Procedure *sradv.End*() simply forwards the call to *sys.End*().

**Correctness.** Here we prove that every adversary, when coupled with Single-response adversary:

### Chapter 3. Sieve

---

- Has a well-defined behavior. An adversary has undefined behavior if, at any point, the sequence of invocations and responses it exchanges with *sradv* is incompatible with every system.
- Is process-sequential, sequential, and Byzantine-counting.

We start by noting that *poisoned* = True if and only if *sys* is poisoned. Indeed, *sradv.Init()* explicitly checks if any correct process has at least  $\hat{E}$  Byzantine processes in its first echo sample.

We distinguish two cases, based on the value of *poisoned*. Let us assume *poisoned* = True. When *sradv.Step()* is invoked, the call is forwarded to *bcadv.Step()* only for the first  $(N - C)C^2$  times. Noting that *bcadv* is an auto-echo adversary, every call to *bcadv.Step()* results in a call to *sradv.Echo(...)*. For the next  $C$  steps, *sradv.Step()* sequentially causes  $\zeta(1), \zeta(2), \dots$  to *pb.Deliver* message 1. Finally, *sradv.Step()* invokes *sys.End()*. Therefore, *sradv* has a well defined behavior and implements a process-sequential adversary. Since it causes only message 1 to be *pb.Delivered*, *sradv* is also trivially sequential.

Let us assume *poisoned* = False. As we proved in Section 3.9, since *sys* is not poisoned, a correct process in *sys* will only deliver a message  $m$  as a result of an invocation to *sys.Deliver*( $\pi, m$ ) for some  $\pi \in \Pi_C$ . Until *cache*  $\neq \emptyset$ , *cache* is updated to *sys.State()* after every call to *sys.Deliver(...)*. Therefore, throughout the first phase, *sradv.State()* is always identical to *sys.State()*. The trace exchanged between *bcadv* and *sradv* is, therefore, trivially compatible with *sys*.

Throughout the second phase, we have *cache*  $\neq \perp$ . Since, throughout the first phase, *cache* is updated after every call to *sys.Deliver(...)*, only one message  $m^*$  exists such that, for some  $\pi^* \in \Pi_C$ ,  $(\pi^*, m^*) \in \text{cache}$ . Noting that *bcadv* is a non-redundant adversary, it will never invoke *sradv.Deliver(...)* on  $m^*$ : indeed, the value returned from *sradv.State()* never changes throughout the second phase. We define a system *sys'* by

$$sys'[\pi][m][n] = \begin{cases} sys[\pi][m][n] & \text{iff } m = m^* \text{ or } sys[\pi][m][n] \in \Pi \setminus \Pi_C \\ \pi^* & \text{otherwise} \end{cases}$$

The trace exchanged between *bcadv* and *sradv* is compatible with *sys'*. Indeed, for every  $\pi \in \Pi_C$ ,  $\pi$ 's sample for  $m^*$  in *sys* is identical to  $\pi$ 's echo sample for  $m^*$  in *sys'*: at any moment,  $\pi$  delivered  $m^*$  in *sys* if and only if  $\pi$  delivered  $m^*$  in *sys'*. For every  $\pi \in \Pi_C$  and  $m \neq m^* \in \mathcal{M}$ ,  $\pi$ 's every correct process in  $\pi$ 's sample for  $m$  is  $\pi^*$ . However,  $\pi^*$  *pb.Delivered*  $m^* \neq m$ . Therefore, since *sys'* is not poisoned, no correct process in *sys'* ever delivers a message other than  $m^*$ .

Every call to  $sradv.Deliver(\dots)$  and  $sradv.Echo(\dots)$  is respectively forwarded to  $sys.Deliver(\dots)$  and  $sys.Echo(\dots)$ . Moreover,  $bcadv$  is process-sequential and sequential. Therefore, if  $poisoned = \text{True}$ ,  $sradv$  is also process-sequential and sequential.

It is immediate to see that Single-response decorator always implements a Byzantine-counting adversary. Indeed, for any  $\pi \in \Pi_C$ ,  $sys.Byzantine(\pi)$  is only invoked from  $|sys.Byzantine(\pi)|$ .

**Single-response.** It is immediate to see that Single-response decorator always implements a single-response adversary. Indeed, when  $sys.State()$  returns a non-empty set for the first time,  $cache$  is set to a non-empty set, and  $sys.State()$  is never invoked again.

**Roadmap.** Let  $\alpha \in \mathcal{A}_{bc}$ , let  $\alpha' = \Delta_{sr}(\alpha)$ . Let  $\sigma$  be a system such that  $\alpha$  compromises the consistency of  $\sigma$ . Let  $\sigma'$  be an identical copy of  $\sigma$ . In order to prove that  $\alpha'$  is as powerful as  $\alpha$ , we prove that  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Poisoned case.** Noting that  $\alpha'$  is an auto-echo adversary, if  $\sigma$  is poisoned we immediately have that  $\alpha'$  compromises the consistency of  $\sigma'$ .

**Trace.** Let us assume that  $\sigma$  is not poisoned. We start by noting that, if we couple Single-response decorator with  $\sigma'$ , we effectively obtain a system instance  $\delta$  with which  $\alpha$  directly exchanges invocations and responses.

We start by defining a boolean sequence  $W$  by setting  $W_n = \text{True}$  if and only if, after the  $n$ -th invocation, two correct processes  $\pi, \pi'$  and two distinct messages  $m, m' \neq m$  exist such that  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$  in  $\sigma$ . Since  $\alpha$  compromises the consistency of  $\sigma$ , for some  $n$  we have  $W_n = \text{True}$ . Let

$$w = \min n \mid W_n = \text{True}$$

Here we show that, for every  $n \leq w$ , the trace  $\tau(\alpha, \sigma)_n$  is identical to the trace  $\tau(\alpha, \delta)_n$ . Intuitively, this means that, until the consistency of  $\sigma$  is compromised,  $\alpha$  has no way of *distinguishing* whether it has been coupled directly with  $\sigma$ , or it has been coupled with  $\sigma'$ , with Single-response decorator acting as an interface. We prove this by induction.

Let us assume

$$\begin{aligned} \tau(\alpha, \sigma) &= ((i_1, r_1), \dots) \\ \tau(\alpha, \delta) &= ((i'_1, r'_1), \dots) \\ i_j = i'_j, r_j = r'_j &\quad \forall j \leq n \end{aligned}$$

### Chapter 3. Sieve

---

We start by noting that, since  $\alpha$  is a deterministic algorithm, we immediately have

$$i_{n+1} = i'_{n+1}$$

and we need to prove that  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{Byzantine}, \pi)$ . Since procedure  $sradv.\text{Byzantine}(\pi)$  forwards the call to  $sys.\text{Byzantine}(\pi)$ ,  $bcadv$  is a Byzantine-counting adversary, and  $\sigma'$  is an identical copy of  $\sigma$ , with a minor abuse of notation we effectively have  $r_{n+1} = r'_{n+1}$ .

Let us assume that  $i_{n+1} = (\text{State})$ . We start by noting that, since all calls to  $sradv.\text{Deliver}(\dots)$  and  $sradv.\text{Echo}(\dots)$  are respectively forwarded to  $sys.\text{Deliver}(\dots)$  and  $sys.\text{Echo}(\dots)$ , a correct process  $\pi$  delivered  $m^*$  in  $\sigma$  if and only if it delivered  $m^*$  in  $\sigma'$  as well. As we proved, throughout the first phase,  $sradv.\text{State}()$  always returns the same value as  $sys.\text{State}()$ . Let us assume  $n > |\eta(\alpha, \sigma)|$ . Let  $m^*$  be the only message that was delivered by at least one correct process in  $\sigma$ . Noting that a correct process delivers a message only as a result of a call to  $sys.\text{Deliver}(\dots)$ , we have  $n < w$ . Therefore, by definition, no correct process in  $\sigma$  delivered a message other than  $m^*$ . Since  $\alpha$  is a non-redundant adversary, it never causes any correct process to  $pb.\text{Deliver } m^*$  throughout the second phase. As a result, no correct process delivers  $m^*$  in  $\sigma$  throughout the second phase. Therefore, all the processes that delivered  $m^*$  in  $\sigma$  are represented in  $cache$ , and no other process delivered a message  $m \neq m^*$ . Consequently,  $r_{n+1} = r'_{n+1}$ .

Noting that procedures  $\text{Deliver}(\dots)$  and  $\text{Echo}(\dots)$  never return a value, we trivially have that if  $i_{n+1} = (\text{Deliver}, \pi, m)$  or  $i_{n+1} = (\text{Echo}, \pi, s, \xi, m)$  then  $r_{n+1} = \perp = r'_{n+1}$ . By induction, we have that, for every  $n \leq w$ ,  $\tau(\alpha, \sigma)_n = \tau(\alpha, \delta)_n$ .

**Consistency of  $\sigma'$ .** We proved that, for all  $n \leq w$ ,  $\tau(\alpha, \sigma)_n = \tau(\alpha, \delta)_n$ . Moreover, we proved that if a correct process  $\pi$  eventually delivers a message  $m$  in  $\sigma$  before the  $w$ -th invocation, then  $\pi$  also delivers  $m$  in  $\sigma'$  before the  $w$ -th invocation.

Since  $\alpha$  compromises the consistency of  $\sigma$  after the  $w$ -th invocation, two correct processes  $\pi$ ,  $\pi'$  and two distinct messages  $m$ ,  $m' \neq m$  exist such that, in  $\sigma$ ,  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$  before the  $w$ -th invocation. Therefore, in  $\sigma'$ ,  $\pi$  delivered  $m$  and  $\pi'$  delivered  $m'$  before the  $w$ -th invocation. Therefore  $\alpha'$  compromises the consistency of  $\sigma'$ .

Consequently, the adversarial power of  $\alpha$  is equal to the adversarial power of  $\alpha' = \Delta_{sr}(\alpha)$ , and the lemma is proved. □



### 3.11.8 Two-phase adversary

**Lemma 27.** *The set of two-phase adversaries  $\mathcal{A}_{tp}$  is optimal.*

*Proof.* We again prove the result using a decorator. Here we show that a decorator  $\Delta_{tp}$  exists such that, for every  $\alpha \in \mathcal{A}_{sm}$ , the adversary  $\alpha' = \Delta_{tp}(\alpha)$  is a two-phase adversary, and more powerful than  $\alpha$ . If this is true, then the lemma is proved: let  $\alpha^*$  be an optimal adversary, then the sequential  $\alpha^+ = \Delta_{tp}(\alpha^*)$  is optimal as well.

**Decorator.** Algorithm 14 implements **Two-phase decorator**, a decorator that transforms a state-polling adversary into a two-phase adversary. Provided with a state-polling adversary  $spadv$ , Two-phase decorator acts as an interface between  $spadv$  and a system  $sys$ , effectively implementing a single-response adversary  $tpadv$ . Two-phase decorator exposes both the adversary and the system interfaces: the underlying adversary  $spadv$  uses  $tpadv$  as its system.

Two-phase decorator works as follows:

- Procedure  $tpadv.Init()$  initializes a *invocations* variable: at any time, *invocations* counts the number of invocations issued by  $spadv$ .
- Procedure  $compatible(invocations)$  returns a set of systems  $\sigma$  that satisfy the following properties:
  - For every correct process  $\pi$ , the number of Byzantine processes in  $\pi$ 's first echo sample is identical in  $\sigma$  and  $sys$ .
  - The length  $|\eta(\alpha, \sigma)|$  of the first phase when  $\alpha$  is coupled with  $\sigma$  is equal to *invocations*.
- Procedure  $tpadv.Step()$  simply forwards the call to  $spadv.Step()$ .
- Procedure  $tpadv.Byzantine(process)$  increments *invocations*, then returns an array of *count* elements, *count* being the number of elements returned from  $sys.Byzantine(process)$ . The array is filled with  $\perp$  values: since  $spadv$  is Byzantine-counting, the content of the array is irrelevant.
- Procedure  $tpadv.State()$  increments *invocations*. It then returns an empty set if  $sys.State()$  is empty. If  $sys.State()$  is not empty, the procedure returns, among all the possible responses that are compatible with the trace exchanged between  $spadv$  and  $tpadv$ , the one that maximizes the probability of  $spadv$  compromising the consistency of  $sys$ . This is achieved as follows:
  - The procedure loops over every system  $\sigma$  in the set  $compatible(invocations)$ . In doing so, the procedure loops over every system  $\sigma$  such that:  $\sigma$  has the same

### Algorithm 14 Two-phase decorator.

---

```
1: Implements:
2:   TwoPhaseAdversary + CobSystem, instance tpadv
3:
4: Uses:
5:   StatePollingAdversary, instance spadv, system tpadv
6:   CobSystem, instance sys
7:
8: procedure tpadv.Init() is
9:   invocations = 0;
10:  spadv.Init();
11:
12: procedure compatible(invocations) is
13:  systems =  $\emptyset$ ;
14:
15:  for all  $\sigma \in \mathcal{S}$  do
16:    match = True;
17:    for all  $\pi \in \Pi_C$  do
18:      if  $|\sigma.\text{Byzantine}(\pi)| \neq |\text{sys}.\text{Byzantine}(\pi)|$  then
19:        match = False;
20:      end if
21:    end for
22:
23:    if match = True and  $|\eta(\alpha, \sigma)| = \text{invocations}$  then
24:      systems  $\leftarrow \text{systems} \cup \{\sigma\}$ ;
25:    end if
26:  end for
27:
28:  return systems;
29:
30: procedure tpadv.Step() is
31:  spadv.Step();
32:
33: procedure tpadv.Byzantine(process) is
34:  invocations  $\leftarrow \text{invocations} + 1$ ;
35:  count = sys.Byzantine(process);
36:  return  $\{\perp\}^{\text{count}}$ ;
37:
```

---

---

```

38: procedure tpadv.State() is
39:   invocations  $\leftarrow$  invocations + 1;
40:
41:   if sys.State()  $\neq$   $\emptyset$  then
42:     outcomes =  $\emptyset$ ;
43:     for all  $\sigma \in$  compatible(invocations) do
44:       (invocation, response) =  $\tau(\alpha, \sigma)_{invocations}$ ;
45:       outcomes  $\leftarrow$  outcomes  $\cup$   $\{(response, \tau(\alpha, \sigma))\}$ ;
46:     end for
47:
48:     best.response =  $\perp$ ; best.compromissions = 0;
49:
50:     for all (response, fulltrace)  $\in$  outcomes do
51:       compromissions = 0;
52:
53:       for all  $\sigma \in$  compatible(invocations) do
54:         if fulltrace  $\setminus$   $\sigma$  then
55:           compromissions  $\leftarrow$  compromissions + 1;
56:         end if
57:       end for
58:
59:       if compromissions > best.compromissions then
60:         best.response  $\leftarrow$  response;
61:         best.compromissions = compromissions;
62:       end if
63:     end for
64:
65:     return best.response;
66:   else
67:     return  $\emptyset$ ;
68:   end if
69:
70: procedure tpadv.Sample(process, message) is
71:   raise error;
72:
73: procedure tpadv.Deliver(process, message) is
74:   invocations  $\leftarrow$  invocations + 1;
75:   sys.Deliver(process, message);
76:
77: procedure tpadv.Echo(process, sample, source, message) is
78:   invocations  $\leftarrow$  invocations + 1;
79:   sys.Echo(process, sample, source, message);
80:
81: procedure tpadv.End() is
82:   sys.End();
83:

```

---

Byzantine count as  $sys$ ; when  $\alpha$  is coupled with  $\sigma$ , it concludes the first phase in exactly  $invocations$  invocations.

- For every process  $\sigma$  in  $compatible(invocations)$ , the procedure stores in a set  $outcome$  a  $(response, fulltrace)$  pair,  $response$  being the  $State()$  of  $\sigma$  at the end of the first phase ( $response$  is extracted from  $\tau(\alpha, \sigma)_{invocations}$ ), and  $fulltrace$  being  $\tau(\alpha, \sigma)$ , the full trace exchanged between  $\alpha$  and  $\sigma$ .
- For every  $(response, fulltrace)$  in  $outcomes$ , the procedure loops over every system  $\sigma$  in  $compatible(invocations)$ , and counts the number of systems whose consistency is compromised by  $fulltrace$ . The procedure returns the value of  $response$  that maximizes the number of systems in  $compatible(invocations)$  whose consistency is compromised by  $fulltrace$ .
- Procedure  $tpadv.Sample(process, message)$  is never called. This is due to the fact that  $spadv$  is sample-blind.
- Procedure  $tpadv.Deliver(process, message)$  increments  $invocations$ , then forwards the call to  $sys.Deliver(process, message)$ .
- Procedure  $tpadv.Echo(process, sample, source, message)$  increments  $invocations$ , then forwards the call to  $sys.Echo(process, sample, source, message)$ .
- Procedure  $tpadv.End()$  simply forwards the call to  $sys.End()$ .

**Correctness.** Here we prove that every adversary, coupled with Two-phase decorator:

- Has a well-defined behavior. An adversary has undefined behavior if, at any point, the sequence of invocations and responses it exchanges with  $tpadv$  is incompatible with every system.
- Is Byzantine-counting and single-response.

We start by noting that, since  $invocations$  is incremented every time  $spadv$  issues an invocation, when  $tpadv.State()$  is invoked and  $sys.State() \neq \emptyset$  we have  $invocations = |\eta(\alpha, \sigma)|$ .

Every invocation of a procedure in  $tpadv$  is always forwarded to the corresponding procedure in  $sys$ , except for  $tpadv.State()$ . Whenever  $sys.State() = \emptyset$ ,  $tpadv.State()$  returns  $\emptyset$  as well. Therefore, up to the  $(|\eta(\alpha, sys)| - 1)$ -th invocation, the trace exchanged between  $spadv$  and  $tpadv$  is trivially compatible with  $sys$ .

Procedure  $compatible(invocations)$  returns all systems  $\sigma$  such that the condition  $|\sigma.Byzantine(\pi)| = |sys.Byzantine(\pi)|$  holds for all  $\pi \in \Pi_C$ , and  $|\eta(\alpha, \sigma)| = invocations = |\eta(\sigma, sys)|$ . It is immediate to see that  $compatible(invocations)$  is non-empty, as it includes  $sys$ . Every system  $\sigma \in compatible(invocations)$  is compatible with the first  $n - 1$  elements

of the trace exchanged between  $spadv$  and  $tpadv$ . Procedure  $tpadv.State()$  then returns a response  $response$ , such that

$$\tau(\alpha, \sigma)_{invocations} = ((State), response)$$

for some  $\sigma \in compatible(invocations)$ . Therefore, the first  $n$  elements of the trace exchanged between  $spadv$  and  $tpadv$  is compatible with  $\sigma$ . Due to Lemma 15, the entire trace exchanged between  $spadv$  and  $tpadv$  is compatible with  $\sigma$ .

It is easy to see that  $tpadv$  always implements a Byzantine-counting and single-response adversary. Indeed: whenever  $tpadv$  invokes  $sys.Byzantine(\pi)$ , it invokes  $|sys.Byzantine(\pi)|$ ;  $tpadv.State()$  returns a non-empty set if and only if  $sys.State()$  returns a non-empty set, and  $spadv$  is a single-response adversary.

**Two-phase.** It is immediate to see that Two-phase decorator always implements a two-phase adversary. Indeed, whenever  $tpadv$  invokes  $sys.State()$ , it invokes  $(sys.State() \neq \emptyset)$ .

**System partitioning.** Let  $\alpha$  be a state-polling adversary, let  $\sigma$  be a system. Let us denote with  $\mathcal{S}^*$  the set of non-poisoned systems. We denote with  $\overset{\alpha}{\sim}$  the two conditions  $\forall \pi \in \Pi_C, \forall m \in \mathcal{M}$ ,

$$|\{n \in 1..E \mid \sigma[\pi][m][n] \in \Pi_C\}| = |\{n \in 1..E \mid \sigma'[\pi][m][n] \in \Pi_C\}|$$

and

$$|\eta(\alpha, \sigma)| = |\eta(\alpha, \sigma')|$$

It is immediate to see that  $\overset{\alpha}{\sim}$  is an equivalence relation, and we can use  $\overset{\alpha}{\sim}$  to partition  $\mathcal{S}^*$ :

$$\mathcal{S}[\alpha]_1, \dots, \mathcal{S}[\alpha]_h = \frac{\mathcal{S}^*}{\overset{\alpha}{\sim}}$$

Let  $i \in 1..h$ . Due to Lemma 14, we have

$$\forall \sigma, \sigma' \in \mathcal{S}[\alpha]_i, \forall n < |\eta(\alpha, \sigma)|, \tau(\alpha, \sigma)_n = \tau(\alpha, \sigma')_n$$

Moreover, since  $\sigma$  is not poisoned,  $\eta(\alpha, \sigma)$  includes at least one call to  $Deliver(\dots)$ . Therefore, for every  $i \in 1..h$ , let  $\sigma \in \mathcal{S}[\alpha]_i$ , we can define a function  $\delta[\alpha]_i : \mathcal{M} \times 1..(|\eta(\alpha, \sigma)|)$  by

$$\pi \in \delta[\alpha]_i(m, n) \stackrel{def}{\iff} \exists j < n \mid \tau(\alpha, \sigma)_j = ((Deliver, \pi), \perp)$$

Intuitively,  $\delta[\alpha]_i(m, n)$  represents the set of correct processes that  $\alpha$  causes to pb.Deliver  $m$  before the  $n$ -th invocation, when  $\alpha$  is coupled with any  $\sigma \in \mathcal{S}[\alpha]_i$ .

### Chapter 3. Sieve

---

We additionally define  $\pi[\alpha]_i: \mathcal{M} \rightarrow \mathbb{P}(\Pi_C)$ ,  $\pi^-[\alpha]_i: \mathcal{M} \rightarrow \mathbb{P}(\Pi_C)$  by, let  $\sigma \in \mathcal{S}[\alpha]_i$ ,

$$\begin{aligned}\pi[\alpha]_i(m) &= \delta[\alpha]_i(m, |\eta(\alpha, \sigma)|) \\ \pi^-[\alpha]_i(m) &= \delta[\alpha]_i(m, |\eta(\alpha, \sigma)| - 1)\end{aligned}$$

Intuitively,  $\pi[\alpha]_i(m)$  represents the set of correct processes that  $\alpha$  causes to pb.Deliver  $m$  throughout the first phase, when  $\alpha$  is coupled with any  $\sigma \in \mathcal{S}[\alpha]_i$ . Noting that  $\alpha$  is a state-polling adversary, and  $\sigma$  is not poisoned, then  $\pi^-[\alpha]_i(m)$  represents the set of correct processes that  $\alpha$  causes to pb.Deliver  $m$  throughout the first phase when  $\alpha$  is coupled with any  $\sigma \in \mathcal{S}[\alpha]_i$ , excluding the last invocation to *Deliver*(..) in  $\eta(\alpha, \sigma)$ .

Finally, we define  $m(\alpha)_i$  by, let  $\sigma \in \mathcal{S}[\alpha]_i$

$$\tau(\alpha, \sigma)_{|\eta(\alpha, \sigma)|} = ((\text{Deliver}, \pi \in \Pi_C, m), \perp)$$

Intuitively,  $m(\alpha)_i$  is the last message that  $\alpha$  causes a correct process to pb.Deliver throughout the first phase, when  $\alpha$  is coupled with any  $\sigma \in \mathcal{S}[\alpha]_i$ . Noting that  $\alpha$  is a state-polling adversary, and that  $\sigma$  is not poisoned,  $m$  is the only message delivered by at least one correct process at the end  $\eta(\alpha, \sigma)$ .

Let  $\sigma, \sigma' \in \mathcal{S}[\alpha]_i$ . We can prove that  $\overset{\alpha}{\sim}$  can be equivalently restated as  $\forall \pi \in \Pi_C, \forall m \in \mathcal{M}$

$$|\{n \in 1..E \mid \sigma[\pi][m][n] \in \Pi_C\}| = |\{n \in 1..E \mid \sigma'[\pi][m][n] \in \Pi_C\}|$$

and

$$\begin{aligned}\nexists \pi \in \Pi_C \mid & \\ & |\{n \in 1..E \mid \sigma[\pi][m(\alpha)_i][n] \in (\pi^-[\alpha]_i(m(\alpha)_i) \cup (\Pi \setminus \Pi_C))\}| \geq \hat{E} \\ \exists \pi \in \Pi_C \mid & \\ & |\{n \in 1..E \mid \sigma[\pi][m(\alpha)_i][n] \in (\pi[\alpha]_i(m(\alpha)_i) \cup (\Pi \setminus \Pi_C))\}| \geq \hat{E} \\ \nexists m \neq m(\alpha)_i, \pi \in \Pi_C \mid & \\ & |\{n \in 1..E \mid \sigma[\pi][m][n] \in (\pi[\alpha]_i(m) \cup (\Pi \setminus \Pi_C))\}| \geq E\end{aligned}$$

Indeed, we are restating the condition  $|\eta(\alpha, \sigma)| = |\eta(\alpha, \sigma')|$  with the following conditions:

- No correct process has, in its echo sample for  $m(\alpha)_i$ , at least  $\hat{E}$  processes that are either Byzantine, or pb.Deliver  $m(\alpha)_i$  as a result of any invocation of *Deliver*(..) in  $\eta(\alpha, \sigma)$  except the last. This encodes the condition that no correct process delivers  $m(\alpha)_i$  before the last invocation of *Deliver*(..) in  $\eta(\alpha, \sigma)$ .
- At least one correct process has, in its echo sample for  $m(\alpha)_i$ , at least  $\hat{E}$  processes that are either Byzantine, or pb.Deliver  $m(\alpha)_i$  throughout the first phase when  $\alpha$  is coupled with  $\sigma$ . This encodes the condition that at least one correct process delivers  $m(\alpha)$  after

the last invocation of *Deliver*(...) in  $\eta(\alpha, \sigma)$ .

- No correct process has, in its echo sample for  $m \neq m(\alpha)_i$ , at least  $\hat{E}$  processes that are either Byzantine, or pb.Deliver  $m$  throughout the first phase when  $\alpha$  is coupled with  $\sigma$ . This encodes the condition that no message is delivered before  $m(\alpha)_i$ .

Let  $\sigma, \sigma' \in \mathcal{S}[\alpha]_i$ . We denote with  $\overset{m}{\sim}$  the condition  $\forall \pi \in \Pi_C$ ,

$$\sigma[\pi][m(\alpha)_i] = \sigma'[\pi][m(\alpha)_i]$$

Again,  $\overset{m}{\sim}$  is an equivalence relation, and can be used to partition  $\mathcal{S}[\alpha]_i$ :

$$\mathcal{S}[\alpha]_1^i, \dots, \mathcal{S}[\alpha]_l^i = \frac{\mathcal{S}[\alpha]_i}{\overset{m}{\sim}}$$

Let  $\bar{\sigma} \in \mathcal{S}[\alpha]_i$ , let  $\tau = \tau(\alpha, \bar{\sigma})$ . For any  $\sigma \in \mathcal{S}[\alpha]_i$ ,  $\tau$  compromises the consistency of  $\sigma$  if  $\tau$  causes at least one correct process to deliver a message  $m' \neq m(\alpha)_i$  throughout the second phase. Since this condition is independent from the echo sample for  $m(\alpha)_i$  of any correct process, we finally have that, for every  $j \in 1..l$ ,

$$\mathcal{P}\left[\tau \searrow \left(\sigma \in \mathcal{S}[\alpha]_j^i\right)\right] = \mathcal{P}\left[\tau \searrow \left(\sigma \in \mathcal{S}[\alpha]_i\right)\right]$$

**Adversarial power.** Here we prove that  $\alpha' = \Delta_{tp}(\alpha)$  is more powerful than  $\alpha$ . Let  $\bar{\sigma}$  denote a random system in  $\mathcal{S}$ .

Let us assume that  $\bar{\sigma}$  is poisoned. Since both  $\alpha$  and  $\alpha'$  are auto-echo adversaries, both compromise  $\bar{\sigma}$ .

Let us assume that  $\bar{\sigma}$  is not poisoned. For some  $i, j$ , we therefore have  $\sigma \in \mathcal{S}[\alpha]_j^i$ . When *tpadv.State*() is invoked and *sys.State*()  $\neq \emptyset$ , the procedure returns a response *best.response* such that the trace  $\tau^*$  that  $\alpha$  issues as a result of *best.response* satisfies

$$\tau^* = \underset{\tau}{\operatorname{argmax}} \mathcal{P}\left[\tau \searrow \left(\sigma \in \mathcal{S}[\alpha]_i\right)\right]$$

As we proved in the previous section, we therefore have

$$\mathcal{P}\left[\tau^* \searrow \left(\sigma \in \mathcal{S}[\alpha]_i\right)\right] \geq \mathcal{P}\left[\tau \searrow \left(\sigma \in \mathcal{S}[\alpha]_i\right)\right] = \mathcal{P}\left[\tau \searrow \left(\sigma \in \mathcal{S}[\alpha]_j^i\right)\right]$$

which proves that, if  $\sigma$  is not poisoned, then the probability of  $\alpha'$  compromising  $\sigma$  is greater or equal to the probability of  $\alpha$  compromising  $\sigma$ .

The adversarial power of  $\alpha'$  is therefore greater or equal to the adversarial power of  $\alpha$ , and the lemma is proved.  $\square$





## 4 Contagion

In this chapter, we present in detail the **probabilistic reliable broadcast** abstraction and discuss its properties. We then present Contagion, an algorithm that implements probabilistic reliable broadcast, and evaluate its **security** and **complexity** as a function of its **parameters**.

The probabilistic reliable broadcast abstraction allows the entire set of correct processes to agree on a single message from a potentially Byzantine designated sender. Probabilistic reliable broadcast is a strictly stronger abstraction than probabilistic consistent broadcast: in the case of a Byzantine sender, while probabilistic consistent broadcast only guarantees that every correct process that delivers a message delivers the same message (**consistency**), probabilistic reliable broadcast also guarantees that either no or every correct process delivers a message (**totality**).

### 4.1 Interface

The **probabilistic reliable broadcast** interface (instance  $prb$ , sender  $\sigma$ ) exposes the following two events:

- **Request:**  $\langle prb.Broadcast \mid m \rangle$ : Broadcasts a message  $m$  to all processes. This is only used by  $\sigma$ .
- **Indication:**  $\langle prb.Deliver \mid m \rangle$ : Delivers a message  $m$  broadcast by process  $\sigma$ .

For any  $\epsilon \in [0, 1]$ , we say that probabilistic reliable broadcast is  $\epsilon$ -secure if:

1. **No duplication:** No correct process delivers more than one message.
2. **Integrity:** If a correct process delivers a message  $m$ , and  $\sigma$  is correct, then  $m$  was previously broadcast by  $\sigma$ .
3.  **$\epsilon$ -Validity:** If  $\sigma$  is correct, and  $\sigma$  broadcasts a message  $m$ , then  $\sigma$  eventually delivers  $m$  with probability at least  $(1 - \epsilon)$ .

4.  $\epsilon$ -**Totality**: If a correct process delivers a message, then every correct process eventually delivers a message with probability at least  $(1 - \epsilon)$ .
5.  $\epsilon$ -**Consistency**: Every correct process that delivers a message delivers the same message with probability at least  $(1 - \epsilon)$ .

### 4.2 Algorithm

Algorithm 15 implements Contagion. Let  $\pi$  be a correct process, let  $m$  be a message. Contagion securely distributes a single message across the system as follows:

- Initially, probabilistic consistent broadcast consistently distributes the same message to a subset of the correct processes.
- $\pi$  can issue a Ready message for more than one message.  $\pi$  issues a Ready message  $m$  when either:
  - $\pi$  receives  $m$  from probabilistic consistent broadcast, or
  - $\pi$  collects enough Ready messages for  $m$  from its *ready sample*.
- $\pi$  delivers  $m$  if  $m$  is the first message for which  $\pi$  collected enough Ready messages from its delivery sample.

A correct process collects Ready messages from two randomly selected samples, the *ready sample* of size  $R$ , and the *delivery sample* of size  $D$ . A correct process issues a Ready message for  $m$  upon collecting  $\hat{R}$  Ready messages for  $m$  from its ready sample, and it delivers  $m$  upon collecting  $\hat{D}$  Ready messages for  $m$  from its delivery sample. We discuss the values of the four parameters of Contagion in Sections 4.4, 4.9 and 4.10.

**Sampling.** Upon initialization (line 12), a correct process randomly selects a **ready sample**  $\mathcal{R}$  of size  $R$ , and a **delivery sample**  $\mathcal{D}$  of size  $D$ . Samples are selected with replacement by repeatedly calling  $\Omega$  (Algorithm 2, line 4).

**Publish-subscribe.** Like Sieve, Contagion uses publish-subscribe to reduce its communication complexity. This is achieved by having each correct process send Ready messages only to its ready subscription set (lines 32 and 50), and accept Ready messages only from its ready and delivery samples (lines 40 and 43).

**Consistent broadcast.** The designated sender  $\sigma$  initially broadcasts its message using probabilistic consistent broadcast (line 27). When message  $m$  is `pcb.Delivered` (correctly signed by  $\sigma$ ) (line 29), a correct process sends a Ready message for  $m$  (line 33) to all the processes in its ready subscription set.

**Algorithm 15 Contagion**


---

```

1: Implements:
2:   ProbabilisticReliableBroadcast, instance prb
3:
4: Uses:
5:   AuthenticatedPointToPointLinks, instance al
6:   ProbabilisticConsistentBroadcast, instance pcb
7:
8: Parameters:
9:    $R$ : ready sample size       $\hat{R}$ : contagion threshold
10:   $D$ : delivery sample size     $\hat{D}$ : delivery threshold
11:
12: upon event  $\langle prb.Init \rangle$  do
13:    $ready = \emptyset$ ;    $delivered = \text{False}$ ;  $\tilde{\mathcal{R}} = \emptyset$ ;
14:
15:    $\mathcal{R} = \text{sample}(\text{ReadySubscribe}, R)$ ;
16:    $\mathcal{D} = \text{sample}(\text{ReadySubscribe}, D)$ ;
17:
18:    $replies.ready = \{\emptyset\}^R$ ;    $replies.delivery = \{\emptyset\}^D$ 
19:
20: upon event  $\langle al.Deliver \mid \pi, [\text{ReadySubscribe}] \rangle$  do
21:   for all  $(message, signature) \in ready$  do
22:     trigger  $\langle al.Send \mid \pi, [\text{Ready}, message, signature] \rangle$ ;
23:   end for
24:    $\tilde{\mathcal{R}} \leftarrow \tilde{\mathcal{R}} \cup \{\pi\}$ ;
25:
26: upon event  $\langle prb.Broadcast \mid message \rangle$  do ▷ only process  $\sigma$ 
27:   trigger  $\langle pcb.Broadcast \mid [\text{Send}, message, sign(message)] \rangle$ ;
28:
29: upon event  $\langle pcb.Deliver \mid [\text{Send}, message, signature] \rangle$  do
30:   if  $verify(\sigma, message, signature)$  then
31:      $ready \leftarrow ready \cup \{(message, signature)\}$ ;
32:     for all  $\rho \in \tilde{\mathcal{R}}$  do
33:       trigger  $\langle al.Send \mid \rho, [\text{Ready}, message, signature] \rangle$ ;
34:     end for
35:   end if
36:

```

---

```

37: upon event  $\langle al.Deliver \mid \pi, [Ready, message, signature] \rangle$  do
38:   if  $verify(\sigma, message, signature)$  then
39:      $reply = (message, signature);$ 
40:     if  $\pi \in \mathcal{R}$  then
41:        $replies.ready[\pi] \leftarrow replies.ready[\pi] \cup \{reply\};$ 
42:     end if
43:     if  $\pi \in \mathcal{D}$  then
44:        $replies.delivery[\pi] \leftarrow replies.delivery[\pi] \cup \{reply\}$ 
45:     end if
46:   end if
47:
48: upon exists message such that  $|\{\rho \in \mathcal{R} \mid (message, signature) \in replies.ready[\rho]\}| \geq \hat{R}$ 
   do
49:    $ready \leftarrow ready \cup \{(message, signature)\};$ 
50:   for all  $\rho \in \tilde{\mathcal{R}}$  do
51:     trigger  $\langle al.Send \mid \rho, [Ready, message, signature] \rangle;$ 
52:   end for
53:
54: upon exists message such that  $|\{\rho \in \mathcal{D} \mid (message, signature) \in$ 
    $replies.delivery[\rho]\}| \geq \hat{D}$  and  $delivered = False$  do
55:    $delivered \leftarrow True;$ 
56:   trigger  $\langle prb.Deliver \mid message \rangle;$ 
57:

```

---

**Contagion.** Upon collecting  $\hat{R}$  Ready messages for a message  $m$  (line 48), a correct process sends a Ready message for  $m$  (line 51) to all the nodes in its ready subscription set.

**Delivery.** Upon collecting  $\hat{D}$  Ready messages for a message  $m$  for the first time, (line 54), a correct process delivers  $m$  (line 56).

### 4.3 No duplication and integrity

We start by verifying that Contagion satisfies both **no duplication** and **integrity**.

**Theorem 10.** *Contagion satisfies no duplication.*

*Proof.* A message is delivered (line 56) only if the variable *delivered* is equal to `False` (line 54). Before any message is delivered, *delivered* is set to `True`. Therefore no more than one message is ever delivered.  $\square$

**Theorem 11.** *Contagion satisfies integrity.*

*Proof.* Upon receiving a Ready message, a correct process checks its signature against the

public key of the designated sender  $\sigma$  (line 38), and the  $(message, signature)$  pair is added to the  $replies.delivery$  variable only if this check succeeds. Moreover, a message is delivered only if it is represented at least  $\hat{D}$  times in  $replies.delivery$  (line 54).

If  $\sigma$  is correct, it only signs  $message$  when broadcasting (line 27). Since we assume that cryptographic signatures cannot be forged, this implies that the message was previously broadcast by  $\sigma$ .  $\square$

#### 4.4 Validity

We now compute, given  $D$  and  $\hat{D}$ , the  $\epsilon$ -**validity** of Contagion. To this end, we prove one preliminary lemma.

**Lemma 28.** *In an execution of Contagion, if  $pcb$  satisfies total validity and the sender has no more than  $D - \hat{D}$  Byzantine processes in its delivery sample, then  $prb$  satisfies validity.*

*Proof.* Let  $m$  be the message broadcast by the correct sender  $\sigma$ . Since  $pcb$  satisfies total validity, every correct process eventually issues a  $Ready(m)$  message (i.e., a Ready message for  $m$ ) (line 33).

By hypothesis,  $\sigma$  has no more than  $D - \hat{D}$  Byzantine processes in its echo sample. Obviously,  $\sigma$  has at least  $\hat{D}$  correct processes in its echo sample. Therefore,  $\sigma$  eventually receives at least  $\hat{D}$   $Ready(m)$  messages (line 37), and delivers  $m$  (line 56).  $\square$

Lemma 28 allows us to bound the  $\epsilon$ -validity of Contagion, given  $D$  and  $\hat{D}$ .

**Theorem 12.** *Contagion satisfies  $\epsilon_v$ -validity, with*

$$\begin{aligned} \epsilon_v &\leq \epsilon_v^{pcb} + (1 - \epsilon_v^{pcb})\epsilon_o \\ \epsilon_o &= \sum_{\bar{F}=D-\hat{D}+1}^D Bin[D, f](\bar{F}) \end{aligned} \tag{4.1}$$

*if the underlying abstraction of  $pcb$  satisfies  $\epsilon_v^{pcb}$ -total validity.*

*Proof.* We compute a bound on  $\epsilon_v$  by assuming that, if the total validity of the underlying  $pcb$  instance is compromised, the validity of  $prb$  is compromised as well. Following from Lemma 28, the validity of  $prb$  can be compromised only if the total validity of  $pcb$  is compromised as well, or if  $\sigma$  has more than  $D - \hat{D}$  Byzantine processes in its delivery sample.

Since procedure *sample* independently picks  $D$  processes with replacement, each element of a correct process' echo sample has an independent probability  $f$  of being Byzantine, i.e., the number of Byzantine processes in a correct delivery sample is binomially distributed.

Therefore,  $\sigma$  has a probability  $\epsilon_o$  of having more than  $D - \hat{D}$  Byzantine processes in its delivery sample.  $\square$

### 4.5 Adversarial execution

In this section, we define the model underlying an adversarial execution of Contagion. Here, a Byzantine adversary is an agent that acts upon a system with the goal to compromise its consistency and / or totality. The main goal of this section is to formalize the information available to the adversary, and the set of actions that it can perform on the system throughout an adversarial execution.

Throughout the rest of this section, we bound the probability of compromising the consistency and totality of Contagion by assuming that, if the consistency of the pcb instance used in Contagion is compromised, then both the consistency and the totality of Contagion are compromised as well. In what follows, therefore, we assume that Sieve satisfies consistency.

#### 4.5.1 Model

Let  $\pi$  be any correct process. We make the following assumptions about an adversarial execution of Contagion:

- As we established in Section 1.3, the adversary does not know which correct processes are in  $\pi$ 's ready or delivery samples. The adversary knows, however, which Byzantine processes are in  $\pi$ 's ready sample, and which Byzantine processes are in  $\pi$ 's delivery sample.
- At any time, the adversary knows the set of messages for which  $\pi$  sent a Ready message.
- At any time, the adversary knows if  $\pi$  delivered a message. If  $\pi$  delivered a message, then the adversary knows which message did  $\pi$  deliver.
- The adversary can arbitrarily cause  $\pi$  to pcb.Deliver a given message  $m^*$ . Since we assume that the underlying pcb instance satisfies consistency, the adversary cannot cause two correct processes to pcb.Deliver two different messages.

Throughout an adversarial execution of Contagion, an adversary performs a sequence of minimal operations on the system. Each operation consists of either of the following:

- Selecting a correct process that did not pcb.Deliver  $m^*$  and causing it to pcb.Deliver  $m^*$ .
- Selecting a Byzantine process and causing it to issue a Ready message to a correct process.

As a result of each operation, zero or more processes send a Ready message and/or deliver a message. The adversary is successful if, at the end of the adversarial execution, either the consistency or the totality of the system is compromised.

## 4.6 Epidemic processes

In the next sections, we compute bounds for the  $\epsilon$ -consistency and  $\epsilon$ -totality of Contagion. In order to do so, in this section we study the *feedback mechanism* produced by Ready messages in an execution of Contagion.

As we discussed in Section 4.2, a correct process issues a Ready message for a message  $m$  after either `pcb.Delivering  $m$`  (line 33) or collecting at least  $\hat{R}$  Ready( $m$ ) messages from its ready sample (line 51). We formalize this observation in the following definition.

**Definition 29** (Ready, E-ready, R-ready). Let  $\pi$  be a correct process, let  $m$  be a message. Throughout an execution of Contagion,  $\pi$  is **E-ready** for  $m$  if  $\pi$  eventually `pcb.Delivers  $m$` ;  $\pi$  is **R-ready** for  $m$  if  $\pi$  eventually receives at least  $\hat{R}$  Ready( $m$ ) messages from its ready sample;  $\pi$  is **ready** for  $m$  if  $\pi$  is either E-ready or R-ready for  $m$ .

We note how a correct process can simultaneously be E-ready and R-ready for the same message.

It is easy to observe that the R-ready condition creates a feedback process: as a result of a correct process being R-ready for a message  $m$ , it issues a Ready( $m$ ) message that might cause other correct processes to become R-ready for  $m$  as well.

Intuitively, this feedback process is designed to have two stable configurations:

- **Few processes are ready:** the fraction of correct processes that are E-ready for a message  $m$  is significantly smaller than  $\hat{R}/R$ . As a result, the probability of a correct process being R-ready for  $m$  becomes very small, and the set of processes that are ready for  $m$  is, with high probability, nearly identical to the set of processes that are E-ready for  $m$ .
- **All processes are ready:** the fraction of correct processes that are E-ready for a message  $m$  is not significantly smaller than  $\hat{R}/R$ . As a result, a correct process that is not E-ready for  $m$  has a significant probability of becoming R-ready for  $m$ . If this happens, the probability of a correct process becoming R-ready for  $m$  further increases, and eventually every correct process is ready for  $m$ .

In this section, we show that the R-ready feedback mechanism is isomorphic to an epidemic process as we define it in Section 4.11. In summary, an epidemic process depends on one parameter (contagion threshold  $\hat{R}$ ) to mimic the spread of a disease in a population:

## Chapter 4. Contagion

---

- A population is represented on the nodes of a directed multigraph, allowing multi-edges and loops. Intuitively, an  $a \rightarrow b$  edge represents the relation  $a$  can infect  $b$ .
- Each member of the population (or node) can be in either of two states: healthy or infected. An infected node stays infected: there is no cure for the infection.
- A set of nodes is initially infected. The epidemic process evolves in steps. At every step, all the nodes that have at least  $\hat{R}$  infected predecessors become infected as well. The process is completed when either all nodes are infected, or no healthy node has at least  $\hat{R}$  infected predecessors.

We refer the reader to Section 4.11 for a more formal discussion of epidemic processes. In this section, we prove the critical result that the R-ready feedback mechanism in Contagion is isomorphic to an epidemic process.

**Definition 30** (Adversarial execution). A **adversarial execution** (or just **execution**) is the sequence of events produced by an execution of Contagion on  $N$  processes, a fraction  $f$  of which are under the control of the adversary described in Section 4.5.1. For the sake of brevity, we omit a more formal definition.

Let  $x, x'$  be executions. We say that  $x$  is equivalent to  $x'$  ( $x = x'$ ) if:

- The sequences of messages exchanged are identical in  $x$  and  $x'$ .
- The values produced by each correct, local source of randomness are identical in  $x$  and  $x'$ .

**Definition 31** (Ready sample matrix). A **ready sample matrix** is an element of the set

$$\mathcal{J} = (\Pi^R)^{\Pi_C}$$

**Definition 32** (Ready sample matrix of an execution). Let  $x$  be an execution, let  $j$  be a ready sample matrix.  $j$  is  $x$ 's **sample matrix** if, for every correct process  $\pi$ ,  $\pi$ 's ready sample in  $x$  is  $j_\pi$ .

**Definition 33** (Random ready sample matrix). A **random ready sample matrix** is a random variable representing the sample matrix of a random execution.

**Lemma 29.** *Random sample matrices are uniformly distributed. More formally, if  $j$  is a random sample matrix, then*

$$\mathcal{P}[j] = \left(\frac{1}{N}\right)^{RC}$$

*Proof.* As we discussed in Section 1.3, the adversary has no control over the local source of randomness of each correct process. Each correct process independently selects with uniform probability  $R$  elements for its ready sample. □



**Lemma 30.** *Let  $j$  be a ready sample matrix. Let  $x, x'$  be executions of Contagion such that:*

- *No Byzantine process issues any Ready message in  $x$  or  $x'$ .*
- *The ready sample matrix of both  $x$  and  $x'$  is  $j$ .*

*Let  $\rho_E, \rho'_E$  denote the set of correct processes that are E-ready for  $m$  in  $x, x'$  respectively. Let  $\rho, \rho'$  denote the set of correct processes that are ready for  $m$  in  $x, x'$  respectively.*

*We have*

$$(\rho_E = \rho'_E) \implies (\rho = \rho')$$

*Proof.* Let us assume  $\rho_E = \rho'_E$ . Let  $\pi$  be a correct process. As we established,  $\pi$  is ready for  $m$  if  $\pi$  is either E-ready or R-ready for  $m$ . Since  $\rho_E = \rho'_E$ , we immediately have that  $\pi$  is E-ready for  $m$  in  $x$  if and only if  $\pi$  is E-ready for  $m$  in  $x'$ .

By definition,  $\pi$  is R-ready for  $m$  in  $x$  ( $x'$ ) if it eventually receives at least  $\hat{R}$  Ready( $m$ ) messages from its ready sample in  $x$  ( $x'$ ). By hypothesis, no Byzantine process issues any Ready message in  $x$  ( $x'$ ). Therefore,  $\pi$  is eventually R-ready for  $m$  in  $x$  ( $x'$ ) if  $\pi$  receives at least  $\hat{R}$  Ready( $m$ ) messages from the correct processes in its ready sample in  $x$  ( $x'$ ).

As we discussed in Section 1.3, we assume that every message is eventually delivered in an unbounded but finite amount of time. Therefore,  $\pi$  is eventually R-ready for  $m$  in  $x$  ( $x'$ ) if at least  $\hat{R}$  correct processes in  $\pi$ 's sample eventually issue a Ready( $m$ ) message in  $x$  ( $x'$ ), i.e., if at least  $\hat{R}$  correct processes in  $\pi$ 's sample are eventually ready for  $m$  in  $x$  ( $x'$ ).

Since the above condition does not depend on the network scheduling, a correct process  $\pi$  is eventually ready for  $m$  in  $x$  if and only if  $\pi$  is also eventually ready for  $m$  in  $x'$ . Therefore,  $\rho = \rho'$ . □

**Lemma 31.** *Let  $x$  be an execution of Contagion where no Byzantine process ever issues any Ready message. Let  $j$  be  $x$ 's ready sample matrix. Let  $m$  be a message, let  $\rho_E$  denote the set of correct processes that are E-ready for  $m$  in  $x$ . Let  $\rho$  denote the set of correct processes that are eventually ready for  $m$  in  $x$ .*

*Let  $s_0 = ((v, e), w_0)$  be a contagion state (as defined in Definition 35), with*

$$\begin{aligned} v &= \Pi_C \\ (\pi, \pi') \in e &\iff \pi \in j_{\pi'} \\ w_0 &= \rho_E \end{aligned}$$

*Let  $s_\infty = ((v, e), w_\infty)$  be the contagion state resulting from the epidemic process with input  $s_0$ . We have*

$$\rho = w_\infty$$

## Chapter 4. Contagion

---

*Proof.* Following from Lemma 30,  $\rho$  does not depend on  $x$ 's network scheduling. Without loss of generality, we can therefore make a synchrony assumption for  $x$ , and assume that every message delay is unitary.

Let  $\rho_t$  denote the set of correct processes that are ready for  $m$  in  $x$  at time  $t$ . We have

$$\rho_0 = w_0 = \rho_E$$

In  $x$ , a correct process that is not ready for  $m$  at time  $t$  becomes ready for  $m$  at time  $t + 1$  if at least  $\hat{R}$  processes in its ready sample are ready for  $m$  at time  $t$ . Therefore

$$\pi \in \rho_{t+1} \iff (\pi \in \rho_t \vee |j_\pi \cap R_t| \geq \hat{R})$$

As we discuss in Section 4.11, at step  $t + 1$ , all the healthy nodes in an epidemic process that have at least  $\hat{R}$  predecessors infected at time  $t$  become infected. Therefore

$$\pi \in w_{t+1} \iff (\pi \in w_t \vee |j_\pi \cap w_t| \geq \hat{R})$$

Therefore, if  $\rho_t = w_t$ , then  $\rho_{t+1} = w_{t+1}$ , and, by induction, for all  $t$ ,  $\rho_t = w_t$ . In Section 4.11, we prove that an epidemic process identically converges in a finite number of steps. Consequently,  $\rho_\infty = w_\infty$ , which proves the lemma.  $\square$

**Lemma 32.** *Let  $m$  be a message. Let  $x$  be an execution of Contagion where every Byzantine process sends a Ready( $m$ ) message to every correct process from which it received a ReadySubscribe message. Let  $j$  be  $x$ 's ready sample matrix. Let  $\rho_E$  denote the set of correct processes that are E-ready for  $m$  in  $x$ . Let  $\rho$  denote the set of correct processes that are eventually ready for  $m$  in  $x$ .*

Let  $s_0 = ((v, e), w_0)$  be a contagion state (as defined in Definition 35), with

$$\begin{aligned} v &= \Pi \\ (\pi, \pi') \in e &\iff (\pi' \in \Pi_C) \wedge (\pi \in j_{\pi'}) \\ w_0 &= \rho_E \cup (\Pi \setminus \Pi_C) \end{aligned}$$

Let  $s_\infty = ((v, e), w_\infty)$  be the contagion state resulting from the epidemic process with input  $s_0$ . We have

$$\rho = w_\infty \setminus (\Pi \setminus \Pi_C)$$

*Proof.* It follows immediately from Lemma 31 and the observation that, in  $x$ , a Byzantine process sends the same Ready messages as a correct process that is E-ready for  $m$ .  $\square$

## 4.7 Threshold contagion (overview)

As we discussed in the previous section, in Section 4.11 we introduce epidemic processes, an abstract model of the feedback mechanism produced by Ready messages in an execution of Contagion. Given the multigraph on which it occurs, an epidemic process is deterministic. In Section 4.11, we also generalize epidemic processes to the probabilistic setting: we introduce and analyze Threshold Contagion, a game where a player infects in rounds arbitrary subsets of a population, causing a sequence of epidemic processes on a random, unknown multigraph.

Threshold Contagion depends on six parameters: node count  $N$ , sample size  $R$ , link probability  $l$ , round count  $K$ , infection batch  $S$ , and contagion threshold  $\hat{R}$ .

In summary, a game of Threshold Contagion is played as follows:

- A random multigraph with  $N$  nodes is generated. The number of predecessors of each node follows a  $\text{Bin}[R, l]$  distribution. Each predecessor of a node is independently picked with uniform probability from the set of nodes.

The topology of the network is not disclosed to the adversary.

- For  $K$  rounds:
  - The player infects an arbitrary set of  $S$  healthy nodes.
  - An epidemic process with contagion threshold  $\hat{R}$  is ran on the resulting contagion state.

We refer the reader to Section 4.11 for a more formal discussion of Threshold Contagion. There we introduce the random variable

$$\gamma(N, R, l, K, S, \hat{R})$$

representing the number of nodes that are infected at the end of a game of Threshold Contagion. We then prove that, by arbitrary choosing which nodes to infect, the adversary has no way to bias  $\gamma$ . Finally, we analitically compute the probability distribution underlying  $\gamma$ .

In this section, we prove the critical result that a game of Threshold Contagion can be used to model two classes of adversarial executions of Contagion.

**Lemma 33.** *Let  $m^*$  be a message. Let  $x$  be an adversarial execution of Contagion where:*

- *No Byzantine process issues any Ready message.*
- *For  $K$  rounds:*
  - *The adversary selects, if possible,  $S$  correct process that are not ready for  $m^*$ , and causes them to pcb.Deliver  $m^*$ .*

## Chapter 4. Contagion

---

– The adversary waits until every resulting Ready message is delivered.

Let  $\rho$  denote the number of correct processes in  $\sigma$  that, at the end of the adversarial execution, are ready for  $m$ . We have

$$\mathcal{P}[\bar{\rho}] = \mathcal{P}[\gamma(C, R, 1 - f, K, S, \hat{R}) = \bar{\rho}]$$

*Proof.* We start by defining a function  $\mathfrak{g} : \mathcal{J} \rightarrow \mathcal{G}$  (see Definition 39 for a definition of  $\mathcal{G}$ ) by

$$\mathfrak{g}(j)_{i,k} = \begin{cases} \zeta^{-1}(j_{\zeta(i),k}) & \text{iff } j_{\zeta(i),k} \in \Pi_C \\ \perp & \text{otherwise} \end{cases}$$

We start by noting that, for every  $\bar{g} \in \mathcal{G}$ ,

$$\mathcal{P}[\bar{g}] = \mathcal{P}[\mathfrak{g}^{-1}(g)]$$

Indeed, following from Lemmas 41 and 42:

$$\begin{aligned} \mathcal{P}[\bar{g}] &= \prod_{i,k} \mathcal{P}[\bar{g}_{i,k}] \\ \mathcal{P}[g_{i,k} = \perp] &= (1 - l) = f \\ \mathcal{P}[g_{i,k} = (\bar{g}_{i,k} \in 1..C)] &= \frac{1}{C} \end{aligned}$$

and following from Definition 33 and Lemma 29 we have

$$\begin{aligned} \mathcal{P}[\bar{j}] &= \prod_{\pi,k} \mathcal{P}[\bar{j}_{\pi,k}] \\ \mathcal{P}[j_{\pi,k} \in \Pi \setminus \Pi_C] &= f \\ \mathcal{P}[j_{\pi,k} = (\bar{\pi}' \in \Pi_C)] &= \frac{1}{C} \end{aligned}$$

We now build from  $x$  a game of Threshold Contagion  $y$ , played on  $\mathfrak{g}(j)$ . At the beginning of each round, if the adversary causes a correct process  $\pi$  to `pcb.Deliver`  $m^*$ , then  $\zeta(\pi)$  is infected.

We can prove that, if  $\pi$  is eventually ready for  $m^*$  in  $x$ , then  $\zeta(\pi)$  is eventually infected in  $y$ . Indeed, following from Lemma 31, if  $\pi$  is ready for  $m^*$  at the end of a round in  $x$ , then  $\zeta(\pi)$  is infected at the end the same round in  $y$ .

Therefore, the following hold true:

- The probability of  $\bar{j}$  is identical to the probability of  $\mathfrak{g}(\bar{j})$ .
- The number of correct processes that are eventually ready for  $m^*$  in  $x$  is identical to the

number of nodes that are eventually infected in  $y$ .

□

**Lemma 34.** *Let  $m$  be a message. Let  $x$  be an adversarial execution of Contagion where:*

- *No correct process  $pcb$ .Delivers  $m$ .*
- *Every Byzantine process sends a Ready( $m$ ) message to every correct process from which it received a ReadySubscribe message.*

*Let  $\rho$  denote the number of correct processes in  $\sigma$  that, at the end of the adversarial execution, are ready for  $m$ . We have*

$$\mathcal{P}[\bar{\rho}] = \mathcal{P}[\gamma(N, R, 0, 1, N - C, \hat{R}) = \bar{\rho} + (N - C)]$$

*Proof.* The proof is similar to the proof of Lemma 33, using Lemma 32 instead of Lemma 31.

□

## 4.8 Preliminary lemmas

In order to compute an upper bound for the probability of the consistency of Contagion being compromised, we will make use of some preliminary lemmas. The statements of those lemmas are independent from the context of Contagion. For the sake of readability, we therefore gather them in this section, and use them throughout the rest of this section.

**Lemma 35.** *Let  $N, K \in \mathbb{N}$  such that  $K \in 0..N$ . Let  $X$  be a random variable defined by*

$$\mathcal{P}[\bar{X}] = \text{Bin}[N, p](\bar{X})$$

*We have that*

$$\mathcal{P}[X \geq K]$$

*is an increasing function of  $p$ .*

*Proof.* We expand

$$\mathcal{P}[X \geq K] = \sum_{\bar{X}=K}^N \text{Bin}[N, p](\bar{X})$$

## Chapter 4. Contagion

---

and take the derivative

$$\begin{aligned}
 & \frac{\partial}{\partial p} \sum_{\bar{X}=K}^N \binom{N}{\bar{X}} p^{\bar{X}} (1-p)^{N-\bar{X}} \\
 &= \sum_{\bar{X}=K}^N \binom{N}{\bar{X}} \left( p^{\bar{X}} (1-p)^{N-\bar{X}} \right) \left( \frac{\bar{X}}{p} - \frac{N-\bar{X}}{1-p} \right) \\
 &= \underbrace{\frac{1}{p(1-p)}}_{\geq 1} \sum_{\bar{X}=K}^N \text{Bin}[N, p](\bar{X})(\bar{X} - pN) \\
 &\geq \sum_{\bar{X}=K}^N \text{Bin}[N, p](\bar{X})(\bar{X} - pN)
 \end{aligned}$$

We now prove that, for every  $K \in [0, N]$ ,

$$\sum_{\bar{X}=K}^N \text{Bin}[N, p](\bar{X})(\bar{X} - pN) \geq 0 \tag{4.2}$$

We start by noting that Equation (4.2) holds true for every  $K > pN$ . Indeed, if  $K > pN$ , then every term of the sum in Equation (4.2) is positive.

We prove that Equation (4.2) holds true for every  $K < pN$  by induction. For  $K = 0$  we have

$$\begin{aligned}
 & \sum_{\bar{X}=0}^N \text{Bin}[N, p](\bar{X})(\bar{X} - pN) \\
 &= \underbrace{\sum_{\bar{X}=0}^N \bar{X} \text{Bin}[N, p](\bar{X})}_{=pN} - pN \underbrace{\sum_{\bar{X}=0}^N \text{Bin}[N, p](\bar{X})}_{=1} \\
 &= 0
 \end{aligned}$$

Let us assume that Equation (4.2) holds true for some  $K < pN$ . We have

$$\begin{aligned}
 & \sum_{\bar{X}=K+1}^N \text{Bin}[N, p](\bar{X})(\bar{X} - pN) \\
 &= \underbrace{\sum_{\bar{X}=K}^N \text{Bin}[N, p](\bar{X})(\bar{X} - pN)}_{\geq 0 \text{ by IH}} - \text{Bin}[N, p](K)(K - pN) \\
 &\geq 0
 \end{aligned}$$

which proves that Equation (4.2) holds true for  $K + 1$  as well. By induction, Equation (4.2) holds true for every  $K < pN$ . This proves that the derivative is positive for all  $p \in [0, 1]$  which proves the lemma.  $\square$

## 4.9 Consistency

In this section, we compute a bound on the  $\epsilon$ -consistency of Contagion. As we discussed in Section 4.5.1, here we bound the probability of compromising the consistency of Contagion by assuming that, if the consistency of the pcb instance used in Contagion is compromised, the consistency of Contagion is compromised as well.

Let  $m^*$  denote the only message that any correct process can pcb.Deliver. We start by noting that, simply by having every Byzantine process behave like a correct process, an adversary can cause any correct process to deliver  $m^*$ : indeed, with  $f = 0$ , Contagion satisfies validity deterministically<sup>1</sup>.

As we discussed in Section 4.2, a correct process can issue a Ready message for an arbitrary number of messages. In other words, causing a correct process to become E-ready for  $m^*$  does not affect its behavior with respect to a message  $m \neq m^*$ .

Therefore, if an adversary can cause at least one correct process  $\pi$  to eventually receive at least  $\hat{R}$  Ready messages for a message  $m \neq m^*$ , it can also compromise the consistency of Contagion.

Indeed, as we discussed in Section 1.3, the adversary has arbitrary control over the network scheduling. Even if  $\pi$  would eventually receive enough Ready( $m^*$ ) messages to deliver  $m^*$ , the adversary can slow those messages down, and cause  $\pi$  to first receive enough Ready( $m$ ) messages to deliver  $m$ . Every other correct process will eventually deliver  $m^*$ , thus compromising the consistency of the system.

We formalize the above intuition in the following lemma.

**Lemma 36.** *Let  $m^*$  denote the only message that any correct process can pcb.Deliver. An optimal adversary causes every Byzantine process to send a Ready( $m$ ) message, for some  $m \neq m^*$ , to every correct process from which it received a ReadySubscribe message.*

*Proof.* Let  $B$  denote the number of Byzantine processes that eventually issue a Ready( $m$ ) message. Let  $\pi$  be a correct process, let  $Q$  denote the number of Ready( $m$ ) messages that  $\pi$  eventually collects. Since  $\pi$  picks each element of its delivery sample independently,  $Q$  is binomially distributed:

$$\mathcal{P}[\bar{Q}] = \text{Bin}\left[D, \frac{B}{N}\right](\bar{Q})$$

Following from Lemma 35,

$$\mathcal{P}[Q \geq \hat{D}] = \sum_{\bar{Q}=\hat{D}}^D \mathcal{P}[\bar{Q}]$$

<sup>1</sup>Here we are slightly abusing the result of Theorem 12, as it only guarantees that a correct sender will eventually deliver its message. The result, however, independently holds for any other correct process as well.

## Chapter 4. Contagion

---

is an increasing function of  $B$ , and maximized by  $B = (N - C)$ . Therefore, the probability of  $\pi$  eventually receiving enough  $\text{Ready}(m)$  messages to deliver  $m$  is maximized if every Byzantine process issues a  $\text{Ready}(m)$  message.

As we previously established, the adversary can cause every correct process to also receive at least  $\hat{D}$   $\text{Ready}(m^*)$  messages. Since the adversary has control over network scheduling, it can cause  $\pi$  to deliver  $m$ , and every other process to deliver  $m^*$ , thus compromising the consistency of the system.  $\square$

**Lemma 37.** *Let  $m^*$  denote the only message that any correct process can pcb.Deliver, let  $m \neq m^*$ . If, throughout an optimal adversarial execution, no correct process eventually collects enough  $\text{Ready}(m)$  messages to deliver  $m$ , then no correct process eventually collects enough  $\text{Ready}(m)$  messages to deliver any message  $m' \neq m$ .*

*Proof.* Following from Lemma 36, the optimal adversary causes every Byzantine process to issue a  $\text{Ready}(m)$  message. In Lemma 32, we use the fact that this strategy makes the Byzantine processes behave identically to correct processes that are E-ready for  $m$  to show that the set of correct processes that are eventually ready for  $m$  only depends on the ready sample matrix of the execution.

Since a correct process does not change its ready or delivery samples throughout an execution, the set of processes that will eventually be ready for  $m'$  is at most the same as the set of processes that will eventually be ready for  $m$ . In turn, this means that if no correct process eventually delivers  $m$ , no correct process eventually delivers  $m'$  either.  $\square$

We can now use Lemma 36 to compute a bound on the  $\epsilon$ -consistency of Contagion.

We introduce the random variable  $\gamma^+$  by

$$\mathcal{P}[\tilde{\gamma}^+] = \mathcal{P}[\gamma(N, R, 0, 1, N - C, \hat{R}) = \tilde{\gamma}^+]$$

Following from Lemmas 34 and 36,  $\gamma^+$  represents the number of processes (Byzantine or correct) that eventually issue a  $\text{Ready}$  message for a message  $m \neq m^*$ , when an optimal adversary is trying to compromise the consistency of the system.

We can finally compute a bound for the  $\epsilon$ -consistency of Contagion. We define

$$\begin{aligned} \mu &= \sum_{\tilde{\gamma}^+ = N - C}^N \left(1 - (1 - \tilde{\mu}(\tilde{\gamma}^+))^C\right) \mathcal{P}[\tilde{\gamma}^+] \\ \tilde{\mu}(\tilde{\gamma}^+) &= \sum_{\tilde{D} = \hat{D}}^D \text{Bin}\left[D, \frac{\tilde{\gamma}^+}{N}\right](\tilde{D}) \end{aligned}$$



**Theorem 13.** *Contagion satisfies  $\epsilon_c$ -consistency, with*

$$\epsilon_c \leq \epsilon_c^{pcb} + (1 - \epsilon_c^{pcb})\mu$$

*if the underlying abstraction of pcb satisfies  $\epsilon_c^{pcb}$ -consistency.*

*Proof.* We start by noting that  $\tilde{\mu}(\bar{\gamma}^+)$  represents the probability that a specific correct process will eventually collect enough Ready( $m$ ) messages to deliver  $m$ , given the number  $\bar{\gamma}^+$  of processes that eventually issue a Ready( $m$ ) message.

Indeed, since every correct process picks its delivery sample independently, each of the  $D$  elements of a correct process' delivery sample has a probability  $\bar{\gamma}^+ / N$  of issuing a Ready( $m$ ) message.

We then note that  $\mu$  represents the probability of any correct process eventually collecting enough Ready( $m$ ) messages to deliver  $m$ .  $\mu$  is obtained by applying the law of total probability to  $\tilde{\mu}(\bar{\gamma}^+)$ .

Finally,  $\epsilon_c$  is obtained by the assumption that, if the consistency of the underlying pcb instance is compromised, the totality of Contagion is compromised as well.  $\square$

## 4.10 Totality

In this section, we compute a bound on the  $\epsilon$ -totality of Contagion. As we discussed in Section 4.5.1, here we bound the probability of compromising the totality of Contagion by assuming that, if the consistency of the pcb instance used in Contagion is compromised, the consistency of Contagion is compromised as well.

### 4.10.1 Minimal operations

Let  $m^*$  be the only message that any correct process can pcb.Deliver. As we discussed in Section 4.5.1, throughout an execution of Contagion, an adversary performs a sequence of minimal operations on the system, i.e., it either causes a correct process to pcb.Deliver  $m^*$ , or it causes a Byzantine process to send an arbitrary Ready( $m$ ) message to a correct process.

We further relax the bound by assuming that, if the adversary can cause any message  $m \neq m^*$  to be delivered by at least one correct process, the totality of Contagion is compromised as well.

Under the assumption that no correct process can eventually collect enough Ready( $m$ ) messages to deliver any message  $m$  different from  $m^*$ , causing a Byzantine process to send a Ready( $m$ ) message has no effect on the totality of the system.

## Chapter 4. Contagion

---

This reduces the set of adversarial operations that have a non-null effect on the totality of the system to:

- Causing an arbitrary correct process to `pcb.Deliver  $m^*$` .
- Causing a Byzantine adversary to send a `Ready( $m^*$ )` message to a correct process.

We now prove a lemma to further reduce the set of minimal operations of an optimal adversary.

**Lemma 38.** *Let  $m^*$  be the only message that any correct process can potentially `pcb.Deliver`. Let  $\pi$  be a correct process, let  $\xi$  be a Byzantine process in  $\pi$ 's ready sample. An optimal adversary never causes  $\xi$  to send a `Ready( $m^*$ )` message to  $\pi$ .*

*Proof.* As a result of receiving a `Ready( $m^*$ )` message from  $\xi$ ,  $\pi$  can either:

- Have collected less than  $\hat{R}$  `Ready( $m^*$ )` messages from its ready sample. The operation has no effect.
- Have collected exactly  $\hat{R}$  `Ready( $m^*$ )` messages from its ready sample. Then  $\pi$  becomes ready for  $m^*$ . However, the same outcome could have been achieved deterministically by causing  $\pi$  to `pcb.Deliver  $m^*$` .

Since every outcome of  $\xi$ 's `Ready( $m^*$ )` message to  $\pi$  can be deterministically emulated by causing  $\pi$  to `pcb.Deliver` (or not `pcb.Deliver`)  $m^*$ , the operation is useless to an optimal adversary.  $\square$

### 4.10.2 Delivery probability

Let  $\gamma^-$  denote the random variable counting the number of correct processes that are eventually ready for  $m^*$ . In this section, we study the probability of totality being compromised, given the value of  $\gamma^-$ .

By definition, totality is compromised if at least one correct process delivers  $m^*$  and one correct process does not deliver  $m^*$ .

Let  $\pi$  be a correct process. We introduce the following events:

- $A_\pi$ : process  $\pi$  delivers  $m^*$ .
- $A$ : all correct processes deliver  $m^*$ .
- $\tilde{A}$ : no correct process delivers  $m^*$ .
- $T$ : the totality of the system is compromised.

Given  $\tilde{\gamma}^-$ , the probability of  $A_\pi$  is bound by

$$\alpha_\pi^-(\tilde{\gamma}^-) \leq \mathcal{P}[A_\pi | \tilde{\gamma}^-] \leq \alpha_\pi^+(\tilde{\gamma}^-)$$

with

$$\begin{aligned} \alpha_\pi^-(\tilde{\gamma}^-) &= \sum_{\tilde{D}=\hat{D}}^D \text{Bin}\left[D, \frac{\tilde{\gamma}^-}{N}\right](\tilde{D}) \\ \alpha_\pi^+(\tilde{\gamma}^-) &= \sum_{\tilde{D}=\hat{D}}^D \text{Bin}\left[D, \frac{\tilde{\gamma}^- + (N-C)}{N}\right](\tilde{D}) \end{aligned}$$

The lower bound is attained when none of the Byzantine processes issue a Ready( $m^*$ ) message, and the upper bound is attained when all Byzantine processes issue a Ready( $m^*$ ) message.

Noting that each correct process independently picks its delivery sample, we can compute, given  $\tilde{\gamma}^-$ , a lower bound for the probability of  $A$ :

$$\mathcal{P}[A | \tilde{\gamma}^-] \geq (\alpha^-(\tilde{\gamma}^-))^C$$

and a lower bound for the probability of  $\tilde{A}$ :

$$\mathcal{P}[\tilde{A} | \tilde{\gamma}^-] \geq (1 - \alpha^+(\tilde{\gamma}^-))^C$$

The above allow us to compute, given  $\tilde{\gamma}^-$ , an upper bound for the probability of  $T$ :

$$\begin{aligned} \mathcal{P}[T | \tilde{\gamma}^-] &= \mathcal{P}[A, \tilde{A} | \tilde{\gamma}^-] \\ &\leq 1 - \mathcal{P}[A | \tilde{\gamma}^-] - \mathcal{P}[\tilde{A} | \tilde{\gamma}^-] \\ &\leq \alpha(\tilde{\gamma}^-) \end{aligned}$$

with

$$\alpha(\tilde{\gamma}^-) = 1 - (\alpha^-(\tilde{\gamma}^-))^C - (1 - \alpha^+(\tilde{\gamma}^-))^C$$

### 4.10.3 C-step Threshold Contagion

Due to Lemma 38, the minimal set of operations for an optimal adversary reduces to

- Causing an arbitrary correct process to pcb.Deliver  $m^*$ .
- Causing a Byzantine process  $\xi$  in the delivery sample of a correct process  $\pi$  to send a Ready( $m^*$ ) message to  $\pi$ .

It is immediate to see that the latter operation has no effect over which correct processes eventually become Ready for  $m^*$ . In the previous section, we computed an upper bound

## Chapter 4. Contagion

---

on the probability of compromising the totality of the system, given the number of correct processes that are eventually ready for  $m^*$ .

In this section, we prove a final constraint on the optimal adversarial strategy, and finally compute a bound on the  $\epsilon$ -totality of Contagion.

**Lemma 39.** *Let  $m^*$  denote the only message that any correct process can `pcb.Deliver`. An optimal adversary executes in  $C$  rounds. At every round, the adversary causes one correct process to `pcb.Deliver`  $m^*$ , then waits until all the resulting `Ready` messages are delivered.*

*Proof.* Due to Lemma 31, the outcome of the execution is not affected by network scheduling: causing one correct process at a time to `pcb.Deliver`  $m^*$  has the same effect, e.g., as causing any set of correct processes to simultaneously `pcb.Deliver`  $m^*$ .  $\square$

Following from Lemma 39, we can intuitively see an adversarial execution whose goal is to compromise the totality of Contagion as a game similar to *blackjack*. The game unfolds in  $C$  rounds. At every round, the adversary causes one more correct process to `pcb.Deliver`  $m^*$ . With high probability, this will have two possible negative outcomes for the player:

- Nothing happens: no correct process is able to deliver  $m^*$ , even if the Byzantine processes in its delivery sample issue a `Ready( $m^*$ )` message. The only possible move is to play again.
- The execution is *busted*: a feedback loop is generated that eventually causes, with high probability, every correct process to deliver  $m^*$ , even if no Byzantine process issues any `Ready( $m^*$ )` message. The adversary fails in compromising the totality of the system.

If the adversary is lucky enough, however, one of the rounds will result in a configuration where no feedback loop occurred, but at least one correct process can deliver  $m^*$ . In that case, the adversary causes that process to deliver  $m^*$ , and stops: totality is compromised.

Following from Lemma 33, the probability distribution underlying the number of correct processes that are ready for  $m^*$  at the end of the  $n$ -th step is

$$\mathcal{P}[\tilde{\gamma}_n^-] = \mathcal{P}[\gamma(C, R, 1 - f, C, 1, \hat{R})]$$

We can finally compute a bound on the  $\epsilon$ -totality of Contagion.

**Theorem 14.** *Contagion satisfies  $\epsilon_t$ -totality, with*

$$\begin{aligned} \epsilon_t &\leq \epsilon_c^{pcb} + \mu + \epsilon_b \\ \epsilon_b &= \sum_{n=0}^C \sum_{\tilde{\gamma}_n^-=0}^C \mathcal{P}[\tilde{\gamma}_n^-] \alpha(\tilde{\gamma}_n^-) \end{aligned}$$

if the underlying abstraction of *pcb* satisfies  $\epsilon_c^{pcb}$ -consistency.

*Proof.* Let  $T_n$  denote the event of totality being compromised at the end of round  $T_n$ .

Under the assumption that the consistency of *pcb* is satisfied, and no message other than  $m^*$  can be delivered by any correct process, the probability of  $T_n$  with  $n > 1$  is

$$\mathcal{P}[T_n] = \sum_{\tilde{\gamma}_n=0}^C \mathcal{P}[T_n | \tilde{\gamma}_n] \mathcal{P}[\tilde{\gamma}_n | T_{n-1}]$$

Indeed, the adversary will proceed to round  $n$  only if round  $n - 1$  was unsuccessful in compromising the totality of the system. We can use the law of total probability to get

$$\begin{aligned} \mathcal{P}[T_n] &\leq \sum_{\tilde{\gamma}_n=0}^C \mathcal{P}[T_n | \tilde{\gamma}_n] (\mathcal{P}[\tilde{\gamma}_n | T_{n-1}] + \mathcal{P}[\tilde{\gamma}_n | T_{n-1}]) \\ &= \sum_{\tilde{\gamma}_n=0}^C \mathcal{P}[T_n | \tilde{\gamma}_n] \mathcal{P}[\tilde{\gamma}_n] \end{aligned}$$

We can use Boole's inequality to get

$$\mathcal{P}[T] \leq \sum_{n=0}^C \mathcal{P}[T_n]$$

and since

$$\mathcal{P}[T_n | \tilde{\gamma}_n] \leq \alpha(\gamma_n^-)$$

we have that  $\epsilon_b$  bounds the probability of compromising totality, if the consistency of *pcb* is satisfied, and no message other than  $m^*$  can be delivered by any correct process.

The value provided for  $\epsilon_t$  follows from applying again Boole's inequality to include  $\epsilon_c^{pcb}$  and  $\mu$  (which, in Section 4.9, we proved to bound the probability of any correct process delivering a message other than  $m^*$ ).  $\square$

## 4.11 Threshold Contagion

In this section we discuss *epidemic processes*, mimicking the spread of a disease in a population, and the Threshold Contagion game, which gives a player the possibility to actively infect parts of a population.

As we discuss in Section 4.2, in Contagion, when a correct process receives enough Ready messages from its *ready sample* for the same message  $m$ , it issues itself a Ready message for  $m$ . This produces a feedback mechanism that, in Section 4.7, we showed to be isomorphic to an epidemic process as we define it below.

## Chapter 4. Contagion

---

Threshold Contagion is a game where a player iteratively applies the epidemic process to chosen inputs. We use Threshold Contagion for modeling and analyzing our Contagion algorithm.

### 4.11.1 Epidemic processes

An epidemic process models the spreading of a disease in a population.

#### Preliminary definitions

**Definition 34** (Directed multigraph). A **directed multigraph** is a pair  $g = (v, e)$ , where  $v$  is a set and  $e : v^2 \rightarrow \mathbb{N}$  is a multiset whose elements are pairs of elements of  $v$ . We call the elements of  $v$  the **vertices** (or **nodes**) of  $g$ . We call the elements of  $e$  the **edges** of  $g$ .

Following from Definition 34, a directed multigraph allows self-loops (let  $a \in v$ ,  $(a, a)$  can be an element of  $e$ ) and multiple edges (let  $a, b \in v$ , the multiplicity of  $(a, b)$  in  $e$  can be greater than one).

#### Contagion state

The spreading of a disease is represented by a **contagion state**.

**Definition 35** (Contagion state). A **contagion state** is a pair  $s = (g, w)$ , where  $g = (v, e)$  is a multigraph and  $w \in \mathbb{P}(v)$ . We call the elements of  $w$  the **infected nodes** of  $s$ .

Let  $s = ((v, e), w)$  be a contagion state. In an epidemic process:

- Each node in  $v$  corresponds to one individual member of the population.
- A node is always in one of two possible states: **healthy** or **infected**. We do not consider any cure—once a node becomes infected, it stays infected forever. The set  $w$  represents the nodes that are infected.
- Edges model interactions between the members of the population. The multiset of edges  $e$  represents the "can infect" relation. Note that this relation is not symmetric. A directed edge  $(a \rightarrow b)$  between nodes  $a$  and  $b$  means that  $a$  can infect  $b$ , but not that  $b$  can infect  $a$ .

**Definition 36** (Predecessors). Let  $g = (v, e)$  be a multigraph, let  $a \in v$ . Then the **predecessors** of  $a$  in  $g$  form the multiset  $p[a] : v \rightarrow \mathbb{N}$  defined by

$$p[a](x \in v) = e(x, a)$$

Following from Definition 36, the predecessors of a node  $a$  in a multigraph  $g$  form the multiset of nodes that have an edge to  $a$ . If a node has multiple edges to  $a$ , then it has a multiplicity greater than one in  $p[a]$ .

### Contagion rule

In an epidemic process, the infection of healthy nodes follows a single rule.

- **Contagion rule:** A healthy node becomes infected if the number of its infected predecessors reaches a critical threshold.

The input to an epidemic process is a contagion state  $s$ . The epidemic process repeatedly applies the contagion rule to  $s$  until either all nodes are infected or no healthy node has enough infected predecessors to become infected itself. The epidemic process outputs the resulting contagion state.

### 4.11.2 Threshold Contagion

Threshold Contagion is a game played on the nodes of a random directed multigraph  $g$ . Threshold Contagion consists of one or more **rounds**. Each round inputs a contagion state  $s$  and outputs a contagion state  $s'$ . The input to the first round is the contagion state  $(g, \emptyset)$ , i.e., the contagion state with no infected nodes whose multigraph is  $g$ . The input to any other round is the output of the previous round.

A round with input  $s$  is played as follows:

- The **player** infects a fixed-size subset of the healthy nodes of the contagion state  $s$ . This results in a contagion state  $s'$ .
- The contagion state  $s'$  is provided as input to an epidemic process. The output  $s''$  of the epidemic process is returned.

### 4.11.3 Rules

In this section, we formally define the rules of Threshold Contagion and introduce its **parameters**.

Threshold Contagion is played on the nodes of a random, directed multigraph  $g = (\nu, e)$ . The in-degree of each node  $n$  in  $\nu$  is independently binomially distributed; each predecessor of  $n$  is uniformly picked with replacement from  $\nu$ .

## Chapter 4. Contagion

---

### Parameters

A game of Threshold Contagion depends on the following numerical parameters:

- **Node count** ( $N \in \mathbb{N}$ ): Represents the number of nodes in the multigraph ( $N = |v|$ ).
- **Sample size** ( $R \in \mathbb{N}$ ): Represents the maximum in-degree of a node in the multigraph.
- **Link probability** ( $l \in [0, 1]$ ): Represents the probability of a predecessor link being successfully established. The in-degree of a node follows the distribution  $\text{Bin}[R, l]$ .
- **Round count** ( $K \in \mathbb{N}_{>0}$ ): Represents the number of rounds in the game.
- **Infection batch** ( $(S < N) \in \mathbb{N}_{>0}$ ): Represents the number of healthy nodes the player infects at the beginning of each round.
- **Contagion threshold** ( $(\hat{R} \leq R) \in \mathbb{N}$ ): Represents the number of infected predecessors that will cause a healthy node to become infected (see Contagion rule).

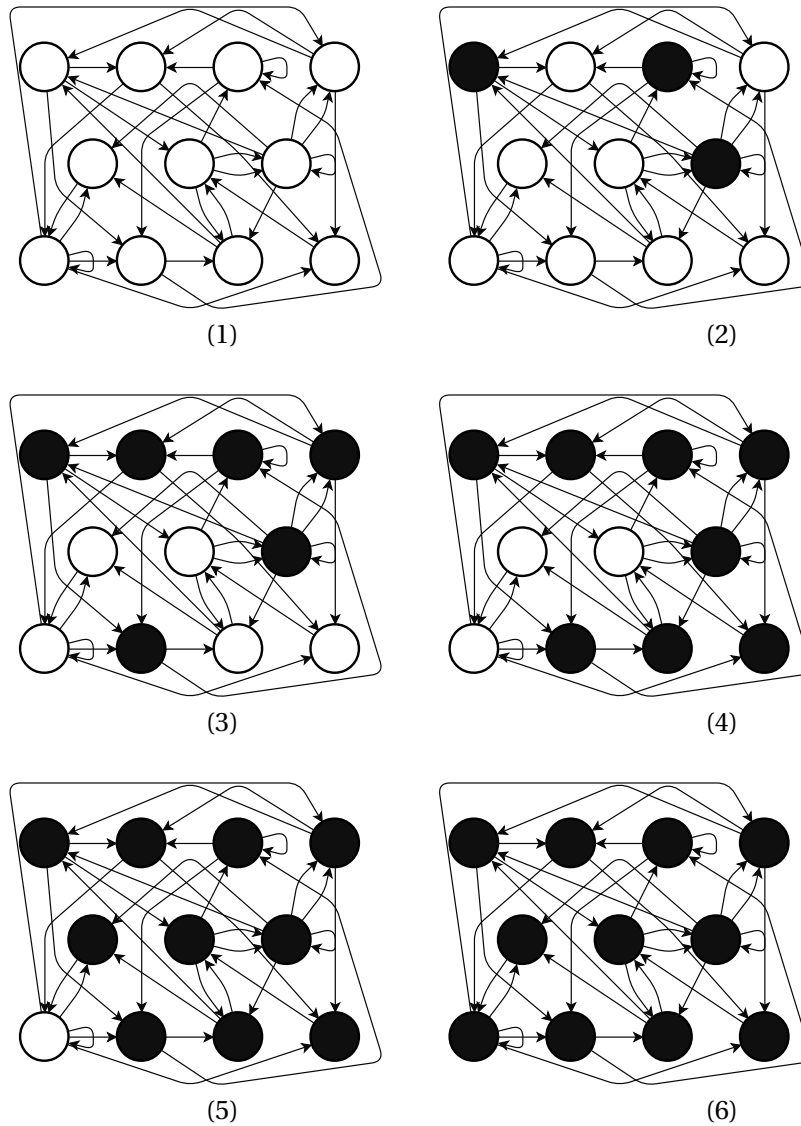
### Game

A game of Threshold Contagion is played as follows:

- A random, directed multigraph  $g = (v, e)$  with  $N$  nodes is built. For every node  $n$  in  $v$ :
  - $R$  times:
    - \* A Bernoulli random variable  $\bar{B} \leftarrow \text{Bern}[l]$  is sampled.
    - \* If  $\bar{B} = 1$ , then a random node  $m$  is selected with uniform probability from  $v$ , and the edge  $m \rightarrow n$  is added to  $e$  (i.e.,  $m$  is added to the predecessors of  $n$ ).
- Let  $s = (g, w = \emptyset)$  be a contagion state. For  $K$  **rounds**:
  - If at least  $S$  nodes in  $v$  are healthy (i.e., they are not in the set of infected nodes  $w$ ), the player selects  $S$  distinct nodes and infects them. The player **cannot** inform this choice with knowledge of the topology of  $g$ .
  - An epidemic process is run on  $s$ : until either every node in  $v$  is infected (i.e.,  $v = w$ ), or no healthy node in  $v$  has at least  $\hat{R}$  infected predecessors, the following **contagion step** is iterated:
    - \* Every node in  $v$  with at least  $\hat{R}$  infected predecessors is infected, i.e., it is added to  $w$ .

Figure 4.1 shows an example game of Threshold Contagion with small parameters.





**Figure 4.1: An example game of Threshold Contagion.** Here  $N = 11$ ,  $l = 1$ ,  $R = 3$ ,  $\hat{R} = 2$ ,  $K = 1$  and  $S = 3$ . Notice how nodes can be linked to themselves, form loops, or be linked more than once. An initial set of  $S$  nodes (1) is infected by the player (2). The game then unfolds in contagion steps (3 to 6): whenever a node has at least  $\hat{R}$  infected predecessors, it becomes infected. This example shows how easily a game of Threshold Contagion can converge to a fully-infected configuration.

#### 4.11.4 Random variables

We introduce the following random variables, which we discuss in more formal detail in the next sections:

- **Infection size**  $N_i^r$ : represents the number of infected nodes at round  $r$  and step  $i$ .
- **Frontier size**  $U_i^r$ : represents the number of nodes that are infected at round  $r$  and step  $i$ , but are not infected at round  $r$  and step  $i - 1$ .
- **Infection status**  $W_i^r[j]$ : represents whether or not the  $j$ -th node is infected at round  $r$  and step  $i$ . We use  $W_i^r[j]$  to signify that the node is infected, and  $\overline{W_i^r[j]}$  to signify that the node is not infected.
- **Infected predecessors count**  $V_i^r[j]$ : represents the number of predecessors of the  $j$ -th node that are infected at round  $r$  and step  $i$ .

**Remark:** for the sake of readability, the round number and/or the node index (for  $W$  and  $V$ ) will be omitted whenever it can be unequivocally inferred from the context.

#### 4.11.5 Goal

The goal of this section is to compute the probability distribution underlying the number of infected nodes at the end of a game of Threshold Contagion.

**Lemma 40.** *For any  $r$ , the random variables  $N_i^r$ ,  $U_i^r$ ,  $W_i^r[j]$ , and  $V_i^r$  converge in a finite number of steps.*

*Proof.* We note the following:

- $N_i^r$  is a non-decreasing function of  $i$ , and  $N_i^r \leq N$ .
- $U_{i>0}^r = N_i^r - N_{i-1}^r$ .
- For any  $j$ ,  $W_i^r[j] \implies W_{i+1}^r[j]$ .
- $V_i^r$  is a non-decreasing function of  $i$  and  $V_i^r \leq R$ .

From the above follows that all random variables converge for  $i \rightarrow \infty$ .

The codomains of  $N$ ,  $U$ ,  $W$  and  $V$  are all finite. Therefore, all random variables converge in a finite number of steps. □

**Corollary 2.** *All rounds terminate in a finite number of contagion steps.*

**Notation 10** (End of round). We use  $N_\infty^r, U_\infty^r, W_\infty^r[j], V_\infty^r$  to denote the values of  $N, U, W, V$  at the end of round  $r$ .

The goal of this section is to compute the probability distribution underlying the random variable

$$\gamma(N, R, l, K, S, \hat{R}) = N_\infty^K \quad (4.3)$$

i.e., the probability of a game of Threshold Contagion resulting in  $N_\infty^K$  nodes ultimately being infected. Lemma 40 proves that  $\Gamma$  is a well defined variable (i.e., the limit exists) and, since  $K$  is finite, can be computed in a finite total number of steps.

### 4.11.6 Sample space

In this section, we define a sample space for Threshold Contagion, i.e., the set of all possible outcomes of a Threshold Contagion game. As we described in Section 4.11.2, the outcome of a game of Threshold Contagion is completely determined by two factors:

1. The topology of the random multigraph  $g$  on which Threshold Contagion is played. The probability distribution underlying  $g$  is known, and we compute it in this section.
2. The player's infection strategy, i.e., the nodes the player chooses to infect at the beginning of each round. The probability distribution underlying the player's choices is **unknown** and arbitrary. In this section, we only formalize their sample space.

Thus, an element of the sample space is a pair consisting of a multigraph (1.) and an infection strategy (2.).

#### Multigraph

As discussed in Section 4.11.3, a game of Threshold Contagion is played on the nodes of a multigraph  $g = (\nu, e)$  allowing multi-edges and loops. Every node in  $\nu$  has at most  $R$  predecessors. Therefore,  $g$  can be represented by a **predecessor matrix** as we define it below.

We start by explicitly labeling the elements of  $\nu$ .

**Notation 11** (Vertices). Let  $g = (\nu, e)$  be a multigraph, with  $|\nu| = N$ . Without loss of generality, we label the elements of  $\nu$  using natural numbers:

$$\nu = 1..N$$

Since every node in  $g$  has at most  $R$  predecessors, for every  $j \in \nu$  we can represent the elements of  $\mathfrak{p}(j)$  as the components of a *predecessor vector*.

**Definition 37** (Predecessor vector). A **predecessor vector** is an element of the set

$$\mathcal{R} = (\{\perp\} \cup \nu)^R$$

In a multigraph  $g = (\nu, e)$ , whose in-degree is bound by  $R$ , we use a predecessor vector to represent the predecessors of a node. Let  $r \in \mathcal{R}$  be the predecessor vector of a node  $j \in \nu$ . If  $r_k = \perp$ , we say that the  $k$ -th predecessor of  $j$  is **missing**.

As discussed in Section 4.11.3, the predecessors of each node in  $\nu$  are generated by independently sampling  $R$  times a value  $\bar{B}$  from a Bernoulli variable; whenever  $\bar{B} = 1$ , an additional predecessor is uniformly picked with replacement from the elements of  $\nu$ . We call a vector of predecessors selected this way a *random predecessor vector*, as formally defined in Definition 38.

**Definition 38** (Random predecessor vector). A **random predecessor vector** is a predecessor vector generated by the procedure described in Section 4.11.3.

Let  $r$  be a random predecessor vector. For every  $k \in \{1, \dots, R\}$ ,  $\bar{B} \leftarrow \text{Bern}[l]$  is independently sampled; if  $\bar{B} = 0$ ,  $r_k$  is set to  $\perp$ , otherwise  $r_k$  is set to an element of  $\nu$ , picked independently with uniform probability.

**Lemma 41.** *Let  $r$  be a random predecessor vector. Then*

$$\begin{aligned} \mathcal{P}[\bar{r}] &= \prod_{k=1}^R \mathcal{P}[\bar{r}_k] \\ \mathcal{P}[r_k = \perp] &= (1 - l) \\ \mathcal{P}[r_k = (\bar{r}_k \in \nu)] &= \frac{l}{N} \end{aligned}$$

*Proof.* Following from Definition 38, each component of  $r$  is independently sampled. Each component has a probability  $(1 - l)$  of being missing. Each non-missing component of  $r$  has an equal probability of being equal to any element of  $\nu$ .  $\square$

As we discussed in Section 4.11.3, the multigraph  $g = (\nu, e)$  is constructed by independently generating the predecessors for each node in  $\nu$ . Therefore, the topology of  $g$  is completely determined by  $N$  predecessor vectors, that can be organized in a *predecessor matrix*.

**Definition 39** (Predecessor matrix). A **predecessor matrix** is an element of the set

$$\mathcal{G} = \mathcal{R}^N$$

**Notation 12** (Predecessor matrix). Since a predecessor matrix uniquely identifies a multigraph, we interchangeably use  $g$  to denote a predecessor matrix and its corresponding multigraph. Let  $g$  be a predecessor matrix defining a multigraph  $(\nu, e)$ , then  $g_j$  is the predecessor vector of node  $j \in \nu$ .

**Definition 40** (Random predecessor matrix). A **random predecessor matrix** is a predecessor matrix representing the outcome of the multigraph generation process described in Section 4.11.3. More formally, a random predecessor matrix consists of  $N$  independent random predecessor vectors.

**Lemma 42.** *Let  $g$  be a random predecessor matrix. Then*

$$\mathcal{P}[\bar{g}] = \prod_{j=1}^N \mathcal{P}[\bar{g}_j] \quad (4.4)$$

*Proof.* It follows immediately from Definition 40. □

### Sub-threshold predecessor set

As discussed in Section 4.11.3, an epidemic process consists of a sequence of contagion steps. Let  $s = ((v, e), w)$  be a contagion state. In a contagion step, a healthy node  $j \in v$  ( $j \in w$ ) becomes infected if at least  $\hat{R}$  of its predecessors are infected, i.e., if

$$|p(j) \cap w| \geq \hat{R}$$

Given the set  $w$ , the set of predecessor vectors that do not satisfy the condition above is uniquely defined. We define *sub-threshold predecessor sets* to capture this notion.

**Definition 41** (Sub-threshold predecessor set). Let  $g$  be a predecessor matrix defining a multigraph  $(v, e)$ . Let  $X \subseteq v$ . The **sub-threshold predecessor set** of  $X$  is the set

$$\tilde{\mathcal{R}}^X = \{r \in \mathcal{R} \mid |\{k \in 1..R \mid r_k \in X\}| < \hat{R}\}$$

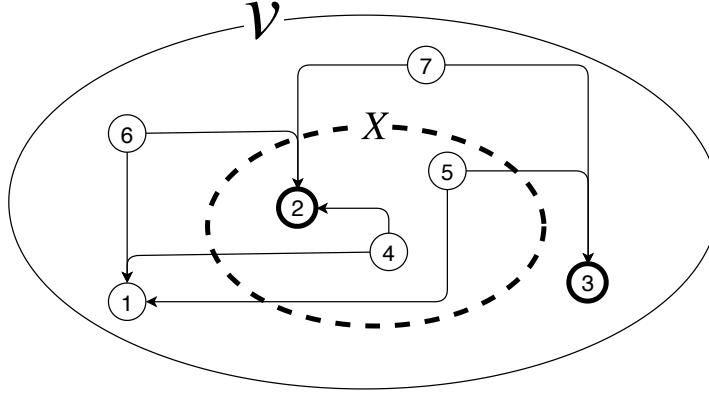
$\tilde{\mathcal{R}}^X$  contains all the predecessor vectors in  $\mathcal{R}$  that have less than  $\hat{R}$  components in  $X$ .

Figure 4.2 shows an example multigraph where the predecessors of three nodes are displayed, two of which are in the sub-threshold predecessor set of a given set  $X$ .

### Player's strategy

As discussed in Section 4.11.3, at the beginning of each round of Threshold Contagion the player selects, if possible,  $S$  distinct healthy nodes and infects them. These are the only  $K$  choices the player makes throughout Threshold Contagion. Moreover, the player has no knowledge of the topology of the multigraph  $g$  on which Threshold Contagion is played.

The player's choices can be expressed in an *infection strategy*, as we formally define it in this section. Together with the topology of the multigraph on which the game is played, an infection strategy uniquely determines the outcome of an instance of Threshold Contagion.



**Figure 4.2: An example multigraph  $g = (v, e)$  with 7 nodes.** A subset  $X \subseteq v$  is highlighted. Numbered dots represent the elements of  $v$ , and the edges to nodes 1, 2 and 3 are displayed. With  $R = 3$  and  $\hat{R} = 2$ , we have  $g_1 \notin \tilde{\mathcal{R}}^X$ ,  $g_2 \in \tilde{\mathcal{R}}^X$ , and  $g_3 \in \tilde{\mathcal{R}}^X$ . Note how the predecessor vector of node 2 is in the sub-threshold predecessor set of  $X$  even if node 2 is in  $X$ . Note how node 3 has one missing predecessor (i.e., one of the elements in  $g_3$  is  $\perp$ ). The nodes whose predecessor vectors are in  $\tilde{\mathcal{R}}^X$  are highlighted.

Let  $g = (v, e)$  be the multigraph on which Threshold Contagion is played. At the beginning of round  $r$ , the player knows the value of  $W_i^{r'}[j]$  for every  $r' < r$ , every  $i \in \mathbb{N}$  and every  $j \in v$ , which we encode in an *infection history*. The player chooses a set of  $S$  of the nodes that are healthy at the beginning of round  $r$ . We model this choice with an *infection function*. We call *infection strategy* the sequence of choices the player makes throughout the game.

**Definition 42** (Infection history). An **infection history** for round  $r > 0$  is an element of the set

$$\mathcal{H}_r = \left( (\{\perp, \top\}^N)^\infty \right)^r$$

An infection history is a table with three indices. The first represents the round, the second represents the step, the third represents the node. Let  $h \in \mathcal{H}_r$ , then  $h_i^{r'}[j] = \top$  signifies that node  $j$  is infected at round  $r'$  and step  $i'$ .

**Notation 13** (Round and step order). Let  $r, r'$  be round numbers, let  $i, i'$  be step numbers. We say that  $(r, i) < (r', i')$  if  $(r, i)$  *temporally precedes*  $(r', i')$ . More formally

$$(r, i) < (r', i') \iff (r' > r) \vee (r' = r \wedge i' > i)$$

**Definition 43** (Valid infection history). A **valid infection history** for round  $r > 0$  is an element of the set

$$\mathcal{H}_r^* = \left\{ h \in \mathcal{H}_r \mid h_i^{r'}[j] = \top \implies h_{i''}^{r''}[j] = \top \forall (r'', i'') > (r', i') \right\}$$

A valid infection history is an infection history where a node is never healed. If a node is

infected at round  $r'$  and step  $i'$ , then it also infected at any subsequent round  $r''$  and step  $i''$ .

**Definition 44** (Incomplete infection history). An **incomplete infection history** for round  $r > 0$  is an element of the set

$$\mathcal{H}_r^+ = \{h \in \mathcal{H}_r^* \mid |\{j \in 1..N \mid h_\infty^{r-1}[j] = \perp\}| \geq S\}$$

An incomplete infection history is a valid infection history with at least  $S$  healthy nodes at the end of round  $r - 1$ .

**Definition 45** (Infection function). An **infection function** for round  $r$  is an element of the set

$$\mathcal{F}_r = \{f : \mathcal{H}_r^+ \rightarrow \mathbb{P}^S(\{1..N\}) \mid \forall x \in f(h), h_\infty^{r-1}[x] = \perp\}$$

An infection function is a function that inputs an incomplete infection history and outputs a set of  $S$  nodes, all of which are healthy at the end of round  $r - 1$ .

**Definition 46** (Infection strategy). An **infection strategy** is an element of the set

$$\mathcal{F} = \mathbb{P}^{1..N}(S) \times \prod_{r=1}^{R-1} \mathcal{F}_r$$

The first element of an infection strategy is a set of  $S$  nodes to infect at the beginning of round 0. Let  $r > 0$ , the  $r$ -th element of an infection strategy is an infection function for round  $r$ .

An infection strategy encodes all the choices a player makes during a game of Threshold Contagion:

- At the beginning of round 0, the player has no information available. All nodes are healthy, and its choice reduces to selecting  $S$  of them to infect.
- At the beginning of round  $r \geq 1$ , the information available to the player is the propagation of the infection throughout all previous rounds. Such information is input to the  $r$ -th infection function, which returns a set of  $S$  healthy nodes to infect.

### Sample space

In Section 4.11.6, we noticed how the outcome of a game of Threshold Contagion is completely determined once both the topology of the multigraph and the strategy of the player are known.

In Section 4.11.6, we showed how a multigraph can be expressed with a predecessor matrix, defined the space of predecessor matrices and derived the probability distribution underlying random predecessor matrices.

## Chapter 4. Contagion

---

In Section 4.11.6, we showed how the choices that a player makes at the beginning of each round in response to the infection history can be encoded in infection strategies. We then defined the space of infection strategies. Unlike random multigraphs, infection strategies are under the control of the player. Therefore, a probability distribution over the space of infection strategies is not available.

As we discussed in Section 4.11.6, an element of the sample space is a pair of a multigraph and an infection strategy.

**Definition 47** (Sample space). The **sample space** for Threshold Contagion is the set  $\Omega = \mathcal{G} \times \mathcal{F}$ .

**Lemma 43.** Let  $\omega = (g, f)$  be a random element of  $\Omega$ . Then  $\mathcal{P}[\bar{g}, \bar{f}] = \mathcal{P}[\bar{g}] \mathcal{P}[\bar{f}]$ , i.e.,  $g$  and  $f$  are independent.

*Proof.* It immediately follows from the fact that the player has no knowledge of the topology of the multigraph  $g$ .  $\square$

### 4.11.7 Random variables as sample functions

In Section 4.11.4 we intuitively defined a set of random variables to capture useful properties of a game of Threshold Contagion. In the next sections, we use those random variables to compute the probability distribution underlying the number of infected nodes at the end of a game.

In Section 4.11.6 we formally defined the sample space of a game of Threshold Contagion. We started by showing that an instance of the game is completely determined once the topology of the multigraph and the strategy of the player are known. We also computed the probability of any specific multigraph topology occurring.

In this section, we rigorously re-define the random variables we defined in Section 4.11.4 by expressing them as functions on the sample space as defined in Section 4.11.6.

#### Infection history

As discussed in Section 4.11.6, an infection function for round  $r$  inputs an incomplete infection history for round  $r$  and outputs a set of  $S$  nodes to infect out of those that are healthy at the end of round  $r - 1$ .

We introduce two useful functions to manipulate infection histories.

**Definition 48** (Sample history function). The **sample history function** for round  $r$  is the function  $\mathfrak{h}_r : \Omega \rightarrow \mathcal{H}_r^*$  defined by

$$(\mathfrak{h}_r(\omega))_i^{r'}[j] = W_i^{r'}[j](\omega)$$



The sample history function for round  $r$  inputs a sample  $\omega$  and outputs the valid infection history for round  $r$  produced by  $\omega$ .

Note how the definition of sample history function relies on the definition of the infection status  $W$ . We introduced  $W$  in Section 4.11.4, and we formally define it in the next section.

**Definition 49** (Sample completion function). The **sample completion function** for round  $r$  is the function  $c_r : \Omega \rightarrow \{\top, \perp\}$  defined by

$$c_r(\omega) = \begin{cases} \perp & \text{iff } \mathfrak{h}_r(\omega) \in \mathcal{H}_r^+ \\ \top & \text{otherwise} \end{cases}$$

The sample completion function for round  $r$  inputs a sample and outputs  $\top$  if the infection history of the sample is complete at round  $r$ , and  $\perp$  otherwise.

### Infection status

As stated in Section 4.11.3, the infection status is defined as follows:

- At the beginning of the game, all the nodes are healthy.
- During the first step of each round, the player selects a set of  $S$  healthy nodes and infects them.
- During every subsequent step, every healthy node that has at least  $\hat{R}$  infected predecessors is infected.
- The infection state at the end of a round is carried without change to the beginning of the next round.

In order to formalize the above in the definition of infection status, we preliminarily define *infection sets*.

**Definition 50** (Infection set). The **infection set** at round  $r$  and step  $i$  is the random variable  $\hat{W}_i^r : \Omega \rightarrow \mathbb{P}(1..N)$  defined by

$$\hat{W}_i^r(\omega) = \{j \in 1..N \mid W_i^r[j](\omega) = \top\}$$

The infection set  $\hat{W}_i^r(\omega)$  represents the set of nodes that are infected in  $\omega$  at round  $r$  and step  $i$ .

Like the sample history function, the definition of infection set relies on the definition of infection status  $W$ , which we can now define by cases.

## Chapter 4. Contagion

---

**Definition 51** (Infection status). Let  $\omega = (g, f) \in \Omega$ . The **infection status** for round  $r$ , step  $i$  and node  $j$  is the random variable  $W_i^r[j] : \Omega \rightarrow \{\top, \perp\}$  defined by

$$W_0^0[j](\omega) = \perp \quad (4.5)$$

$$W_0^{r>0}[j](\omega) = W_\infty^{r-1}[j](\omega) \quad (4.6)$$

$$W_1^0[j](\omega) = \begin{cases} \top & \text{iff } j \in f_0 \\ W_0^0[j](\omega) & \text{otherwise} \end{cases} \quad (4.7)$$

$$W_1^{r>0}[j](\omega) = \begin{cases} \top & \text{iff } \mathfrak{c}_r(\omega) = \perp \wedge j \in f_r(\mathfrak{h}_r(\omega)) \\ W_0^r[j](\omega) & \text{otherwise} \end{cases} \quad (4.8)$$

$$W_{i>1}^r[j](\omega) = \begin{cases} W_{i-1}^r[j](\omega) & \text{iff } g_j \in \tilde{\mathcal{R}}^{\hat{W}_{i-1}^r(\omega)} \\ \top & \text{otherwise} \end{cases} \quad (4.9)$$

The above equations encode the following properties:

- At the beginning of the game (Equation (4.5)), all nodes are healthy.
- The infection status at the beginning of round  $r > 0$  (Equation (4.6)) is equal to the infection status at the end of round  $r - 1$ .
- During step 1 of round 0 (Equation (4.7)), all the nodes in  $f_0$  are infected. Intuitively, the player selects  $S$  nodes and infects them. Note how this choice is not informed by any history (following from Definition 46,  $f_0$  is a set and not a function).
- During step 1 of round  $r > 0$  (Equation (4.8)), if  $\omega$  is not complete (i.e., there are at least  $S$  healthy nodes at the beginning of round  $r$ ), all the nodes in  $f_r(\mathfrak{h}_r(\omega))$  are infected. Intuitively, the player selects  $S$  healthy nodes and infects them. This choice is informed by the infection history for round  $r$  (see Definition 48).
- During step  $i > 0$  of any round  $r$  (Equation (4.9)), all the nodes whose predecessor vector is not in the sub-threshold predecessor set (see Definition 41) of the infection set at step  $i - 1$  are infected. In other words, the contagion rule (see Section 4.11.1) is applied, and all the nodes that have at least  $\hat{R}$  infected predecessors are infected.

Following from Definition 51, we prove that nodes are never healed in a game of Threshold Contagion.

**Lemma 44.** *Let  $j \in 1..N$ ,  $r, r' \in 1..K$ ,  $i, i' \in \mathbb{N}$ , let  $\omega \in \Omega$ . If  $(r', i') \geq (r, i)$ , then*

$$W_i^r[j](\omega) = \top \implies W_{i'}^{r'}[j](\omega)$$

*Proof.* Let  $r'' \in 1..K$ ,  $i'' \in \mathbb{N}$ . Following from Equations (4.5) to (4.9), we have

$$W_{i''+1}^{r''}[j](\omega) \neq W_{i''}^{r''}[j](\omega) \implies W_{i''+1}^{r''}[j](\omega) = \top \quad (4.10)$$

The lemma is proved by induction on Equations (4.6) and (4.10).  $\square$

**Corollary 3.** *The infection set  $\hat{W}_i^r(\omega)$  is non-decreasing in  $(r, i)$ .*

### Infection size, frontier size and infected predecessors count

In Section 4.11.7, we defined the infection status  $W_i^r[j]$  as a function on the sample space (see Definition 51). We also defined the infection set  $\hat{W}_i^r$  as the set of nodes for which  $W_i^r = \top$  (see Definition 50).

As stated in Section 4.11.4, the infection size  $N_i^r$  represents the number of infected nodes at round  $r$  and step  $i$ , and the frontier size  $U_{i>0}^r$  represents the number of nodes that are infected at round  $r$  and step  $i$ , but not at step  $i - 1$ . We can formalize the above in the following definitions.

**Definition 52** (Infection size). The **infection size** for round  $r$  and step  $i$  is the random variable  $N_i^r : \Omega \rightarrow 0..N$  defined by

$$N_i^r(\omega) = |\hat{W}_i^r(\omega)|$$

The infection size counts the infected nodes at step  $(r, i)$ .

**Definition 53** (Frontier size). The **frontier size** for round  $r$  and step  $i$  is the random variable  $U_{i>0}^r : \Omega \rightarrow 0..N$  defined by

$$U_{i>0}^r(\omega) = N_i^r(\omega) - N_{i-1}^r(\omega)$$

The infection size counts the nodes that are infected at step  $(r, i)$ , but not at step  $(r, i - 1)$ .

As stated in Section 4.11.4, the infected predecessors count of node  $j$  for round  $r$  and step  $i$  represents the number of predecessors of node  $j$  that are infected at round  $r$  and step  $i$ . We can formalize this definition in the following.

**Definition 54** (Infected predecessors count). Let  $\omega = (g, f) \in \Omega$ . The **infected predecessors count** of node  $j$  for round  $r$  and step  $i$  is the random variable  $V_i^r[j] : \Omega \rightarrow 0..R$  defined by

$$V_i^r[j](\omega) = |\{k \mid (g_{j,k}) \in \hat{W}_i^r(\omega)\}|$$

The infected predecessors count counts the number of predecessors of node  $j$  that are infected at step  $(r, i)$ .

**Lemma 45.** *Let  $\omega = (g, f) \in \Omega$ , let  $j \in 1..N$ ,  $r \in 1..K$ ,  $i \in \mathbb{N}$ . Then*

$$g \in \tilde{\mathcal{R}}^{\hat{W}_i^r[j]} \iff V_i^r[j](\omega) \leq \hat{R}$$

*Proof.* It follows immediately from Definitions 41 and 54.  $\square$

### 4.11.8 Contagion step

In Section 4.11.7, we expressed the random variables we introduced in Section 4.11.4 as functions over the elements of the sample space we defined in Section 4.11.6. As we established in Section 4.11.5, the goal of this section is to compute the distribution underlying  $N_\infty^K$  (see Equation (4.3)).

Here we focus on the contagion steps of a round of Threshold Contagion. As per Equation (4.9), at every step  $(r, i)$  such that  $i > 1$ , all the healthy nodes that have at least  $\hat{R}$  infected predecessors become infected.

In this section, we show that a contagion step defines a Markov chain with states  $(\bar{N}_i^r, \bar{U}_i^r)$ . More formally, we show that a transition matrix  $\mathcal{M}$  exists such that, for every  $(\bar{N}, \bar{U})$ ,  $(\bar{N}', \bar{U}')$  and for every  $r \in 1..K$ ,  $i \geq 1$ ,

$$\mathcal{M}_{\bar{N}, \bar{U}}^{\bar{N}', \bar{U}'} = \mathcal{P}[N_{i+1}^r = \bar{N}', U_{i+1}^r = \bar{U}' \mid N_i^r = \bar{N}, U_i^r = \bar{U}] \quad (4.11)$$

Intuitively, this means that, once the infection size and the frontier size at step  $(r, i)$  are determined, no other knowledge is needed to compute the probability distribution underlying the frontier size at step  $(r, i + 1)$ . This means, in particular, that the player's infection strategy does not affect the end result of the game. This result is somewhat unsurprising: since the player has no knowledge of the multigraph on which Threshold Contagion is played, the player has no way to meaningfully distinguish two nodes by the effect that their infection will have on the system. Since the number of infected nodes per round is determined, every choice of the player can be shown to be effectively equivalent to the infection of  $S$  random healthy nodes.

#### Roadmap

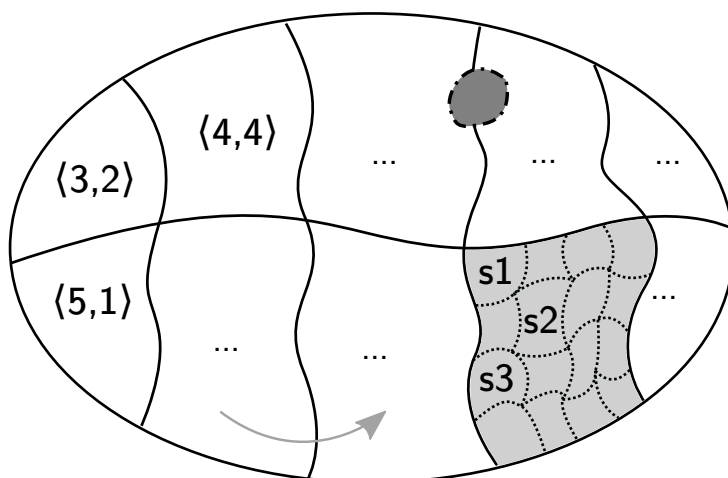
**Notation 14** (Markov states). We use  $\langle \bar{N}_i^r, \bar{U}_i^r \rangle$  to denote the subset of the sample space  $\Omega$  that satisfies  $N_i^r(\omega \in \Omega) = \bar{N}_i^r, U_i^r(\omega \in \Omega) = \bar{U}_i^r$ .

Equivalently,

$$\langle \bar{N}_i^r, \bar{U}_i^r \rangle = (N_i^r)^{-1}(\bar{N}_i^r) \cap (U_i^r)^{-1}(\bar{U}_i^r)$$

In order to show that a infection step defines a Markov chain with states  $(\bar{N}_i^r, \bar{U}_i^r)$ , we:

- Define a set of *partition functions*  $\mathcal{S}_i^r : \Omega \rightarrow \mathbb{P}(\Omega)$  that map elements of  $\Omega$  into well-structured subsets of  $\Omega$ . Intuitively,  $\mathcal{S}_i^r$  maps a sample  $\omega = (g, f)$  to a set of samples that are *similar to it* (by a notion of similarity that we define later).
- Let  $\omega' \in \mathcal{S}_i^r(\omega)$ . We show that  $\omega$  and  $\omega'$  result in the same infection history up to step  $(r, i)$ .



**Figure 4.3: An illustration of sample space and the steps needed to show that a contagion step defines a Markov chain.** The grey arrow represents a transition from a state to another. One of the states is further partitioned by  $\mathcal{S}_i^r$ . The dark grey area represents a case that we prove won't happen.

- We show that  $\mathcal{S}_i^r$  can be used to define an equivalence relation on the sample space  $\Omega$ .
- Let  $\omega$  be equivalent to  $\omega'$  through  $\mathcal{S}_i^r$ . We show that, since  $N_i^r(\omega) = N_i^r(\omega')$  and  $U_i^r(\omega) = U_i^r(\omega')$ , then  $\mathcal{S}_i^r$  can be used to quotient  $\langle \tilde{N}_i^r, \tilde{U}_i^r \rangle$ .
- Let  $r \in 1..K$ ,  $i > 1$ . We use  $\mathcal{S}_i^r$  to partition  $\langle \tilde{N}_i^r, \tilde{U}_i^r \rangle$  in  $s_1, \dots, s_q$ . We show that the probability of  $\omega$  being in  $\langle \tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r \rangle$  given that  $\omega$  is in  $s_h$  is analytically computable and independent of  $h$ .
- We use the independence across partitions to compute the terms of  $\mathcal{M}_{\tilde{N}, \tilde{U}}^{\tilde{N}', \tilde{U}'}$ .

### Partition functions

We start by defining a set of *partition functions*  $\mathcal{S}_i^r : \Omega \rightarrow \mathbb{P}(\Omega)$  that map elements of  $\Omega$  into subsets of  $\Omega$ . Intuitively, a partition function maps a sample to a set of samples that are *similar* to it.

Let  $\omega = (g, f) \in \Omega$ , let  $\omega' = (g', f') \in \mathcal{S}_i^r(\omega)$ . We define  $\mathcal{S}_i^r$  such that the following hold:

- $f' = f$ , i.e., the player's strategy is left unchanged by  $\mathcal{S}_i^r$ .
- Let  $j \in 1..N$  be a node. If  $j$  is infected in  $\omega$  at step  $(r, i)$ , then  $g'_j = g_j$ . In other words, the predecessors of a node that is infected at step  $(r, i)$  in  $\omega$  are left unchanged by  $\mathcal{S}_i^r$ .
- Let  $j \in 1..N$  be a node. If  $j$  is not infected in  $\omega$  at step  $(r, i)$ , then  $g'_j$  is an element of the sub-threshold predecessor set of  $\hat{W}_{i-1}^r(\omega)$ . In other words, the predecessors of a node

## Chapter 4. Contagion

---

that is not infected at step  $(r, i)$  in  $\omega$  can be changed by  $\mathcal{S}_i^r$ , as long as no more than  $\hat{R}$  of them are infected in  $\omega'$  at step  $(r, i-1)$ . Intuitively, we allow the predecessors of  $j$  to change in a way that does not make it infected in  $\omega$  at step  $(r, i)$ .

We formalize the above in the following definition.

**Definition 55** (Partition function). Let  $r \in 1..K$ ,  $i \geq 1$ , let  $\omega = (g, f) \in \Omega$ . The **partition function** for round  $r$  and step  $i$  is the function  $\mathcal{S}_i^r : \Omega \rightarrow \mathbb{P}(\Omega)$  defined by

$$\mathcal{S}_i^r(\omega) = \left( \prod_{j=1}^N \mathcal{S}_i^r[j](\omega) \right) \times \{f\} \quad (4.12)$$

$$\mathcal{S}_i^r[j](\omega) = \begin{cases} \{g_j\} & \text{iff } W_i^r[j](\omega) = \top \\ \hat{R}\hat{W}_{i-1}^r(\omega) & \text{otherwise} \end{cases} \quad (4.13)$$

### Infection history

In Section 4.11.8 we defined a set of partitions functions that map a sample  $\omega \in \Omega$  to a set of samples that are *similar* to  $\omega$ .

Let  $\omega = (g, f) \in \Omega$ . We designed  $\mathcal{S}_i^r$  to leave unchanged the player's strategy and the predecessors of every node that is infected in  $\omega$  at step  $(r, i)$ . The predecessors of the nodes that are not infected in  $\omega$  at step  $(r, i)$  can change, as long as less than  $\hat{R}$  of them are among the nodes that are infected in  $\omega$  at step  $(r, i-1)$ .

Intuitively,  $\mathcal{S}_i^r$  is designed to alter the topology of  $g$  in a way that does not affect its infection history: since the predecessors of the nodes that are not infected in  $\omega$  at step  $(r, i)$  are not changed, they will still be infected in  $\omega'$ . Similarly, if a node is not infected in  $\omega$  at step  $(r, i)$ , its predecessors are not changed in a way that makes it infected in  $\omega'$  at step  $(r, i)$ .

In this section, we formally prove this intuitive result.

**Lemma 46.** Let  $j \in 1..N$ , let  $\omega, \omega' \in \Omega$ . If  $\omega' \in \mathcal{S}_i^r(\omega)$ , then for every  $(r', i') \leq (r, i)$

$$W_{i'}^{r'}[j](\omega') = W_{i'}^{r'}[j](\omega)$$

*Proof.* Let  $\omega = (g, f)$  and  $\omega' = (g', f')$ . We prove the lemma by induction. We start by noting that, following from Equation (4.5),

$$W_0^0[j](\omega') = \perp = W_0^0[j](\omega)$$

Now, assume that  $(r', i') < (r, i)$  and, for all  $j \in 1..N$ ,  $\hat{W}_{i'}^{r'}[j](\omega') = \hat{W}_{i'}^{r'}[j](\omega)$ .

If  $r' = 0$  and  $i' = 0$ , then from Equation (4.7) it follows that, if  $j \in (f_0 = f'_0)$ ,

$$W_1^0[j](\omega') = \top = W_1^0[j](\omega)$$

and, otherwise,

$$W_1^0[j](\omega') = W_0^0[j](\omega') = W_0^0[j](\omega) = W_1^0[j](\omega)$$

If  $r' > 0$  and  $i' = 0$ , then  $h_r(\omega') = h_r(\omega)$ . Following from Equation (4.8), if  $(c_r(\omega) = c_r(\omega')) = \perp$  and  $j \in (f_r(h_r(\omega)) = f'_r(h_r(\omega')))$ , then

$$W_1^{r'}[j](\omega') = \top = W_1^{r'}[j](\omega)$$

and otherwise

$$W_1^{r'}[j](\omega') = W_0^{r'}[j](\omega') = W_0^{r'}[j](\omega) = W_1^{r'}[j](\omega)$$

We now consider the case  $i' \geq 1$ . We start by noting that, since  $\hat{W}_{i'}^{r'}(\omega') = \hat{W}_{i'}^{r'}(\omega)$ , then  $\tilde{\mathcal{R}}\hat{W}_{i'}^{r'}(\omega') = \tilde{\mathcal{R}}\hat{W}_{i'}^{r'}(\omega)$ .

If  $W_i^r[j](\omega) = \top$ , then  $g_j = g'_j$ . Following from Equation (4.9), if  $(g'_j = g_j) \in (\tilde{\mathcal{R}}\hat{W}_{i'}^{r'}(\omega') = \tilde{\mathcal{R}}\hat{W}_{i'}^{r'}(\omega))$ , then

$$W_{i'+1}^{r'}[j](\omega') = W_{i'}^{r'}[j](\omega') = W_{i'}^{r'}[j](\omega) = W_{i'+1}^{r'}[j](\omega)$$

and otherwise

$$W_{i'+1}^{r'}[j](\omega') = \top = W_{i'+1}^{r'}[j](\omega)$$

If  $W_i^r[j](\omega) = \perp$ , then  $g'_j \in \tilde{\mathcal{R}}\hat{W}_{i-1}^r(\omega)$ . Noting that  $(r', i') \leq (r, i - 1)$ , from Lemma 44 it follows that  $\hat{W}_{i'}^{r'}(\omega) \subseteq \hat{W}_{i-1}^r(\omega)$ , and consequently, from Equation (4.9), we have

$$g'_j \in (\tilde{\mathcal{R}}\hat{W}_{i-1}^r(\omega) \subseteq \tilde{\mathcal{R}}\hat{W}_{i'}^{r'}(\omega) = \tilde{\mathcal{R}}\hat{W}_{i'}^{r'}(\omega'))$$

Moreover, since  $W_i^r[j](\omega) = \perp$ , from Lemma 44 it follows

$$W_{i'+1}^{r'}(\omega) = W_{i'}^{r'}(\omega) = W_{i'}^{r'}(\omega') = \perp$$

and therefore

$$W_{i'+1}^{r'}(\omega') = W_{i'}^{r'}(\omega') = W_{i'+1}^{r'}(\omega)$$

Finally, if  $i = \infty$ , then following from Equation (4.6) we have

$$W_0^{r'+1}[j](\omega') = W_\infty^{r'}[j](\omega') = W_\infty^{r'}[j](\omega) = W_0^{r'+1}[j](\omega)$$

□

**Corollary 4.** *Let  $\omega, \omega' \in \Omega$ . If  $\omega' \in \mathcal{S}_i^r(\omega)$ , then*

$$\begin{aligned} N_i^r(\omega') &= N_i^r(\omega) \\ U_i^r(\omega') &= U_i^r(\omega) \end{aligned}$$

### Equivalence relation

In Section 4.11.8, we introduced a set of functions  $\mathcal{S}_i^r : \Omega \rightarrow \mathbb{P}(\Omega)$  that map a sample into a set of *similar* samples. In Section 4.11.8, we proved that, if  $\omega \in \Omega$  and  $\omega' \in \mathcal{S}_i^r(\omega)$ , then  $\omega$  and  $\omega'$  produce the same infection history (i.e., the same values for  $\hat{W}_i^r$ ) up to round  $r$  and step  $i$ .

In this section, we show that  $\mathcal{S}_i^r$  can be used to define an equivalence relation on  $\Omega$ .

**Lemma 47.** *Let  $\omega, \omega' \in \Omega$ . If  $\omega' \in \mathcal{S}_i^r(\omega)$ , then  $\mathcal{S}_i^r(\omega') = \mathcal{S}_i^r(\omega)$ .*

*Proof.* Let  $\omega = (g, f)$ ,  $\omega' = (g', f')$ . Following from Lemma 46, for every  $j$  we have

$$\begin{aligned} W_i^r[j](\omega') &= W_i^r[j](\omega) \\ W_{i-1}^r[j](\omega') &= W_{i-1}^r[j](\omega) \end{aligned}$$

Following from Definition 55, if  $W_i^r[j](\omega) = \top$ , then  $g'_j = g_j$ . Consequently,

$$\mathcal{S}_i^r[j](\omega') = \{g'_j\} = \{g_j\} = \mathcal{S}_i^r[j](\omega)$$

If  $W_i^r[j](\omega) = \perp$ , then

$$\mathcal{S}_i^r[j](\omega') = \tilde{\mathcal{R}}^{\hat{W}_{i-1}^r[j](\omega')} = \tilde{\mathcal{R}}^{\hat{W}_{i-1}^r[j](\omega)} = \mathcal{S}_i^r[j](\omega)$$

Therefore,

$$\mathcal{S}_i^r(\omega') = \prod_{j=1}^N \mathcal{S}_i^r[j](\omega') = \prod_{j=1}^N \mathcal{S}_i^r[j](\omega) = \mathcal{S}_i^r(\omega)$$

□

**Definition 56** (Partition relation). Let  $\omega, \omega' \in \Omega$ . If  $\omega' \in \mathcal{S}_i^r(\omega)$ , then  $\omega'$  has a **partition relation** with  $\omega$  at round  $r$  and step  $i$ :

$$\omega' \stackrel{(r,i)}{\sim} \omega$$

**Lemma 48.**  $\stackrel{(r,i)}{\sim}$  is an equivalence relation.



*Proof.* Let  $j \in 1..N$ , let  $\omega \in \Omega$ . Following from Equation (4.9), if  $W_i^r[j](\omega) = \perp$ , then  $g_j \in \tilde{\mathcal{R}}^{\hat{W}_{i-1}^r(\omega)}$ . Consequently, following from Definition 55, if  $W_i^r[j](\omega) = \top$ , then

$$g_j \in \{g_j\} = \mathcal{S}_i^r[j](\omega)$$

and if  $W_i^r[j](\omega) = \perp$ , then

$$g_j \in \tilde{\mathcal{R}}^{\hat{W}_{i-1}^r[j](\omega)} = \mathcal{S}_i^r[j](\omega)$$

Therefore,  $\omega \in \mathcal{S}_i^r(\omega)$ , and

$$\omega \stackrel{(r,i)}{\sim} \omega$$

therefore  $\stackrel{(r,i)}{\sim}$  is reflexive.

Let  $\omega' \in \mathcal{S}_i^r(\omega)$ . By Lemma 47,  $\mathcal{S}_i^r(\omega') = \mathcal{S}_i^r(\omega)$ . Consequently

$$\omega \in (\mathcal{S}_i^r(\omega) = \mathcal{S}_i^r(\omega'))$$

and

$$\omega' \stackrel{(r,i)}{\sim} \omega \implies \omega \stackrel{(r,i)}{\sim} \omega'$$

therefore  $\stackrel{(r,i)}{\sim}$  is symmetric.

Let  $\omega'' \in \mathcal{S}_i^r(\omega')$ . Again by Lemma 47,

$$\omega'' \in (\mathcal{S}_i^r(\omega') = \mathcal{S}_i^r(\omega))$$

and

$$\omega' \stackrel{(r,i)}{\sim} \omega, \omega'' \stackrel{(r,i)}{\sim} \omega' \implies \omega'' \stackrel{(r,i)}{\sim} \omega$$

therefore,  $\stackrel{(r,i)}{\sim}$  is transitive. □

### Transition probabilities

In Section 4.11.8 we showed that the partition function we introduced in Section 4.11.8 can be used to induce an equivalence relation on the sample space  $\Omega$ .

In this section, we use this result to show that a contagion step defines a Markov chain with states  $(\tilde{N}_i^r, \tilde{U}_i^r)$ , and compute the values of its associated transition matrix  $\mathcal{M}$ .

More formally, let  $r \in 1..K$ ,  $i \geq 1$ . In this section, we compute

$$\mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r \mid \tilde{N}_i^r, \tilde{U}_i^r]$$

and we show that its value is independent of the player's strategy.

As we established in Lemma 48,  $\stackrel{(r,i)}{\sim}$  is an equivalence relation on  $\Omega$ . Moreover, let  $\omega \in \Omega$ , by

## Chapter 4. Contagion

---

Corollary 4 we have  $\mathcal{S}_i^r(\omega) \subseteq \langle \bar{N}_i^r, \bar{U}_i^r \rangle$ .

We can therefore use  $(\simeq^{r,i})$  to partition  $\langle \bar{N}_i^r, \bar{U}_i^r \rangle$ :

$$\{s_1, \dots, s_q\} = \frac{\langle \bar{N}_i^r, \bar{U}_i^r \rangle}{(\simeq^{r,i})}$$

By the law of total probability,

$$\begin{aligned} \mathcal{P}[\bar{N}_{i+1}^r, \bar{U}_{i+1}^r | \bar{N}_i^r, \bar{U}_i^r] &= \mathcal{P}[\bar{N}_{i+1}^r, \bar{U}_{i+1}^r | \langle \bar{N}_i^r, \bar{U}_i^r \rangle] \\ &= \sum_{l=1}^q \mathcal{P}[\bar{N}_{i+1}^r, \bar{U}_{i+1}^r | s_l] \mathcal{P}[s_l | \langle \bar{N}_i^r, \bar{U}_i^r \rangle] \end{aligned}$$

Note how  $\mathcal{P}[s_l | \langle \bar{N}_i^r, \bar{U}_i^r \rangle]$  is unknown, as it depends on the probability distribution underlying the player's strategy. For a given  $h$ , we instead focus on computing  $\mathcal{P}[\bar{N}_{i+1}^r, \bar{U}_{i+1}^r | s_h]$ .

**Roadmap.** In order to compute  $\mathcal{P}[\bar{N}_{i+1}^r, \bar{U}_{i+1}^r | s_h]$ , we compute the probability for a node that is not infected in  $s_h$  at step  $(r, i)$  to become infected at time  $(r, i + 1)$ . Let  $j$  be a node that is not infected in  $s_h$  at step  $(r, i)$ . We compute the probability of it becoming infected at step  $(r, i + 1)$  by first computing the probability distribution underlying  $V_{i-1}^r[j]$ . Given  $\bar{V}_{i-1}^r[j]$ , we then compute the probability distribution underlying  $V_i^r[j]$ , and threshold it with  $\hat{R}$  to compute the probability of  $j$  becoming infected at step  $i + 1$ .

**Notation 15** (Kronecker delta). We use  $\delta$  to denote the **Kronecker delta**. Let  $i, j \in \mathbb{N}$ , then

$$\delta_{i,j} = I(i = j)$$

Let  $\bar{\omega} = (\bar{g}, \bar{f}) \in s_h$  be an example of  $s_h$ . Let  $W, \mathcal{W}$  denote the set of nodes that are infected and not infected in  $\bar{\omega}$  at step  $(r, i)$ , respectively:

$$\begin{aligned} W = \{w_1, \dots, w_n\} &= \hat{W}_r^i(\bar{\omega}) \\ \mathcal{W} = \{\mathcal{w}_1, \dots, \mathcal{w}_m\} &= 1..N \setminus \hat{W}_r^i(\bar{\omega}) \end{aligned}$$

with  $n = N_r^i(\bar{\omega})$  and  $m = N - n$ . Let  $\omega = (g, f)$ , following from Lemma 45 we have

$$(\omega \in s_h) \iff (g_{w_1} = \bar{g}_{w_1}, \dots, g_{w_n} = \bar{g}_{w_n}, V_{i-1}^r[\mathcal{w}_1] \leq \hat{R}, \dots, V_{i-1}^r[\mathcal{w}_m] \leq \hat{R})$$

Let  $j \in \mathcal{W}$ , i.e.,  $W_i^r[j] = \perp$ . Using the independence of the distribution of each predecessor vector in  $s_h$  (see Equation (4.4) and Definition 55), we can compute the probability distribution

underlying  $V_{i-1}^r[j]$  in  $s_h$ :

$$\begin{aligned}
 & \mathcal{P}[\bar{V}_{i-1}^r[j] | \cancel{W_i^r[j]}, s_h] \\
 &= \mathcal{P}[\bar{V}_{i-1}^r[j] | \cancel{W_i^r[j]}, \bar{g}_{w_1}, \dots, \bar{g}_{w_n}, r_{i-1}[\omega_1] < \hat{R}, \dots, V_{i-1}^r[\omega_m] < \hat{R}] \\
 &= \mathcal{P}[\bar{V}_{i-1}^r[j] | \cancel{W_i^r[j]}, V_{i-1}^r[\omega_1] < \hat{R}, \dots, V_{i-1}^r[\omega_m] < \hat{R}] \\
 &= \mathcal{P}[\bar{V}_{i-1}^r[j] | V_{i-1}^r[j] < \hat{R}]
 \end{aligned}$$

Using Bayes' theorem we get

$$\mathcal{P}[\bar{V}_{i-1}^r | V_{i-1}^r < \hat{R}] = \frac{\mathcal{P}[V_{i-1}^r < \hat{R} | \bar{V}_{i-1}^r] \mathcal{P}[\bar{V}_{i-1}^r]}{\mathcal{P}[V_{i-1}^r < \hat{R}]} \quad (4.14)$$

Following from Lemma 41, each predecessor of  $j$  is independently selected with uniform probability. Given  $\bar{N}_{i-1}^r$ , each predecessor of  $j$  has a probability  $l(\bar{N}_{i-1}^r/N)$  of being in  $\hat{W}_{i-1}^r$ . The unconditioned number of infected predecessors of  $j$  is therefore binomially distributed:

$$\mathcal{P}[\bar{V}_{i-1}^r] = \text{Bin} \left[ E, l \frac{\bar{N}_{i-1}^r}{N} \right] (\bar{V}_{i-1}^r) \quad (4.15)$$

Plugging Equation (4.15) in Equation (4.14) and noting that  $N_{i-1}^r = N_i^r - U_i^r$  we get

$$\mathcal{P}[\bar{V}_{i-1}^r | V_{i-1}^r < \hat{R}] = \frac{I(\bar{V}_{i-1}^r < \hat{R}) \text{Bin} \left[ R, l \frac{\bar{N}_{i-1}^r - \bar{U}_i^r}{N} \right] (\bar{V}_{i-1}^r)}{\sum_{\bar{V}=0}^{\hat{R}-1} \text{Bin} \left[ R, l \frac{\bar{N}_{i-1}^r - \bar{U}_i^r}{N} \right] (\bar{V})}$$

We now compute the distribution underlying  $V_i^r$ , given  $\bar{V}_{i-1}^r$ ,  $\cancel{W_i^r}$  and  $s_h$ . Given  $\bar{V}_{i-1}^r$ ,  $\cancel{W_i^r}$  and  $s_h$ ,  $j$  has  $E - \bar{V}_{i-1}^r$  predecessors that are not in  $\hat{W}_{i-1}^r$ . Let  $g_{j,k}$  be a predecessor of  $j$  that is not in  $\hat{W}_{i-1}^r$ , we have

$$\begin{aligned}
 \mathcal{P}[g_{j,k} \in \hat{W}_i^r | g_{j,k} \notin \hat{W}_{i-1}^r] &= \frac{\mathcal{P}[g_{j,k} \in \hat{W}_i^r, g_{j,k} \notin \hat{W}_{i-1}^r]}{\mathcal{P}[g_{j,k} \notin \hat{W}_{i-1}^r]} \\
 &= \frac{l \frac{\bar{U}_i^r}{N}}{1 - l \frac{\bar{N}_{i-1}^r - \bar{U}_i^r}{N}}
 \end{aligned}$$

Following from Equation (4.4), each predecessor of  $j$  that is not in  $\hat{W}_{i-1}^r$  has an independent chance of being in  $\hat{W}_i^r$ . Therefore, the number of newly infected predecessors for  $j$  at step  $i$  is

## Chapter 4. Contagion

---

binomially distributed:

$$\mathcal{P}[\tilde{V}_i^r | \tilde{V}_{i-1}^r, \mathcal{W}_i^r, s_h] = \text{Bin} \left[ R - \tilde{V}_{i-1}^r, \frac{l \frac{\tilde{U}_i}{N}}{1 - l \frac{\tilde{N}_i - \tilde{U}_i}{N}} \right] (\tilde{V}_i^r - \tilde{V}_{i-1}^r) \quad (4.16)$$

Using the law of total probability, we can now use Equations (4.14) and (4.16) to compute the probability distribution underlying  $V_i^r[j]$ , given  $\mathcal{W}_i^r$  and  $s_h$ :

$$\mathcal{P}[\tilde{V}_i^r | \mathcal{W}_i^r, s_h] = \sum_{\tilde{V}_{i-1}^r=0}^{\hat{R}-1} \mathcal{P}[\tilde{V}_i^r | \tilde{V}_{i-1}^r, \mathcal{W}_i^r, s_h] \mathcal{P}[\tilde{V}_{i-1}^r | \mathcal{W}_i^r, s_h]$$

Finally, following from Lemma 45, we get the probability of  $W_i^r[j]$ , given  $\mathcal{W}_i^r$  and  $s_h$ :

$$\mathcal{P}[W_i^r | \mathcal{W}_i^r, s_h] = \sum_{\tilde{V}_i^r=\hat{R}}^R \mathcal{P}[\tilde{V}_i^r | \mathcal{W}_i^r, s_h]$$

Since each of the  $N - \tilde{N}_i$  nodes in  $\mathcal{W}$  has an independent probability of becoming infected at round  $i + 1$ , the frontier size at step  $i + 1$ , given  $s_h$  is binomially distributed:

$$\mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r | s_h] = \text{Bin}[N - \tilde{N}_i, \mathcal{P}[W_{i+1}^r | \mathcal{W}_i^r, s_h]] (\tilde{U}_{i+1}^r) \delta_{\tilde{N}_{i+1}^r - \tilde{N}_i^r, \tilde{U}_{i+1}^r}$$

We can now note how, when computing  $\mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r | s_h]$ , the condition on  $s_h$  reduces only to a condition on the values of  $\tilde{N}_i^r$  and  $\tilde{U}_i^r$ . Since  $s_1, \dots, s_q$  share the same values of  $(\tilde{N}_i^r, \tilde{U}_i^r)$ , the transition probability for the Markov chain underlying a contagion step reduces to

$$\begin{aligned} \mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r | \tilde{N}_i^r, \tilde{U}_i^r] &= \sum_{l=1}^q \mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r | s_l] \mathcal{P}[s_l | \langle \tilde{N}_i^r, \tilde{U}_i^r \rangle] \\ &= \mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r | s_h] \sum_{l=1}^q \mathcal{P}[s_l | \langle \tilde{N}_i^r, \tilde{U}_i^r \rangle] \\ &= \mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r | s_h] \end{aligned} \quad (4.17)$$

### 4.11.9 Final infection size

In Section 4.11.8, we showed that a contagion step defines a Markov chain with states  $(\tilde{N}_i^r, \tilde{U}_i^r)$ , and we computed the values of its associated transition matrix  $\mathcal{M}$ . In this section, we use this result to achieve our goal to compute the probability distribution underlying the infection size at the end of a game of Threshold Contagion.

As we established in Section 4.11.8, provided with  $\mathcal{P}[\tilde{N}_i^r, \tilde{U}_i^r]$ , we can compute  $\mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r]$ . Moreover, following from Corollary 2, every configuration  $\mathcal{P}[\tilde{N}_i^r, \tilde{U}_i^r]$  converges in a finite

number of steps  $i^*$  to satisfy

$$\begin{aligned}\mathcal{P}[\tilde{N}_i^r, \tilde{U}_i^r] &= \mathcal{P}[\tilde{N}_{i+1}^r, \tilde{U}_{i+1}^r] & \forall i \geq i^* \\ \mathcal{P}[U_i^r > 0] &= 0 & \forall i \geq i^*\end{aligned}$$

It is easy to see that the first step of each round (where the player selects  $S$  healthy node and infects them) also defines a Markov chain that deterministically increases, if possible, the infection size by  $S$ .

Specifically, the transition probabilities from step 0 to step 1 in each round are defined by:

$$\mathcal{P}[\tilde{N}_1^r, \tilde{U}_1^r] = \begin{cases} \mathcal{P}[N_0^r = \tilde{N}_1^r - S] & \text{iff } \tilde{N}_1^r \geq S, \tilde{U}_1^r = S \\ \mathcal{P}[N_0^r = \tilde{N}_1^r] & \text{iff } \tilde{N}_1^r > (N - S), \tilde{U}_1^r = 0 \\ 0 & \text{otherwise} \end{cases} \quad (4.18)$$

The distribution underlying the final infection size can be computed as follows:

- The distribution underlying the first step of the game is known:

$$\mathcal{P}[\tilde{N}_0^0, \tilde{U}_0^0] = \delta_{\tilde{N}_0^0, 0} \delta_{\tilde{U}_0^0, 0}$$

- For  $K$  rounds:

- If  $r > 0$ , then  $\mathcal{P}[\tilde{N}_0^r, \tilde{U}_0^r] = \mathcal{P}[\tilde{N}_\infty^{r-1}, \tilde{U}_\infty^{r-1}]$ .
- Apply Equation (4.18) to compute  $\mathcal{P}[\tilde{N}_1^r, \tilde{U}_1^r]$ .
- Until convergence:
  - \* Apply Equation (4.17) to compute  $\mathcal{P}[\tilde{N}_i^r, \tilde{U}_i^r]$ .



**Oracular Byzantine Reliable Broadcast** **Part II**





# 5 Overview

## 5.1 Introduction

Byzantine reliable broadcast (BRB) is one of the most fundamental and versatile building blocks in distributed computing, powering a variety of Byzantine fault-tolerant (BFT) systems [37, 61]. The BRB abstraction has recently been shown to be strong enough to process payments, enabling cryptocurrency deployments in an asynchronous environment [80]. Originally introduced by Bracha [31] to allow a set of processes to agree on a single message from a designated sender, BRB naturally generalizes to the multi-shot case, enabling higher-level abstractions such as Byzantine FIFO [34, 132] and causal [11, 24] broadcast. We study a practical, multi-shot variant of BRB whose interface is split between broadcasting *clients* and delivering *servers*. We call this abstraction Client-Server Byzantine Reliable Broadcast (CSB).

**CSB in brief.** Clients broadcast, and servers deliver, *payloads* composed by a *context* and a *message*. This interface allows, for example, Alice to announce her wedding as well as will her fortune by respectively broadcasting

$$\left( \underbrace{\text{"My wife is", "Carla"}}_{\text{context } c_w}, \underbrace{\text{"Carla"}}_{\text{message } m_w} \right) \left( \underbrace{\text{"All my riches go to", "Bob"}}_{\text{context } c_r}, \underbrace{\text{"Bob"}}_{\text{message } m_r} \right)$$

CSB guarantees that: (*Consistency*) no two correct servers deliver different messages for the same client and context; (*Totality*) either all correct servers deliver a message for a given client and context, or no correct server does; (*Integrity*) if a correct server delivers a payload from a correct client, then the client has broadcast that payload; and (*Validity*) a payload broadcast by a correct client is delivered by at least one correct server. Following from the above example, Carla being Alice's wife does not conflict with Bob being her sole heir (indeed,  $c_w \neq c_r$ ), but Alice would not be able to convince two correct servers that she married Carla and Diana, respectively. Higher-level broadcast abstractions can be easily built on top of CSB. For example, using integer sequence numbers as contexts and adding a reordering layer yields

## Chapter 5. Overview

---

Client-Server Byzantine FIFO Broadcast. For the sake of CSB, however, it is not important for contexts to be integers, or satisfy any property other than comparability. Throughout the remainder of this Part, the reader can picture contexts as opaque binary blobs. Lastly, while the set of servers is known, CSB as presented does not assume any client to be known a priori. The set of clients can be permissionless, with servers discovering new clients throughout the execution.

**A utopian model.** Real-world BRB implementations are often bottlenecked either by expensive signature verifications [52] or by communication overhead [32, 110, 111]. With the goal of broadening those bottlenecks, simplified, more trustful models are useful to establish a (sometimes grossly unreachable) bound on the efficiency that an algorithm can attain in the Byzantine setting. For example, in a utopian model where any agreed-upon process can be trusted to never fail (let us call it an *oracle*), CSB can easily be implemented with great efficiency. Upon initialization, the oracle organizes all clients in a list, which it disseminates to all servers. For simplicity, let us call *id* a client's position in the list. To broadcast a payload  $p$ , a client with *id*  $i$  simply sends  $p$  to the oracle: the oracle checks  $p$  for equivocation (thus ensuring consistency), then forwards  $(i, p)$  to all servers (thus ensuring validity and totality). Upon receiving  $(i, p)$ , a server blindly trusts the oracle to uphold all CSB properties, and delivers  $(i, p)$ . Oracle-CSB is clearly very efficient. On the one hand, because the oracle can be trusted not to attribute spurious payloads to correct clients, integrity can be guaranteed without any server-side signature verification. On the other, in order to deliver  $(i, p)$ , a server needs to receive just  $(\lceil \log_2(c) \rceil + |p|)$  bits, where  $c$  denotes the total number of clients, and  $|p|$  measures  $p$ 's length in bits. This is optimal assuming the rate at which clients broadcast is unknown<sup>1</sup> or uniform<sup>2</sup> [51].

**Matching the oracle.** Due to its reliance on a single infallible process, Oracle-CSB is not a fault-tolerant distributed algorithm: shifting back to the Byzantine setting, a single failure would be sufficient to compromise all CSB properties. Common sense suggests that Byzantine resilience will necessarily come at some cost: protocol messages must be exchanged to preserve consistency and totality, signatures must be produced and verified to uphold integrity and, lacking the totally-ordering power that only consensus can provide, *ids* cannot be assigned in an optimally dense way. However, this Part proves the counter-intuitive result that an asynchronous, optimally-resilient, Byzantine implementation of CSB can asymptotically match the efficiency of Oracle-CSB. This is not just up to a constant, but identically. In a synchronous execution, free from Byzantine misbehaviour, and as the number of concurrently

---

<sup>1</sup>Lacking an assumption on broadcasting rates, an adversarial scheduler could have all messages broadcast by the client with the longest *id*, which we cannot guarantee to be shorter than  $\lceil \log_2(c) \rceil$  bits.

<sup>2</sup>Should some clients be expected to broadcast more frequently than others, we could further optimize Oracle-CSB by assigning smaller *ids* to more active clients, possibly at the cost of having less active clients have *ids* whose length exceeds  $\lceil \log_2(c) \rceil$ . Doing so, however, is beyond the scope of this Part.

broadcasting clients goes to infinity (we call these conditions the *batching limit*<sup>3</sup>), our CSB implementation Draft delivers a payload  $p$  at an asymptotic<sup>4</sup>, amortized cost of 0 signature verifications<sup>5</sup> and  $(\lceil \log_2(c) \rceil + |p|)$  bits exchanged per server, the same as in Oracle-CSB (we say that Draft achieves *oracular efficiency*). At the batching limit a Draft server is dispensed from nearly all signature verifications, as well as nearly all traffic that would be normally required to convey protocol messages, signatures, or client public keys. Network is the limit: payloads are delivered as quickly as they can be received.

**CSB's common bottlenecks.** To achieve oracular efficiency, we focus on three types of server overhead that commonly affect a real-world implementation of CSB:

- *Protocol overhead.* Safekeeping consistency and totality typically requires some form of communication among servers. This communication can be direct (as in Bracha's original, all-to-all BRB implementation) or happen through an intermediary (as in Bracha's signed, one-to-all-to-one BRB variant), usually employing signatures to establish authenticated, intra-server communication channels through a (potentially Byzantine) relay.
- *Signature overhead.* Upholding integrity usually requires clients to authenticate their messages using signatures. For servers, this entails both a computation and a communication overhead. On the one hand, even using well-optimized schemes, signature verification is often CPU-heavy enough to dominate a server's computational budget, dwarfing in particular the CPU footprint of much lighter, symmetric cryptographic primitives such as hashes and ciphers. On the other hand, transmitting signatures results in a fixed communication overhead per payload delivered. While the size of a signature usually ranges from a few tens to a few hundreds of bytes, this overhead is non-negligible in a context where many clients broadcast small messages. This is especially true in the case of payments, where a message reduces to the identifier of a target account and an integer to denote the amount of money to transfer.
- *Identifier overhead.* CSB's multi-shot nature calls for a sender identifier to be attached to each broadcast payload. Classically, the client's public key is used as identifier. This is convenient for two reasons. First, knowing a client's identifier is sufficient to authenticate its payloads. Second, asymmetric keypairs have very low probability of collision. As such, clients can create identities in the system without any need for coordination: locally generating a keypair is sufficient to begin broadcasting messages. By cryptographic design, however, public keys are sparse, and their size does not change with the

<sup>3</sup>The batching limit includes other easily achievable, more technical conditions that we omit in this section for the sake of brevity. We formally define the batching limit in Section 6.4.2

<sup>4</sup>The asymptotic costs are reached quite fast, at rates comparable to  $C^{-1}$  or  $\log(C) \cdot C^{-1}$ .

<sup>5</sup>This does not mean that batches are processed in constant time: hashes and signature aggregations, for example, still scale linearly in the size of a batch. The real-world computational cost of such simple operations, however, is several orders of magnitude lower than that of signature verification.

## Chapter 5. Overview

---

number of clients. This translates to tens to hundreds of bytes being invested to identify a client from a set that can realistically be enumerated by a few tens of bits. Again, this communication overhead is heavier on systems where broadcasts are frequent and brief.

On the way to matching Oracle-CSB's performance, we develop techniques to negate all three types of overhead: at the batching limit, a Draft server delivers a payload wasting 0 bits to protocol overhead, performing 0 signature verifications, and exchanging  $\lceil \log_2(c) \rceil$  bits of identifier, the minimum required to enumerate the set of clients. We outline our contributions below, organized in three (plus one) take-home messages (T-HMs).

**T-HM1: The effectiveness of batching goes beyond total order.** In the totally ordered setting, batching is famously effective at amortizing protocol overhead [8, 133]. Instead of disseminating its message to all servers, a client hands it over to (one or more)<sup>6</sup> batching processes. Upon collecting a large enough set of messages, a batching process organizes all messages in a batch, which it then disseminates to the servers. Having done so, the batching process submits the batch's hash to the system's totally-ordering primitive. Because hashes are constant in length, the cost of totally ordering a batch does not depend on its size. Once batches are totally ordered, so too are messages (messages within a batch can be ordered by any deterministic function), and equivocations can be handled at the application layer (for example, in the context of a cryptocurrency, the second request to transfer the same asset can be ignored by all correct servers, with no need for additional coordination). At the limit of infinitely large batches, the relative overhead of the ordering protocol becomes vanishingly small, and a server can allocate virtually all of its bandwidth to receiving batches. This strategy, however, does not naturally generalize to CSB, where batches lack total order. As payloads from multiple clients are bundled in the same batch, a correct server might detect equivocation for only a subset of the payloads in the batch. Entirely accepting or entirely rejecting a partially equivocated batch is not an option. In the first case, consistency could be violated. In the second case, a single Byzantine client could single-handedly "poison" the batches assembled by every correct batching process with equivocated payloads, thus violating validity. In Draft, a server can partially reject a batch, acknowledging all but some of its payloads. Along with its partial acknowledgement, a server provides a proof of equivocation to justify each exception. Having collected a quorum of appropriately justified partial acknowledgements, a batching process has servers deliver only those payloads that were not excepted by any server. Because proofs of equivocations cannot be forged for correct clients, a correct client handing over its payload to a correct batching process is guaranteed to have that payload delivered. In the common case where batches have little to no equivocations, servers exchange either empty or small lists of exceptions, whose size does not scale with that of the batch. This extends the protocol-amortizing power of batching to CSB and, we conjecture, other non-totally ordered abstractions.

---

<sup>6</sup>In most real-world implementations, a client optimistically entrusts its payload to a single process, extending its request to larger portions of the system upon expiration of a suitable timeout.

**T-HM2: Interactive multi-signing can slash signature overhead.** Traditionally, batching protocols are non-interactive on the side of clients. Having offloaded its message to a correct batching process, a correct client does not need to interact further for its message to be delivered: the batching process collects an arbitrary set of independently signed messages and turns to the servers to get each signature verified, and the batch delivered. This approach is versatile (messages are not tied to the batch they belong to) and reliable (a client crashing does not affect a batch's progress) but expensive (the cost of verifying each signature is high and independent of the batch's size). In Draft, batching processes engage in an interactive protocol with clients to replace, in the good case, all individual signatures in a batch with a single, batch-wide *multi-signature*. In brief, multi-signature schemes extend traditional signatures with a mechanism to *aggregate* signatures and public keys: an arbitrarily large set of signatures for the same message<sup>7</sup> can be aggregated into a single, constant-sized signature; similarly, a set of public keys can be aggregated into a single, constant-sized public key. The aggregation of a set of signatures can be verified in constant time against the aggregation of all corresponding public keys. Unlike verification, aggregation is a cheap operation, reducing in some schemes to a single multiplication on a suitable field. Multi-signature schemes open a possibility to turn expensive signature verification into a once-per-batch operation. Intuitively, if each client contributing to a batch could multi-sign the entire batch instead of its individual payload, all multi-signatures could be aggregated, allowing servers to authenticate all payloads at once. However, as clients cannot predict how their payloads will be batched, this must be achieved by means of an interactive protocol. Having collected a set of individually-signed payloads in a batch, a Draft batching process shows to each contributing client that its payload was included in the batch. In response, clients produce their multi-signatures for the batch's hash, which the batching process aggregates. Clients that fail to engage in this interactive protocol (e.g., because they are faulty or slow) do not lose liveness, as their original signature can still be attached to the batch to authenticate their individual payload. In the good case, all clients reply in a timely fashion, and each server has to verify a single multi-signature per batch. At the limit of infinitely large batches, this results in each payload being delivered at an amortized cost of 0 signature verifications. The usefulness of this interactive protocol naturally extends beyond CSB to all multi-shot broadcast abstractions whose properties include integrity.

**T-HM3: Dense id assignment can be achieved without consensus.** In order to efficiently convey payload senders, Oracle-CSB's oracle organizes all clients in a list, attaching to each client a successive integral identifier. Once the list is disseminated to all servers, the oracle can identify each client by its identifier, sparing servers the cost of receiving larger, more sparse, client-generated public keys. Id-assignment strategies similar to that of Oracle-CSB can be developed, in the distributed setting, building on top of classical algorithms that identify clients by their full public keys (we call such algorithms *id-free*, as opposed to algorithms such as Draft, which are *id-optimized*). In a setting where consensus can be achieved, the identifier

<sup>7</sup>Some multi-signature schemes also allow the aggregation of signatures on heterogeneous messages. In that case, however, aggregation is usually as expensive as signature verification. Given our goal to reduce CPU complexity for servers, this Part entirely disregards heterogeneous aggregation schemes.

## Chapter 5. Overview

---

density of Oracle-CSB is easily matched. Upon initialization, each client submits its public key to an id-free implementation of Atomic Broadcast. Upon delivery of a public key, every correct process agrees on its position within the common, totally-ordered log. As in Oracle-CSB, each client can then use its position in the list as identifier within some faster, id-optimized broadcast implementation. In a consensus-less setting, achieving a totally-ordered list of public keys is famously impossible [62]. This Part, however, proves the counter-intuitive result that, when batching is used, the density of ids assigned by a consensus-less abstraction can asymptotically match that of those produced by Oracle-CSB or consensus. In Dibs, our consensus-less id-assigning algorithm, a client requests an id from every server. Each server uses an id-free implementation of FIFO Broadcast to order the client's public key within its own log. Having observed its public key appear in at least one log, the client publicly elects the server in charge of that log to be its *assigner*. Having done so, the client obtains an id composed of the assigner's public key and the client's position within the assigner's log. We call the two components of an id *domain* and *index*, respectively. Because the set of servers is known to (and can be enumerated by) all processes, an id's domain can be represented in  $\lceil \log_2(n) \rceil$  bits, where  $n$  denotes the total number of servers. Because at most  $c$  distinct clients can appear in the FIFO log of any server, indices are at most  $\lceil \log_2(c) \rceil$  bits long. In summary, Dibs assigns ids to clients without consensus, at an additional cost of  $\lceil \log_2(n) \rceil$  bits per id. Interestingly, even this additional complexity can be amortized by batching. Having assembled a batch, a Draft batching process represents senders not as a list of ids, but as a map, associating to each of the  $n$  domains the indices of all ids in the batch under that domain. At the limit of infinitely large batches ( $C \gg N$ ), the bits required to represent the map's keys are entirely amortized by those required to represent its values. This means that, while  $(\lceil \log_2(n) \rceil + \lceil \log_2(c) \rceil)$  bits are required to identify a client in isolation,  $\lceil \log_2(c) \rceil$  bits are sufficient if the client is batched: even without consensus, Draft asymptotically matches the id efficiency of Oracle-CSB.

**Bonus T-HM: Untrusted processes can carry the system.** In THM1, we outlined how batching can be generalized to the consensus-less case, and discussed its role in removing protocol overhead. In THM2, we sketched how an interactive protocol between clients and batching processes can eliminate signature overhead. In employing these techniques, we shifted most of the communication and computation complexity of our algorithms from servers to batching processes. Batching processes verify all client signatures, create batches, verify and aggregate all client multi-signatures, then communicate with servers in an expensive one-to-all pattern, engaging server resources (at the batching limit) as little as an oracle would. Our last contribution is to observe that a batching process plays no role in upholding CSB's safety. As we discuss in detail throughout the remainder of this Part, a malicious batching process cannot compromise consistency (it would need to collect two conflicting quorums of acknowledgements), totality (any server delivering a batch has enough information to convince all others to do the same) or integrity (batches are still signed, and forged or improperly aggregated multi-signatures are guaranteed to be detected). Intuitively, the only damage a batching pro-

cess can do to the system is to refuse to process client payloads<sup>8</sup>. This means that a batching process does not need to satisfy the same security properties as a server. CSB’s properties cannot be upheld if a third of the servers are faulty. Conversely, Draft has both liveness and safety as long as *a single* batching process is correct. This observation has profound practical implications. In the real world, scaling the resources of a permissioned, security-critical set of servers can be hard. On the one hand, reputable, dependable institutions partaking in the system might not have the resources to keep up with its demands. On the other, more trusted hardware translates to a larger security cross-section. Trustless processes, however, are plentiful to the point that permissionless cryptocurrencies traditionally waste their resources, making them compete against each other in expensive proofs of Sybil-resistance [122]. In this Part, we extend the classical client-server model with *brokers*, a permissionless, scalable set of processes whose only purpose is to alleviate server complexity. Unlike servers, more than two-thirds of which we assume to be correct, all brokers but one can be faulty. In Draft, brokers act as an intermediary between clients and servers, taking upon themselves the batching of payloads, verification and aggregation of signatures, the dissemination of batches, and the transmission of protocol messages.

**Roadmap.** We discuss related work in Section 5.2. We state our model and recall useful cryptographic background in Section 5.3. We discuss our CSB implementation Draft in Chapter 6: we formally define the Client-Server Byzantine Reliable Broadcast abstraction in Section 6.1; we present Draft in Section 6.2; we prove Draft’s correctness in Section 6.3; we analyze Draft’s complexity at the batching limit in Section 6.4. We discuss our Directory implementation Dibs in Chapter 7: we introduce the Directory abstraction in Section 7.1; we present Dibs’s pseudocode in Section 7.2; we prove the correctness of Dibs in Section 7.3.

## 5.2 Related Work

Byzantine Reliable Broadcast (BRB) is a classical primitive of distributed computing, with widespread practical applications such as in State Machine Replication (SMR) [33, 37, 118], Byzantine agreement [44, 92, 100, 123, 141], blockchains [8, 53, 55], and online payments [50, 80, 101]. In classical BRB, a system of  $n$  processes agree on a single message from a single *source* (one of the  $n$  processes), while tolerating up to  $f$  Byzantine failures ( $f$  of the  $n$  processes can behave arbitrarily). A well known solution to asynchronous BRB with provably optimal resilience ( $f < n/3$ ) was first proposed by Bracha [30, 31] who introduced the problem. Bracha’s broadcast reaches  $O(n^2)$  message complexity, and  $O(n^2L)$  communication complexity (total number of transmitted bits between correct processes [154]), where  $L$  is the length of the message. Since  $O(n^2)$  message complexity is provably optimal [59], the main focus of BRB-related research has been on reducing its communication complexity. The best lower bound

<sup>8</sup>Or cause servers to waste resources, e.g., by transmitting improperly signed batches. Simple accountability measures, we conjecture, would be sufficient to mitigate these attacks in Draft. A full discussion of Denial of Service, however, is beyond the scope of this Part.

for communication complexity is  $\Omega(nL + n^2)$ , although it is unknown whether it is tight. The  $nL$  term comes from all processes having to receive the message (length  $L$ ), while the  $n^2$  term comes from each of the  $n$  processes having to receive  $\Omega(n)$  protocol messages to ensure agreement in the presence of  $f = \Theta(n)$  failures [59]. One line of research focuses on worst-case complexity, predominantly using error correcting codes [17, 130] or erasure codes [5, 38, 85, 126], and has produced various BRB protocols with improved complexity [5, 35, 38, 57, 123], many of them quite recently. The work of Das, Xiang and Ren [57] achieves  $O(nL + kn^2)$  communication complexity (specifically,  $7nL + 2kn^2$ ), where  $k$  is the security parameter (e.g., the length of a hash, typically 256 bits). As the authors note, the value of hidden constants (and  $k$ , which is sometimes considered as a constant in literature) is particularly important when considering practical implementations of these protocols. Another line of research focuses on optimizing the good case performance of BRB, i.e., when the network behaves synchronously and no process misbehaves [1, 35, 44, 100, 129]. As the good case is usually the common case, in practice, the real-world communication complexity of these optimistic protocols matches that of the good case. A simple and widely-used hash-based BRB protocol is given by Cachin *et al.* [35]. It replaces the echo and ready phase messages in Bracha's protocol with hashes, achieving  $O(nL + kn^2)$  in the good case (specifically,  $nL + 2kn^2$ ), and  $O(n^2L)$  in the worst-case. Considering practical throughput, some protocols also focus on the *amortized* complexity per source message [44, 111, 129]. Combining techniques such as *batching* [44] and threshold signatures [137], at the limit (of batch size), BRB protocols reach  $O(nL)$  amortized communication complexity in the good case [129]. At this point, the remaining problem lies in the hidden constants. In the authenticated setting, batching-based protocols rely on digital signatures to validate (source) messages before agreeing to deliver them [129]. In reality, each source message in a batch includes its content, an identifier of the source (e.g., a  $k$ -sized public key), a sequence id (identifying the message), and a  $k$ -sized signature. When considering systems where  $L$  is small (e.g., online payments), these can take up a large fraction of the communication. To be precise, the good-case amortized communication complexity would be  $O(nL + kn)$ . In fact, message signatures (the  $kn$  factor) are by far the main bottleneck in practical applications of BRB today [55, 141], both in terms of communication and computation (signature verification), leading to various attempts at reducing or amortizing their cost [53, 111]. For example, Crain *et al.* [53] propose *verification sharding*, in which only  $f + 1$  processes have to receive and verify all message signatures in the good case, which is a 3-fold improvement over previous systems (on the  $kn$  factor) where all  $n$  processes verify all signatures. However, by itself, this does not improve on the amortized cost of  $O(nL + kn)$  per message. When contrasting theoretical research with practical systems, it is interesting to note the gap that can surge between the theoretical model and reality. The recent work of Abraham *et al.* [1], focused on the good-case latency of Byzantine broadcast, expands on some of these mismatches and argues about the practical limitations of focusing on the worst-case. Another apparent mismatch lies in the classical model of Byzantine broadcast. In many of the applications of BRB mentioned previously (e.g., SMR, permissioned blockchains, online payments), there is usually a set of *servers* ( $n$ , up to  $f$  of which are faulty), and a set of external clients ( $X$ ) which are the true sources of messages. The usual transformation



from BRB's classical model into these practical settings maps the set of  $n$  servers as the  $n$  processes and simply excludes clients as system entities, e.g., assuming their messages are relayed through one of the servers. Since the number of clients can be very large ( $|X| \gg n$ ), clients are untrusted (which can limit their usefulness), and the focus is on the communication complexity of the *servers*, this transformation seems reasonable and simplifies the problem. However, it can also limit the search for more practical solutions. In this Part, in contrast with the classical model of BRB, we explicitly include the set of clients  $X$  in our system while focusing on the communication complexity surrounding the servers (i.e., the bottleneck). Furthermore, we introduce *brokers*, an untrusted set  $B$  of processes, only one of which is assumed to be correct, whose goal is to assist servers in their operation. By doing this, we can leverage brokers to achieve a good-case, amortized communication complexity (for servers, information received or sent) of  $nL + o(nL)$ .

### 5.3 Model & background

#### 5.3.1 Model

**System and adversary.** We assume an asynchronous message-passing system where the set  $\Pi$  of processes is the distinct union of three sets: **servers** ( $\Sigma$ ), **brokers** ( $B$ ), and **clients** ( $X$ ). We use  $n = |\Sigma|$ ,  $k = |B|$  and  $c = |X|$ . Any two processes can communicate via reliable, FIFO, point-to-point links (messages are delivered in the order they are sent). Faulty processes are Byzantine, i.e., they may fail arbitrarily. Byzantine processes know each other, and may collude and coordinate their actions. At most  $f$  servers are Byzantine, with  $n = 3f + 1$ . At least one broker is correct. All clients may be faulty. We use  $\Pi_C$  and  $\Pi_F$  to respectively identify the set of correct and faulty processes. The adversary cannot subvert cryptographic primitives (e.g., forge signatures). Servers and brokers<sup>9</sup> are permissioned (every process knows  $\Sigma$  and  $B$ ), clients are permissionless (no correct process knows  $X$  a priori). We call *certificate* a statement signed by either a plurality ( $f + 1$ ) or a quorum ( $2f + 1$ ) of servers. Since every process knows  $\Sigma$ , any process can verify a certificate.

**Good case.** The algorithms presented in this Part are designed to uphold all their properties in the model above. Draft, however, achieves oracular efficiency only in the *good case*. In the good case, links are synchronous (messages are delivered at most one time unit after they are sent), all processes are correct, and the set of brokers contains only one element. To take advantage of the good case, Draft makes use of timers (which is uncommon for purely asynchronous algorithms). A timer with timeout  $\delta$  set at time  $t$  rings: after time  $(t + \delta)$ , if the system is synchronous; after time  $t$ , otherwise. Intuitively, in the non-synchronous case, timers disregard their timeout entirely, and are guaranteed to ring only eventually.

---

<sup>9</sup>The assumption that brokers are permissioned is made for simplicity, and can be easily relaxed to the requirement that every correct process knows at least one correct broker.

### 5.3.2 Background

Besides commonly used hashes and signatures, the algorithms presented in this Part make use of two less often used cryptographic primitives, namely, multi-signatures and Merkle trees. We briefly outline their use below. An in-depth discussion of their inner workings, however is beyond the scope of this Part.

**Multi-signatures.** Like traditional signatures, multi-signatures [15] are used to publicly authenticate messages: a public / secret keypair  $(p, r)$  is generated locally;  $r$  is used to produce a signature  $s$  for a message  $m$ ;  $s$  is publicly verified against  $p$  and  $m$ . Unlike traditional signatures, however, multi-signatures for the same message can be *aggregated*. Let  $(p_1, r_1), \dots, (p_n, r_n)$  be a set of keypairs, let  $m$  be a message, and let  $s_i$  be  $r_i$ 's signature for  $m$ .  $(p_1, \dots, p_n)$  and  $(s_1, \dots, s_n)$  can be respectively aggregated into a constant-sized public key  $\hat{p}$  and a constant-sized signature  $\hat{s}$ . As with individually-generated multi-signatures,  $\hat{s}$  can be verified in constant time against  $\hat{p}$  and  $m$ . Aggregation is cheap and non-interactive: provided with  $(p_1, \dots, p_n)$  (resp.,  $(s_1, \dots, s_n)$ ) any process can compute  $\hat{p}$  (resp.,  $\hat{s}$ ).

**Merkle trees.** Merkle trees [117] extend traditional hashes with compact *proofs of inclusion*. As with hashes, a sequence  $(x_1, \dots, x_n)$  of values can be hashed into a preimage and collision-resistant digest (or *root*)  $r$ . Unlike hashes, however, a proof  $p_i$  can be produced from  $(x_1, \dots, x_n)$  to attest that the  $i$ -th element of the sequence whose root is  $r$  is indeed  $x_i$ . In other words, provided with  $r$ ,  $p_i$  and  $x_i$ , any process can verify that the  $i$ -th element of  $(x_1, \dots, x_n)$  is indeed  $x_i$ , without having to learn  $(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$ . The size of a proof of inclusion for a sequence of  $n$  elements is logarithmic in  $n$ .

# 6 Draft

In this chapter, we present in detail the **Client-Server Byzantine Reliable Broadcast** abstraction and discuss its properties. We then present Draft, an algorithm that implements Client-Server Byzantine Reliable Broadcast, and evaluate its **security** and **complexity**, showing that it achieves oracularity at the batching limit.

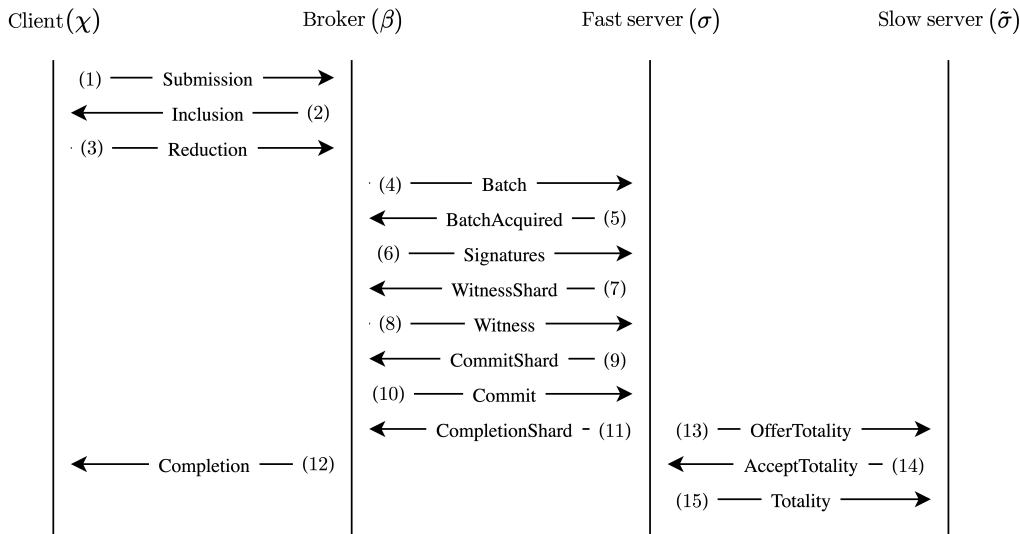
## 6.1 Interface

A **Client-Server Byzantine Reliable Broadcast (CSB) system** offers two interfaces, CSB Client (instance  $cl$ ) and CSB Server (instance  $sr$ ), exposing the following events:

- **Request:**  $\langle cl.Broadcast \mid context, message \rangle$ : Broadcasts a message  $message$  for context  $context$  to all servers.
- **Indication**  $\langle sr.Deliver \mid client, context, message \rangle$ : Delivers the message  $message$  broadcast by client  $client$  for context  $context$ .

A Client-Server Byzantine Reliable Broadcast system satisfies the following properties:

1. **No duplication:** No correct server delivers more than one message for the same client and context.
2. **Integrity:** If a correct server delivers a message  $m$  for context  $c$  from a correct client  $\chi$ , then  $\chi$  previously broadcast  $m$  for  $c$ .
3. **Consistency:** No two correct servers deliver different messages for the same client and context.
4. **Validity:** If a correct client  $\chi$  broadcasts a message for context  $c$ , then eventually a correct server delivers a message for  $c$  from  $\chi$ .



**Figure 6.1: Draft’s protocol.** Having collected a batch of client payloads, a broker engages in an interactive protocol with clients to reduce the batch, replacing (most of) its individual payload signatures with a single, batch-wide multi-signature. The broker then disseminates the batch to all servers, successively gathering a witness for its correctness and a certificate to commit (some of) its payloads. Having had a plurality of servers deliver the batch, the broker notifies all clients with a suitable certificate. In the bad case, servers can ensure totality without any help from the broker, propagating batches and commit certificates in an all-to-all fashion.

5. **Totality:** If a correct server delivers a message for context  $c$  from a client  $\chi$ , then eventually every correct server delivers a message for  $c$  from  $\chi$ .

## 6.2 Algorithm

In this section, we present our CSB implementation Draft. We outline Draft’s protocol in Section 6.2.1, providing qualitative arguments for its security. We qualitatively analyze Draft’s complexity in Section 6.2.2. We present Draft’s full pseudocode in Sections 6.2.3 (client), 6.2.4 (broker) and 6.2.5 (server). The remainder of Chapter 6 proves Draft’s security and complexity to the fullest extent of formal detail.

### 6.2.1 Protocol & correctness overview

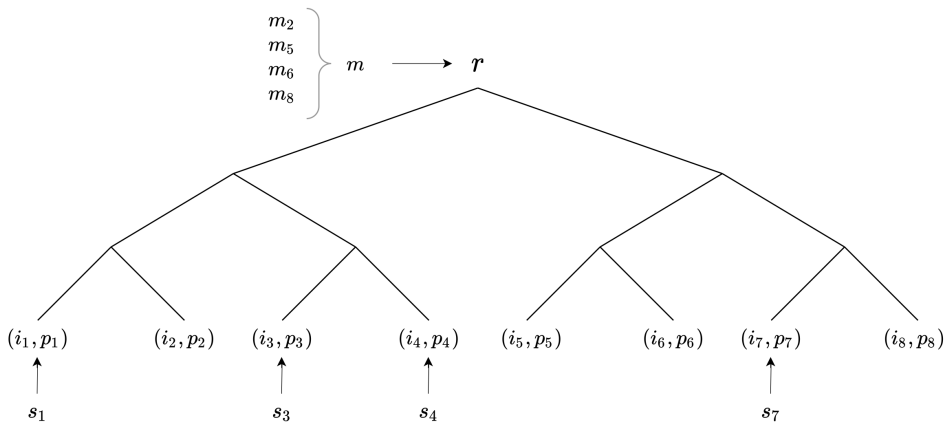
**Dramatis personae.** The goal of this section is to provide an intuitive understanding of Draft’s protocol. In order to do this, we focus on four processes: a correct client  $\chi$ , a correct broker  $\beta$ , a correct and fast server  $\sigma$ , and a correct but slow server  $\bar{\sigma}$ . We follow the messages exchanged between  $\chi$ ,  $\beta$ ,  $\sigma$  and  $\bar{\sigma}$  as the protocol unfolds, as captured by Figure 6.1.

**The setting.**  $\chi$ 's goal is to broadcast a payload  $p$ .  $\chi$  has already used Draft's underlying Directory abstraction (DIR) to obtain an id  $i$ . In brief, DIR guarantees that  $i$  is assigned to  $\chi$  only, and provides  $\chi$  with an *assignment certificate*  $a$ , which  $\chi$  can use to prove that its id is indeed  $i$ . As we discussed in Section 5.1, Draft uses DIR-assigned ids to identify payload senders. This is essential to Draft's performance, as DIR guarantees *density*: as we outline in Section 6.2.2,  $\lceil \log_2(c) \rceil$  bits are asymptotically sufficient to represent each id in an infinitely large batch. We discuss the details of the DIR abstraction in Section 7, along with our DIR implementation, Dibs. Throughout the remainder of this Part, we say that a process  $\pi$  *knows* an id  $\hat{i}$  iff  $\pi$  knows the public keys to which  $\hat{i}$  is assigned.

**Building a batch.** In order to broadcast its payload  $p$ ,  $\chi$  produces a signature  $s$  for  $p$ , and then sends a Submission message to  $\beta$  (fig. 6.1, step 1). The Submission message contains  $p$ ,  $s$ , and  $\chi$ 's assignment certificate  $a$ . Upon receiving the Submission message,  $\beta$  learns  $\chi$ 's id  $i$  from  $a$ , then verifies  $s$  against  $p$ . Having done so,  $\beta$  stores  $(i, p, s)$  in its *submission pool*. For a configurable amount of time,  $\beta$  fills its pool with submissions from other clients, before *flushing* it into a *batch*. Let us use  $(i_1, p_1, s_1), \dots, (i_b, p_b, s_b)$  to enumerate the elements  $\beta$  flushes from the submission pool (for some  $n$ , we clearly have  $(i, p, s) = (i_n, p_n, s_n)$ ). For convenience, we will also use  $\chi_j$  to identify the sender of  $p_j$  (owner of  $i_j$ ). Importantly,  $\beta$  flushes the pool in such a way that  $i_j \neq i_k$  for all  $j \neq k$ : for safety reasons that will soon be clear, Draft's protocol prevents a client from having more than one payload in any specific batch. Because of this constraint, some payloads might linger in  $\beta$ 's pool. This is not an issue:  $\beta$  will simply flush those payloads to a different batch at a later time. When building the batch,  $\beta$  splits submissions and signatures, storing  $(i_1, p_1), \dots, (i_b, p_b)$  separately from  $s_1, \dots, s_b$ .

**Reducing the batch.** Having flushed submissions  $(i_1, p_1), \dots, (i_b, p_b)$  and signatures  $s_1, \dots, s_b$ ,  $\beta$  moves on to *reduce* the batch, as exemplified in Figure 6.2. In an attempt to minimize signature overhead for servers,  $\beta$  engages in an interactive protocol with clients  $\chi_1, \dots, \chi_b$  to replace as many signatures as possible with a single, batch-wide multi-signature. In order to do so,  $\beta$  organizes  $(i_1, p_1), \dots, (i_b, p_b)$  in a Merkle tree with root  $r$  (for brevity, we call  $r$  the batch's root).  $\beta$  then sends an Inclusion message to each  $\chi_j$  (fig. 6.1, step 2). Each Inclusion message contains  $r$ , along with a proof of inclusion  $q_j$  for  $(i_j, p_j)$ . Upon receiving its Inclusion message,  $\chi$  checks  $q_n$  against  $r$ . In doing so,  $\chi$  comes to two conclusions. First,  $\chi$ 's submission  $(i, p) = (i_n, p_n)$  is part of a batch whose root is  $r$ . Second, because no Draft batch can contain multiple payloads from the same client, that batch does not attribute  $\chi$  any payload other than  $p$ . In other words,  $\chi$  can be certain that  $\beta$  will not broadcast some spurious payload  $p' \neq p$  in  $\chi$ 's name: should  $\beta$  attempt to do that, the batch would be verifiably malformed, and immediately discarded. This means  $\chi$  can safely produce a multi-signature  $m$  for  $r$ : as far as  $\chi$  is concerned, the batch with root  $r$  upholds integrity. Having signed  $r$ ,  $\chi$  sends  $m$  to  $\beta$  by means of a Reduction message (fig. 6.1, step 3). Upon receiving  $\chi_j$ 's Reduction message,  $\beta$  checks  $\chi_j$ 's multi-signature  $m_j$  against  $r$ . Having done so,  $\beta$  discards  $\chi_j$ 's original signature  $s_j$ . Intuitively, with  $m_j$ ,  $\chi_j$  attested its agreement with whatever payload the batch

attributes to  $\chi_j$ . Because this is equivalent to individually authenticating  $p_j$ ,  $s_j$  is redundant and can be dropped. Upon expiration of a suitable timeout,  $\beta$  stops collecting `Reduction` messages: clearly, if  $\beta$  waited for every  $\chi_j$  to produce  $m_j$ , a single Byzantine client could prevent the protocol from moving forward by refusing to send its `Reduction` message.  $\beta$  aggregates all the multi-signatures it collected for  $r$  into a single, batch-wide multi-signature  $m$ . In the good case, every  $\chi_j$  is correct and timely. If so,  $\beta$  drops all individual signatures, and the entire batch is authenticated by  $m$  alone.



**Figure 6.2: An example of partially reduced batch.**  $B = 8$  submissions are organized on the leaves of a Merkle tree with root  $r$ . Each submission  $(i_j, p_j)$  is originally authenticated by an individual signature  $s_j$ . Upon collecting a multi-signature  $m_j$  for  $r$ , the broker drops  $s_j$ . Here the broker collected multi-signatures  $m_2, m_5, m_6$  and  $m_8$ , leaving a *straggler set*  $S = \{(i_1, s_1), (i_3, s_3), (i_4, s_4), (i_7, s_7)\}$ . Upon expiration of a suitable timeout, the broker aggregates  $m_2, m_5, m_6$  and  $m_8$  into a single multi-signature  $m$ . As such, every payload in the batch is authenticated either by  $m$  or by  $S$ .

**The perks of a reduced batch.** Having reduced the batch,  $\beta$  is left with a sequence of submissions  $(i_1, p_1), \dots, (i_b, p_b)$ , a multisignature  $m$  on the Merkle root  $r$  of  $(i_1, p_1), \dots, (i_b, p_b)$ , and a *straggler set*  $S$  holding the individual signatures that  $\beta$  failed to reduce. More precisely,  $S$  contains  $(i_j, s_j)$  iff  $\beta$  did not receive a valid `Reduction` message from  $\chi_j$  before the reduction timeout expired. We recall that  $m$ 's size is constant, and  $S$  is empty in the good case. Once reduced, the batch is cheap to authenticate: it is sufficient to verify the batch's multi-signature against the batch's root, and each straggler signature against its individual payload. More precisely, let  $T$  denote the set of *timely* clients ( $\chi_j$  is in  $T$  iff  $(i_j, ..)$  is not in  $S$ ). Let  $t$  denote the aggregation of  $T$ 's public keys. Provided with  $(i_1, p_1), \dots, (i_b, p_b)$ ,  $m$  and  $S$ , any process that knows  $i_1, \dots, i_b$  can verify that the batch upholds integrity by: (1) computing  $r$  and  $t$  from  $(i_1, p_1), \dots, (i_b, p_b)$  and  $S$ ; (2) using  $t$  to verify  $m$  against  $r$ ; and (3) verifying each  $s_j$  in  $S$  against  $p_j$ . In the good case, authenticating the batch reduces to verifying a single multi-signature. This is regardless of the batch's size.

**The pitfalls of a reduced batch.** As we discussed in the previous paragraph, reducing a batch makes it cheaper to verify its integrity. Reduction, however, hides a subtle trade-off: once reduced, a batch gets easier to authenticate *as whole*. Its *individual payloads*, however, become harder to authenticate. For the sake of simplicity, let us imagine that  $\beta$  successfully dropped all the individual signatures it originally gathered from  $\chi_1, \dots, \chi_b$ . In order to prove that  $(\chi = \chi_n)$  broadcast  $(p = p_n)$ ,  $\beta$  could naively produce the batch's root  $r$ ,  $(i_n, p_n)$ 's proof of inclusion  $q_n$ , and the batch's multi-signature  $m$  for  $r$ . This, however, would not be sufficient to authenticate  $p$ : because the multi-signature  $m_n$  that  $\chi$  produced for  $r$  was aggregated with all others,  $m$  can only be verified by the aggregation of *all*  $\chi_1, \dots, \chi_b$ 's public keys. This makes authenticating  $p$  as expensive as authenticating the entire batch: in order to verify  $m$ , *all*  $(i_j, p_j)$  must be produced and checked against  $r$ , so that all corresponding public keys can be safely aggregated.

**Witnessing the batch.** As we discuss next, proving the integrity of individual payloads is fundamental to ensure Draft's validity. In brief, to prove that some  $\chi_k$  equivocated its payload  $p_k = (c_k, l_k)$ , a server must prove to  $\beta$  that  $\chi_k$  also issued some payload  $p'_k = (c_k, l'_k \neq l_k)$ . Lacking this proof, a single Byzantine server could, for example, claim without basis that  $\chi$  equivocated  $p$ . This could trick  $\beta$  into excluding  $p$ , thus compromising Draft's validity. As we discussed in the previous paragraph, however, proving the integrity of an individual payload in a reduced batch is difficult. While we conjecture that purely cryptographic solutions to this impasse might be achievable in some schemes<sup>1</sup>, Draft has  $\beta$  engage in a simple protocol to further simplify the batch's authentication, replacing all client-issued (multi-)signatures with a single, server-issued certificate. Having collected and reduced the batch,  $\beta$  sends a Batch message to all servers (fig. 6.1, step 4). The Batch message only contains  $(i_1, p_1), \dots, (i_b, p_b)$ . Upon receiving the Batch message,  $\sigma$  collects in a set  $U_\sigma$  all the ids it does not know ( $i_j$  is in  $U_\sigma$  iff  $\sigma$  does not know  $i_j$ ), and sends  $U_\sigma$  back to  $\beta$  by means of a BatchAcquired message (fig. 6.1, step 5). Upon receiving  $\sigma$ 's BatchAcquired message,  $\beta$  builds a set  $A_\sigma$  containing all id assignments that  $\sigma$  is missing ( $a_j$  is in  $A_\sigma$  iff  $i_j$  is in  $U_\sigma$ ). Having done so,  $\beta$  sends a Signatures message to  $\sigma$  (fig. 6.1, step 6). The Signatures message contains the batch's multi-signature  $m$ , the straggler set  $S$ , and  $A_\sigma$ . We underline the importance of sending id assignments upon request only. Thinking to shave one round-trip off the protocol,  $\beta$  could naively package in a single message all submissions, all (multi-)signatures, and all assignments relevant to the batch. In doing so, however,  $\beta$  would force each server to receive one assignment per submission, immediately forfeiting Draft's oracular efficiency. As we discuss in Section 6.4.2, at the batching limit we assume that all servers already know all broadcasting clients. In that case, both  $U_\sigma$  and  $A_\sigma$  are constant-sized, empty sets, adding

<sup>1</sup>For example, using BLS,  $\beta$  could aggregate the public keys of  $\chi_1, \dots, \chi_{n-1}, \chi_{n+1}, \dots, \chi_b$  into a public key  $\tilde{t}_n$ , then show that the aggregation of  $\tilde{t}_n$  with  $\chi$ 's public key correctly verifies  $m$  against  $r$ . Doing so, however, would additionally require  $\beta$  to exhibit a proof that  $\tilde{t}_n$  is not a rogue public key, i.e., that  $\tilde{t}_n$  indeed results from the aggregation of client public keys. This could be achieved by additionally having  $\chi_1, \dots, \chi_b$  multi-sign some hard-coded statement to prove that they are not rogues.  $\beta$  could aggregate such signatures on the fly, producing a rogue-resistance proof for  $\tilde{t}_n$  that can be transmitted and verified in constant time. This, however, is expensive (and, frankly, at the limit of our cryptographic expertise).

only a vanishing amount of communication complexity to the protocol. Upon receiving the Signatures message,  $\sigma$  verifies and learns all assignments in  $A_\sigma$ . Having done so,  $\sigma$  knows  $i_1, \dots, i_b$ . As we outlined above,  $\sigma$  can now efficiently authenticate the whole batch, verifying  $m$  against the batch's root  $r$ , and each  $i_j$  in  $S$  against  $p_j$ . Having established the integrity of the whole batch,  $\sigma$  produces a *witness shard* for the batch, i.e., a multi-signature  $w_\sigma$  for  $[\text{Witness}, r]$ , effectively affirming to have successfully authenticated the batch.  $\sigma$  sends  $w_\sigma$  back to  $\beta$  by means of a WitnessShard message (fig. 6.1, step 7). Having received a valid WitnessShard message from  $f + 1$  servers,  $\beta$  aggregates all witness shards into a *witness*  $w$ . Because  $w$  is a plurality ( $f + 1$ ) certificate, at least one correct server necessarily produced a witness shard for the batch. This means that at least one correct server has successfully authenticated the batch by means of client (multi-)signatures. Because  $w$  could not have been gathered if the batch was not properly authenticated,  $w$  itself is sufficient to authenticate the batch, and  $\beta$  can drop all (now redundant) client-generated (multi-)signatures for the batch. Unlike  $m$ ,  $w$  is easy to verify, as it is signed by only  $f + 1$ , globally known servers. Like  $m$ ,  $w$  authenticates  $r$ . As such, any  $p_j$  can now be authenticated just by producing  $w$ , and  $(i_j, p_j)$ 's proof of inclusion  $q_j$ .

**Gathering a commit certificate.** Having successfully gathered a witness  $w$  for the batch,  $\beta$  sends  $w$  to all servers by means of a Witness message (fig. 6.1, step 8). Upon receiving the Witness message,  $\sigma$  moves on to check  $(i_1, p_1), \dots, (i_b, p_b)$  for equivocations. More precisely,  $\sigma$  builds a set of *exceptions*  $E_\sigma$  containing the ids of all equivocating submissions in the batch ( $i_j$  is in  $E_\sigma$  iff  $\sigma$  previously observed  $\chi_j$  submit a payload  $p'_j$  that conflicts with  $p_j$ ; we recall that  $p_j$  and  $p'_j$  conflict if their contexts are the same, but their messages are different).  $\sigma$  then produces a *commit shard* for the batch, i.e., a multi-signature  $c_\sigma$  for  $[\text{Commit}, r, E_\sigma]$ , effectively affirming that  $\sigma$  has found all submissions in the batch to be non-equivocated, *except* for those in  $E_\sigma$ . In the good case, every client is correct and  $E_\sigma$  is empty. Having produced  $c_\sigma$ ,  $\sigma$  moves on to build a set  $Q_\sigma$  containing a *proof of equivocation* for every element in  $E_\sigma$ . Let us assume that  $\sigma$  previously received from some  $\chi_k$  a payload  $p'_k$  that conflicts with  $p_k$ .  $\sigma$  must have received  $p'_k$  as part of some witnessed batch. Let  $r'_k$  identify the root of  $p'_k$ 's batch, let  $w'_k$  identify  $r'_k$ 's witness, let  $q'_k$  be  $(i_k, p'_k)$ 's proof of inclusion in  $r'_k$ . By exhibiting  $(r'_k, w'_k, p'_k)$ ,  $\sigma$  can prove to  $\beta$  that  $\chi_k$  equivocated:  $p_k$  conflicts with  $p'_k$ , and  $(i_k, p'_k)$  is provably part of a batch whose integrity was witnessed by at least one correct server. Furthermore, because correct clients never equivocate,  $(r'_k, w'_k, p'_k)$  is sufficient to convince  $\beta$  that  $\chi_k$  is Byzantine. For each  $i_j$  in  $E_\sigma$ ,  $\sigma$  collects in  $Q_\sigma$  a proof of equivocation  $(r'_j, w'_j, p'_j)$ . Finally,  $\sigma$  sends a CommitShard message back to  $\beta$  (fig. 6.1, step 9). The CommitShard message contains  $c_\sigma$ ,  $E_\sigma$  and  $Q_\sigma$ . Upon receiving  $\sigma$ 's CommitShard message,  $\beta$  verifies  $c_\sigma$  against  $r$  and  $E_\sigma$ , then checks all proofs in  $Q_\sigma$ . Having collected valid CommitShard messages from a quorum of servers  $\sigma_1, \dots, \sigma_{2f+1}$ ,  $\beta$  aggregates all commit shards into a *commit certificate*  $c$ . We underline that each  $\sigma_j$  signed the same root  $r$ , but a potentially different set of exceptions  $E_{\sigma_j}$ . Let  $E$  denote the union of  $E_{\sigma_1}, \dots, E_{\sigma_{2f+1}}$ . We call  $E$  the batch's *exclusion set*. Because a proof of equivocation cannot be produced against a correct client,  $\beta$  knows that all clients identified by  $E$  are necessarily



Byzantine. In particular, because  $\chi$  is correct,  $(i = i_n)$  is guaranteed to not be in  $E$ .

**Committing the batch.** Having collected a commit certificate  $c$  for the batch,  $\beta$  sends  $c$  to all servers by means of a `Commit` message (fig. 6.1, step 10). Upon receiving the `Commit` message,  $\sigma$  verifies  $c$ , computes the exclusion set  $E$ , then delivers every payload  $p_j$  whose id  $i_j$  is not in  $E$ . Recalling that  $c$  is assembled from a quorum of commit shards, at least  $f + 1$  correct servers contributed to  $c$ . This means that, if some  $i_k$  is not in  $E$ , then at least  $f + 1$  correct servers found  $p_k$  not to be equivocated. As in most BRB implementations [31], this guarantees that no two commit certificates can be gathered for equivocating payloads: Draft’s consistency is upheld.

**The role of equivocation proofs.** As the reader might have noticed,  $\beta$  does not attach any proof of equivocation to its `Commit` message. Having received  $\beta$ ’s commit certificate  $c$ ,  $\sigma$  trusts  $\beta$ ’s exclusion set  $E$ , ignoring every payload whose id is in  $E$ . This is not because  $\sigma$  can trust  $\beta$  to uphold validity. On the contrary,  $\sigma$  has *no way* to determine that  $\beta$  is not maliciously excluding the payload of a correct client. Indeed, even if  $\sigma$  were to verify a proof of exclusion for every element in  $E$ , a malicious  $\beta$  could still censor a correct client simply by ignoring its `Submit` message in the first place. Equivocation proofs are fundamental to Draft’s validity not because they *force* malicious brokers to uphold validity, but because they *enable* correct brokers to do the same. Thanks to equivocation proofs, a malicious server cannot trick a correct broker into excluding the payload of a correct client. This is enough to guarantee validity. As we discuss below,  $\chi$  successively submits  $p$  to all brokers until it receives a certificate attesting that  $p$  was delivered by at least one correct server. Because we assume at least one broker to be correct,  $\chi$  is eventually guaranteed to succeed.

**Notifying the clients.** Having delivered every payload whose id is not in the exclusion set  $E$ ,  $\sigma$  produces a *completion shard* for the batch, i.e., a multi-signature  $z_\sigma$  for  $[\text{Completion}, r, E]$ , effectively affirming that  $\sigma$  has delivered all submissions in the batch whose id is not in  $E$ .  $\sigma$  sends  $z_\sigma$  to  $\beta$  by means of a `CompletionShard` message (fig. 6.1, step 11). Upon receiving  $f + 1$  valid `CompletionShard` messages,  $\beta$  assembles all completion shards into a *completion certificate*  $z$ . Finally,  $\beta$  sends a `Completion` message to  $\chi_1, \dots, \chi_b$  (fig. 6.1, step 12). The `Completion` message contains  $z$  and  $E$ . Upon receiving the `Completion` message,  $\chi$  verifies  $z$  against  $E$ , then checks that  $i$  is not in  $E$ . Because at least one correct server contributed a completion shard to  $z$ , at least one correct process delivered all payloads that  $E$  did not exclude, including  $p$ . Having succeeded in broadcasting  $p$ ,  $\chi$  does not need to engage further, and can stop successively submitting  $p$  to all brokers.

**No one is left behind.** As we discussed above, upon receiving the commit certificate  $c$ ,  $\sigma$  delivers every payload in the batch whose id is not in the exclusion set  $E$ . Having gotten at least one correct server to deliver the batch,  $\beta$  is free to disengage, and moves on to assembling

and brokering its next batch. In a moment of asynchrony, however, all communications between  $\beta$  and  $\tilde{\sigma}$  might be arbitrarily delayed. This means that  $\tilde{\sigma}$  has no way of telling whether or not it will eventually receive batch and commit certificate: a malicious  $\beta$  might have deliberately left  $\tilde{\sigma}$  out of the protocol. Server-to-server communication is thus required to guarantee totality. Having delivered the batch,  $\sigma$  waits for an interval of time long enough for all correct servers to deliver batch and commit certificate, *should the network be synchronous and  $\beta$  correct*.  $\sigma$  then sends to all servers an OfferTotality message (fig. 6.1, step 13). The OfferTotality message contains the batch's root  $r$ , and the exclusion set  $E$ . In the good case, upon receiving  $\sigma$ 's OfferTotality message, every server has delivered the batch and ignores the offer. This, however, is not the case for slow  $\hat{\sigma}$ , which replies to  $\sigma$  with an AcceptTotality message (fig. 6.1, step 14). Upon receiving  $\hat{\sigma}$ 's AcceptTotality message,  $\sigma$  sends back to  $\tilde{\sigma}$  a Totality message (fig. 6.1, step 15). The Totality message contains all submissions  $(i_1, p_1), \dots, (i_b, p_b)$ , id assignments for  $i_1, \dots, i_b$ , and the commit certificate  $c$ . Upon delivering  $\sigma$ 's Totality message,  $\tilde{\sigma}$  computes  $r$  from  $(i_1, p_1), \dots, (i_b, p_b)$ , checks  $c$  against  $r$ , computes  $E$  from  $c$ , and delivers every payload  $p_j$  whose id  $i_j$  is not in  $E$ . This guarantees totality and concludes the protocol.

### 6.2.2 Complexity overview

**Directory density.** As we introduced in Section 6.2.1, Draft uses ids assigned by its underlying Directory (DIR) abstraction to identify payload senders. A DIR-assigned id is composed of two parts: a *domain* and an *index*. Domains form a finite set  $\mathbb{D}$  whose size does not increase with the number of clients, indices are natural numbers. Along with safety (e.g., no two processes have the same id) and liveness (e.g., every correct client that requests an id eventually obtains an id), DIR guarantees *density*: the index part of any id is always smaller than the total number of clients  $c$  (i.e., each id index is between 0 and  $(c - 1)$ ). Intuitively, this echoes the (stronger) density guarantee provided by Oracle-CSB, the oracle-based implementation of CSB we introduced in Section 5.1 to bound Draft's performance. In Oracle-CSB, the oracle organizes all clients in a list, effectively labeling each client with an integer between 0 and  $(c - 1)$ . In a setting where consensus cannot be achieved, agreeing on a totally-ordered list of clients is famously impossible: a consensus-less DIR implementation cannot assign ids if  $|\mathbb{D}| = 1$ . As we show in Section 7, however, DIR can be implemented without consensus if servers are used as domains ( $\mathbb{D} = \Sigma$ ). In our DIR implementation Dibs, each server maintains an independent list of public keys. In order to obtain an id, a client  $\chi$  has each server add its public key to its list, then selects a server  $\sigma$  to be its *assigner*. In doing so,  $\chi$  obtains an id  $(\sigma, n)$ , where  $n \in 0..(c - 1)$  is  $\chi$ 's position in  $\sigma$ 's log. In summary, a consensus-less implementation of DIR still guarantees that indices will be smaller than  $c$ , at the cost of a non-trivial domain component for each id. This inflates the size of each individual id by  $\lceil \log_2(|\mathbb{D}|) \rceil$  bits.

**Batching ids.** While DIR-assigned ids come with a non-trivial domain component, the size overhead due to domains vanishes when infinitely many ids are organized into a batch. This

is because domains are constant in the number of clients. Intuitively, as infinitely many ids are batched together, repeated domains become compressible. When building a batch, a Draft broker represents the set  $I$  of sender ids not as a list, but as a map  $\tilde{i}$ . To each domain,  $\tilde{i}$  associates all ids in  $I$  under that domain ( $n$  is in  $\tilde{i}[d]$  iff  $(d, n)$  is in  $I$ ). Because  $\tilde{i}$ 's keys are fixed, as the size of  $I$  goes to infinity, the bits required to represent  $\tilde{i}$ 's keys are completely amortized by those required to represent  $\tilde{i}$ 's values. At the batching limit, the cost of representing each id in  $\tilde{i}$  converges to that of representing its index only,  $\lceil \log_2(c) \rceil$ .

**Protocol cost.** As we discuss in Section 6.4.2, at the batching limit we assume a good-case execution: links are synchronous, all processes are correct, and the set of brokers contains only one element. We additionally assume that infinitely many clients broadcast concurrently. Finally, we assume all servers to already know all broadcasting clients. Let  $\beta$  denote the only broker. As all broadcasting clients submit their payloads to  $\beta$  within a suitably narrow time window,  $\beta$  organizes all submissions into a single batch with root  $r$ . Because links are synchronous and all clients are correct, every broadcasting client submits its multi-signature for  $r$  in time. Having removed all individual signatures from the batch,  $\beta$  is left with a single, aggregated multi-signature  $m$  and an empty straggler set  $S$ .  $\beta$  compresses the sender ids and disseminates the batch to all servers. As  $m$  and  $S$  are constant-sized, the amortized cost for a server to receive each payload  $p$  is  $(\lceil \log_2(c) \rceil + |p|)$  bits. As  $m$  authenticates the entire batch, a server authenticates each payload at an amortized cost of 0 signature verifications. The remainder of the protocol unfolds as a sequence of constant-sized messages: because all broadcasting clients are known to all servers, no server requests any id assignment; witnesses are always constant-sized; and because all processes are correct, no client equivocates and all exception sets are empty. Finally, again by the synchrony of links, all offers of totality are ignored. In summary, at the batching limit a server delivers a payload at an amortized cost of 0 signature verifications and  $(\lceil \log_2(c) \rceil + |p|)$  bits exchanged.

**Latency.** As depicted in Figure 6.1, the latency of Draft is 10 message delays in the synchronous case (fast servers deliver upon receiving the broker's Commit message), and at most 13 message delays in the asynchronous case (slow servers deliver upon receiving other servers' Totality messages). By comparison, the latency of the optimistic reliable broadcast algorithm by Cachin *et al.* [35] is respectively 4 message delays (synchronous case) and 6 message delays (asynchronous case). Effectively, Draft trades oracular efficiency for a constant latency overhead.

**Worst-case complexity.** In the worst case, a Draft server delivers a  $b$ -bits payload by exchanging  $O((\log(c) + b)kn)$  bits, where  $c$ ,  $k$  and  $n$  respectively denote the number of clients, brokers and servers. In brief, the same id, payload and signature is included by each broker in a different batch (hence the  $k$  term) and propagated in an all-to-all fashion (carried by Totality messages) across correct servers (hence the  $n$  term). By comparison, the worst-case

## Chapter 6. Draft

---

communication complexity of Cachin *et al.*'s optimistic reliable broadcast is  $O(ln)$  per server, where  $l$  is the length of the broadcast payload. A direct batched generalization of the same algorithm, however, would raise the worst-case communication to  $O(ln^2)$  per server, similar to that of Draft when  $n \sim k$ . Both batched Bracha and Draft can be optimized by polynomial encoding, reducing their per-server worst-case complexity to  $O(ln)$  and  $O((\log(c) + b)k)$  respectively. Doing so for Draft, however, is beyond the scope of this Part.

### 6.2.3 Pseudocode (Client)

```
1 implements:
2   CSBClient, instance cl
3
4
5 uses:
6   Directory, instance dir
7   PerfectPointToPointLinks, instance pl
8
9
10 parameters:
11   b: Interval // Batching window
12
13
14 struct Submission:
15   message: Message,
16   signature: Signature,
17   submitted_to: {Broker},
18   included_in: {Root}
19
20
21 upon <cl.Init>:
22   trigger <dir.Signup>;
23   await <dir.SignupComplete | id>;
24
25   submissions: {Context: Submission} = {};
26   completed: {Root} = {};
27
28
29 upon <cl.Broadcast | context, message>:
30   signature = sign([Message, context, message]);
31
32   submissions[context] = Submission {
33     message, signature,
```

```

    submitted_to: {}, included_in: {}
  };
33
34  submit(context):
35
36
37  procedure submit(context):
38    if let submission = submissions[context]
39      and let  $\beta$  in ( $B \setminus$  submission.submitted_to):
40      assignment = dir.export(id);
41      trigger <pl.Send |  $\beta$ ,
42        [Submission, assignment, (context, submission.message,
43          submission.signature)]>;
44
45      submission.submitted_to.add( $\beta$ );
46      trigger <timer.Set | [Submit, context], 13 + b>;
47
48
49
50  upon <timer.Ring | [Submit, context]>:
51    submit(context);
52
53
54
55  upon <pl.Deliver |  $\beta$ , [Inclusion, context, root, proof]>:
56    if let submission = submissions[context]:
57      if verify(root, proof, (id, context, submission.message)):
58        submission.included_in.add(root);
59        multisignature = multisign([Reduction, root]);
60        trigger <pl.Send |  $\beta$ , [Reduction, root, multisignature]>;
61
62
63
64  upon <pl.Deliver |  $\beta$ , [Completion, root, exclusions, certificate]>:
65    if verify_plurality(certificate, [Completion, root, exclusions])
66      and id not in exclusions:
67      completed.add(root);
68
69
70
71  upon exists (context, submission) in submissions such that
72    (submission.included_in  $\cap$  completed)  $\neq$  {}:
73    submissions.remove(context);

```

### 6.2.4 Pseudocode (Broker)

```
1 implements:
2     CSBBroker, instance bk
3
4
5 uses:
6     Directory, instance dir
7     PerfectPointToPointLinks, instance pl
8
9
10 parameters:
11     b: Interval // Batching window
12
13
14 struct Submission:
15     context: Context,
16     message: Message,
17     signature: Signature
18
19
20 enum Batch:
21     payloads: {Id: (Context, Message)},
22     signatures: {Id: Signature},
23     reductions: {Id: MultiSignature},
24
25     commit_to: {Server},
26     committable: bool,
27
28     variant Reducing:
29         (empty)
30
31     variant Witnessing:
32         witnesses: {Server: MultiSignature}
33
34     variant Committing:
35         commits: {Server: ({Id}, MultiSignature)}
36
37     variant Completing:
38         exclusions: {Id},
39         completions: {Server: MultiSignature}
40
41
```

```

42 upon <bk.Init>:
43     pending: {Id: [Submission]} (default []) = {};
44     pool: {Id: Submission} = {};
45     collecting: bool = false;
46
47     batches: {Root: Batch} = {};
48
49
50 upon <pl.Deliver |  $\chi$ , [Submission, assignment, (context, message,
    signature)]>:
51     dir.import(assignment);
52
53     if  $\chi$ .verify(signature, [Message, context, message]) and let id =
    dir[ $\chi$ ]:
54         pending[id].push_back(Submission {context, message, signature
    });
55
56
57 upon exists id in pending such that pending[id] != [] and id not in
    pool:
58     submission = pending[id].pop_front();
59     pool[id] = submission;
60
61
62 upon pool != {} and collecting = false:
63     collecting = true;
64     trigger <timer.Set | [Flush], b + 1>;
65
66
67 upon <timer.Ring | [Flush]>:
68     collecting = false;
69
70     submissions = pool;
71     pool = {};
72
73     leaves = [(id, context, message )
        for (id, Submission {context, message, ..}) in submissions];
74
75     tree = merkle_tree(leaves);
76     root = tree.root();
77
78     for (id, Submission {context, message, ..}) in submissions:

```

## Chapter 6. Draft

---

```
79      $\chi$  = dir[id];
80     proof = tree.prove((id, context, message));
81     trigger <pl.Send |  $\chi$ , [Inclusion, context, root, proof]>;
82
83     payloads = {id: (context, message)
84               for (id, Submission {context, message, ..}) in submissions};
85
86     signatures = {id: signature for (id, Submission {signature, ..})
87                 in submissions};
88
89     batches[root] = Reducing {
90         payloads, signatures,
91         reductions: {},
92         commit_to: {},
93         committable: false
94     };
95
96     trigger <timer.Set | [Reduce, root], 2>;
97
98
99     upon <pl.Deliver |  $\chi$ , [Reduction, root, multisignature]>:
100         if let batch alias batches[root] and batch is Reducing
101             and let id = dir[ $\chi$ ] and id in batch.payloads:
102             if  $\chi$ .multiverify(multisignature, [Reduction, root]):
103                 batch.reductions[id] = multisignature;
104                 batch.signatures.remove(id);
105
106
107     upon <timer.Ring | [Reduce, root]>:
108         batch alias batches[root];
109         compressed_ids = compress(batch.payloads.keys());
110         payloads = batch.payloads.values()
111
112         for  $\sigma$  in  $\Sigma$ :
113             trigger <pl.Send | [Batch, compressed_ids, payloads]>;
114
115         batch = Witnessing {witnesses: {}, ..batch};
116         trigger <timer.Set | [Committable, root], 4>;
117
118
119     upon event <timer.Ring | [Committable, root]>:
120         if let batch = batches[root]:
```



```

113     batch.committable = true;
114
115
116 upon <pl.Deliver |  $\sigma$ , [BatchAcquired, root, unknowns]>:
117     if let batch = batches[root]:
118         if exists unknown in unknowns such that unknown not in dir:
119             return;
120
121         assignments = dir.export(unknowns..);
122
123         multisignature = aggregate(batch.reductions.values());
124         signatures = batch.signatures;
125
126         trigger <pl.Send |  $\sigma$ ,
            [Signatures, root, assignments, multisignature, signatures
            ]>;
127
128
129 upon <pl.Deliver |  $\sigma$ , [WitnessShard, root, shard]>:
130     if let batch alias batches[root]:
131         batch.commit_to.add( $\sigma$ );
132
133         if batch is Witnessing:
134             if  $\sigma$ .multiverify(shard, [Witness, root]):
135                 batch.witnesses[ $\sigma$ ] = shard;
136
137
138 upon exists root in batches such that batches[root] is Witnessing
            and |batches[root].witnesses| >= f + 1:
139     batch alias batches[root];
140     certificate = aggregate(batch.witnesses);
141
142     for  $\sigma$  in  $\Sigma$ :
143         trigger <pl.Send | [Witness, root, certificate]>;
144
145     batch = Committing {commits: {}, ..batch};
146
147
148 upon <pl.Deliver |  $\sigma$ , [CommitShard, root, conflicts, commit]>:
149     if let batch alias batches[root] and batch is Committing:
150         if ! $\sigma$ .multiverify(commit, [Commit, root, conflicts.keys()]):
151             return;

```

```
152
153     for (id, (conflict_root, conflict_witness, conflict_proof,
conflict_message))
154         in conflicts:
155             if id not in batch.payloads:
156                 return;
157
158             (context, message) = batch.payloads[id];
159
160             if not verify_plurality(conflict_witness, [Witness,
conflict_root]):
161                 return;
162
163             if not verify(conflict_root, conflict_proof,
(id, context, conflict_message)):
164                 return;
165
166             if message = conflict_message:
167                 return;
168
169             batch.commits[ $\sigma$ ] = (conflicts.keys(), commit);
170
171 upon exists root in batches such that batches[root] is Committing
and batches[root].committable and |batches[root].commits| >= 2f +
1:
172
173     batch alias batches[root];
174     patches: {{Id}: {MultiSignature}} (default {}) = {};
175
176     for (_, (conflicts, commit)) in batch.commits:
177         patches[conflicts].add(commit);
178
179     patches = {ids: aggregate(signatures) for (ids, signatures) in
patches};
180
181     for  $\sigma$  in batch.commit_to:
182         trigger <pl.Send |  $\sigma$ , [Commit, root, patches]>;
183
184     exclusions = union(patches.keys());
185     batch = Completing {exclusions, ..batch};
186
```

```

187
188 upon <pl.Deliver |  $\sigma$ , [CompletionShard, root, completion]>:
189     if let batch alias batches[root] and batch in Completing:
190         if multiverify(completion, [Completion, root, batch.exclusions
191             ]):
192             batch.completions[ $\sigma$ ] = completion
193
194 upon exists root in batches such that batches[root] is Completing
195     and |batches[root].completions| >= f + 1:
196     batch alias batches[root];
197     certificate = aggregate(batch.completions);
198
199     for id in payloads:
200          $\chi$  = dir[id];
201         trigger <pl.Send |  $\chi$ , [Completion, root, batch.exclusions,
202             certificate]>;

```

### 6.2.5 Pseudocode (Server)

```

1 implements:
2     CSBServer, instance sr
3
4
5 uses:
6     Directory, instance dir
7     PerfectPointToPointLinks, instance pl
8
9
10 struct Batch:
11     ids: [Id],
12     payloads: [(Context, Message)],
13     tree: MerkleTree
14
15
16 upon event <sr.Init>:
17     batches: {Root: Batch} = {};
18     witnesses: {Root: Certificate} = {};
19     commits: {(Root, {Id}): {{Id}: Certificate}} = {};
20

```

## Chapter 6. Draft

---

```
21     messages: {(Id, Context): (Message, Root)} = {};
22     delivered: {(Id, Context)} = {};
23
24
25     procedure compress(ids):
26         compressed_ids: {Domain: {Id}} (default {}) = {};
27
28         for (domain, index) in ids:
29             compressed_ids[domain].add(index)
30
31         return compressed_ids
32
33     procedure expand(compressed_ids):
34         ids = {};
35
36         for (domain, indices) in compressed_ids:
37             for index in indices:
38                 ids.push_back((domain, index));
39
40         ids.sort();
41         return ids;
42
43
44     procedure join(ids, payloads):
45         return [(id, context, message)
46                 for (id, (context, message)) in zip(ids, payloads)];
47
48     procedure handle_batch(compressed_ids, payloads):
49         ids = expand(compressed_ids);
50
51         if ids.has_duplicates() or |ids| != |payloads|:
52             return ⊥;
53
54         unknowns = {id for id in ids such that id not in dir};
55
56         leaves = join(ids, payloads);
57         tree = merkle_tree(leaves);
58         root = tree.root();
59
60         batches[root] = Batch {ids, payloads, tree};
61         return [BatchAcquired, root, unknowns];
```

```

62
63
64 upon event <pl.Deliver |  $\pi$ , [Batch, compressed_ids, payloads]>:
65     response = handle_batch(compressed_ids, payloads);
66
67     if response !=  $\perp$  and  $\pi$  in  $B$ :
68         trigger <pl.Send |  $\pi$ , response>;
69
70
71 procedure handle_signatures(root, assignments, multisignature,
    signatures):
72     dir.import(assignments..);
73
74     if root not in batches:
75         return  $\perp$ ;
76
77     Batch {ids, payloads, tree} = batches[root];
78
79     if exists id in ids such that id not in dir:
80         return  $\perp$ ;
81
82     for (id, signature) in signatures:
83         if id not in ids:
84             return  $\perp$ ;
85
86         (context, message) = payloads[ids.index_of(id)];
87
88         if not dir[id].verify(signature, [Message, context, message]):
89
90             return  $\perp$ ;
91
92     multisigners = ids \ signatures.keys();
93
94     if not dir[multisigners..].multiverify(multisignature, [Reduction,
    tree.root()]):
95         return  $\perp$ ;
96
97     shard = multisign([Witness, root]);
98     return [WitnessShard, root, shard];
99
100 upon <pl.Deliver |  $\beta$ , [Signatures, root, assignments, multisignature,

```

```
    signatures]>:
101     response = handle_signatures(root, assignments, multisignature,
    signatures);
102
103     if response != ⊥:
104         trigger <pl.Send | β, response>;
105
106
107 procedure handle_witness(root, certificate):
108     if root not in batches:
109         return ⊥;
110
111     Batch {ids, payloads, tree} = batches[root];
112
113     if !verify_plurality(certificate, [Witness, root]):
114         return ⊥;
115
116     witnesses[root] = certificate;
117     conflicts: {Id: (Root, Certificate, MerkleProof, Message)} = {};
118
119     for (id, context, message) in join(ids, payloads):
120         if let (original_message, original_root) = messages[(id,
    context)]
121             and original_message != message:
122                 original_witness = witnesses[original_root];
123                 original_batch = batches[original_root];
124
125                 Batch {
126                     ids: original_ids,
127                     payloads: original_payloads,
128                     tree: original_tree
129                 } = original_batch;
130
131                 original_leaf = (id, context, original_message);
132
133                 conflicts[id] = (original_root,
134                                 original_witness,
135                                 original_tree.prove(original_leaf),
136                                 original_message
137                             );
138     else:
139         messages[(id, context)] = (message, root);
```

```

131
132     commit = multisign([Commit, root, conflicts.keys()]);
133     return [CommitShard, root, conflicts, commit];
134
135
136 upon <pl.Deliver |  $\beta$ , [Witness, root, certificate]>:
137     response = handle_witness(root, certificate);
138
139     if response !=  $\perp$ :
140         trigger <pl.Send |  $\beta$ , response>;
141
142
143 procedure handle_commit(root, patches):
144     if root not in batches:
145         return  $\perp$ ;
146
147     Batch {ids, payloads, ..} = batches[root];
148
149     if exists id in ids such that id not in dir:
150         return  $\perp$ ;
151
152     signers: {Server} = {};
153
154     for (conflicts, certificate) in patches:
155         if not certificate.verify([Commit, root, conflicts]):
156             return  $\perp$ ;
157
158         signers.extend(certificate.signers());
159
160     if |signers| < 2f + 1:
161         return  $\perp$ ;
162
163     exclusions = union(patches.keys());
164     commits[(root, exclusions)] = patches;
165
166     for (id, context, message) in join(ids, payloads):
167         keycard = dir[id];
168
169         if id not in exclusions
170             and (keycard, context) not in delivered:
171             delivered.add((keycard, context));
172             trigger <sr.Deliver| keycard, context, message>;

```

```
172
173     trigger <timer.Set | [OfferTotality, root, exclusions], 7>;
174
175     shard = multisign([Completion, root, exclusions]);
176     return [CompletionShard, root, shard];
177
178
179 upon <pl.Deliver |  $\pi$ , [Commit, root, patches]>:
180     response = handle_commit(root, patches);
181
182     if response !=  $\perp$  and  $\pi$  in B:
183         trigger <pl.Send |  $\pi$ , response>;
184
185
186 upon <timer.Ring | [OfferTotality, root, exclusions]>:
187     for  $\sigma$  in  $\Sigma$ :
188         trigger <pl.Send |  $\sigma$ , [OfferTotality, root, exclusions]>;
189
190
191 upon <pl.Deliver |  $\sigma$ , [OfferTotality, root, exclusions]>:
192     if (root, exclusions) not in commits:
193         trigger <pl.Send |  $\sigma$ , [AcceptTotality, root, exclusions]>;
194
195
196 upon <pl.Deliver |  $\sigma$ , [AcceptTotality, root, exclusions]>:
197     if let batch = batches[root]
198         and let patches = commits[(root, exclusions)]:
199         Batch {ids, payloads, ..} = batch;
200
201         assignments = dir.export(ids..);
202         compressed_ids = compress(ids);
203
204         trigger <pl.Send |  $\sigma$ , [Totality, root, assignments,
205             (compressed_ids, payloads), patches]>;
206
207
208 upon <pl.Deliver |  $\sigma$ , [Totality, root, assignments, (ids, payloads),
209     patches]>:
210     dir.import(assignments..);
211
212     handle_batch(ids, payloads);
213     handle_commit(root, patches)
```



## 6.3 Correctness

In this section, we prove to the fullest extent of formal detail that Draft implements a Client-Server Byzantine Reliable Broadcast system.

### 6.3.1 No duplication

In this section, we prove that Draft satisfies no duplication.

**Lemma 49.** *Let  $\sigma$  be a correct server, let  $\chi$  be a client, let  $c$  be a context. If  $(\chi, c) \notin \text{delivered}$  at  $\sigma$ , then server did not deliver a message from  $\chi$  for  $c$ .*

*Proof.* Upon initialization,  $\text{delivered}$  is empty at  $\sigma$  (line 22). Moreover,  $\sigma$  adds  $(\chi, c)$  to  $\text{delivered}$  only by executing line 170. Immediately after doing so,  $\sigma$  delivers a message from  $\chi$  for  $c$  (line 171), and because  $\sigma$  never removes elements from  $\text{delivered}$ , the lemma is proved.  $\square$

**Theorem 15.** *Draft satisfies no duplication.*

*Proof.* Let  $\sigma$  be a correct server, let  $\chi$  be a client, let  $c$  be a context.  $\sigma$  delivers a message for  $c$  from  $\chi$  only by executing line 171.  $\sigma$  does so only if  $(\chi, c) \notin \text{delivered}$  (line 169). By Lemma 49, we then have that  $\sigma$  never delivers a message for  $c$  from  $\chi$  more than once, and the theorem is proved.  $\square$

### 6.3.2 Consistency

**Notation 16** (Contexts, messages). We use  $\mathbb{C}$  and  $\mathbb{M}$  to respectively denote the set of contexts and messages broadcast in a Client-Server Byzantine Reliable Broadcast system.

**Definition 57** (Ideal accumulator). d An **ideal accumulator** is a tuple  $(\mathcal{U}, \mathcal{R}, \mathcal{P}, \rho, \zeta, \nu)$  composed by

- An **universe**  $\mathcal{U}$ ;
- A set of **roots**  $\mathcal{R}$ ;
- A set of **proofs**  $\mathcal{P}$ ;
- A **root function**  $\rho : \mathcal{U}^{<\infty} \rightarrow \mathcal{R}$ ;
- A **proof function**  $\zeta : \mathcal{U}^{<\infty} \times \mathbb{N} \rightarrow \mathcal{P}$ ;
- A **verification function**  $\nu : \mathcal{R} \times \mathcal{P} \times \mathbb{N} \times \mathcal{U} \rightarrow \{\text{True}, \text{False}\}$ ;

## Chapter 6. Draft

---

such that

$$\begin{aligned} \forall z \in \mathcal{U}^{<\infty}, \forall n \leq |z|, v(\rho(z), \zeta(z, n), n, z_n) &= \text{True} \\ \forall z \in \mathcal{U}^{<\infty}, \forall n \leq |z|, \forall q \neq z_n, \forall p \in \mathcal{P}, v(\rho(z), p, n, q) &= \text{False} \end{aligned}$$

Let  $r \in \mathcal{R}$ , let  $p \in \mathcal{P}$ , let  $n \in \mathbb{N}$ , let  $x \in \mathcal{U}$  such that  $v(r, p, n, x) = \text{True}$ . We say that  $p$  is a proof for  $x$  from  $r$ .

**Lemma 50.** *Let  $(\mathcal{U}, \mathcal{R}, \mathcal{P}, \rho, \zeta, v)$  be an ideal accumulator. We have that  $\rho$  is injective.*

*Proof.* Let us assume by contradiction that  $z, z' \in \mathcal{U}^{<\infty}$  exist such that  $z \neq z'$  and  $\rho(z) = \rho(z')$ . Let  $n \in \mathbb{N}$  such that  $z_n \neq z'_n$ . By Definition 57 we have

$$\begin{aligned} \text{True} &= v(\rho(z), \zeta(z, n), n, z_n) = \\ &= v(\rho(z'), \zeta(z, n), n, z_n) = \text{False} \end{aligned}$$

which immediately proves the lemma. □

Throughout the remainder of this document, we use Merkle hash-trees as a cryptographic approximation for an ideal accumulator

$$((\mathbb{I} \times \mathbb{C} \times \mathbb{M}), \mathcal{R}, \mathcal{P}, \rho, \zeta, v)$$

on the universe of triplets containing an id, a context and a message. In brief: while a successfully verifiable, corrupted proof can theoretically be forged for a Merkle tree (Merkle tree roots are highly non-injective), doing so is unfeasible for a computationally bounded adversary. An in-depth discussion of Merkle trees is beyond the scope of this Part.

**Lemma 51.** *Let  $\sigma$  be a correct server. Let  $(r, b) \in \text{batches}$  at  $\sigma$ , let*

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

*We have  $|i| = |p|$ .*

*Proof.* We start by noting that, upon initialization,  $\text{batches}$  is empty at  $\sigma$  (line 17). Moreover,  $\sigma$  adds  $\text{Batch}\{ids: i, payloads: p, ..\}$  to  $\text{batches}$  only by executing line 60:  $\sigma$  does so only if  $|i| = |p|$  (lines 51 and 52). □

**Lemma 52.** *Let  $\sigma$  be a correct server. Let  $(r, b) \in \text{batches}$  at  $\sigma$ , let*

$$i = b.ids$$

*The elements of  $i$  are all distinct.*

*Proof.* The proof of this lemma is identical to that of Lemma 51, and we omit it for the sake of brevity.  $\square$

**Definition 58** (Join function). The **join function**

$$j : \{(i, p) \mid i \in \mathbb{I}^{<\infty}, p \in (\mathbb{C} \times \mathbb{M})^{<\infty}, |i| = |p|\} \rightarrow (\mathbb{I} \times \mathbb{C} \times \mathbb{M})^{<\infty}$$

is defined by

$$\begin{aligned} |j(i, p)| &= (|i| = |p|) \\ j(i, p)_n &= (i_n, c_n, m_n), \text{ where } (c_n, m_n) = p_n \end{aligned}$$

**Lemma 53.** *Let  $\sigma$  be a correct server. Let  $(r, b) \in \text{batches}$  at  $\sigma$ , let  $l = j(b.ids, b.payloads)$ . We have  $r = \rho(l)$ .*

*Proof.* We underline that, because by Lemma 51 we have  $|b.ids| = |b.payloads|$ , the definition of  $l$  is well-formed. We start by noting that, upon initialization,  $\text{batches}$  is empty at  $\sigma$  (line 17). Moreover,  $\sigma$  sets

$$\text{batches}[r] = \text{Batch}\{ids: i, payloads: p, ..\}$$

only by executing line 60. Immediately before doing so,  $\sigma$  computes  $r = \rho(l)$  (lines 56, 57 and 58).  $\square$

**Lemma 54.** *Let  $\sigma, \sigma'$  be correct servers. Let  $(r, b) \in \text{batches}$  at  $\sigma$ , let  $(r, b') \in \text{batches}$  at  $\sigma'$ . We have  $b = b'$ .*

*Proof.* Upon initialization,  $\text{batches}$  is empty at both  $\sigma$  and  $\sigma'$ . Moreover,  $\sigma$  and  $\sigma'$  respectively set

$$\begin{aligned} \text{batches}[r] &= \text{Batch}\{ids: i, payloads: p, tree: t\} \\ \text{batches}[r] &= \text{Batch}\{ids: i', payloads: p', tree: t'\} \end{aligned}$$

only by executing line 60. By Lemma 53, we have

$$\rho(j(i, p)) = r = \rho(j(i', p'))$$

which, by Lemma 50 and the injectiveness of  $j$ , proves  $i = i'$  and  $p = p'$ . Because  $\sigma$  (resp.,  $\sigma'$ ) computes  $t$  (resp.  $t'$ ) as a pure function of  $i$  and  $p$  (resp.,  $i' = i$  and  $p' = p$ ) (line 57), we have  $t = t'$ . This proves  $b = b'$  and concludes the lemma.  $\square$

**Lemma 55.** *Let  $\sigma$  be a correct server, let  $i$  be an id. Whenever  $\sigma$  invokes  $\text{dir}[i]$ ,  $\sigma$  knows  $i$ .*

## Chapter 6. Draft

---

*Proof.*  $\sigma$  invokes  $dir[i]$  only by executing lines 88, 93 or 167. If  $\sigma$  invokes  $dir[i]$  by executing lines 88 or 93, then  $i \in I$  (lines 83 and 84, line 91 respectively) for some  $I \subseteq \mathbb{I}$  such that, for all  $i' \in I$ ,  $\sigma$  knows  $i'$  (lines 79 and 80). Similarly, if  $\sigma$  invokes  $dir[i]$  by executing line 167, then  $i \in I$  (line 166) for some  $I \subseteq \mathbb{I}$  such, for all  $i' \in I$ ,  $\sigma$  knows  $i'$  (lines 149 and 150).  $\square$

**Lemma 56.** *Let  $\sigma$  be a correct server. Let  $\chi$  be a client, let  $c$  be a context, let  $m$  be a message such that  $\sigma$  delivers  $(\chi, c, m)$ . We have that  $\sigma$  knows  $\chi$ .*

*Proof.*  $\sigma$  delivers  $(\chi, c, m)$  only by executing line 171. Immediately before doing so,  $\sigma$  invokes  $dir[i]$  for some  $i \in \mathbb{I}$  to obtain  $\chi$  (line 167). By Lemma 55 we then have that  $\sigma$  knows  $i$  which, by Definition 66, proves that  $\sigma$  knows  $k$  as well, and concludes the lemma.  $\square$

**Lemma 57.** *Let  $\sigma$  be a correct server. Let  $\chi$  be a client, let  $c$  be a context, let  $m$  be a message such that  $\sigma$  delivers  $(\chi, c, m)$ . Let  $i = D(\chi)$ . Some  $l \in (\mathbb{I} \times \mathbb{C} \times \mathbb{M})^{<\infty}$ ,  $\sigma_1, \dots, \sigma_{2f+1} \in \Sigma$ ,  $\epsilon_1, \dots, \epsilon_{2f+1} \subseteq \mathbb{I}$  exist such that:*

- For some  $k$ ,  $l_k = (i, c, m)$ ;
- For all  $n$ ,  $i \notin \epsilon_n$ ;
- For all  $n$ ,  $\sigma_n$  signed  $[Commit, \rho(l), \epsilon_n]$ .

*Proof.* We underline the soundness of the lemma's statement. Because  $\sigma$  delivers  $(\chi, ..)$ , by Lemma 56  $\sigma$  knows  $\chi$ , and by Notation 30  $D(\chi)$  is well-defined.

We start by noting that some  $b = batches[r]$  exists (lines 144, 145 and 147) such that, for some  $k$ , we have  $i = b.ids[k]$  and  $(c, m) = b.payloads[k]$  (line 166). Let  $l = j(b.ids, b.payloads)$ . By Definition 58 we immediately have  $l_k = (i, c, m)$ . Moreover, by Lemma 53, we have  $r = \rho(l)$ .

Moreover, some  $S \subseteq \Sigma$  exists such that:

- $|S| \geq 2f + 1$  (lines 160 and 161);
- For all  $\sigma' \in S$ , some  $\epsilon(\sigma')$  exists such that  $\sigma'$  signed  $[Commit, r, \epsilon(\sigma')]$  (lines 152, 155, 156 and 158:  $S$  is obtained by iteratively extending an initially empty set with certificate signers);

Let  $\sigma_1, \dots, \sigma_{2f+1}$  be distinct elements of  $S$ , let  $\epsilon_n = \epsilon(\sigma_n)$ .

Finally, some  $\epsilon \supseteq \epsilon_1 \cup \dots \cup \epsilon_{2f+1}$  exists (lines 154 and 163) such that  $i \notin \epsilon$  (line 169). For all  $n$  we then obviously have  $i \notin \epsilon_n$ .

In summary:  $l_k = (i, c, m)$ ; for all  $n$ ,  $\sigma_n$  signed  $[Commit, \rho(l) = r, \epsilon_n]$ ; for all  $n$ ,  $i \notin \epsilon_n$ . The lemma is therefore proved.  $\square$

**Lemma 58.** *Let  $\sigma$  be a correct server, let  $i$  be an id, let  $c$  be a context, let  $m$  be a message. Let  $t \in \mathbb{R}$  such that, at time  $t$ ,  $messages[(i, c)] = (m, \_)$ . At any time  $t' > t$ , we have  $messages[(i, c)] = (m, \_)$  as well.*

*Proof.* After setting  $messages[(i, c)] = (m, \_)$ ,  $\sigma$  sets  $messages[(i, c)] = (m, \_)$  only by executing line 130. It does so only if  $m = m'$  (line 120).  $\square$

**Lemma 59.** *Let  $\sigma$  be a correct server. Let  $i$  be an id, let  $c$  be a context, let  $m$  be a message. Let  $l \in (\mathbb{I} \times \mathbb{C} \times \mathbb{M})^{<\infty}$ , let  $\epsilon \subseteq \mathbb{I}$  such that:*

- For some  $k$ ,  $l_k = (i, c, m)$ ;
- $i \notin \epsilon$ ;
- $\sigma$  signs  $[Commit, \rho(l), \epsilon]$ .

*We have that  $\sigma$  sets  $messages[(i, c)] = (m, \_)$ .*

*Proof.* Prior to signing  $[Commit, \rho(l), \epsilon]$  (line 132),  $\sigma$  retrieves  $b' = batches[\rho(l)]$  (lines 108 and 109 and 111). Let  $l' = j(b'.ids, b'.payloads)$ . By Lemma 53 we have  $\rho(l') = \rho(l)$  which, by Lemma 50, proves  $l = l'$ .

Subsequently,  $\sigma$  loops through all elements of  $l = l'$  (line 119). For each  $(i_j, c_j, m_j)$  in  $l$ ,  $\sigma$  either adds  $i_j$  to  $\epsilon$  (line 128), or sets  $messages[i_j, c_j] = (m_j, \_)$  (line 130). Because  $\sigma$  does so in particular for  $j = k$ , and  $i \notin \epsilon$ ,  $\sigma$  sets  $messages[(i, c)] = (m, \_)$ , and the lemma is proved.  $\square$

**Theorem 16.** *Draft satisfies consistency.*

*Proof.* Let  $\sigma, \sigma'$  be two correct servers. Let  $\chi$  be a client, let  $c$  be a context, let  $m, m'$  be messages such that  $\sigma$  and  $\sigma'$  deliver  $(\chi, c, m)$  and  $(\chi, c, m')$  respectively. By Lemma 56,  $\sigma$  knows  $\chi$ . Let  $i = D(\chi)$ . By Lemma 57, and noting that at most  $f$  servers are Byzantine, some  $l \in (\mathbb{I} \times \mathbb{C} \times \mathbb{M})^{<\infty}$ ,  $\sigma_1, \dots, \sigma_{f+1} \in \Sigma$ , and  $\epsilon_1, \dots, \epsilon_{f+1}$  exist such that: for some  $k$ ,  $l_k = (i, c, m)$ ; for all  $n$ ,  $\sigma_n$  is correct; for all  $n$ ,  $i \notin \epsilon_n$ ; for all  $n$ ,  $\sigma_n$  signed  $[Commit, \rho(l), \epsilon_n]$ . By Lemmas 59 and 58,  $messages[(i, c)]$  is permanently set to  $(m, \_)$  at  $\sigma_1, \dots, \sigma_{f+1}$ .

By an identical reasoning, some  $\sigma'_1, \dots, \sigma'_{f+1}$  exist such that for all  $n$ ,  $\sigma'_n$  is correct and permanently sets  $messages[(i, c)] = (m', \_)$ . Because  $(\sigma_1, \dots, \sigma_{f+1})$  and  $(\sigma'_1, \dots, \sigma'_{f+1})$  intersect in at least one server, we have  $m = m'$ , and the theorem is proved.  $\square$

### 6.3.3 Totality

**Notation 17** (Ordering). Let  $X$  be a set endowed with a total order relationship. We use  $X' \subset X^{<\infty}$  to denote the set of finite, non-decreasing sequences on  $X$ . Let  $Y \subseteq X$ . We use  $\mathcal{S}(Y)$  to denote **sorted**  $Y$ , i.e., the sequence containing all elements of  $Y$  in ascending order.

**Definition 59** (Id compression and expansion). The **id compression function**

$$c: \mathbb{N}' \rightarrow (\mathbb{D} \rightarrow \mathbb{P}(\mathbb{N}))$$

is defined by

$$n \in c(i)(d) \iff \exists j \mid i_j = (d, n)$$

The **id expansion function**

$$e: (\mathbb{D} \rightarrow \mathbb{P}(\mathbb{N})) \rightarrow \mathbb{N}'$$

is defined by

$$e(f) = \mathcal{S}(\{(d, n) \mid n \in f(d)\})$$

**Lemma 60.** *Let  $i \in \mathbb{N}'$ . We have  $e(c(i)) = i$ .*

*Proof.* It follows immediately from Definition 59. □

**Lemma 61.** *Let  $\sigma, \sigma'$  be correct processes, let  $(r, b) \in \text{batches}$  at  $\sigma$ , let*

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

*Upon evaluating  $\text{handle\_batch}(c(i), p)$ ,  $\sigma'$  sets  $\text{batches}[r] = b$ .*

*Proof.* We start by noting that, upon initialization,  $\text{batches}$  is empty at  $\sigma$  (line 17). Moreover  $\sigma$  adds  $(r, b)$  to  $\text{batches}$  only by executing line 60.  $\sigma$  does so only if  $i$  has no duplicates, and  $|i| = |p|$  (line 51).

Upon invoking  $\text{handle\_batch}(c(i), p)$ ,  $\sigma'$  verifies that  $i = e(c(i))$  (line 49, see Lemma 60) has no duplicates and satisfies  $|i| = |p|$  (line 51). Having done so,  $\sigma'$  sets  $\text{batches}[r] = b'$  for some  $b'$ . However, by Lemma 54 we have  $b = b'$ , and the lemma is proved. □

**Lemma 62.** *Let  $\sigma, \sigma'$  be correct processes, let  $((r, e), \zeta) \in \text{commits}$  at  $\sigma$ , let  $(r, b) \in \text{batches}$  at  $\sigma'$  such that for all  $i \in b.ids$ ,  $\sigma'$  knows  $i$ . Upon evaluating  $\text{handle\_commit}(r, \zeta)$ ,  $\sigma'$  sets  $\text{commits}[(r, e)] = \zeta$ .*

*Proof.* We start by noting that, upon initialization,  $\text{commits}$  is empty at  $\sigma$  (lines 19 and 156). Moreover  $\sigma$  adds  $((r, e), \zeta)$  to  $\text{commits}$  only by executing line 164.  $\sigma$  does so only if, for all

$(c, z) \in \zeta$  (line 154),  $z$  verifies correctly against  $[\text{Commit}, r, c]$  (line 155) and

$$\left| \bigcup_{(c,z) \in \zeta} z.\text{signers}() \right| \geq 2f + 1$$

(lines 160 and 161). We additionally have

$$\epsilon = \bigcup_{(c,_) \in \zeta} c$$

(line 163).

Upon invoking  $\text{handle\_commits}(r, \zeta)$ , by hypothesis  $\sigma'$  verifies that  $r \in \text{batches}$  (line 144) and, for all  $i \in b.\text{ids}$  (line 147),  $i \in \text{dir}$  (line 149). Moreover, because  $\text{verify}$  is a pure procedure,  $\sigma'$  passes the checks at lines 155 and 160. Next,  $\sigma'$  computes

$$\epsilon' = \bigcup_{(c,_) \in \zeta} c = \epsilon$$

(line 163) and sets  $\text{commits}[(r, \epsilon')] = \zeta$  (line 164), concluding the lemma.  $\square$

**Lemma 63.** *Let  $\sigma$  be a correct process, let  $(r, \_) \in \text{commits}$  at  $\sigma$ . We have  $r \in \text{batches}$  at  $\sigma$ .*

*Proof.* Upon initialization,  $\text{commits}$  is empty at  $\sigma$  (line 19). Moreover,  $\sigma$  adds  $(r, \_)$  to  $\text{commits}$  only by executing line 164.  $\sigma$  does so only if  $r \in \text{batches}$  (lines 144 and 145).  $\square$

**Lemma 64.** *Let  $\sigma$  be a correct process, let  $(r, \_) \in \text{commits}$  at  $\sigma$ . Let  $b = \text{batches}[r]$  at server. For all  $i \in b.\text{ids}$ ,  $\sigma$  knows  $i$ .*

*Proof.* We start by noting that, by Lemma 63, the definition of  $b$  is well-formed. Upon initialization,  $\text{commits}$  is empty at  $\sigma$  (line 19). Moreover,  $\sigma$  adds  $(r, \_)$  to  $\text{commits}$  only by executing line 164.  $\sigma$  does so only if, for all  $i$  in  $b.\text{ids}$ ,  $\sigma$  knows  $i$  (lines 147 and 149 and 150).  $\square$

**Lemma 65.** *Let  $\sigma$  be a correct process, let  $(r, \epsilon) \in \text{commits}$  at  $\sigma$ . Let  $b = \text{batches}[r]$  at  $\sigma$ , let*

$$\begin{aligned} i &= b.\text{ids} \\ p &= b.\text{payloads} \end{aligned}$$

*For all  $j$ , let  $(c_j, \_) = p_j$ . For all  $j$  such that  $i_j \notin \epsilon$ ,  $\sigma$  delivered a message from  $D(i_j)$  for  $c_j$ .*

*Proof.* We underline that, by Lemma 64,  $D(i_j)$  is well-defined for all  $j$ . Upon initialization,  $\text{commits}$  is empty at  $\sigma$  (line 19). Moreover,  $\sigma$  adds  $(r, \epsilon)$  to  $\text{commits}$  only by executing line 164. Upon doing so,  $\sigma$  retrieves  $b$  from  $\text{batches}$  (line 147) and delivers a message from  $D(i_j)$  for  $c_j$  (lines 167 and 171) for all  $j$  such that:  $i_j \notin \epsilon$ ; and  $\sigma$  did not previously deliver a message from  $D(i_j)$  for  $c_j$  (line 169, see Lemma 49).  $\square$

**Theorem 17.** *Draft satisfies totality.*

*Proof.* Let  $\sigma, \sigma'$  be correct servers. Let  $\chi$  be a client, let  $c$  be a context such that  $\sigma$  delivers a message from  $\chi$  for  $c$ .  $\sigma$  delivers a message from  $\chi$  for  $c$  only by executing line 171. Immediately before doing so,  $\sigma$  sets  $commits[r, \epsilon] = \zeta$  for some  $r, \epsilon$  and  $\zeta$  such that  $D(\chi) \notin \epsilon$  (lines 164, 167 and 169). Let  $b = batches[r]$  at  $\sigma$  (by Lemma 63,  $r \in batches$  at  $\sigma$ ), let

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

For all  $j$ , let  $(c_j, \_) = p_j$ . For some  $k$  we have  $i_k = D(\chi)$  and  $c_k = c$  (lines 147, 166 and 167).

Immediately after delivering a message from  $\chi$  for  $c$ ,  $\sigma$  sets a timer for  $[OfferTotality, r, \epsilon]$  (line 173). When the timer eventually rings (line 186),  $\sigma$  sends an  $[OfferTotality, r, \epsilon]$  message to all servers, including  $\sigma'$  (lines 187 and 188). Upon eventually delivering  $[OfferTotality, r, \epsilon]$  (line 191),  $\sigma'$  checks if  $(r, \epsilon) \in commits$  (line 192). Let us assume that  $(r, \epsilon) \in commits$  at  $\sigma'$ . Let  $b' = batches[r]$  at  $\sigma'$  (by Lemma 63,  $r \in batches$  at  $\sigma'$ ). By Lemma 54, we have  $b' = b$ . As a result, by Lemma 65,  $\sigma'$  already delivered a message from  $D(i_k) = \chi$  for  $c_k = c$ . Throughout the remainder of this proof, we assume  $(r, \epsilon) \notin commits$  at  $\sigma'$ .

Upon verifying that  $(r, \epsilon) \notin commits$  (line 192),  $\sigma'$  sends back to  $\sigma$  an  $[AcceptTotality, r, \epsilon]$  message (line 193). Upon delivering  $[AcceptTotality, r, \epsilon]$  (line 196),  $\sigma$  verifies that  $r \in batches$  and  $(r, \epsilon) \in commits$  (line 197), exports an array  $a$  of assignments for all elements of  $i$  (line 200) (by Lemma 64  $\sigma$  knows all elements of  $i$ ), and sends  $[Totality, r, a, (c(i), p), \zeta]$  to  $\sigma'$  (line 203)

Upon delivering  $[Totality, r, a, (c(i), p), \zeta]$  (line 206),  $\sigma'$  imports all elements of  $a$  (line 207). Having done so,  $\sigma'$  knows all elements of  $i$ . Then,  $\sigma'$  invokes  $handle\_batch(c(i), p)$  (line 209). By Lemma 61, having done so  $\sigma'$  sets  $batches[r] = b$ . Finally,  $\sigma'$  invokes  $handle\_commit(r, \zeta)$  (line 210). By Lemma 62, having done so  $\sigma'$  sets  $commits[(r, \epsilon)] = \zeta$ . By Lemma 65,  $\sigma'$  delivers a message from  $D(\chi) = i_k$  for  $c = c_k$ .

In summary, if  $\sigma$  delivers a message from  $\chi$  for  $c$ , then  $\sigma'$  eventually delivers a message from  $\chi$  for  $c$  as well, and the theorem is proved.  $\square$

### 6.3.4 Integrity

In this section, we prove that Draft satisfies integrity.

**Lemma 66.** *Let  $\sigma$  be a correct server. Let  $\chi$  be a client, let  $c$  be a context, let  $m$  be a message such that  $\sigma$  delivers  $m$  from  $\chi$  for  $c$ . Some correct server  $\sigma'$  and some  $(r, b) \in batches$  at  $\sigma'$  exists such*



that  $\sigma'$  signed  $[\text{Witness}, r]$  and, with

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

some  $k$  exists such that  $i_k = D(\chi)$  and  $p_k = (c, m)$ .

*Proof.* We start by noting that  $\sigma$  delivers  $m$  from  $\chi$  for  $c$  only by executing line 171. When  $\sigma$  does so, some  $(r, b) \in \text{batches}$  exists at  $\sigma$  (line 147) such that, with

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

some  $k$  exists such that  $i_k = D(\chi)$  and  $p_k = (c, m)$  (lines 166 and 167). Moreover, a set  $S \subseteq \Sigma$  exists such that  $|S| \geq 2f + 1$  (lines 160 and 161) and, for all  $\tilde{\sigma} \in S$ ,  $\tilde{\sigma}$  signed a  $[\text{Commit}, r, \_]$  message (lines 152, 155, 156 and 158). Noting that at most  $f$  processes are Byzantine, at least one element  $\sigma^*$  of  $S$  is correct.

$\sigma^*$  signs  $[\text{Commit}, r, \_]$  only upon executing line 132.  $\sigma^*$  does so only if at least  $f + 1$  servers signed  $[\text{Witness}, r]$  (lines 113 and 114). Noting that at most  $f$  processes are Byzantine, at least one correct server  $\sigma'$  signed  $[\text{Witness}, r]$ .

$\sigma'$  signs  $[\text{Witness}, r]$  only by executing line 96. It does so only if some  $(r, b') \in \text{batches}$  exists at  $\sigma'$  (lines 74, 75). By Lemma 54, we have  $b' = b$ .

In summary, some correct server  $\sigma'$  exists such that  $\sigma'$  signed  $[\text{Witness}, r]$  and  $(r, b) \in \text{batches}$  at  $\sigma'$ , with  $i_k = D(\chi)$  and  $p_k = (c, m)$ : the lemma is proved.  $\square$

**Lemma 67.** *Let  $\chi$  be a correct client. Let  $(c, s) \in \text{submissions}$  at  $\chi$ , let*

$$m = s.message$$

*We have that  $\chi$  broadcast  $m$  for  $c$ .*

*Proof.* Upon initialization,  $\text{submissions}$  is empty at  $\chi$  (line 25);  $\chi$  adds  $(c, s)$  to  $\text{submissions}$  (line 32) only upon broadcasting  $m$  for  $c$  (line 29).  $\square$

**Lemma 68.** *Let  $\sigma$  be a correct server, let  $l \in (\mathbb{I} \times \mathbb{C} \times \mathbb{M})'$  such that  $\sigma$  signs  $[\text{Witness}, \rho(l)]$ . For all  $j$ , let  $(i_j, c_j, m_j) = l_j$ , let  $\chi_j = D(i_j)$ . For all  $k$  such that  $\chi_k$  is correct,  $\chi_k$  broadcast  $m_k$  for  $c_k$ .*

*Proof.* We start by noting that  $\sigma$  signs  $[\text{Witness}, \rho(l)]$  only by executing line 96. Immediately before doing so,  $\sigma$  loads a batch  $b = \text{batches}[\rho(l)]$  (line 77). By Lemmas 53 and 50, we have

$$l = j(b.ids, b.payloads)$$

## Chapter 6. Draft

---

Moreover, by Lemma 52, all elements of  $i$  are distinct.

Next,  $\sigma$  uses a set  $H \subseteq \mathbb{I}$  to partition the elements of  $i$ . For each  $h \in H$  (line 82),  $\sigma$  verifies that: for some  $j$ , we have  $h = i_j$  (lines 83 and 89); and  $\chi_j$  signed  $[\text{Message}, c_j, m_j]$  (lines 86, 88 and 89). Finally, for all  $j$  such that  $i_j \notin H$  (line 91),  $\sigma$  verifies that  $\chi_j$  signed  $[\text{Reduction}, \rho(l)]$  (lines 93 and 94).

In summary, for all  $j$ , we have that either  $\chi_j$  signed  $[\text{Message}, c_j, m_j]$ , or  $\chi_j$  signed  $[\text{Reduction}, \rho(l)]$ .

Let  $k$  such that  $\chi_k$  is correct. Let us assume that  $\chi_k$  signed  $[\text{Message}, c_k, m_k]$ .  $\chi_k$  signs  $[\text{Message}, c_k, m_k]$  only by executing line 30. It does so only upon broadcasting  $m_k$  for  $c_k$  (line 29). Throughout the remainder of this proof, we assume that  $\chi_k$  signed  $[\text{Reduction}, \rho(l)]$ .

$\chi_k$  signs  $[\text{Reduction}, \rho(l)]$  only by executing line 55. It does so only if, for some  $p \in \mathcal{P}$ ,  $n \in \mathbb{N}$ ,  $c' \in \mathbb{C}$ ,  $m' \in \mathbb{M}$ , we have

$$v(\rho(l), p, n, (i_k, c', m')) = \text{True} \quad (6.1)$$

(line 52), and  $\chi_k$  previously broadcast  $m'$  for  $c'$  (line 51, see Lemma 67). By Equation 6.1 and Definition 57, we must have  $l_n = (i_k, \dots)$ . As a result, because  $i$  is non-repeating, we immediately have  $n = k$ ,  $c' = c_k$  and  $m' = m_k$ . We then have that  $\chi_k$  previously broadcast  $m_k$  for  $c_k$ , and the lemma is proved.  $\square$

**Theorem 18.** *Draft satisfies integrity.*

*Proof.* Let  $\sigma$  be a correct server. Let  $\chi$  be a correct client, let  $c$  be a context, let  $m$  be a message such that  $\sigma$  delivers  $m$  from  $\chi$  for  $c$ . By Lemma 66, some correct server  $\sigma'$  and some  $(r, b) \in \text{batches}$  at  $\sigma'$  exist such that  $\sigma'$  signed  $[\text{Witness}, r]$  and, with

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

some  $k$  exists such that  $i_k = D(\chi)$  and  $p_k = (c, m)$ . Let  $l = j(i, p)$ . We trivially have  $l_k = (D(\chi), c, m)$ . By Lemma 53 we have  $r = \rho(l)$ . For all  $j$ , let  $(c_j, m_j) = p_j$ . By Lemma 68, we have that  $(\chi = D(D(\chi))) = D(i_k)$  broadcast  $(m = m_k)$  for  $(c = c_k)$ , and the theorem is proved.  $\square$

### 6.3.5 Validity

In this section, we prove that Draft satisfies validity.

**Lemma 69.** *Let  $\chi$  be a correct client. Upon initialization,  $\chi$  knows  $\chi$ .*

*Proof.* It follows immediately from lines 22 and 23, and the signup validity of Directory.  $\square$

**Lemma 70.** *Let  $\chi$  be a correct client, let  $c$  be a context such that  $\chi$  broadcast a message for  $c$ . If  $c \notin \text{submissions}$  at  $\chi$ , then at least one correct server delivered a message from  $\chi$  for  $c$ .*

*Proof.* Let  $i = D(\chi)$  (by Lemma 69,  $i$  is well-defined). Let  $m$  be the message  $\chi$  broadcast for  $c$ . We start by noting that  $\chi$  adds  $c$  to  $\text{submissions}$  (line 32) upon broadcasting a message for  $c$  (line 29). As a result,  $\chi$  satisfies  $c \notin \text{submissions}$  only upon removing  $(c, s)$  from  $\text{submissions}$  (line 64 only).  $\chi$  does so only if

$$s.\text{included\_in} \cap \text{completed} \neq \emptyset$$

Let  $r \in s.\text{included\_in} \cap \text{completed}$ .

Noting that  $s.\text{included\_in}$  is initially empty at  $\chi$  (line 32),  $\chi$  must have added  $r$  to  $s.\text{included\_in}$  (line 53 only). Noting that  $\chi$  permanently sets  $s.\text{message}$  to  $m$  upon initialization (line 32 only, see line 29),  $\chi$  adds  $r$  to  $s.\text{included\_in}$  (line 53) only if, for some  $p \in \mathcal{P}$  and  $n \in \mathbb{N}$ , we have

$$v(r, p, n, (i, c, m)) = \text{True} \tag{6.2}$$

(line 52).

Noting that  $\text{completed}$  is initially empty at  $\chi$  (line 26),  $\chi$  must have added  $r$  to  $\text{completed}$  (line 60 only).  $\chi$  does so only if some  $\epsilon$  exists such that  $i \notin \epsilon$  and at least  $f + 1$  servers signed  $[\text{Completion}, r, \epsilon]$  (line 59).

Noting that at most  $f$  servers are Byzantine, some correct server  $\sigma$  exists such that  $\sigma$  signed  $[\text{Completion}, r, \epsilon]$  (line 175 only).  $\sigma$  does so only if  $r \in \text{batches}$  (lines 144 and 145). Let  $b = \text{batches}[r]$  at  $\sigma$ , let  $l = j(b.\text{ids}, b.\text{payloads})$ , let  $(i_j, c_j, m_j) = l_j$ , let  $\chi_j = D(i_j)$ . By Lemma 53, we have  $r = \rho(l)$ , and by Equation 6.2 and Definition 57 we have  $l_n = (i, c, m)$ .

Immediately before signing  $[\text{Completion}, r, \epsilon]$ ,  $\sigma$  loops through  $l$  (line 166) and delivers  $m_j$  from  $\chi_j$  for  $c_j$  (lines 167 and 171) if and only if  $i_j \notin \epsilon$  and  $\sigma$  did not previously deliver a message from  $\chi_j$  for  $c_j$  (line 169, see Lemma 49). Recalling that  $(i = i_n) \notin \epsilon$ ,  $\sigma$  delivers a message from  $(\chi = \chi_j)$  for  $(c = c_j)$  before or upon signing  $[\text{Completion}, r, \epsilon]$ , and the lemma is proved.  $\square$

**Lemma 71.** *Let  $\chi$  be a correct client, let  $c$  be a context, let  $m$  be a message such that  $\chi$  broadcast  $m$  for  $c$ . Eventually, either a correct server delivers a message from  $\chi$  for  $c$ , or a correct broker delivers a  $[\text{Submission}, a, (c, m, s)]$  message, where  $a$  is an assignment for  $\chi$  and  $s$  is  $\chi$ 's signature for  $[\text{Message}, c, m]$ .*

*Proof.* Upon broadcasting  $m$  for  $c$  (line 29),  $\chi$  sets

$$\text{submissions}[c] = \text{Submission}\{\text{message}: m, \text{signature}: s, \text{submitted\_to}: \emptyset, \dots\}$$

## Chapter 6. Draft

---

where  $s$  is a signature for  $[\text{Message}, c, m]$  (lines 30 and 32).  $\chi$  never updates  $\text{submissions}[c].\text{message}$  or  $\text{submissions}[c].\text{signature}$ . Subsequently,  $\chi$  invokes  $\text{submit}(c)$  (line 34).

Upon executing  $\text{submit}(c)$  (line 37),  $\chi$  immediately returns if and only if  $c \notin \text{submissions}$ , or no  $\beta \in (B \setminus \text{submissions}[c].\text{submitted\_to})$  exists (line 38). Otherwise,  $\chi$  sends a  $[\text{Submission}, a, (c, m, s)]$  message to  $\beta$  (where  $a$  is an assignment for  $\chi$ ) (lines 39 and 40), then adds  $\beta$  to  $\text{submissions}[c].\text{submitted\_to}$  (line 42), and finally schedules  $\text{submit}(c)$  for eventual re-execution (lines 43, 46 and 47).

Noting that  $\chi$  updates  $\text{submissions}[c].\text{submitted\_to}$  only by executing line 42,  $\chi$  keeps re-executing  $\text{submit}(c)$  until either

- $c \notin \text{submissions}$ : by Lemma 70, at least one correct server delivered a message from  $\chi$  for  $c$ ; or
- $\chi$  sent to all brokers a  $[\text{Submission}, a, (c, m, s)]$  message, where  $a$  is an assignment for  $\chi$

The lemma follows immediately from the observation that  $B$  is finite and the assumption that at least one broker is correct.  $\square$

**Lemma 72.** *Let  $\beta$  be a correct broker. If  $\text{pool} \neq \{\}$  at  $\beta$ , then eventually  $\text{pool} = \{\}$  at  $\beta$ .*

*Proof.* We start by noting that the  $[\text{Flush}]$  timer is pending at  $\beta$  if and only if  $\text{collecting} = \text{True}$  at  $\beta$ . Indeed, upon initialization, we have  $\text{collecting} = \text{False}$  at  $\beta$ . Moreover,  $\beta$  sets  $\text{collecting} = \text{True}$  (line 63) only upon setting the  $[\text{Flush}]$  timer (line 64 only). Finally,  $\beta$  resets  $\text{collecting} = \text{False}$  (line 68 only) only upon ringing  $[\text{Flush}]$  (line 67).

Upon ringing  $[\text{Flush}]$  (line 67),  $\beta$  resets  $\text{pool} = \{\}$  (line 71). As a result, whenever  $\text{pool} \neq \{\}$  at  $\beta$ , we either have:

- $\text{collecting} = \text{True}$  at  $\beta$ :  $[\text{Flush}]$  is pending at  $\beta$  and  $\beta$  will eventually reset  $\text{pool} = \{\}$ ; or
- $\text{collecting} = \text{False}$  at  $\beta$ :  $\beta$  will eventually detect that  $\text{pool} \neq \{\}$  and  $\text{collecting} = \text{False}$  (line 62) and set  $\text{collecting} = \text{true}$  - the above will then apply.

$\square$

**Lemma 73.** *Let  $\beta$  be a correct broker, let  $i$  be an id such that*

$$\text{pending}[i] = [s_1, \dots, s_S]$$

at  $\beta$ , with  $S \geq 1$ . We eventually have

$$\begin{aligned} pool[i] &= s_1 \\ pending[i] &= [s_2, \dots, s_S, \dots] \end{aligned}$$

at  $\beta$ .

*Proof.* We start by noting that, by Lemma 72, if  $i \in pool$  then eventually  $pool = \{\}$  at  $\beta$ . Noting that  $\beta$  adds elements to  $pool$  only by executing line 59, we then have that eventually  $\beta$  detects  $pending[i] \neq []$  and  $i \notin pool$  (line 57). Upon doing so,  $\beta$  sets  $pending[i] = [s_2, \dots, s_S, \dots]$  (line 58) and  $pool[i] = s_1$  (line 59).  $\square$

**Lemma 74.** *Let  $\beta$  be a correct broker, let  $i$  be an id, let  $s$  be a submission such that  $\beta$  pushes  $s$  to  $pending[i]$ . We eventually have  $pool[i] = s$  at  $\beta$ .*

*Proof.* It follows immediately by induction on Lemma 73, and the observation that  $\beta$  only pushes elements to the back of  $pending$ .  $\square$

**Lemma 75.** *Let  $\beta$  be a correct broker. Let  $\chi$  be a client, let  $a$  be an assignment for  $\chi$ , let  $c$  be a context, let  $m$  be a message, let  $s$  be  $\chi$ 's signature for  $[Message, c, m]$ . If  $\beta$  delivers a  $[Submission, a, (c, m, s)]$ , then  $\beta$  eventually sets*

$$pool[i] = Submission\{context : c, message : m, ..\}$$

where  $i = D(\chi)$ .

*Proof.* Upon delivering  $[Submission, a, (c, m, s)]$  (line 50),  $\beta$  imports  $a$  (line 51). Because  $s$  successfully verifies against  $[Message, c, m]$  and  $\beta$  knows  $\chi$  (line 53),  $\beta$  pushes  $Submission\{context : c, message : m, ..\}$  to  $pending[i]$  (line 54). The lemma immediately follows from Lemma 74.  $\square$

**Definition 60** (Map). Let  $X, Y$  be sets. A **map**  $m : X \rightarrow Y$  is a subset of  $(X \times Y)$  such that

$$\forall (x, y), (x', y') \neq (x, y) \in m, x \neq x'$$

We use

$$\begin{aligned} \mathcal{D}(m) &= \{x \in X \mid (x, \_) \in m\} \\ \mathcal{C}(m) &= \{y \in Y \mid (\_, y) \in m\} \end{aligned}$$

to respectively denote the **domain** and **codomain** of  $m$ .

**Notation 18** (Sorting maps). Let  $X, Y$  be sets such that  $X$  is endowed with a total order relationship, let  $m : X \rightarrow Y$ . With a slight abuse of notation, we use  $\mathcal{S}(\mathcal{C}(m))$  to denote the

codomain of  $m$ , sorted by  $m$ 's domain. For example, let  $m = \{(1, \square), (2, \circ), (3, \Delta)\}$ , we have  $\mathcal{S}(\mathcal{C}(m)) = (\square, \circ, \Delta)$ .

**Definition 61** (Broker variables). Let  $\beta$  be a correct broker, let  $t \in \mathbb{R}$ , let  $r$  be a root.

- We use  $P_\beta(t)[r] : \mathbb{I} \rightarrow (\mathbb{M} \times \mathbb{C})$  to denote, if it exists, the value of  $batches[r].payloads$  at  $\beta$  at time  $t$ ; we use  $P_\beta(t)[r] = \perp$  otherwise.
- We define  $C_\beta(t)[r] : \mathbb{I} \rightarrow \mathbb{C}$  and  $M_\beta(t)[r] : \mathbb{I} \rightarrow \mathbb{M}$  by

$$(C_\beta(t)[r][i], M_\beta(t)[r][i]) = P_\beta(t)[r][i]$$

if  $P_\beta(t)[r][i] \neq \perp$ ; we use  $C_\beta(t)[r][i] = M_\beta(t)[r][i] = \perp$  otherwise.

- We define  $L_\beta(t)[r] \in (\mathbb{I} \times \mathbb{C} \times \mathbb{M})^\sphericalangle$  by

$$L_\beta(t)[r] = \mathcal{S}(\{(i, c, m) \mid (i, (c, m)) \in P_\beta(t)[r]\})$$

if  $P_\beta(t)[r] \neq \perp$ ; we use  $L_\beta(t)[r] = \perp$  otherwise.

- We use  $S_\beta(t)[r] : \mathbb{I} \rightarrow \mathbb{S}^1$  and  $Q_\beta(t)[r] : \mathbb{I} \rightarrow \mathbb{S}^+$  to respectively denote, if they exist, the values of  $batches[r].signatures$  and  $batches[r].reductions$  at  $\beta$  at time  $t$ ; we use  $S_\beta(t)[r] = Q_\beta(t)[r] = \perp$  otherwise.
- We use  $Y_\beta(t)[r] \subseteq \Sigma$  and  $Z_\beta(t)[r] \in \{\text{True}, \text{False}\}$  to respectively denote, if they exist, the values of  $batches[r].commit\_to$  and  $batches[r].committable$  at  $\beta$  at time  $t$ ; we use  $Y_\beta(t)[r] = Z_\beta(t)[r] = \perp$  otherwise.
- We use  $W_\beta(t)[r] : \Sigma \rightarrow \mathbb{S}^+$  to denote, if it exists, the value of  $batches[r].witnesses$  at  $\beta$  at time  $t$ ; we use  $W_\beta(t)[r] = \perp$  otherwise. We underline that  $batches[r].witnesses$  exists only if  $batches[r]$  exists and is Witnessing.
- We use  $X_\beta(t)[r] : \Sigma \rightarrow (\mathbb{P}(\mathbb{I}) \times \mathbb{S}^+)$  to denote, if it exists, the value of  $batches[r].commits$  at  $\beta$  at time  $t$ ; we use  $X_\beta(t)[r] = \perp$  otherwise. We underline that  $batches[r].commits$  exists only if  $batches[r]$  exists and is Committing.
- We use  $E_\beta(t)[r] \in \mathbb{P}(\mathbb{I})$  and  $T_\beta(t)[r] : \Sigma \rightarrow \mathbb{S}^+$  to respectively denote, if they exist, the values of  $batches[r].exclusions$  and  $batches[r].completions$  at  $\beta$  at time  $t$ ; we use  $E_\beta(t)[r] = T_\beta(t)[r] = \perp$  otherwise. We underline that  $batches[r].exclusions$  and  $batches[r].completions$  exist only if  $batches[r]$  exists and is Completing.

**Notation 19** (Broker variables). Let  $r$  be a root. Wherever it can be unequivocally inferred from context, we omit the correct broker and the time from  $P[r]$ ,  $C[r]$ ,  $M[r]$ ,  $L[r]$ ,  $S[r]$ ,  $Q[r]$ ,  $Y[r]$ ,  $Z[r]$ ,  $W[r]$ ,  $X[r]$ ,  $E[r]$  and  $T[r]$ .

**Lemma 76.** *Let  $\beta$  be a correct broker, let  $r$  be a root. We have*

$$\begin{aligned} & (P[r] = \perp) \iff (C[r] = \perp) \iff (M[r] = \perp) \iff \\ & \iff (Y[r] = \perp) \iff (Z[r] = \perp) \iff (L[r] = \perp) \iff \\ & \iff (S[r] = \perp) \iff (Q[r] = \perp) \end{aligned}$$

and

$$(P[r] = \perp) \implies (W[r] = X[r] = E[r] = T[r] = \perp)$$

*Proof.* It follows immediately from Definition 61.  $\square$

**Lemma 77.** *Let  $\beta$  be a correct broker, let  $i$  be an id, let  $c$  be a context, let  $m$  be a message such that*

$$pool[i] = Submission\{context : c, message : m, \dots\}$$

at  $\beta$ . Eventually some root  $r$  exists such that  $P[r][i] = (c, m)$ .

*Proof.* By Lemma 72, eventually  $pool = \{\}$  at  $\beta$ . Moreover,  $\beta$  resets  $pool = \{\}$  only by executing line 71. Upon first doing so,  $\beta$  stores the original value of  $pool$  in a variable  $u$  (line 70), which it uses to define a variable  $p$  by

$$((i, (c, m)) \in p) \iff ((i, Submission\{c, m, \dots\}) \in u)$$

(line 83). We immediately have  $p[i] = (c, m)$ . Finally, for some  $r$ ,  $\beta$  adds  $(r, b)$  to  $batches$ , with  $b.payloads = p$  (line 87). The lemma follows immediately from Definition 61.  $\square$

**Lemma 78.** *Let  $\chi$  be a correct client, let  $i = D(\chi)$ . Let  $c$  be a context, let  $m$  be a message such that  $\chi$  broadcast  $m$  for  $c$ . Eventually, either a correct server delivers a message from  $\chi$  for  $c$ , or some correct broker  $\beta$  and some root  $r$  exist such that eventually  $P[r][i] = (c, m)$ .*

*Proof.* The lemma immediately follows from Lemmas 71, 75 and 77.  $\square$

**Lemma 79.** *Let  $\beta$  be a correct broker, let  $i$  be an id. If  $i \in pool$  at  $\beta$ , then  $\beta$  knows  $i$ .*

*Proof.* We start by noting that, upon initialization,  $pool$  is empty at  $\beta$  (line 44). Moreover,  $\beta$  adds  $(i, \_)$  to  $pool$  only by executing line 59.  $\beta$  does so only if  $pending[i]$  is not empty (line 57). Upon initialization, all values of  $pending$  are also empty at  $\beta$  (line 43). Finally,  $\beta$  adds elements to  $pending[i]$  only by executing line 54. Because  $\beta$  does so only if  $\beta$  knows  $i$  (line 53), the lemma is proved.  $\square$

**Lemma 80.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $P[r] \neq \perp$ . For all  $i \in \mathcal{D}(P[r])$ ,  $\beta$  knows  $i$ .*

*Proof.* We start by noting that, upon initialization,  $batches$  is empty at  $\beta$  (line 47). Moreover,  $\beta$  adds  $(r, b)$  to  $batches$  only by executing line 87. Upon doing so,  $\beta$  satisfies  $\mathcal{D}(b.payloads) = \mathcal{D}(u)$  (line 83) for some  $u$  initialized to  $pool$  (line 70). The lemma immediately follows from Lemma 79.  $\square$

**Lemma 81.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $(P[r] \neq \perp) \iff (S[r] \neq \perp)$ . We have  $\mathcal{D}(P[r]) = \mathcal{D}(S[r]) \cup \mathcal{D}(Q[r])$  and  $\mathcal{D}(S[r]) \cap \mathcal{D}(Q[r]) = \emptyset$*

*Proof.* We start by noting that, upon initialization,  $batches$  is empty at  $\beta$  (line 47). Moreover,  $\beta$  adds  $(r, b)$  to  $batches$  only by executing line 87. Upon doing so,  $\beta$  satisfies  $\mathcal{D}(b.payloads) = \mathcal{D}(b.signatures)$  (lines 83 and 85). Subsequently,  $\beta$  updates  $batches[r].signatures$  and  $batches[r].reductions$  only concurrently, by executing lines 95 and 96. Upon doing so,  $\beta$  shifts an id  $i$  in  $batches[r].payloads$  (line 93) from  $\mathcal{D}(batches[r].signatures)$  to  $\mathcal{D}(batches[r].reductions)$ .  $\square$

**Lemma 82.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $S[r] \neq \perp$ . Let  $i \in \mathcal{D}(S[r])$ , let  $\chi = D(i)$ , let  $(c, m) = P[r][i]$ , let  $s = S[r][i]$ . We have that  $s$  is  $\chi$ 's signature for  $[Message, c, m]$ .*

*Proof.* We underline that  $c$  and  $m$  are well-defined by Lemma 81, and  $i$  is well-defined by Lemmas 81 and 80. We start by noting that, upon initialization,  $batches$  is empty at  $\beta$  (line 47). Moreover,  $\beta$  adds  $(r, b)$  to  $batches$  only by executing line 87. Upon last doing so,  $\beta$  satisfies

$$(c, m) = b.payloads[i] = (w.context, w.message)$$

$$s = b.signatures[i] = w.signature$$

with  $w = u[i]$  for some copy  $u$  of  $pool$  (lines 70, 83 and 85).

Upon initialization,  $pool$  is empty at  $\beta$  as well (line 44). Moreover,  $\beta$  sets  $pool[i] = w$  only by executing line 59. It does so only if  $w$  was in  $pending[i]$  (lines 57 and 58). Finally,  $pending$  is initially empty at  $\beta$  (line 43), and  $\beta$  adds  $w$  to  $pending[i]$  only by executing line 54.  $\beta$  does so only if  $w.signature$  is  $D(i)$ 's signature for  $[Message, w.context, w.message]$  (line 53). This proves that  $s$  is  $\chi$ 's signature for  $[Message, c, m]$  and concludes the lemma.  $\square$

**Lemma 83.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $Q[r] \neq \perp$ . Let  $i \in \mathcal{D}(Q[r])$ , let  $\chi = D(i)$ , let  $q = Q[r][i]$ . We have that  $q$  is  $\chi$ 's multisignature for  $[Reduction, r]$ .*

*Proof.* We start by noting that, upon initialization,  $batches$  is empty at  $\beta$  (line 47). Moreover,  $\beta$  adds  $(r, b)$  to  $batches$  only by executing line 87. Upon doing so,  $\beta$  satisfies  $b.reductions = \{\}$ . Finally,  $\beta$  adds  $(i, q)$  to  $batches[r].reductions$  only by executing line 95.  $\beta$  does so only if  $q$  is  $(\chi = D(i))$ 's multisignature for  $[Reduction, r]$  (lines 93 and 94).  $\square$

**Lemma 84.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $P[r] \neq \perp$ . We have  $r = \rho(L[r])$ .*



*Proof.* We start by noting that  $\beta$  adds elements to  $batches$  only by executing line 87, and  $\beta$  never updates  $batches[r].payloads$ . Immediately before adding  $(r, b)$  to  $batches$  (line 87),  $\beta$  computes  $r = \rho(\mathcal{S}(l))$  (lines 75 and 76), with

$$\begin{aligned} ((i, c, m) \in l) &\iff ((i, Submission\{c, m, \dots\}) \in u) \\ &\iff ((i, (c, m)) \in b.payloads) \end{aligned}$$

for some  $u$  (lines 73 and line 83). The lemma immediately follows from Definition 61.  $\square$

**Lemma 85.** *Let  $\beta$  be a correct broker, let  $r$  be a root, let  $t, t' \in \mathbb{R}$  such that  $P_\beta(t)[r] \neq \perp$  and  $P_\beta(t')[r] \neq \perp$ . We have  $P_\beta(t)[r] = P_\beta(t')[r]$ .*

*Proof.* It follows immediately from Lemmas 84 and 50 and Definition 61.  $\square$

**Lemma 86.** *Let  $\beta$  be a correct broker, let  $\sigma$  be a correct server, let  $r$  be a root such that  $P[r] \neq \perp$  and  $r \in batches$  at  $\sigma$ . Let  $b = batches[r]$  at  $\sigma$ , let*

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

*We have*

$$\begin{aligned} i &= \mathcal{S}(\mathcal{D}(P[r])) \\ p &= \mathcal{S}(\mathcal{C}(P[r])) \end{aligned}$$

*Proof.* It follows immediately from Definition 61, Lemmas 84 and 53, and Lemma 50.  $\square$

**Lemma 87.** *Let  $\beta$  be a correct broker, let  $r$  be a root, let  $t, t'' \in \mathbb{R}$  such that  $t'' > t$ ,  $W_\beta(t)[r] \neq \perp$  and  $W_\beta(t'')[r] = \perp$ . For some  $t' \in [t, t'']$  we have  $|W_\beta(t')[r]| \geq f + 1$ .*

*Proof.* We start by noting that, by Definition 61,  $batches[r]$  is Witnessing at  $\beta$  at time  $t$ . Let  $t'$  identify the first moment after  $t$  when  $\beta$  updates  $batches[r]$ 's variant.  $\beta$  does so only by executing line 145, and only if  $|batches[r].witnesses| \geq f + 1$  (line 138). We then have  $|W_\beta(t')[r]| \geq f + 1$ , and the lemma is proved.  $\square$

**Notation 20** (Sequence elements). Let  $X$  be a set, let  $z \in X^{<\infty}$ . We use

$$\{z\} = \{z_n \mid n \leq |z|\} \tag{6.3}$$

to denote the **elements** of  $z$ .

**Notation 21** (Sequence indexing). Let  $X$  be a set, let  $z \in X^{<\infty}$  such that all elements of  $z$  are distinct. Let  $x \in \{z\}$ . We use

$$(z_x) = (n \iff z_n = x)$$

to identify the **index** of  $x$  in  $z$ .

**Lemma 88.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $W[r] \neq \perp$ .  $\beta$  has sent a  $[Batch, \tilde{i}, p]$  message to all servers, with*

$$\begin{aligned}\tilde{i} &= c(i) \\ i &= \mathcal{S}(\mathcal{D}(P[r])) \\ p &= \mathcal{S}(\mathcal{C}(P[r]))\end{aligned}$$

*Proof.*  $\beta$  sets  $batches[r]$ 's variant to Witnessing only by executing line 107. The lemma immediately follows from lines 100, 101, 102, 104 and 105 and Definition 61.  $\square$

**Lemma 89.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $P[r] \neq \perp$ , let*

$$\begin{aligned}\tilde{i} &= c(i) \\ i &= \mathcal{S}(\mathcal{D}(P[r])) \\ p &= \mathcal{S}(\mathcal{C}(P[r]))\end{aligned}$$

*Let  $\sigma$  be a correct server. Upon delivering a  $[Batch, \tilde{i}, p]$  message from  $\beta$ ,  $\sigma$  sets*

$$batches[r] = Batch\{ids: i, payloads: p, \dots\}$$

*and sends a  $[BatchAcquired, r, u]$  message back to  $\beta$ , with  $u \subseteq \mathcal{D}(P[r])$  such that  $\sigma$  knows all elements of  $(\mathcal{D}(P[r]) \setminus u)$ .*

*Proof.* Noting that  $i$  and  $p$  respectively list the domain and codomain of the same map, we have that  $i$  has no duplicates and  $|i| = |p|$ . Moreover, by Definitions 58 and 61, we have  $j(i, p) = L[r]$ . By Lemma 84, this proves  $r = \rho(j(i, p))$ .

Upon delivering  $[Batch, \tilde{i}, p]$  (line 64),  $\sigma$  expands  $\tilde{i}$  back into  $i = e(\tilde{i})$  (line 49) and verifies that  $i$  has no duplicates and  $|i| = |p|$  (line 51).  $\sigma$  then computes the set  $u$  of ids in  $\{i\}$  that  $\sigma_n$  does not know (line 54). We immediately have  $u \subseteq \mathcal{D}(P[r])$ , and  $\sigma$  knows all elements of  $(\mathcal{D}(P[r]) \setminus u)$ . Next,  $\sigma_n$  computes  $r = \rho(j(i, p))$  (lines 56, 57 and 58). Finally,  $\sigma_n$  sets  $batches[r] = Batch\{ids: i, payloads: p, \dots\}$  (line 60), and sends a  $[BatchAcquired, r, u]$  message back to  $\beta$  (lines 61, 67 and 68).  $\square$

**Lemma 90.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that, at some point in time, we have  $P[r] \neq \perp$ . At some point in time we have  $|W[r]| \geq f + 1$ .*

*Proof.* Upon initialization,  $batches$  is empty at  $\beta$  (line 47). Moreover, upon adding  $(r, \_)$  to  $batches$  (line 87 only),  $\beta$  sets a  $[Reduce, r]$  timer (line 64). When  $[Reduce, r]$  eventually rings

at  $\beta$  (line 99),  $\beta$  updates  $batches[r]$ 's variant to *Witnessing*. By Lemma 88, before doing so  $\beta$  sends a  $[Batch, \tilde{i}, p]$  message to all servers, with

$$\begin{aligned}\tilde{i} &= c(i) \\ i &= \mathcal{S}(\mathcal{D}(P[r])) \\ p &= \mathcal{S}(\mathcal{C}(P[r]))\end{aligned}$$

Let  $\sigma_1, \dots, \sigma_{f+1}$  be distinct correct servers (noting that at most  $f$  servers are Byzantine,  $\sigma_1, \dots, \sigma_{f+1}$  are guaranteed to exist). Let  $n \leq f+1$ . By Lemma 89, upon delivering  $[Batch, \tilde{i}, p]$   $\sigma_n$  sets  $batches[r] = Batch\{ids: i, payloads: p, \dots\}$ , and sends a  $[BatchAcquired, r, u]$  message back to  $\beta$ , with  $u \subseteq \mathcal{D}(P[r])$  such that  $\sigma$  knows all elements of  $(\mathcal{D}(P[r]) \setminus u)$ .

Let us assume that, upon delivering  $[BatchAcquired, r, u]$  from  $\sigma_n$ , we have  $r \notin batches$  at  $\beta$ . By Definition 61 we have  $P[r] = \perp$  which, by Lemma 76, implies  $W[r] = \perp$ . As a result, by Lemma 87, at some point in time we must have had  $|W[r]| \geq f+1$ . Throughout the remainder of this proof we assume that, upon delivering  $[BatchAcquired, r, u]$  from  $\sigma_n$ , we have  $r \in batches$  at  $\beta$ .

Upon delivering  $[BatchAcquired, r, u]$  from  $\sigma_n$  (line 116),  $\beta$  verifies that  $b \in batches$  (line 117) and, because  $\beta$  knows all elements of  $u$  (line 118),  $\beta$  maps  $u$  into a corresponding set of assignments  $a$  (line 121). Next,  $\beta$  aggregates all elements of  $\mathcal{C}(Q[r])$  into a multisignature  $q$  (line 123) and copies  $S[r]$  into a map  $s$  (line 124).

By Lemmas 85 and 81 we have

$$\{\hat{i}\} \setminus \mathcal{D}(s) = \mathcal{D}(P[r]) \setminus \mathcal{D}(S[r]) = \mathcal{D}(Q[r]) \quad (6.4)$$

By Lemma 82, for all  $(\hat{i}, \hat{s}) \in s$ ,  $(\hat{s} = s[\hat{i}] = S[r][i])$  is  $D(\hat{i})$ 's signature for  $[Message, c, m]$ , with

$$(c, m) = p_{(i_i)} = P[r][\hat{i}]$$

By Lemma 83 and Equation 6.4,  $q$  is  $(\{\hat{i}\} \setminus \mathcal{D}(s) = \mathcal{D}(Q[r]))$ 's multisignature for  $[Reduction, r]$ .

Having computed  $a$ ,  $q$  and  $s$ ,  $\beta$  sends a  $[Signatures, r, a, q, s]$  message back to  $\sigma_n$  (line 126).

Upon delivering  $[Signatures, r, a, q, s]$  from  $\beta$  (line 100),  $\sigma_n$  imports all elements of  $a$  (line 72). Noting that  $a$  contains assignments for all elements of  $u$ , and any element of  $(\{\hat{i}\} \setminus u)$  was known to  $\sigma_n$  upon delivering  $[Batch, \dots]$ ,  $\sigma_n$  knows all elements of  $i$ . Noting that  $\sigma_n$  never modifies or removes elements of  $batches$ ,  $\sigma_n$  successfully retrieves  $i$  and  $p$  from  $batches[r]$  (lines 74 and 77).  $\sigma_n$  then verifies to know all elements of  $i$  (line 79).

Next,  $\sigma_n$  loops through all elements of  $s$ . For each  $(\hat{i}, \hat{s})$  in  $s$ ,  $\sigma_n$  verifies that  $\hat{i} \in i$  (line 83, see Lemma 81), then verifies that  $\hat{s}$  is  $D(\hat{i})$ 's signature for  $[\text{Message}, c, m]$ , with  $(c, m) = p_{(i_i)}$  (lines 86 and 88). Subsequently,  $\sigma_n$  verifies that  $q$  is  $(i \setminus \mathcal{D}(s))$ 's multisignature for  $[\text{Reduction}, r]$  (lines 91 and 93). Finally,  $\sigma$  produces a multisignature  $w$  for  $[\text{Witness}, r]$  (line 96) and sends a  $[\text{WitnessShard}, r, w]$  message back to  $\beta$  (lines 97, 103 and 104).

Let us assume that, upon delivering  $[\text{WitnessShard}, r, w]$  from  $\sigma_n$ , we have  $r \notin \text{batches}$  or  $\text{batches}[r]$  not Witnessing at  $\beta$ . By Definition 61 we have  $W[r] = \perp$ . As a result, by Lemma 87, at some point in time we must have had  $|W[r]| \geq f + 1$ . Throughout the remainder of this proof we assume that, upon delivering  $[\text{WitnessShard}, r, w]$ , we have  $r \in \text{batches}$  and  $\text{batches}[r]$  Witnessing at  $\beta$ .

Upon delivering  $[\text{WitnessShard}, r, w]$  from  $\sigma_n$  (line 129),  $\beta$  verifies that  $r \in \text{batches}$  and  $\text{batches}[r]$  is Witnessing (lines 130 and 133), then verifies that  $w$  is  $\sigma_n$ 's multisignature for  $[\text{Witness}, r]$  (line 134) and finally adds  $(\sigma_n, w)$  to  $\text{batches}[r].\text{witnesses}$  (line 135).

In summary, for all  $n \leq f + 1$ , either  $|W[r]| \geq f + 1$ , or  $\beta$  adds a distinct element to  $\text{batches}[r].\text{witnesses}$ . Noting that  $\beta$  never removes elements from  $\text{batches}[r].\text{witnesses}$ , this trivially reduces to  $\beta$  eventually satisfying  $|W[r]| \geq f + 1$ , and the lemma is proved.  $\square$

**Lemma 91.** *Let  $\sigma$  be a correct server, let  $(r, w) \in \text{witnesses}$  at  $\sigma$ . We have that  $w$  is a plurality certificate for  $[\text{Witness}, r]$ .*

*Proof.* Upon initialization,  $\text{witnesses}$  is empty at  $\sigma$  (line 18). Moreover,  $\sigma$  adds  $(r, w)$  to  $\text{witnesses}$  only by executing line 116. Immediately before doing so,  $\sigma$  verifies that  $w$  is a plurality certificate for  $[\text{Witness}, r]$  (lines 113 and 114).  $\square$

**Lemma 92.** *Let  $\sigma$  be a correct server, let  $t \in \mathbb{R}$ , let  $(r, b) \in \text{batches}$  at  $\sigma$  at time  $t$ . For all  $t' \geq t$ ,  $(r, b) \in \text{batches}$  at  $\sigma$  at time  $t'$ .*

*Proof.* Let  $t' \geq t$ . Because  $\sigma$  never removes elements from  $\text{batches}$ , for some  $b'$  we have  $(r, b') \in \text{batches}$  at  $\sigma$  at time  $t'$ . By Lemmas 53 and 50 we have  $b' = b$ , and the lemma is proved.  $\square$

**Lemma 93.** *Let  $\sigma$  be a correct server, let  $((i, c), (\_, r)) \in \text{messages}$  at  $\sigma$ . We have  $r \in \text{batches}$  at  $\sigma$ .*

*Proof.* We start by noting that, upon initialization,  $\text{messages}$  is empty at  $\sigma$  (line 21). Moreover,  $\sigma$  adds  $((i, c), (\_, r))$  to  $\text{messages}$  only by executing line 130.  $\sigma$  does so only if  $r \in \text{batches}$  (lines 108 and 109). The lemma immediately follows from Lemma 92.  $\square$

**Lemma 94.** *Let  $\sigma$  be a correct server, let  $((i, c), (m, r)) \in \text{messages}$  at  $\sigma$ . Let  $b = \text{batches}[r]$  at  $\sigma$ , let*

$$\begin{aligned} i &= b.ids \\ p &= b.payloads \end{aligned}$$

*let  $l = j(i, p)$ . We have  $(i, c, m) \in \{l\}$ .*

*Proof.* We underline that, by Lemma 93,  $b$  is well-defined. We start by noting that, upon initialization,  $\text{messages}$  is empty at  $\sigma$  (line 21). Moreover,  $\sigma$  adds  $((i, c), (m, r))$  to  $\text{messages}$  only by executing line 130.  $\sigma$  does so only if (see Lemma 92)  $(i, c, m) \in \{l\}$  (lines 111, 119).  $\square$

**Lemma 95.** *Let  $\beta$  be a correct broker, let  $r$  be a root, let  $t, t'' \in \mathbb{R}$  such that  $t'' > t$ ,  $X_\beta(t)[r] \neq \perp$  and  $X_\beta(t'')[r] = \perp$ . For some  $t' \in [t, t'']$  we have  $|X_\beta(t')[r]| \geq 2f + 1$ .*

*Proof.* We start by noting that, by Definition 61,  $\text{batches}[r]$  is *Committing* at  $\beta$  at time  $t$ . Let  $t'$  identify the first moment after  $t$  when  $\beta$  updates  $\text{batches}[r]$ 's variant.  $\beta$  does so only by executing line 185, and only if  $|\text{batches}[r].\text{commits}| \geq 2f + 1$  (line 171). We then have  $|X_\beta(t')[r]| \geq 2f + 1$ , and the lemma is proved.  $\square$

**Lemma 96.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $|W[r]| \geq f + 1$ . At some point in time we have  $|X[r]| \geq 2f + 1$ .*

*Proof.* We start by noting that, if  $\text{batches}[r]$  is *Witnessing* at broker,  $\beta$  never removes elements from  $\text{batches}[r].\text{witnesses}$ , and  $\beta$  updates  $\text{batches}[r]$ 's variant only by executing line 145. As a result,  $\beta$  is eventually guaranteed to detect that  $\text{batches}[r]$  is *Witnessing* and that  $|\text{batches}[r].\text{witnesses}| \geq f + 1$  (line 138). Upon doing so,  $\beta$  aggregates  $\text{batch}[r].\text{witnesses}$  into a certificate  $c$  (lines 139 and 140) and sends a  $[\text{Witness}, r, c]$  message to every server (lines 142 and 143).

Let  $\sigma_1, \dots, \sigma_{2f+1}$  be distinct correct servers (noting that at most  $f$  servers are Byzantine,  $\sigma_1, \dots, \sigma_{2f+1}$  are guaranteed to exist). Let  $n \leq 2f + 1$ . Noting that  $W[r] \neq \perp$ , by Lemma 88, the source-order delivery of perfect links and Lemmas 89 and 92, by the time  $\sigma_n$  handles the delivery of  $[\text{Witness}, r, c]$ ,  $\sigma_n$  satisfies

$$\text{batches}[r] = \text{Batch}\{ids: i, payloads: p, \dots\}$$

with

$$\begin{aligned} i &= \mathcal{S}(\mathcal{D}(\mathcal{P}[r])) \\ p &= \mathcal{S}(\mathcal{C}(\mathcal{P}[r])) \end{aligned}$$

## Chapter 6. Draft

---

Upon delivering  $[\text{Witness}, r, c]$  (line 136),  $\sigma_n$  verifies that  $r \in \text{batches}$  (line 108) and retrieves  $i$  and  $p$  from  $\text{batches}[r]$  (line 111).  $\sigma_n$  then initializes an empty map  $f : \emptyset \rightarrow (\mathcal{R}, \mathbb{S}^+, \mathcal{P}, \mathbb{M})$  (line 117), and loops through all elements of  $j(i, p)$  (line 119).

For all  $j \leq (|i| = |p|)$ , let  $(c_j, m_j) = p_j$ . Noting that  $\sigma_n$  adds to  $f$  only keys that belong to  $i$  (line 128), let  $k \leq (|i| = |p|)$  such that  $\sigma_n$  sets

$$f[i_k] = (r'_k, w'_k, p'_k, m'_k)$$

By line 120 we immediately have  $m'_k \neq m_k$ . Additionally, because  $w'_k = \text{witnesses}[r'_k]$  at  $\sigma$  (line 121), by Lemma 91  $w'_k$  is a plurality certificate for  $[\text{Witness}, r'_k]$ . Finally, by line 120 we have  $((i_k, c_k), (m'_k, r'_k)) \in \text{messages}$  at  $\sigma_n$ . Let  $b'_k = \text{batches}[r'_k]$  at  $\sigma_n$ , let

$$l'_k = j(b'_k.\text{ids}, b'_k.\text{payloads})$$

By Lemma 94 we have  $(i_k, c_k, m'_k) \in l'_k$ . As a result, by lines 122, 124, 126 and 128, we have that  $p'_k$  is a proof for  $(i_k, c_k, m'_k)$  from  $r'_k$ .

Having looped over all elements of  $j(i, p)$ ,  $\sigma_n$  produces a signature  $s$  for  $[\text{Commit}, r, \mathcal{D}(f)]$  (line 132) and sends a  $[\text{CommitShard}, r, f, s]$  message back to  $\beta$  (lines 133, 139 and 140).

Let us assume that, upon delivering  $[\text{CommitShard}, r, f, s]$  from  $\sigma_n$ , we have  $r \notin \text{batches}$  or  $\text{batches}[r]$  not  $\text{Witnessing}$  at  $\beta$ . By Definition 61 we have  $X[r] = \perp$ . As a result, by Lemma 95, at some point in time we must have had  $|X[r]| \geq 2f + 1$ . Throughout the remainder of this proof we assume that, upon delivering  $[\text{CommitShard}, r, f, s]$ , we have  $r \in \text{batches}$  and  $\text{batches}[r]$   $\text{Witnessing}$  at  $\beta$ .

Upon delivering  $[\text{CommitShard}, r, f, s]$  from  $\sigma_n$  (line 148),  $\beta$  verifies that  $r \in \text{batches}$  and  $\text{batches}[r]$  is  $\text{Witnessing}$  (line 149).  $\beta$  then verifies that  $s$  is  $\sigma_n$ 's multisignature for  $[\text{Commit}, r, \mathcal{D}(f)]$  (line 150). Recalling that  $\mathcal{D}(f) \subseteq \{i\}$  and, by Lemma 85, we still have

$$i = \mathcal{S}(\mathcal{D}(\mathcal{P}(r)))$$

$$p = \mathcal{S}(\mathcal{C}(\mathcal{P}(r)))$$

for all  $k \leq |i|$  such that  $(i_k, (r'_k, w'_k, p'_k, m'_k)) \in f$  (line 153),  $\beta$  successfully verifies that:  $i_k \in \{i\}$  (line 154);  $w'_k$  is a plurality certificate for  $[\text{Witness}, r'_k]$  (line 159);  $p'_k$  is a proof for  $(i_k, c_k, m'_k)$  from  $r'_k$  (lines 157 and 162); and  $m'_k \neq m_k$  (lines 157 and 165). Having done so,  $\beta$  adds  $(\sigma_n, \_)$  to  $\text{batches}[r].\text{commits}$  (lines 149 and 168).

In summary, for all  $n \leq 2f + 1$ , either  $|X[r]| \geq 2f + 1$ , or  $\beta$  adds a distinct element to  $\text{batches}[r].\text{commits}$ . Noting that  $\beta$  never removes elements from  $\text{batches}[r].\text{commits}$ , this trivially reduces to  $\beta$  eventually satisfying  $|X[r]| \geq 2f + 1$ , and the lemma is proved.  $\square$

**Lemma 97.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $X[r] \neq \perp$ . Let  $(\sigma, (i, \_)) \in X[r]$ . We have  $i \subseteq \mathcal{D}(P[r])$ .*

*Proof.* Upon first setting  $batches[r]$ 's variant to `Committing` (line 145 only),  $\beta$  sets  $batches[r].commits = \{\}$ . Moreover,  $\beta$  adds  $(\sigma, (i, \_))$  to  $batches[r].commits$  only by executing line 168. Before doing so,  $\beta$  retrieves  $b = batches[r]$  (line 149), loops through every element  $\hat{i}$  of  $i$  (line 153) and verifies  $\hat{i} \in b.payloads$  (line 154). The lemma immediately follows from Definition 61.  $\square$

**Lemma 98.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $P[r] \neq \perp$ , let  $(i, (c, m)) \in P[r]$  such that  $\chi = D(i)$  is correct. We have that  $\chi$  broadcast  $m$  for  $c$ .*

*Proof.* Upon initialization,  $batches$  is empty at  $\beta$  (line 47). Moreover,  $\beta$  adds  $(r, b)$  to  $batches$  only by executing line 87. Upon doing so,  $\beta$  sets  $b.reductions = \{\}$ . By Definition 61 and Lemma 81 we then have  $\mathcal{D}(S[r]) = \mathcal{D}(P[r])$ . Therefore, by Lemma 82,  $\chi$  signed  $[Message, c, m]$ . Because  $\chi$  does so (line 30 only) only upon broadcasting  $m$  for  $c$  (line 29), the lemma is proved.  $\square$

**Lemma 99.** *Let  $\sigma$  be a correct server, let  $r$  be a root such that  $\sigma$  signs  $[Witness, r]$ . Some  $l \in (\mathbb{I} \times \mathbb{C} \times \mathbb{M})^\nearrow$  exists such that  $r = \rho(l)$ .*

*Proof.*  $\sigma$  signs  $[Witness, r]$  only by executing line 96.  $\sigma$  does so only if  $r \in batches$  (lines 74 and 75). Upon initialization,  $batches$  is empty at  $\sigma$  (line 17). Moreover,  $\sigma$  adds  $(r, \_)$  to  $batches$  (line 60 only) only if, for some sorted  $l$  (lines 49 and 56), we have  $r = \rho(l)$  (lines 57 and 58).  $\square$

**Lemma 100.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $X[r] \neq \perp$ . Let  $(\_, (i, \_)) \in X[r]$ , let  $\hat{i} \in i$ , let  $\chi = D(\hat{i})$ . We have that  $\chi$  is Byzantine.*

*Proof.* We underline that, by Lemmas 97 and 80,  $\chi$  is well-defined. Upon first setting  $batches[r]$ 's variant to `Committing` (line 145 only),  $\beta$  sets  $batches[r].commits = \{\}$ . Moreover,  $\beta$  adds  $(\_, (i, \_))$  to  $batches[r].commits$  only by executing line 168. Before doing so,  $\beta$  retrieves  $b = batches[r]$  (line 149), then loops through all elements of  $i$  (line 153). Upon looping over  $\hat{i}$ ,  $\beta$  retrieves  $(c, m) = b.payloads[r]$ .  $\beta$  then verifies that, for some root  $r'$ , at least  $f + 1$  servers signed  $[Witness, r']$  (line 159). Next,  $\beta$  verifies, that for some proof  $p$ , index  $n$  and message  $m'$ , we have

$$v(r', p, n, (\hat{i}, c, m')) = \text{True} \quad (6.5)$$

(line 162). Finally,  $\beta$  verifies that  $m \neq m'$  (line 165).

Let us assume by contradiction that  $\chi$  is correct. By Lemma 98 we immediately have that  $\chi$  broadcast  $m$  for  $c$ . Moreover, noting that at most  $f$  servers are Byzantine, at least one correct server signed  $[Witness, r']$ . As a result, by Lemma 99, some  $l' \in \mathbb{I} \times \mathbb{C} \times \mathbb{M}^\nearrow$  exists such that  $r' = \rho(l')$ . Moreover, by Equation 6.5 and Definition 57, we have  $l'_n = (i, c, m')$ . By Lemma

68, this proves that  $\chi$  broadcast  $m'$  for  $c$ . In summary,  $\chi$  broadcast  $m$  and  $m' \neq m$  for  $c$ . This contradicts  $\chi$  being correct and proves the lemma.  $\square$

**Lemma 101.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $X[r] \neq \perp$ . Let  $(\sigma, (\epsilon, s)) \in X[r]$ . We have that  $s$  is  $\sigma$ 's multisignature for  $[\text{Commit}, r, \epsilon]$ .*

*Proof.* Upon setting  $batches[r]$ 's variant to `Committing` (line 145 only),  $\beta$  initializes  $batches[r].commits$  to an empty map. Moreover,  $\beta$  adds  $(\sigma, (\epsilon, s))$  to  $batches[r].commits$  only by executing line 168.  $\beta$  does so only if  $s$  is  $\sigma$ 's multisignature for  $[\text{Commit}, r, \epsilon]$  (lines 150, 151).  $\square$

**Lemma 102.** *Let  $\beta$  be a correct broker, let  $r$  be a root, let  $\sigma \in Y[r]$  be a correct server. We have that  $r \in batches$  at server, and  $\sigma$  knows all elements of  $\mathcal{D}(P[r])$ .*

*Proof.* Upon initialization of  $batches[r]$ , we have  $batches[r].commit\_to = \emptyset$  at  $\beta$  (line 87 only). Moreover,  $\beta$  adds  $\sigma$  to  $batches[r].commit\_to$  only by executing line 131.  $\beta$  does so only upon receiving a `[WitnessShard, r]` message from  $\sigma$  (line 129). In turn,  $\sigma$  sends a `[WitnessShard, r]` message to  $\beta$  only by executing line 104.  $\sigma$  does so only if  $r \in batches$  at  $\sigma$  (lines 74 and 75) and  $\beta$  knows all elements of  $batches[r].ids$  (lines 88 and 89). By Lemma 86, however, we have  $batches[r].ids = \mathcal{S}(\mathcal{D}(P[r]))$ , and the lemma is proved.  $\square$

**Lemma 103.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $X[r] \neq \perp$ . We have  $|Y[r]| \geq f + 1$ .*

*Proof.* We start by noting that  $\beta$  updates  $batches[r]$ 's variant to `Committing` only by executing line 145.  $\beta$  does so only if  $batches[r]$  is `Witnessing` and  $|batches[r].witnesses| \geq f + 1$  (line 138). Upon setting  $batches[r]$ 's variant to `Witnessing` (line 107 only),  $\beta$  initializes  $batches[r].witnesses$  to an empty map. Finally, whenever  $\beta$  adds  $(\sigma, \_)$  to  $batches[r].witnesses$  (line 135 only),  $\beta$  also adds  $\sigma$  to  $batches[r].commit\_to$  (line 131). The lemma immediately follows from the observation that  $\beta$  never removes elements from  $batches[r].commit\_to$ .  $\square$

**Lemma 104.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $X[r] \neq \perp$ . We eventually have  $Z[r] = \text{True}$ .*

*Proof.* We start by noting that  $\beta$  removes  $r$  from  $batches$  (line 202 only) only if  $batches[r]$  is `Completing` (line 194). Moreover,  $\beta$  sets  $batches[r]$ 's variant to `Completing` (line 185 only) only if  $batches[r].committable = \text{True}$ . In other words, if  $X[r] \neq \perp$ , then  $\beta$  removes  $r$  from  $batches$  only if  $Z[r] = \text{True}$ .

$\beta$  updates  $batches[r]$ 's variant to `Committing` (line 145 only) only if  $batches[r]$ 's variant is `Witnessing`. Immediately after setting  $batches[r]$ 's variant to `Witnessing` (line 107 only),  $\beta$  sets a `[Committable, r]` timer (line 108). When `[Committable, r]` eventually rings (line 111),  $\beta$  checks if  $r \in batches$  (line 112). If so,  $\beta$  sets  $batches[r].committable$  to `True`. Otherwise, as



we proved above, we previously had  $batches[r].committable = \text{True}$  at  $\beta$ , and the lemma is proved.  $\square$

**Lemma 105.** *Let  $\beta$  be a correct broker, let  $r$  be a root such that  $|X[r]| \geq 2f + 1$ . Let*

$$E = \bigcup_{(\_, \epsilon, \_) \in X[r]} \epsilon$$

*Let  $\hat{i} \in (\mathcal{D}(P[r]) \setminus E)$ , let  $\hat{\chi} = D(\hat{i})$ , let  $(\hat{c}, \_) = P[r][\hat{i}]$ . Some correct server  $\sigma$  exists such that  $\sigma$  eventually delivers a message from  $\hat{\chi}$  for  $\hat{c}$ .*

*Proof.* We start by noting that, if  $batches[r]$  is *Committing* at  $\beta$ ,  $\beta$  never removes elements from  $batches[r].commits$ , and  $\beta$  updates  $batches[r]$ 's variant only by executing line 185. As a result, by Lemma 104,  $\beta$  is eventually guaranteed to detect that:  $batches[r]$  is *Committing*;  $batches[r].committable = \text{true}$ ; and  $|batches[r].commits| \geq 2f + 1$  (line 171). Upon doing so,  $\beta$  builds a map  $p : \mathbb{P}(\mathbb{I}) \rightarrow \mathbb{P}(\mathbb{S}^+)$  such that

$$((\epsilon, \_) \in p) \iff ((\_, \epsilon, \_) \in X[r]) \quad (6.6)$$

and

$$p[\epsilon] = \bigcup_{(\_, \epsilon, s) \in X[r]} s \quad (6.7)$$

(lines 174, 176, 177). From  $p$ ,  $\beta$  builds a map  $q : \mathbb{P}(\mathbb{I}) \rightarrow \mathbb{S}^+$  that to each  $\epsilon$  in  $p$  associates the aggregation of  $p[\epsilon]$  (line 179). By Equation 6.6 we immediately have

$$\bigcup_{(\epsilon, \_) \in q} \epsilon = E$$

For all  $(\_, c) \in q$ , let  $|c|$  denote the number of signers of  $c$ . By Equations 6.6 and 6.7 we have

$$\sum_{(\_, c) \in q} |c| = \sum_{(\_, s) \in p} |s| = |X[r]| \geq 2f + 1$$

Finally, by Lemma 101, for every  $(\epsilon, c) \in q$ ,  $c$  certifies  $[\text{Commit}, r, \epsilon]$ . Having computed  $q$ ,  $\beta$  sends a  $[\text{Commit}, r, q]$  message to all servers in  $Y[r]$  (lines 181 and 182).

By Lemma 103, we have  $|Y[r]| \geq f + 1$ . Let  $\sigma \in X[r]$  be a correct server. Noting that at most  $f$  servers are Byzantine,  $\sigma$  is guaranteed to exist. By Lemma 102, upon receiving a  $[\text{Commit}, r, q]$  message from  $\beta$  (line 179),  $\sigma$  verifies that  $r \in batches$  (line 144). By Lemma 86,  $\sigma$  then verifies to know all elements ( $\{batches[r].ids\} = \mathcal{D}(P[r])$ ) (line 149). Next,  $\sigma$  verifies that, for every  $(\epsilon, c) \in q$ ,  $c$  certifies  $[\text{Commit}, r, \epsilon]$  (lines 154 and 155). Finally,  $\sigma$  verifies that

$$\sum_{(\_, c)} |c| \geq 2f + 1$$

(lines 152, 154, 158 and 160). Having done so,  $\sigma$  computes  $E$  (line 163). Next, by Lemma

86  $\sigma$  loops through all elements of  $P[r]$  (line 166). For every  $(i, (c, m)) \in P[r]$  such that  $i \notin E$  (line 169),  $\sigma$  either delivers  $m$  from  $D(i)$  for  $c$  (lines 167 and 171), or  $\sigma$  has already delivered a message from  $D(i)$  for  $c$  (line 169, see Lemma 49). This proves in particular that  $\sigma$  delivers a message from  $\hat{\chi}$  for  $\hat{c}$ , and concludes the lemma.  $\square$

**Theorem 19.** *Draft satisfies validity.*

*Proof.* Let  $\chi$  be a correct client, let  $c$  be a context, let  $m$  be a message such that  $\chi$  broadcasts  $m$  for  $c$ . By Lemma 78, either a correct server delivers message from  $\chi$  for  $c$ , or some correct broker  $\beta$  and some root  $r$  exist such that eventually  $P[r][i] = (c, m)$ . Throughout the remainder of this proof, we assume the existence of  $\beta$  and  $r$ . By Lemma 90, at some point in time we have  $|W[r]| \geq f + 1$ . Therefore, by Lemma 96, at some point in time we have  $|X[r]| \geq 2f + 1$ . Moreover, by Lemma 100, for all  $(\_, (c, \_)) \in X[r]$ , we have  $D(\chi) \notin c$ . As a result, by Lemma 105, some correct server delivers a message from  $\chi$  for  $c$ , and the theorem is proved.  $\square$

## 6.4 Complexity

In this section, we prove to the fullest extent of formal detail the good-case signature and communication complexity of Draft.

### 6.4.1 Auxiliary results

In this section we gather definitions and lemmas that we will use to prove, in the next sections, the good-case signature and communication complexity of Draft. The results presented in this section hold independently of Draft itself, and could be applicable to a broader spectrum of analyses.

**Notation 22** (Bit strings). We use  $\mathbb{S} = \{0, 1\}^{<\infty}$  to denote all finite strings of bits. We use  $\epsilon$  to denote the empty sequence of bits. We use  $\mathbb{S}^{<b} = \{0, 1\}^{<b}$ ,  $\mathbb{S}^{\leq b} = \{0, 1\}^{\leq b}$ ,  $\mathbb{S}^b = \{0, 1\}^b$ ,  $\mathbb{S}^{\geq b} = \{0, 1\}^{\geq b}$  and  $\mathbb{S}^{>b} = \{0, 1\}^{>b}$ . We use  $\mathbb{S}^{even}$  and  $\mathbb{S}^{odd}$  to denote the sets of strings with an even and odd number of bits, respectively. We use programming notation when indexing a string of bits: for all  $s \in \mathbb{S}$  we use  $s = (s_0, s_1, \dots)$ .

**Notation 23** (Cropping). Let  $n \in \mathbb{N}$ , let  $s \in \mathbb{S}$ . We use the following **cropping** notation:

$$\begin{aligned} (s|_{\leq n}) &= (s_0, \dots, s_n) \\ (s|_{\geq n}) &= (s_n, \dots, s_{|s|-1}) \end{aligned}$$

We also use  $(s|_{<n}) = (s|_{\leq n-1})$  and  $(s|_{>n}) = (s|_{\geq n+1})$ .

**Definition 62** (Integer encoding). For all  $b \in \mathbb{N}$ , the  **$b$ -bits integer representation**  $\tilde{t}_b :$

$(0..(2^b - 1)) \leftrightarrow \mathbb{S}^b$  is defined by

$$\begin{aligned}\tilde{t}_b(n)_i &= \lfloor \frac{n}{2^i} \rfloor \bmod 2 \\ \tilde{t}_b^{-1}(s) &= \sum_{i=0}^b 2^i s_i\end{aligned}$$

The  $b$ -bits integer encoding  $\iota_b : \mathbb{S} \times (0..(2^b - 1)) \leftrightarrow \mathbb{S}^{\geq b}$  is defined by

$$\begin{aligned}\iota_b(s, n) &= \tilde{t}_b(n) \frown s \\ \iota_b^{-1}(s) &= ((s|_{\geq b}), \tilde{t}_b^{-1}((s|_{< b})))\end{aligned}$$

**Lemma 106.** *Let  $b \in \mathbb{N}$ .  $\tilde{t}_b^{-1}$  is injective.*

*Proof.* Let  $s, s' \in \mathbb{S}$  such that  $s \neq s'$ . Let

$$k = \max i \mid s_i \neq s'_i$$

Noting that  $s \neq s'$ ,  $k$  is guaranteed to exist. Without loss of generality, let us assume  $s_k = 1$ ,  $s'_k = 0$ . We by Definition 62 have

$$\begin{aligned}\tilde{t}_b^{-1}(s) - \tilde{t}_b^{-1}(s') &= \sum_{i=0}^b 2^i (s_i - s'_i) = \\ &= \sum_{i=0}^{k-1} 2^i (s_i - s'_i) + 2^k \geq \\ &\geq \sum_{i=0}^{k-1} -2^i + 2^k \geq 1\end{aligned}$$

which proves  $\tilde{t}_b^{-1}(s) \neq \tilde{t}_b^{-1}(s')$  and concludes the lemma □

**Lemma 107.** *Let  $b \in \mathbb{N}$ , let  $s \in \mathbb{S}^b$ . We have  $\tilde{t}_b(\tilde{t}_b^{-1}(s)) = s$ .*

*Proof.* By Definition 62 we have

$$\begin{aligned}
 \tilde{t}_b(\tilde{t}_b^{-1}(s))_i &= \lfloor \frac{\sum_{j=0}^b 2^j s_j}{2^i} \rfloor \bmod 2 = \\
 &= \lfloor \frac{\sum_{j=0}^{i-1} 2^j s_j}{2^i} + s_i + \frac{\sum_{j=i+1}^b 2^j s_j}{2^i} \rfloor \bmod 2 = \\
 &= \underbrace{\lfloor \frac{\sum_{j=0}^{i-1} 2^j s_j}{2^i} \rfloor}_{<1} + s_i + 2 \underbrace{\sum_{j=0}^{b-i-1} 2^j s_{j+i+1}}_{\in \mathbb{N}} \bmod 2 = \\
 &= \left( s_i + 2 \underbrace{\sum_{j=0}^{b-i-1} 2^j s_{j+i+1}}_{\in \mathbb{N}} \right) \bmod 2 = s_i
 \end{aligned}$$

□

**Lemma 108.** Let  $b \in \mathbb{N}$ .  $\tilde{t}_b$  is a bijection.

*Proof.* It follows immediately from Lemmas 106 and 107.

□

**Lemma 109.** Let  $b \in \mathbb{N}$ .  $\iota_b$  is a bijection.

*Proof.* It follows immediately from Definition 62 and Lemma 108.

□

**Definition 63** (Varint encoding). The **varint representation**  $\tilde{v} : \mathbb{N}^+ \leftrightarrow \mathbb{S}^{even}$  is defined by

$$\begin{aligned}
 |\tilde{v}(n)| &= 2 \lceil \log_2(n+1) \rceil \\
 \tilde{v}(n)_i &= \begin{cases} 1 & \text{iff } i \bmod 2 = 0, i < 2 \lceil \log_2(n+1) \rceil - 2 \\ 0 & \text{iff } i = 2 \lceil \log_2(n+1) \rceil - 2 \\ \tilde{t}_{\lceil \log_2(n+1) \rceil}(n)_{\frac{i-1}{2}} & \text{otherwise} \end{cases} \\
 \tilde{v}^{-1}(s) &= \tilde{t}_{\frac{|s|}{2}}^{-1}(s_1, s_3, s_5, \dots, s_{|s|-1})
 \end{aligned}$$

The set of **varint-parseable strings** is the set

$$\mathbb{S}^v = \{s \in \mathbb{S} \mid (\exists k \in \mathbb{N} \mid s_{2k} = 0)\}$$

The **varint encoding**  $v : (\mathbb{S} \times \mathbb{N}^+ \leftrightarrow \mathbb{S}^v)$  is defined by

$$\begin{aligned}
 v(s, n) &= \tilde{v}(n) \smallfrown s \\
 v^{-1}(s) &= ((s|_{\geq \lambda(s)}), \tilde{v}^{-1}((s|_{< \lambda(s)})))
 \end{aligned}$$

with  $\lambda : \mathbb{S} \rightarrow \mathbb{N}$  defined by

$$\lambda(s) = 2(\min k \mid s_{2k} = 0) + 2$$

Definition 63 contains a slight abuse of notation: we prove below that  $\tilde{v}^{-1}$  inverts  $\tilde{v}$ , but we do not prove that  $\tilde{v}$  inverts  $\tilde{v}^{-1}$  (it does not). Similarly we prove that  $v^{-1}$  inverts  $v$ , but not that  $v$  inverts  $v^{-1}$ .

**Lemma 110.** *Let  $n \in \mathbb{N}^+$ . We have  $\tilde{v}^{-1}(\tilde{v}(n)) = n$ .*

*Proof.* By Definition 63 and Lemma 109 we have

$$\begin{aligned} \tilde{v}^{-1}(\tilde{v}(n)) &= \tilde{t}_{\frac{2\lceil \log_2(n+1) \rceil}{2}}^{-1} \left( \tilde{t}_{\lceil \log_2(n+1) \rceil}(n)_0, \tilde{t}_{\lceil \log_2(n+1) \rceil}(n)_1, \dots, \tilde{t}_{\lceil \log_2(n+1) \rceil}(n)_{\lceil \log_2(n+1) \rceil - 1} \right) = \\ &= \tilde{t}_{\lceil \log_2(n+1) \rceil}^{-1}(\tilde{t}_{\lceil \log_2(n+1) \rceil}(n)) = n \end{aligned}$$

□

**Lemma 111.** *Let  $s \in \mathbb{S}$ , let  $n \in \mathbb{N}^+$ . We have  $v^{-1}(v(s, n)) = (s, n)$ .*

*Proof.* All derivations in this lemma follow directly from Definition 63. We have

$$v(s, n) = \tilde{v}(n) \frown s$$

Moreover, we have

$$\forall k < \lceil \log_2 n + 1 \rceil - 1, v(s, n)_{2k} = \tilde{v}(n)_{2k} = 1$$

and

$$v(s, n)_{2(\lceil \log_2(n+1) \rceil - 1)} = \tilde{v}^{-1}(n)_{2(\lceil \log_2(n+1) \rceil - 1)} = 0$$

which proves

$$\lambda(v(s, n)) = 2\lceil \log_2(n+1) \rceil$$

This implies

$$(v(s, n)|_{<\lambda(v(s, n))}) = (\tilde{v}(n) \frown s|_{<2\lceil \log_2(n+1) \rceil}) = \tilde{v}(n)$$

and similarly

$$(v(s, n)|_{\geq \lambda(v(s, n))}) = s$$

By Lemma 110 we then have

$$v^{-1}(v(s, n)) = (s, \tilde{v}^{-1}(\tilde{v}(n))) = (s, n)$$

and the lemma is proved. □

**Notation 24** (Integer partition). Let  $X$  be a finite set, let  $\mu : X \rightarrow \mathbb{P}^{<\infty}(\mathbb{N})$ . We call  $\mu$  a **integer partition** on  $X$ . We use

$$|\mu| = \sum_{x \in X} |\mu(x)|$$

$$\max \mu = \max_{x \in X} \max \mu(x)$$

**Notation 25** (Enumeration). Let  $X = \{x_1, \dots, x_{|X|}\}$  be a finite set. We use  $\mathcal{E}(X) = (x_1, \dots, x_{|X|})$  to denote any specific **enumeration** of  $X$ .

**Definition 64** (Partition representation). Let  $X$  be a finite set. The **partition representation** on  $X$  is the function  $\rho : (X \rightarrow \mathbb{P}^{<\infty}(\mathbb{N})) \leftrightarrow \mathbb{S}$  defined, with

$$(x_1, \dots, x_{|X|}) = \mathcal{E}(X)$$

by

$$w(\mu) = \lceil \log_2(\max \mu + 1) \rceil$$

$$d(\mu) = \lceil \log_2(|\mu| + 1) \rceil$$

$$y_n(\mu) = |\mu(x_n)|$$

$$(z_1(\mu), \dots, z_{|\mu|}(\mu)) = \mathcal{E}(\mu(x_1)) \frown \dots \frown \mathcal{E}(\mu(x_{|X|}))$$

$$v_0(\mu) = \emptyset$$

$$v_n(\mu) = \iota_{w(\mu)}(v_{n-1}(\mu), z_n(\mu))$$

$$v(\mu) = v_{|\mu|}(\mu)$$

$$k_0(\mu) = v(\mu)$$

$$k_n(\mu) = \iota_{d(\mu)}(k_{n-1}(\mu), y_n(\mu))$$

$$k(\mu) = k_{|X|}(\mu)$$

$$g(\mu) = v(k(\mu), w(\mu))$$

$$h(\mu) = v(g(\mu), d(\mu))$$

$$\rho(\mu) = h(\mu)$$

and

$$\begin{aligned}
 \mathbf{h}(s) &= s \\
 (\mathbf{g}(s), d(s)) &= v^{-1}(\mathbf{h}(s)) \\
 (\mathbf{k}(s), w(s)) &= v^{-1}(\mathbf{g}(s)) \\
 \mathbf{k}_{|X|}(s) &= \mathbf{k}(s) \\
 (\mathbf{k}_{n-1}(s), y_n(s)) &= \iota_{d(s)}^{-1}(\mathbf{k}_n(s)) \\
 t(s) &= \sum_{n=1}^{|X|} y_n(s) \\
 \mathbf{v}(s) &= \mathbf{k}_0(s) \\
 \mathbf{v}_{t(s)}(s) &= \mathbf{v}(s) \\
 (\mathbf{v}_{n-1}(s), z_n(s)) &= \iota_{w(s)}^{-1}(\mathbf{v}_n(s)) \\
 p_n(s) &= \sum_{i=1}^n y_n(s) \\
 \rho^{-1}(s)(x_n) &= \{z_{p_{n-1}(s)+1}(s), \dots, z_{p_n(s)}(s)\}
 \end{aligned}$$

Definition 64 contains a slight abuse of notation: we prove below that  $\rho^{-1}$  inverts  $\rho$ , but  $\rho$  is obviously not defined on all of  $\mathbb{S}$  ( $\checkmark$  being a trivial counterexample).

**Lemma 112.** *Let  $X$  be a finite set, let  $\mu: X \rightarrow \mathbb{P}^{<\infty}(\mathbb{N})$ . We have  $\rho^{-1}(\rho(\mu)) = \mu$ .*

*Proof.* All derivations in this lemma follow from Notation 24, Definition 64 and Lemmas 109 and 111. Let  $s = \rho(\mu)$ . We trivially have

$$\mathbf{h}(s) = s = \rho(\mu) = \mathbf{h}(\mu)$$

In steps, we can derive

$$\begin{aligned}
 (\mathbf{g}(s), d(s)) &= v^{-1}(\mathbf{h}(s)) = v^{-1}(\mathbf{h}(\mu)) = v^{-1}(v(\mathbf{g}(\mu), d(\mu))) = (\mathbf{g}(\mu), d(\mu)) \\
 (\mathbf{k}(s), w(s)) &= v^{-1}(\mathbf{g}(s)) = v^{-1}(\mathbf{g}(\mu)) = v^{-1}(v(\mathbf{k}(\mu), w(\mu))) = (\mathbf{k}(\mu), w(\mu))
 \end{aligned}$$

and

$$\mathbf{k}_{|X|}(s) = \mathbf{k}(s) = \mathbf{k}(\mu) = \mathbf{k}_{|X|}(\mu)$$

By backwards induction, for all  $n \in 1..|X|$  we then have

$$\begin{aligned}
 (\mathbf{k}_{n-1}(s), y_n(s)) &= \iota_{d(s)}^{-1}(\mathbf{k}_n(s)) = \iota_{d(\mu)}^{-1}(\mathbf{k}_n(\mu)) = \\
 &= \iota_{d(\mu)}^{-1}(\iota_{d(\mu)}(\mathbf{k}_{n-1}(\mu), y_n(\mu))) = (\mathbf{k}_{n-1}(\mu), y_n(\mu))
 \end{aligned}$$

## Chapter 6. Draft

---

from which follows

$$t(s) = \sum_{n=1}^{|X|} y_n(s) = \sum_{n=1}^{|X|} y_n(\mu) = \sum_{n=1}^{|X|} |\mu(x_n)| = |\mu|$$

In steps, we can then derive

$$\begin{aligned} v(s) &= \mathbf{k}_0(s) = \mathbf{k}_0(\mu) = v(\mu) \\ v_{|\mu|}(s) &= v_{t(s)}(s) = v(s) = v(\mu) = v_{|\mu|}(\mu) \end{aligned}$$

Again by backwards induction, for all  $n \in 1..|\mu|$  we then have

$$\begin{aligned} (v_{n-1}(s), z_n(s)) &= \iota_{w(s)}^{-1}(v_n(s)) = \iota_{w(\mu)}^{-1}(v_n(\mu)) \\ &= \iota_{w(\mu)}^{-1}(\iota_{w(\mu)}(v_{n-1}(\mu), z_n(\mu))) = (v_{n-1}(\mu), z_n(\mu)) \end{aligned}$$

In summary we have

$$\begin{aligned} \forall n \in 1..|X|, y_n(s) &= y_n(\mu) \\ \forall n \in 1..|\mu|, z_n(s) &= z_n(\mu) \end{aligned}$$

from which finally follows

$$\begin{aligned} \rho^{-1}(s)(x_n) &= \{z_{p_{n-1}(s)+1}(s), \dots, z_{p_n(s)}(s)\} = \\ &= \{z_{p_{n-1}(s)+1}(\mu), \dots, z_{p_n(s)}(\mu)\} = \\ &= \{z_{(\sum_{i=1}^{n-1} y_i(s)+1)}(\mu), \dots, z_{(\sum_{i=1}^n y_i(s))}(\mu)\} = \\ &= \{z_{(\sum_{i=1}^{n-1} y_i(\mu)+1)}(\mu), \dots, z_{(\sum_{i=1}^n y_i(\mu))}(\mu)\} = \\ &= \{z_{(\sum_{i=1}^{n-1} |\mu(x_i)|+1)}(\mu), \dots, z_{(\sum_{i=1}^n |\mu(x_i)|)}(\mu)\} = \\ &= \{z_{(\sum_{i=1}^{n-1} |\mathcal{E}(\mu(x_i))|+1)}(\mu), \dots, z_{(\sum_{i=1}^n |\mathcal{E}(\mu(x_i))|)}(\mu)\} = \\ &= \{\mathcal{E}(\mu(x_n))_1, \dots, \mathcal{E}(\mu(x_n))_{|\mu(x_n)|}\} = \mu(x_n) \end{aligned}$$

which proves  $\rho^{-1}(\rho(\mu)) = \rho^{-1}(s) = \mu$ , and concludes the lemma.  $\square$

**Lemma 113.** *Let  $X$  be a finite set, let  $\mu : \mathbb{N} \rightarrow (X \rightarrow \mathbb{P}^{<\infty}(\mathbb{N}))$  be a sequence of partition representations such that*

$$\begin{aligned} \lim_{n \rightarrow \infty} |\mu_n| &= \infty \\ \lim_{n \rightarrow \infty} \frac{\log_2(\log_2(\max \mu_n))}{|\mu_n|} &= 0 \end{aligned}$$

We have

$$\lim_{n \rightarrow \infty} \frac{|\rho(\mu_n)|}{|\mu_n|} = \lceil \log_2(\max \mu_n + 1) \rceil$$



*Proof.* All derivations in this proof follow from Definitions 64, 63 and 62. Let  $n \in \mathbb{N}$ . We have

$$|v_0(\mu_n)| = 0$$

and, for all  $k \in 1..|\mu|$ ,

$$|v_k(\mu_n)| = |v_{k-1}(\mu_n)| + w(\mu_n) = |v_{k-1}(\mu_n)| + \lceil \log_2(\max \mu_n + 1) \rceil$$

from which, by induction, follows

$$|v(\mu_n)| = |\mu_n| \lceil \log_2(\max \mu_n + 1) \rceil$$

Similarly we have

$$|k_0(\mu_n)| = |\mu_n| \lceil \log_2(\max \mu + 1) \rceil$$

and, for all  $k \in 1..|X|$ ,

$$|k_k(\mu_n)| = |k_{k-1}(\mu_n)| + d(\mu_n) = |k_{k-1}(\mu_n)| + \lceil \log_2(|\mu_n| + 1) \rceil$$

from which, by induction, follows

$$|k(\mu_n)| = |\mu_n| \lceil \log_2(\max \mu_n + 1) \rceil + |X| \lceil \log_2(|\mu_n| + 1) \rceil$$

In steps we can then derive

$$\begin{aligned} |g(\mu_n)| &= |k(\mu_n)| + 2 \lceil \log_2(w(\mu_n) + 1) \rceil = \\ &= |\mu_n| \lceil \log_2(\max \mu_n + 1) \rceil + |X| \lceil \log_2(|\mu_n| + 1) \rceil + 2 \lceil \log_2(\lceil \log_2(\max \mu_n + 1) \rceil + 1) \rceil \end{aligned}$$

and finally

$$\begin{aligned} |\rho(\mu_n)| &= |h(\mu_n)| = |g(\mu_n)| + 2 \lceil \log_2(d(\mu_n)) \rceil = \\ &= |\mu_n| \lceil \log_2(\max \mu_n + 1) \rceil + |X| \lceil \log_2(|\mu_n| + 1) \rceil + 2 \lceil \log_2(\lceil \log_2(\max \mu_n + 1) \rceil + 1) \rceil \\ &\quad + 2 \lceil \log_2(\lceil \log_2(|\mu_n| + 1) \rceil + 1) \rceil \end{aligned}$$

The above holds for any  $n \in \mathbb{N}$ . Moreover, by hypothesis we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{|X| \lceil \log_2(|\mu_n| + 1) \rceil}{|\mu_n|} &= 0 \\ \lim_{n \rightarrow \infty} \frac{2 \lceil \log_2(\lceil \log_2(|\mu_n| + 1) \rceil + 1) \rceil}{|\mu_n|} &= 0 \end{aligned}$$

and

$$\lim_{n \rightarrow \infty} \frac{2 \lceil \log_2 (\lceil \log_2 (\max \mu_n + 1) \rceil + 1) \rceil}{|\mu_n|} =$$

$$\lim_{n \rightarrow \infty} \frac{2 \lceil \log_2 (\lceil \log_2 (\max \mu_n) \rceil) \rceil}{|\mu_n|} \frac{\lceil \log_2 (\lceil \log_2 (\max \mu_n + 1) \rceil + 1) \rceil}{\lceil \log_2 (\lceil \log_2 (\max \mu_n) \rceil) \rceil} = 0$$

This proves that

$$\lim_{n \rightarrow \infty} \frac{|\rho(\mu_n)|}{|\mu_n|} = \lceil \log_2 (\max \mu_n + 1) \rceil$$

and concludes the lemma. □

### 6.4.2 Batching limit

As we discussed in Section 5.1, Draft is designed to asymptotically match, in the good case, the signature and communication complexity of Oracle-CSB, a toy implementation of CSB that relies on an infallible oracle to uphold all CSB properties. As we discussed, Oracle-CSB achieves optimal signature complexity (as it requires no signature verification) and optimal communication complexity (if the frequency at which clients broadcast their payloads is uniform or unknown). Below we establish and discuss Draft's **batching limit**, i.e., the collection of assumptions and limits at which Draft matches Oracle-CSB's complexity:

- **Assumption: good-case execution.** As we discussed in Section 5.3.1, in the good case: *links are synchronous* (messages are delivered at most one time unit after they are sent); *all processes are correct*; and *the set of brokers contains only one element*. **Discussion.** We assume only one broker for pedagogical reasons: as long as every broker is exposed to a high enough rate of submissions, all derivations in this section still hold true. In the real world, the assumptions of synchrony and correctness are not strict for clients: if a small fraction of broadcasting clients is slow or Byzantine, Draft still achieves near-oracular efficiency, linearly degrading its performance as more client fail to engage with the broker to reduce its batch<sup>2</sup>.
- **Assumption: steady-state directory.** We assume that *all servers know all broadcasting clients*. **Discussion.** This assumption is naturally satisfied if all broadcasting clients have already broadcast at least one message. In a real-world, long-lived system, most clients can safely be assumed to broadcast more than one message (this is especially true in the context of a cryptocurrency). Similarly to client synchrony and correctness, Draft's performance degrades linearly in the number of unknown, broadcasting clients.
- **Assumption: concurrent broadcasts.** We assume that *all clients broadcast their message within  $b$  time units of each other*, where  $b$  is Draft's batching window parameter (see Sections 6.2.3 and 6.2.4). **Discussion.** Similarly to broker count, this assumption is made

---

<sup>2</sup>Proving this result is beyond the scope of this Part, and left as an exercise to the interested reader.

for pedagogical reasons: as long as the rate at which payloads are submitted to each broker is high enough, all derivations in this section still hold true.

- **Limit: infinite broadcasts.** We derive Draft's complexity at the limit of *infinitely many broadcasting clients*. **Discussion.** All limits derived in this section converge approximately inversely with the number of broadcasting clients. This means that even a finite, real-world implementation of Draft can achieve near-oracular efficiency.
- **Limit: sub-double-exponential clients.** We assume that *the number of clients is infinitely small with respect to the exponential of the exponential of the number of broadcasting clients*. **Discussion.** This limit is only technical, and trivially satisfied by any realistic number of broadcasting clients.

### 6.4.3 Protocol analysis

In this section, we establish Draft's signature and communication complexity at the batching limit.

**Theorem 20.** *At the batching limit, a Draft server delivers a payload  $p$  by performing 0 signature verifications and exchanging at most  $(\lceil \log_2(c) \rceil + |p|)$  bits.*

*Proof.* All derivations in this proof follow from the batching limit's assumptions and limits (see Section 6.4.2). Let  $\beta$  denote the only broker in the system. Let  $\chi_1, \dots, \chi_M$  denote the set of broadcasting clients. We have  $M \rightarrow \infty$ . For all  $j$ , let  $i_j = \mathcal{D}(\chi_j)$ , let  $p_j = (c_j, m_j)$  identify the payload broadcast by  $\chi_j$ , let  $a_j$  be  $i_j$ 's id assignment certificate. Without loss of generality we assume that for all  $j < j'$  we have  $i_j < i_{j'}$ . The goal of this proof is to compute, for all  $j \in 1..M$ ,  $p_j$ 's maximum amortized signature complexity  $\psi_j$  and  $p_j$ 's maximum amortized bit complexity  $\chi_j$ .

Let  $n \in 1..M$ . Without loss of generality, we assume that  $\chi_n$  triggers `<cl.Broadcast>` (line 29) between time 0 and  $b$ . Upon doing so,  $\chi$  produces a signature  $s_n$  for `[Message,  $c_n, m_n$ ]` (line 30), then sends a `[Submission,  $a_n, (c_n, m_n, s_n)$ ]` message to  $\beta$  (line 40).

$\beta$  receives `[Submission,  $a_n, (c_n, m_n, s_n)$ ]` (line 50) between time 0 and time  $B + 1$ : indeed,  $\chi_n$ 's `Submission` message takes between 0 and 1 time units to reach  $\beta$ . Upon delivering  $\chi_n$ 's `Submission` message,  $\beta$  pushes

$$y_n = \text{Submission}\{\text{context}: c_n, \text{message}: m_n, \text{signature}: s_n\}$$

to `pending[ $i_n$ ]` (line 54), detects that `pending[ $i_n$ ]` is not empty (line 57), empties `pending[ $i_n$ ]` (line 58) and adds  $(i_n, y_n)$  to `pool` (line 59). Initially, we have `collecting = False` and `pool =  $\emptyset$`  at  $\beta$ . Moreover,  $\beta$  adds its first element to `pool` no earlier than time 0. Upon doing so (line 62),  $\beta$  sets a `[Flush]` timer to ring at some time  $t_f \geq B + 1$  (line 64).

In summary, for all  $j \in 1..M$ ,  $\beta$  adds  $(i_j, y_j)$  to *pool* by time  $B + 1$ , and [Flush] rings at time  $t_f \geq B + 1$ . When [Flush] rings (line 67),  $\beta$  takes  $(i_1, p_1), \dots, (i_M, p_M)$  from *pool* (line 70) and computes  $r = \rho(j(i, p))$  (line 76). For each  $j \in 1..M$ ,  $\beta$  produces an inclusion proof  $q_j$  for  $(i_j, c_j, m_j)$  in  $r$  (line 80), then sends an [Inclusion,  $c_j, r, q_j$ ] message to  $\chi_j$ . Next,  $\beta$  sets

$$batches[r] = Reducing\{payloads: p, signatures: s, \dots\}$$

(line 87). Finally,  $\beta$  sets a [Reduce,  $r$ ] to ring at some time  $t_r \geq t_f + 2$  (line 89).

$\chi_n$  receives [Inclusion,  $c_n, r, q_n$ ] by time  $t_f + 1$  (line 50). Upon doing so,  $\chi_n$  produces a multi-signature  $g_n$  for [Reduction,  $r$ ] (line 54), then sends a [Reduction,  $r, g_n$ ] message to  $\beta$  (line 55).

$\beta$  receives [Reduction,  $r, g_n$ ] by time  $t_r$  (line 92). Because [Reduce,  $r$ ] has not yet rung, *batches[r]* is still Reducing at  $\beta$  (line 93).  $\beta$  adds  $(i_n, g_n)$  to *batches[r].reductions* (line 95), then removes  $(i_n, s_n)$  from *batches[r].signatures* (line 96).

In summary, before [Reduce,  $r$ ] rings at time  $t_r$ , for all  $j \in 1..M$ ,  $\beta$  added  $(i_j, g_j)$  to *batches[r].reductions* and removed  $(i_j, s_j)$  from *batches[r].signatures*. This means in particular that, when [Reduce,  $r$ ], *batches[r].signatures* is empty at  $\beta$ . Upon ringing [Reduce,  $r$ ] (line 99),  $\beta$  builds the integer partition  $\tilde{i}: \mathbb{D} \rightarrow \mathbb{P}^{<\infty}(\mathbb{N})$  defined by

$$((d, k) \in \tilde{i}) \iff ((d, k) \in i)$$

(line 101). We have  $|\tilde{i}| = |i| = M$ . Moreover, by the density of Directory, we have  $\max \tilde{i} \leq c - 1$ . Recalling that  $M \rightarrow \infty$  and  $\log(\log(c))/M \rightarrow \infty$ , by Lemmas 112 and 113,  $\tilde{i}$  can be represented by a string of bits  $\hat{i}$  such that

$$|\hat{i}| = M \lceil \log_2(c) \rceil + o(M)$$

To each server,  $\beta$  sends a

$$x^{(b)} = [\text{Batch}, \tilde{i}, p]$$

message (line 105). We have

$$\begin{aligned} |x^{(b)}| &= M \lceil \log_2(c) \rceil + \sum_{j=1}^M |p_j| + o(M) = \\ &= \sum_{j=1}^M (\lceil \log_2(c) \rceil + |p_j| + o(1)) \end{aligned}$$

Consequently, for all  $j \in 1..M$ ,  $x^{(b)}$ 's amortized size for  $p_j$  is

$$|x^{(b)}|_j = (\lceil \log_2(c) \rceil + |p_j|)$$

Having sent  $x^{(b)}$  to all servers,  $\beta$  updates *batches[r]* to Witnessing (line 107).

Let  $\sigma$  be a server. Upon receiving  $[\text{Batch}, \tilde{l}, p]$  (line 64),  $\sigma$  observes to know all ids in  $i_1, \dots, i_M$  (line 54), computes  $r$  (line 58) then sends back to  $\beta$  a

$$x^{(ba)} = [\text{BatchAcquired}, r, \emptyset]$$

message (line 68). Noting that

$$|x^{(ba)}| = O(1) = \sum_{j=1}^M o(1)$$

for all  $j \in 1..M$ ,  $x^{(ba)}$ 's amortized size for  $p_j$  is

$$|x^{(ba)}|_j = 0$$

Upon receiving  $\sigma$ 's  $[\text{BatchAcquired}, r, \emptyset]$  (line 116),  $\beta$  exports no assignment (line 121), aggregates all  $g_1, \dots, g_M$  into a single  $g$  (line 123), observes  $batches[r].signatures$  to be empty, then sends back to  $\sigma$  a

$$x^{(s)} = [\text{Signatures}, r, \emptyset, g, \emptyset]$$

message (line 126). Noting that  $|x^{(s)}| = O(1)$ ,  $x^{(s)}$ 's amortized size for  $p_j$  is

$$|x^{(s)}|_j = 0$$

Upon receiving  $[\text{Signatures}, r, \emptyset, g, \emptyset]$  (line 100),  $\sigma$  verifies only  $g$  (line 91, the loop at line 82 does no iterations). Noting that  $\sigma$  performs signature verifications only upon receiving a  $\text{Signatures}$  message, and noting that

$$1 = \sum_1^M o(1)$$

we can immediately derive

$$\psi_j = 0$$

Next,  $\sigma$  produces a multi-signature  $w_\sigma$  for  $[\text{Witness}, r]$  (line 96), then sends back to  $\beta$  a

$$x_\sigma^{(ws)} = [\text{WitnessShard}, r, w_\sigma]$$

message (line 104). Noting that  $|x_\sigma^{(ws)}| = O(1)$ ,  $x_\sigma^{(ws)}$ 's amortized size for  $p_j$  is

$$|x_\sigma^{(ws)}|_j = 0$$

Upon receiving  $f + 1$   $\text{WitnessShard}$  messages (lines 129, 135, 138),  $\beta$  aggregates all witnesses

into a certificate  $w$  (line 140), then sends a

$$x^{(w)} = [\text{Witness}, r, w]$$

message to all servers (line 143). Noting that  $|x^{(w)}| = O(1)$ ,  $x^{(w)}$ 's amortized size for  $p_j$  is

$$|x^{(w)}|_j = 0$$

Upon receiving  $[\text{Witness}, r, w]$  (line 136),  $\sigma$  observes that, for all  $j$ ,  $p_j$  is not equivocated (lines 119 and 120) (because all clients are correct, no client equivocates), produces a multi-signature  $c_\sigma$  for  $[\text{Commit}, r, \emptyset]$  (line 132), then sends a

$$x_\sigma^{(cs)} = [\text{CommitShard}, r, \emptyset, c_\sigma]$$

message (line 140). Noting that  $|x_\sigma^{(cs)}| = O(1)$ ,  $x_\sigma^{(cs)}$ 's amortized size for  $p_j$  is

$$|x_\sigma^{(cs)}|_j = 0$$

Upon receiving  $2f + 1$  CommitShard messages (lines 148, 168 and 171),  $\beta$  builds a map  $h : \mathbb{P}(\mathbb{I}) \mapsto \mathbb{S}^+$  with only the key  $\emptyset$  (lines 176, 177, 179) (we recall that every server produced an empty exception set for  $r$ ).  $\beta$  then sends a

$$x^{(c)} = [\text{Commit}, r, h]$$

message to all servers (line 182). Noting that  $|x^{(c)}| = O(1)$ ,  $x^{(c)}$ 's amortized size for  $p_j$  is

$$|x^{(c)}|_j = 0$$

Upon delivering  $[\text{Commit}, r, h]$  (line 179),  $\sigma$  delivers  $p_1, \dots, p_M$  (line 171), produces a multi-signature  $z_\sigma$  for  $[\text{Completion}, r, \emptyset]$  (line 175), then sends a

$$x_\sigma^{(zs)} = [\text{CompletionShard}, r, z_\sigma]$$

message to  $\beta$  (line 183). Noting that  $|x_\sigma^{(zs)}| = O(1)$ ,  $x_\sigma^{(zs)}$ 's amortized size for  $p_j$  is

$$|x_\sigma^{(zs)}|_j = 0$$

Finally,  $\sigma$  sets an  $[\text{OfferTotality}]$  timer to ring after 7 time units (173). Summarizing the scheduling of messages we have that: every server delivers the Batch message between time  $t_r$  and  $t_r + 1$ ;  $\beta$  delivers all BatchAcquired messages between time  $t_r$  and  $t_r + 2$ ; every server delivers a Signatures message between time  $t_r$  and  $t_r + 3$ ;  $\beta$  delivers all WitnessShard messages between time  $t_r$  and  $t_r + 4$ ; every server delivers the Witness message between

time  $t_r$  and  $t_r + 5$ ;  $\beta$  delivers all CommitShard messages between time  $t_r$  and  $t_r + 6$ ; all servers deliver the Commit message between time  $t_r$  and  $t_r + 7$ . As a result, [OfferTotality,  $r$ ] rings at  $\sigma$  after all servers delivered  $p_1, \dots, p_M$ . Upon ringing [OfferTotality,  $r$ ] (line 186),  $\sigma$  sends a

$$x^{(ot)} = [\text{OfferTotality}, r]$$

message to all servers (line 188).

Let  $\sigma'$  be a server. Upon receiving  $\sigma$ 's [OfferTotality,  $r$ ] message,  $\sigma'$  verifies to have already delivered a batch with root  $r$  and no exclusions (line 192) and ignores the message.

In summary, noting that  $\sigma$  exchanges  $2n$  OfferTotality messages ( $N$  outgoing,  $N$  incoming),  $\sigma$ 's amortized bit complexity for  $p_j$  is

$$\psi_j^\sigma = |x^{(b)}|_j + |x^{(ba)}|_j + |x^{(s)}|_j + |x_\sigma^{(ws)}|_j + |x^{(w)}|_j + |x_\sigma^{(cs)}|_j + |x^{(c)}|_j + |x_\sigma^{(zs)}|_j + 2n|x^{(ot)}|_j$$

from which immediately follows

$$\psi_j^\sigma = (\lceil \log_2(c) \rceil + |p_j|)$$

and

$$\psi_j = \max_\sigma \psi_j^\sigma = (\lceil \log_2(c) \rceil + |p_j|)$$

which concludes the theorem. □





# 7 Dibs

In this chapter, we present in detail the **Directory** abstraction and discuss its properties. We then present Dibs, an algorithm that implements Directory, and prove its **security**.

## 7.1 Interface

**Notation 26** (Signatures and multi-signatures). We use  $\mathbb{S}^1$  and  $\mathbb{S}^+$  to respectively denote the set of signatures and multisignatures.

**Definition 65** (Id). An **id** is an element of  $(\mathbb{I} = \mathbb{D} \times \mathbb{N})$ , where  $\mathbb{D}$  is a finite set of domains. Let  $i = (d, n)$  be an id, we call  $d$  and  $n$  the **domain** and **index** of  $i$ , respectively.

The **Directory** interface (instance  $dir$ ) exposes the following procedures and events:

- **Request**  $\langle dir.Signup \rangle$ : requests that an id be assigned to the local process.
- **Indication**  $\langle dir.SignupComplete \rangle$ : indicates that an id was successfully assigned to the local process.
- **Getter**  $dir[id]$ : returns the process associated with id  $id$ , if known. Otherwise, returns  $\perp$ .
- **Getter**  $dir[process]$ : returns the id associated with process  $process$ , if known. Otherwise, returns  $\perp$ .
- **Getter**  $dir.export(id)$ : returns the assignment for id  $id$ , if known. Otherwise, returns  $\perp$ .
- **Setter**  $dir.import(assignment)$ : imports assignment  $assignment$ .

**Notation 27** (Bijective relations). Let  $X, Y$  be sets. We use  $(X \leftrightarrow Y)$  to denote the set of **bijective relations** between  $X$  and  $Y$ .

## Chapter 7. Dibs

---

**Notation 28** (Tuple binding). Let  $t = (t_1, \dots, t_n)$  be a tuple. When binding  $t$ , we use the **any** symbol  $\_$  to mark which elements of  $t$  are discarded from the binding. For example,

$$(x, \_, \dots, \_, y, \_) = t$$

binds  $x = t_1$ , and  $y = t_{n-1}$ . We use the **tail** symbol  $\dots$  to indicate that all subsequent elements of  $t$  are discarded from the binding. For example,

$$(x', y', \dots) = t$$

binds  $x' = t_1$  and  $y' = t_2$ , regardless of  $n$ . Let  $X$  be a set. We use the tuple binding notation to **filter** the elements  $X$ . For example, we use  $\{(x'', \_, y'', \dots)\}$  to identify the set

$$\{t \in X \mid |t| \geq 3, (x'', \_, y'', \dots) = t\}$$

of tuples in  $X$  whose first element is  $x''$  and whose third element is  $y''$ .

**Definition 66** (Directory record). Let  $D : (\Pi_C \times \mathbb{R}) \rightarrow (\mathbb{I} \leftrightarrow \Pi)$ .  $D$  is a **directory record** if and only if:

- For all  $\pi \in \Pi_C$ ,  $D_\pi(t)$  is non-decreasing in  $t$ .
- For all  $\pi, \pi' \in \Pi_C$ ,  $t, t' \in \mathbb{R}$ ,  $D_\pi(t) \cup D_{\pi'}(t')$  is a bijective relation.

Let  $\pi$  be a correct process, let  $t$  be a time, let  $i$  be an id, let  $\rho$  be a process.  $\pi$  **associates**  $i$  and  $\rho$  by time  $t$  if and only if  $(i, \rho) \in D_\pi(t)$ .  $\pi$  **knows**  $i$  (resp.,  $\rho$ ) by time  $t$  if and only if  $(i, \_) \in D_\pi(t)$  (resp.,  $(\_, \rho) \in D_\pi(t)$ ).

**Notation 29** (Directory record). Let  $\pi$  be a correct process. Wherever it can be unequivocally inferred from context, we omit the time from the directory record  $D_\pi$ .

A Directory satisfies the following properties:

- **Correctness:** Some directory record  $D$  exists such that, for any id  $i$  and process  $\rho$ , if a process  $\pi$  invokes  $dir[i]$  (resp.,  $dir[\rho]$ ),  $\pi$  obtains  $\rho$  (resp.,  $i$ ) if and only if  $\pi$  associates  $i$  and  $\rho$ .
- **Signup Integrity:** A correct process never triggers  $\langle dir.SignupComplete \rangle$  before triggering  $\langle dir.Signup \rangle$ .
- **Signup Validity:** If a correct process triggers  $\langle dir.Signup \rangle$ , it eventually triggers  $\langle dir.SignupComplete \rangle$ .
- **Self-knowledge:** Upon triggering  $\langle dir.SignupComplete \rangle$ , a correct process knows itself.

- **Transferability:** If a correct process invokes *dir.export(i)* to obtain an assignment *a*, then any correct process knows *i* upon invoking *dir.import(a)*.
- **Density:** For every correct process  $\pi$ , for every time  $t$ , for every  $(\_, n, \_) \in D_\pi(t)$ , we have  $n < |\Pi|$ .

**Notation 30** (Directory mapping). Let  $i$  be an id, let  $\rho$  be a process such that, for some  $\pi \in \Pi_C$ ,  $t \in \mathbb{R}$ , we have  $(i, \rho) \in D_\pi(t)$ . We say that  $i$  and  $\rho$  are **known**, and we use

$$i = D(\rho) \quad \rho = D(i)$$

We underline the soundness of Notation 30: by Definition 66, if  $(i, \rho) \in D_\pi(t)$ , then no  $i' \neq i$  or  $\rho' \neq \rho$  exist such that, for some  $\pi', t'$ , we have  $(i, \rho') \in D_{\pi'}(t')$  or  $(i', \rho) \in D_{\pi'}(t')$ . In other words, no  $i' \neq i$  or  $\rho' \neq \rho$  exist such that  $i' = D(\rho)$  or  $\rho' = D(i)$ .

## 7.2 Algorithm

In this section, we present Dibs's pseudocode. The remainder of Chapter 7 proves Dibs's security to the fullest extent of formal detail.

### 7.2.1 Pseudocode (Client)

```

1 implements:
2   Directory, instance dir
3
4
5 uses:
6   AuthenticatedPointToPointLinks, instance al
7
8
9 struct Assignment:
10  id: Id,
11  process: Process,
12  certificate: Certificate
13
14
15 enum Status(Outsider, SigningUp, SignedUp)
16
17
18 upon <dir.Init>:
19   rankings: {Server: {Server}} (default {}) = {};
20   assigner: (Server or  $\perp$ ) =  $\perp$ ;

```

## Chapter 7. Dibs

---

```
21   assignments: {Integer: {Server: MultiSignature}} (default {}) =
    {};
22   status: Status = Outsider;
23
24   directory: {(Id, Process)} = {};
25   certificates: {Id: Certificate} = {};
26
27
28   upon <dir.Signup>:
29     status = SigningUp;
30
31     for  $\sigma$  in  $\Sigma$ :
32       trigger <al.Send |  $\sigma$ , [Signup]>;
33
34
35   upon <al.Deliver |  $\sigma$ , [Ranked, source]>:
36     rankings[source].add( $\sigma$ );
37
38
39   upon exists source such that |rankings[source]| >= f + 1 and assigner
    =  $\perp$ :
40     assigner = source;
41
42     for  $\sigma$  in  $\Sigma$ :
43       trigger <al.Send |  $\sigma$ , [Assigner, assigner]>;
44
45
46   upon event <al.Deliver |  $\sigma$ , [Assignment, index, signature]>:
47     if  $\sigma$ .multiverify(signature, [Assignment, (assigner, index), self])
    :
48       assignments[index][ $\sigma$ ] = signature;
49
50
51   upon exists index such that |assignments[index]| >= 2f+1 and status =
    SigningUp:
52     status = SignedUp;
53
54     id = (assigner, index);
55     certificate = aggregate(assignments[index]);
56
57     dir.import(Assignment {id, process: self, certificate});
58     trigger <dir.SignupComplete>;
```

```

59
60
61 procedure dir.[] (query):
62     if query is Id:
63         if exists process such that (query, process) in directory:
64             return process;
65     else if query is Process:
66         if exists id such that (id, query) in directory:
67             return id;
68
69     return ⊥;
70
71
72 procedure dir.export(id):
73     if let certificate = certificates[id]:
74         process = dir[id];
75         return Assignment {id, process, certificate};
76     else:
77         return ⊥;
78
79
80 procedure dir.import(assignment):
81     if assignment.certificate.verify_quorum(
82         [Assignment, assignment.id, assignment.process]
83     ):
84         directory.add((assignment.id, assignment.process));
85         certificates[assignment.id] = assignment.certificate;

```

### 7.2.2 Pseudocode (Server)

```

1 implements:
2     DirectoryServer, instance dsr
3
4
5 uses:
6     FifoBroadcast, instance fb
7     AuthenticatedPointToPointLinks, instance al
8
9
10 upon <dsr.Init>:
11     rankings: {Server: [KeyCard]} (default []) = {};
12     assigners: {KeyCard: Server} = {};

```

```
13   certified: {KeyCard} = {};
14
15
16 upon <al.Deliver |  $\pi$ , [Signup]>:
17   trigger <fb.Broadcast | [Rank,  $\pi$ ]>;
18
19
20 upon <fb.Deliver |  $\sigma$ , [Rank, process]>:
21   if process not in rankings[ $\sigma$ ]:
22     rankings[ $\sigma$ ].push_back(process);
23     trigger <al.Send | process, [Ranked,  $\sigma$ ]>
24
25
26 upon <al.Deliver |  $\pi$ , [Assigner, assigner]>:
27   if  $\pi$  not in assigners:
28     assigners[ $\pi$ ] = assigner;
29
30
31 upon exists process such that (process in assigners)
   and (process in rankings[assigners[process]])
   and (process not in certified):
32   certified.add(process);
33
34   assigner = assigners[process];
35   index = rankings[assigner].index_of(process);
36
37   signature = multisign([Assignment, (assigner, index), process]);
38   trigger <al.Send | process, [Assignment, index, signature]>;
```

### 7.3 Correctness

In this section, we prove to the fullest extent of formal detail that Dibs implements a Directory.

#### 7.3.1 Correctness

In this section, we prove that Dibs satisfies correctness.

**Lemma 114.** *Let  $\sigma$  be a correct server, let  $\hat{\sigma}$  be a server. All elements of  $\text{ranking}[\hat{\sigma}]$  at  $\sigma$  are distinct.*

*Proof.* Upon initialization,  $\text{rankings}[\hat{\sigma}]$  is empty at  $\sigma$  (line 11). Moreover,  $\sigma$  appends  $\pi$  to

$rankings[\hat{\sigma}]$  only by executing line 22. Because  $\sigma$  does so only if  $\pi \notin rankings[\hat{\sigma}]$  (line 21), the lemma is proved.  $\square$

**Lemma 115.** *Let  $\sigma, \sigma'$  be correct servers, let  $\hat{\sigma}$  be a server, let  $n \in \mathbb{N}$ , let  $\pi$  be a process. If  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma$ , then eventually  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma'$  as well.*

*Proof.* Let  $[Rank, \pi_1], \dots, [Rank, \pi_R]$  denote the sequence of  $[Rank, \_]$  messages  $\sigma$  FIFO-delivered from  $\hat{\sigma}$ . By the totality and FIFO properties of FIFO broadcast,  $\sigma'$  eventually delivers  $[Rank, \pi_1], \dots, [Rank, \pi_R]$  as well. Upon initialization,  $rankings[\hat{\sigma}]$  is empty at both  $\sigma$  and  $\sigma'$  (line 11). Moreover,  $\sigma$  (resp.,  $\sigma'$ ) adds some process  $\hat{\pi}$  to  $rankings[\hat{\sigma}]$  (line 22 only) only upon FIFO-delivering a  $[Rank, \hat{\pi}]$  message from  $\hat{\sigma}$  (line 20), and only if  $rankings[\hat{\sigma}]$  satisfies a deterministic condition on  $\hat{\pi}$  (line 21). As a result, noting  $\sigma$  (resp.,  $\sigma'$ ) never removes elements from  $rankings[\hat{\sigma}]$ , if  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma$  as a result of  $\sigma$  delivering  $[Rank, \pi_1], \dots, [Rank, \pi_R]$ , then eventually  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma'$  as well, as it also eventually delivers  $[Rank, \pi_1], \dots, [Rank, \pi_R]$ .  $\square$

**Lemma 116.** *Let  $\sigma, \sigma'$  be correct servers. Let  $\hat{\sigma}$  be a server, let  $n \in \mathbb{N}$ , let  $\pi, \pi'$  be processes such that  $\sigma$  and  $\sigma'$  respectively sign  $[Assignment, (\hat{\sigma}, n), \pi]$  and  $[Assignment, (\hat{\sigma}, n), \pi']$ . We have  $\pi = \pi'$ .*

*Proof.*  $\sigma$  (resp.,  $\sigma'$ ) signs  $[Assignment, (\hat{\sigma}, n), \pi]$  (resp.,  $[Assignment, (\hat{\sigma}, n), \pi']$ ) only by executing line 37.  $\sigma$  (resp.,  $\sigma'$ ) does so only if  $n$  is the index of  $\pi$  (resp.,  $\pi'$ ) in  $rankings[\hat{\sigma}]$  at  $\sigma$  (resp.,  $\sigma'$ ) (line 35). We underline that, by Lemma 114, at most one element of  $rankings[\hat{\sigma}]$  is  $\pi$  (resp.,  $\pi'$ ) at  $\sigma$  (resp.,  $\sigma'$ ), hence  $n$  is well-defined. By Lemma 115, however, if  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma$ , then  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma'$  as well. This proves  $\pi = \pi'$  and concludes the lemma.  $\square$

**Lemma 117.** *Let  $\hat{\sigma}$  be a server, let  $n \in \mathbb{N}$ , let  $\pi, \pi' \neq \pi$  be processes. If a quorum certificate for  $[Assignment, (\hat{\sigma}, n), \pi]$  exists, then no quorum certificate for  $[Assignment, (\hat{\sigma}, n), \pi']$  exists.*

*Proof.* Let us assume by contradiction that quorum certificates exist for both  $[Assignment, (\hat{\sigma}, n), \pi]$  and  $[Assignment, (\hat{\sigma}, n), \pi']$ . Noting that at most  $f$  servers are Byzantine, at least one correct server  $\sigma$  (resp.,  $\sigma'$ ) signed  $[Assignment, (\hat{\sigma}, n), \pi]$  (resp.,  $[Assignment, (\hat{\sigma}, n), \pi']$ ). By Lemma 116 we then have  $\pi = \pi'$ , which contradicts  $\pi \neq \pi'$  and proves the lemma.  $\square$

**Lemma 118.** *Let  $\sigma, \sigma'$  be correct servers. Let  $\hat{\sigma}$  be a server, let  $n, n' \in \mathbb{N}$ , let  $\pi$  be a process such that  $\sigma$  and  $\sigma'$  respectively sign  $[Assignment, (\hat{\sigma}, n), \pi]$  and  $[Assignment, (\hat{\sigma}, n'), \pi]$ . We have  $n = n'$ .*

*Proof.*  $\sigma$  (resp.,  $\sigma'$ ) signs  $[Assignment, (\hat{\sigma}, n), \pi]$  (resp.,  $[Assignment, (\hat{\sigma}, n'), \pi]$ ) only by executing line 37.  $\sigma$  (resp.,  $\sigma'$ ) does so only if  $n$  (resp.,  $n'$ ) is the index of  $\pi$  in  $rankings[\hat{\sigma}]$  at  $\sigma$  (resp.,  $\sigma'$ ) (line 35). We underline that, by Lemma 114, at most one element of

## Chapter 7. Dibs

---

$rankings[\hat{\sigma}]$  is  $\pi$  at  $\sigma$  (resp.,  $\sigma'$ ), hence  $n$  (resp.,  $n'$ ) is well-defined. By Lemma 115, however, if  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma$ , then  $rankings[\hat{\sigma}][n] = \pi$  at  $\sigma'$  as well. Again by Lemma 114, this proves  $n = n'$  and concludes the lemma.  $\square$

**Lemma 119.** *Let  $\sigma$  be a correct server. Let  $\hat{\sigma}, \hat{\sigma}'$  be servers, let  $\pi$  be a process such that  $\sigma$  signs both an  $[Assignment, (\hat{\sigma}, \_), \pi]$  and an  $[Assignment, (\hat{\sigma}', \_), \pi]$  message. We have  $\hat{\sigma} = \hat{\sigma}'$ .*

*Proof.*  $\sigma$  signs an  $[Assignment, (\hat{\sigma}, \_), \pi]$  (resp.,  $[Assignment, (\hat{\sigma}', \_), \pi]$ ) message only by executing 37.  $\sigma$  does so only if  $\hat{\sigma} = assigners[\pi]$  (resp.,  $\hat{\sigma}' = assigners[\pi]$ ) (line 34). The lemma follows immediately from the observation that  $\sigma$  sets  $assigners[\pi]$  (line 28 only) at most once (line 27).  $\square$

**Lemma 120.** *Let  $\hat{\sigma}, \hat{\sigma}'$  be servers, let  $n, n' \in \mathbb{N}$  such that  $(\hat{\sigma}, n) \neq (\hat{\sigma}', n')$ , let  $\pi$  be a process. If a quorum certificate for  $[Assignment, (\hat{\sigma}, n), \pi]$  exists, then no quorum certificate for  $[Assignment, (\hat{\sigma}', n'), \pi]$  exists.*

*Proof.* Let us assume that  $\hat{\sigma} = \hat{\sigma}'$ . Because  $(\hat{\sigma}, n) \neq (\hat{\sigma}', n')$ , we have  $n \neq n'$ . Let us assume by contradiction that quorum certificates exist for both  $[Assignment, (\hat{\sigma}, n), \pi]$  and  $[Assignment, ((\hat{\sigma} = \hat{\sigma}'), n'), \pi]$ . Noting that at most  $f$  servers are Byzantine, at least one correct server  $\sigma$  (resp.  $\sigma'$ ) signed an  $[Assignment, (\hat{\sigma}, n), \pi]$  (resp.,  $[Assignment, (\hat{\sigma}, n'), \pi]$ ) message. By Lemma 118 we then have  $n = n'$ , which contradicts  $n \neq n'$ .

Let us assume that  $\hat{\sigma} \neq \hat{\sigma}'$ . We have that at least  $f + 1$  correct processes signed  $[Assignment, (\hat{\sigma}, \_), \pi]$  (resp.,  $[Assignment, (\hat{\sigma}', \_), \pi]$ ). Because any two sets of  $f + 1$  correct processes intersect in at least one element, some correct server  $\sigma^*$  exists that signed both  $[Assignment, (\hat{\sigma}, \_), \pi]$  and  $[Assignment, (\hat{\sigma}', \_), \pi]$ . By Lemma 119 we then have  $\hat{\sigma} = \hat{\sigma}'$  which contradicts  $\hat{\sigma} \neq \hat{\sigma}'$  and proves the lemma.  $\square$

**Definition 67** (Certifiable assignments). Let  $i$  be an id, let  $\pi$  be a process. The set  $D^\infty$  of **certifiable assignments** contains  $(i, \pi)$  if and only if a quorum certificate ever exists for  $[Assignment, i, \pi]$ .

**Lemma 121.**  $D^\infty$  is a bijection.

*Proof.* It follows immediately from Definition 67 and Lemmas 117 and 120.  $\square$

**Notation 31** (Directory record). Let  $\pi$  be a correct process, let  $t \in \mathbb{R}$ . We use  $D_\pi(t)$  to denote the value of *directory* at  $\pi$  at time  $t$ .

As we immediately prove,  $D$  is a directory record, which makes Notation 31 compatible with Definition 66.

**Lemma 122.** *Let  $\pi$  be a correct process, let  $t, t' \in \mathbb{R}$  such that  $t' \geq t$ . We have  $D_\pi(t') \supseteq D_\pi(t)$ .*



*Proof.* The lemma immediately follows from Notation 31 and the observation that  $\pi$  never removes elements from *directory*.  $\square$

**Lemma 123.** *Let  $\pi$  be a correct process, let  $t \in \mathbb{R}$ . We have  $D_\pi(t) \subseteq D^\infty$ .*

*Proof.* Upon initialization, *directory* is empty at  $\pi$  (line 24). Moreover,  $\pi$  adds  $(\hat{i}, \hat{\pi})$  only by executing line 82.  $\pi$  does so only upon verifying a quorum certificate for  $[\text{Assignment}, \hat{i}, \hat{\pi}]$  (line 81). The lemma immediately follows from Definition 67 and Notation 31.  $\square$

**Lemma 124.**  *$D$  is a directory record.*

*Proof.* By Lemma 122 we immediately have that, for every correct process  $\pi$ ,  $D_\pi(t)$  is non-decreasing in  $t$ . Let  $\pi, \pi'$  be correct processes, let  $t, t' \in \mathbb{R}$ . By Lemma 123 we have

$$D_\pi(t) \cup D_{\pi'}(t') \subseteq D^\infty$$

By Lemma 121, this proves that  $D_\pi(t) \cup D_{\pi'}(t')$  is bijective. The lemma immediately follows from Definition 66.  $\square$

**Theorem 21.** *Dibs satisfies correctness.*

*Proof.* Let  $\pi$  be a correct process, let  $\hat{i}$  be an id, let  $\hat{\pi}$  be a process. By Notation 31, upon invoking  $\text{dir}[\hat{i}]$  (resp.,  $\text{dir}[\hat{\pi}]$ ) (line 61),  $\pi$  obtains  $\hat{\pi}$  (resp.,  $\hat{i}$ ) (line 64, resp., line 67) if and only if  $(\hat{i}, \hat{\pi}) \in D_\pi$  (line 63, resp., line 66). The theorem follows immediately from Lemma 124.  $\square$

### 7.3.2 Signup integrity

In this section, we prove that Dibs satisfies signup integrity.

**Theorem 22.** *Dibs satisfies signup integrity.*

*Proof.* Let  $\pi$  be a correct process. Upon initialization, we have  $\text{status} = \text{Outsider}$  at  $\pi$  (line 22). Moreover,  $\pi$  triggers  $\langle \text{dir.SignupComplete} \rangle$  only by executing line 58.  $\pi$  does so only if  $\text{status} = \text{SigningUp}$  (line 51). The theorem follows immediately from the observation that  $\pi$  sets  $\text{status} = \text{SigningUp}$  (line 29 only) only upon triggering  $\langle \text{dir.Signup} \rangle$  (line 28).  $\square$

### 7.3.3 Signup validity

In this section, we prove that Dibs satisfies signup validity.

**Lemma 125.** *Let  $\pi$  be a correct process that triggers  $\langle \text{dir.Signup} \rangle$ . We eventually have  $\text{assigner} \neq \perp$  at  $\pi$ .*

## Chapter 7. Dibs

---

*Proof.* Upon triggering  $\langle dir.Signup \rangle$  (line 28),  $\pi$  sends a  $[Signup]$  message to all servers (lines 31 and 32). Let  $\hat{\sigma}$  be a correct server. Noting that at most  $f$  servers are Byzantine,  $\hat{\sigma}$  is guaranteed to exist.

Upon delivering  $[Signup]$  from  $\pi$  (line 16),  $\hat{\sigma}$  FIFO-broadcasts a  $[Rank, \pi]$  message (line 17). Let  $\sigma_1, \dots, \sigma_{f+1}$  be distinct correct servers. Noting that at most  $f$  servers are Byzantine,  $\sigma_1, \dots, \sigma_{f+1}$  are guaranteed to exist.

Let  $n \leq f + 1$ . We start by noting that, upon initialization,  $rankings[\hat{\sigma}]$  is empty at  $\sigma_n$  (line 11). Moreover,  $\sigma_n$  adds  $\pi$  to  $rankings[\hat{\sigma}]$  (line 22) only upon delivering a  $[Rank, \pi]$  message from  $\hat{\sigma}$  (line 20). By the validity and totality of FIFO broadcast,  $\sigma_n$  is eventually guaranteed to FIFO-deliver  $[Rank, \pi]$  from  $\hat{\sigma}$ . Upon doing so (line 20),  $\sigma_n$  verifies that  $\pi \notin rankings[\hat{\sigma}]$  (line 22), then sends a  $[Ranked, \hat{\sigma}]$  message to  $\pi$ .

Upon delivering a  $[Ranked, \hat{\sigma}]$  message from  $\sigma_n$  (line 35),  $\pi$  adds  $\sigma_n$  to  $rankings[\hat{\sigma}]$  (line 36). In summary, recalling that the above holds true for any  $n \leq f + 1$ ,  $\pi$  adds  $\sigma_1, \dots, \sigma_{f+1}$  to  $rankings[\hat{\sigma}]$ . Noting that  $\pi$  never removes elements from  $rankings[\hat{\sigma}]$ , that initially  $assigner = \perp$  at  $\pi$  (line 20), and that  $\pi$  updates  $assigner$  only by executing line 40,  $\pi$  is eventually guaranteed to detect that, for some  $\sigma^*$ ,  $|rankings[\sigma^*]| \geq f + 1$  and  $assigner = \perp$  (line 39). Upon doing so,  $\pi$  sets  $assigner = \sigma^*$ , and the lemma is proved.  $\square$

**Lemma 126.** *Let  $\pi$  be a correct process, let  $\hat{\sigma}$  be a server such that  $assigner = \hat{\sigma}$  at  $\pi$ . We eventually have  $\pi \in rankings[\hat{\sigma}]$  at all correct processes.*

*Proof.* We start by noting that, upon initialization, we have  $assigner = \perp$  at  $\pi$  (20). Moreover,  $\pi$  sets  $assigner = \hat{\sigma}$  (line 40 only) only if  $|rankings[\hat{\sigma}]| \geq f + 1$  at  $\pi$  (line 39). Upon initialization,  $rankings[\hat{\sigma}]$  is empty at  $\pi$  (line 19). Moreover,  $\pi$  adds  $\sigma$  to  $rankings[\hat{\sigma}]$  only upon delivering a  $[Ranked, \hat{\sigma}]$  message from  $\sigma$ . In summary, at least  $f + 1$  servers sent a  $[Ranked, \hat{\sigma}]$  message to  $\pi$ . Let  $\sigma$  be a correct server that sent a  $[Ranked, \hat{\sigma}]$  message to  $\pi$ . Noting that at most  $f$  servers are Byzantine,  $\sigma$  is guaranteed to exist.

$\sigma$  sends a  $[Ranked, \hat{\sigma}]$  to  $\pi$  (line 23 only) only upon FIFO-delivering a  $[Rank, \pi]$  message from  $\hat{\sigma}$ . By the totality of FIFO broadcast, eventually every correct process delivers  $[Rank, \pi]$ . Let  $\sigma'$  be a correct process. Upon FIFO-delivering  $[Rank, \pi]$  from  $\hat{\sigma}$  (line 20),  $\sigma'$  either already satisfies  $\pi \in rankings[\hat{\sigma}]$  (line 21) or adds  $\pi$  to  $rankings[\hat{\sigma}]$ . Noting that no correct process ever removes elements from  $rankings[\hat{\sigma}]$ , eventually  $\pi \in rankings[\hat{\sigma}]$  at all correct processes, and the lemma is proved.  $\square$

**Lemma 127.** *Let  $\pi$  be a correct process, let  $\hat{\sigma}$  be a server such that  $assigner = \hat{\sigma}$  at  $\pi$ , let  $r \in \mathbb{N}$  such that, eventually,  $rankings[\hat{\sigma}][r] = \pi$  at all correct servers. We eventually have  $|assignments[r]| \geq 2f + 1$  at  $\pi$ .*

*Proof.* We start by noting that, upon initialization, we have  $assigner = \perp$  at  $\pi$  (line 20). Moreover,  $\pi$  updates  $assigner$  to  $\hat{\sigma}$  only by executing line 40. Upon doing so,  $\pi$  sends a  $[Assigner, \hat{\sigma}]$  to all servers (lines 42 and 43). Moreover, because  $\pi$  updates  $assigner$  to a non- $\perp$  value (line 40) immediately before disseminating  $[Assigner, \hat{\sigma}]$ , and  $\pi$  never resets  $assigner$  back to  $\perp$ ,  $\pi$  never issues any  $[Assigner, (\hat{\sigma} \neq \hat{\sigma})]$  message (see line 39). Let  $\sigma_1, \dots, \sigma_{2f+1}$  be distinct correct servers. Noting that at most  $f$  servers are Byzantine,  $\sigma_1, \dots, \sigma_{2f+1}$  are guaranteed to exist.

Let  $n \leq 2f + 1$ . Upon initialization,  $assigners$  is empty at  $\sigma_n$  (line 12). Moreover,  $\sigma_n$  adds  $\pi$  to  $assigners$  (line 28 only) only upon delivering a  $[Assigner, \_]$  message from  $\pi$  (line 26). As a result, upon delivering  $[Assigner, \hat{\sigma}]$  from  $\pi$  (line 26),  $\sigma_n$  verifies  $\pi \notin assigners$  (line 27) and sets  $assigners[\pi] = \hat{\sigma}$  (line 28). Upon initialization,  $certified$  is empty at  $\sigma_n$  (line 13). Moreover,  $\sigma_n$  adds  $\pi$  to  $certified$  only by executing line 32. As a result,  $\sigma_n$  is eventually guaranteed to observe  $assigners[\pi] = \hat{\sigma}$ ,  $rankings[\hat{\sigma}] = r$ , and  $\pi \in certified$  (line 31). Upon doing so,  $\sigma_n$  produces a signature  $s_n$  for  $[Assignment, (\hat{\sigma}, r), \pi]$  (lines 34, 35 and 37) and sends an  $[Assignment, r, s_n]$  message back to  $\pi$  (line 38).

Upon delivering  $[Assignment, r, s_n]$  from  $\sigma_n$  (line 46),  $\pi$  verifies  $s_n$  against  $[Assignment, (\hat{\sigma}, r), \pi]$  (line 47), then adds  $\sigma_n$  to  $assignments[r]$ . In summary, recalling that the above holds true for all  $n \leq 2f + 1$ ,  $assignments[r]$  eventually contains  $\sigma_1, \dots, \sigma_{2f+1}$  at  $\pi$ , and the lemma is proved. □

**Theorem 23.** *Dibs satisfies signup validity.*

*Proof.* Let  $\pi$  be a correct process that triggers  $\langle dir.Signup \rangle$ . By Lemma 125, for some  $\hat{\sigma}$  we eventually have  $assigner = \hat{\sigma}$  at  $\pi$ . By Lemmas 126 and 115, some  $r \in \mathbb{N}$  exists such that, eventually,  $rankings[\hat{\sigma}][r] = \pi$  at all correct processes. By Lemma 127, we eventually have  $|assignments[r]| \geq 2f + 1$  at  $\pi$ . Upon triggering  $\langle dir.Signup \rangle$  (line 28),  $\pi$  sets  $status = SigningUp$  (line 29). Moreover, if  $status = SigningUp$ ,  $\pi$  updates  $status$  only by executing line 52. As a result,  $\pi$  is eventually guaranteed to detect that  $|assignments[r]| \geq 2f + 1$  and  $status = SigningUp$  (line 51). Upon doing so,  $\pi$  triggers  $\langle dir.SignupComplete \rangle$  (line 58), and the theorem is proved. □

### 7.3.4 Self-knowledge

In this section, we prove that Dibs satisfies self-knowledge.

**Lemma 128.** *Let  $\pi$  be a correct process. Let  $\hat{\sigma}$  be a server such that  $assigner = \hat{\sigma}$  at  $\pi$ . Let  $r \in \mathbb{N}$ , let  $\sigma$  be a server such that, for some  $s$ , we have  $assignments[r][\sigma] = s$  at  $\pi$ . We have that  $s$  is  $\sigma$ 's signature for  $[Assignment, (\hat{\sigma}, r), \pi]$ .*

## Chapter 7. Dibs

---

*Proof.* We start by noting that  $\pi$  updates *assigner* (line 40 only) only if *assigner* =  $\perp$  (line 46). Consequently, at all times we either have *assigner* =  $\perp$  or *assigner* =  $\hat{\sigma}$  at  $\pi$ . Upon initialization, *assignments* is empty at  $\pi$  (line 21). Moreover,  $\pi$  sets *assignments*[ $r$ ][ $\sigma$ ] =  $s$  only by executing line 48.  $\pi$  does so only if  $s$  is  $\sigma$ 's multisignature for  $[\text{Assignment}, (\hat{\sigma}, r), \pi]$  (line 47).  $\square$

**Theorem 24.** *Dibs satisfies self-knowledge.*

*Proof.* Let  $\pi$  be a correct process that triggers  $\langle \text{dir.SignupComplete} \rangle$ .  $\pi$  does so (line 58 only) only if, for some  $r$ , we have  $|\text{assignments}[r]| \geq 2f + 1$  at  $\pi$  (line 51). Let  $\hat{\sigma}$  denote the value of *assigner* at  $\pi$ , let  $i = (\hat{\sigma}, r)$ . Immediately before triggering  $\langle \text{dir.SignupComplete} \rangle$ ,  $\pi$  aggregates *assignments*[ $r$ ] into a certificate  $t$  (line 55). By Lemma 128,  $t$  is a quorum certificate for  $[\text{Assignment}, i, \pi]$ .  $\pi$  then invokes

$$\text{dir.import}(\text{Assignment}\{id : i, process : \pi, certificate : t\})$$

(line 57). Upon doing so,  $\pi$  verifies that  $t$  is indeed a certificate for  $[\text{Assignment}, i, \pi]$  (line 81), and adds  $(i, \pi)$  to *directory* (line 82). By Notation 31 and Definition 66,  $\pi$  knows itself, and the theorem is proved.  $\square$

### 7.3.5 Transferability

In this section, we prove that Dibs satisfies transferability.

**Lemma 129.** *Let  $\pi$  be a correct process, let  $i$  be an id, be a process such that  $(i, \_) \in \text{directory}$  at  $\pi$ . We have  $i \in \text{certificates}$  at  $\pi$ .*

*Proof.* Upon initialization, *directory* is empty at  $\pi$  (line 24). Moreover,  $\pi$  adds  $(i, \_)$  to *directory* only by executing line 82. Immediately after doing so,  $\pi$  adds  $i$  to *certificates* (line 83).  $\square$

**Lemma 130.** *Let  $\pi$  be a correct process, let  $(\hat{i}, \hat{\pi}) \in \text{directory}$  at  $\pi$ . Let  $\hat{c} = \text{certificates}[\hat{i}]$  at  $\pi$ . We have that  $\hat{c}$  is a quorum certificate for  $[\text{Assignment}, \hat{i}, \hat{\pi}]$ .*

*Proof.* We underline that, by Lemma 129,  $\hat{c}$  is guaranteed to exist. By Notation 31 and Lemmas 123 and 121, for all  $\hat{\pi}' \neq \hat{\pi}$  we always have  $(\hat{i}, \hat{\pi}') \notin \text{directory}$  at  $\pi$ . Upon initialization, *certificates* is empty at  $\pi$  (line 25). Moreover,  $\pi$  adds  $(\hat{i}, \hat{c})$  to *certificates* only by executing line 83.  $\pi$  does so only if  $\hat{c}$  is a quorum certificate for  $[\text{Assignment}, \hat{i}, \hat{\pi}]$  for some  $\hat{\pi}$  (line 81). Immediately before adding  $(\hat{i}, \hat{c})$  to *certificates*, however,  $\pi$  adds  $(\hat{i}, \hat{\pi})$  to *directory* (line 82). This proves that  $\hat{\pi} = \hat{\pi}$  and concludes the lemma.  $\square$

**Theorem 25.** *Dibs satisfies transferability.*

*Proof.* Let  $\pi, \pi'$  be correct processes, let  $\hat{i}$  be an id such that  $a = \text{dir.export}(\hat{i})$  at  $\pi$  is an assignment. Let

$$\begin{aligned}\tilde{i} &= a.id \\ \hat{\pi} &= a.process \\ \hat{c} &= a.certificate\end{aligned}$$

By lines 75 and 72, we immediately have  $\tilde{i} = \hat{i}$ . Moreover, by line 73 we have  $\hat{c} = \text{certificates}[\hat{i}]$  at  $\pi$ . Finally, by line 74 we have  $(\hat{i}, \hat{\pi}) \in \text{directory}$  at  $\pi$ . By Lemma 130, we then have that  $\hat{c}$  is a quorum certificate for  $[\text{Assignment}, \hat{i}, \hat{\pi}]$ .

Upon invoking  $\text{dir.import}(a)$ ,  $\pi'$  verifies that  $\hat{c}$  is indeed a quorum certificate for  $[\text{Assignment}, \hat{i}, \hat{\pi}]$  (line 81), then adds  $(i, \_)$  to  $\text{directory}$  (line 82). By Notation 31 and Definition 66,  $\pi'$  knows  $\hat{i}$ , and the theorem is proved.  $\square$

### 7.3.6 Density

In this section, we prove that Dibs satisfies density.

**Lemma 131.** *Let  $\sigma$  be a correct server, let  $\hat{\sigma}$  be a server. We have  $|\text{rankings}[\hat{\sigma}]| < |\Pi|$  at  $\sigma$ .*

*Proof.* Upon initialization,  $\text{rankings}[\hat{\sigma}]$  is empty at  $\sigma$  (line 19). Moreover,  $\sigma$  adds a process  $\pi$  to  $\text{rankings}[\hat{\sigma}]$  (line 22 only) only if  $\pi$  is not already in  $\text{rankings}[\hat{\sigma}]$  (line 21). This proves that, at  $\sigma$ , the elements of  $\text{rankings}[\hat{\sigma}]$  are distinct processes. We then have  $|\text{rankings}[\hat{\sigma}]| < |\Pi|$  at  $\sigma$ , and the lemma is concluded.  $\square$

**Lemma 132.** *Let  $(d, n)$  be an id, let  $\pi$  be a process. If a quorum certificate exists for  $[\text{Assignment}, (d, n), \pi]$ , then  $n < |\Pi|$ .*

*Proof.* Let us assume that a quorum certificate exists for  $[\text{Assignment}, (d, n), \pi]$ . Noting that at most  $f$  servers are Byzantine, at least one correct server  $\sigma$  signed  $[\text{Assignment}, (d, n), \pi]$ .  $\sigma$  does so (line 37 only) only if  $n < |\text{rankings}[\hat{\sigma}]|$ , for some server  $\hat{\sigma}$  (line 35). By Lemma 131 we then have  $n < |\Pi|$ , and the lemma is proved.  $\square$

**Theorem 26.** *Dibs satisfies density.*

*Proof.* Let  $\pi$  be a correct process, let  $t \in \mathbb{R}$ , let  $((\hat{d}, \hat{n}), \hat{\pi}) \in D_\pi(t)$ . By Lemma 123 we have  $((\hat{d}, \hat{n}), \hat{\pi}) \in D^\infty$ . By Definition 67 and Lemma 132 we then have  $n < |\Pi|$ , and the theorem is proved.  $\square$



# **Byzantine Atomic Broadcast to the Network Limit**

**Part III**





# 8 Chop Chop

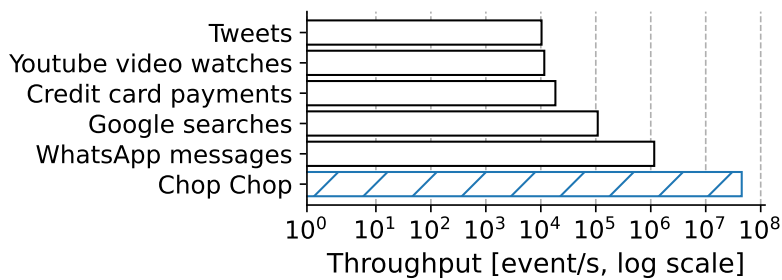
## 8.1 Overview

This Part focuses on Atomic Broadcast, the powerful consensus-equivalent primitive that has a set of processes agree on the (total) order of client-issued messages. Unlike Reliable Broadcast, which can only be applied to a relatively narrow class of useful problems (such as Asset Transfer [80]), Atomic Broadcast powers state machine replication (SMR) [62, 135], enabling decentralized universal computation in the face of arbitrary failures [104, 138].

Despite its versatility, however, Atomic Broadcast comes with fundamental bounds [59] and constraints [68], hindering its real-world performance despite decades of extensive research [10, 22, 36, 43, 49, 99, 108, 121, 155, 156, 159] and attention from industry, where SMR powers a myriad of blockchains and ledgers [8, 78, 98, 102, 107, 116, 143, 144, 145, 152, 153].

When deployed globally, seminal Atomic Broadcast implementations, such as BFT-SMaRt [22] and HotStuff [156], fail to stretch past a few thousand messages per second, three orders of magnitude short of the millions of requests per second collectively handled by the Internet's largest, centralized services (Figure 8.1). As we discussed in the introduction to this thesis, bringing Atomic Broadcast into the tens of millions of messages per second seems a necessary stepping stone towards planetary-scale, decentralized computation.

**Towards line rate.** While slow and expensive, ordering messages in Atomic Broadcast is amenable to *batching* [43]: order once, deliver in bulk. This observation motivated the development of *memory pool* (mempool) protocols [55, 70, 140], as initiated by Narwhal [55], designed to amortize ordering. This strategy proved effective, e.g., Bullshark [140] delivers in the order of 380,000 messages per second when accelerated by Narwhal. Despite this improvement, however, state-of-the-art batching still falls short of achieving *line rate*, i.e., matching the communication complexity of a protocol that does not ensure any ordering, authentication, or Byzantine resilience. In such a simplified setting, a server could simply deliver a sequence of application messages as it receives them from the network:  $b$  bits



**Figure 8.1:** Throughput of Internet-scale services.

received,  $b$  bits delivered. Modern connections have enough bandwidth to receive tens of millions of application messages per second:<sup>1</sup> 2.5 orders of magnitude of gap still exist between Atomic Broadcast and unordered, unauthenticated dissemination. It is natural to ask if such a large gap is inherent to atomicity's unavoidable cost of ordering, authenticating and deduplicating messages. In this Part, we answer in the negative, accelerating Atomic Broadcast by a further two orders of magnitude with a system that performs within 8% of line-rate<sup>2</sup>, even when handling 40 million requests per second.

**Chop Chop.** We present Chop Chop, a Byzantine Atomic Broadcast system. As a mempool, Chop Chop uses an underlying instance of Atomic Broadcast to order batches, aggregating messages to amortize costs. Classic methods of batching, however, fail to also amortize authentication and deduplication: each payload in a batch still has to carry an individual public key, signature and sequence number.

Chop Chop addresses this shortcoming with a new form of batches: *distilled batches*. Unlike a classic batch, a distilled batch contains condensed information that allows to authenticate and deduplicate its messages in bulk, much faster than in existing schemes. Distilled batches leverage the strong ordering of Atomic Broadcast to minimize redundant information.

**Trustless brokers.** Chop Chop produces distilled batches using a novel interactive protocol involving *brokers*, a layer of facilitating processes between clients and servers. Counterintuitively, distilled batches are faster for servers to receive and process, but expensive for brokers to produce: distillation is interactive and relies on expensive cryptographic operations.

Importantly, however, incorrectly distilled batches are visibly malformed. As such, brokers are *untrusted*: good brokers take load off the servers; bad ones cannot compromise the system's safety. Servers are exposed to every message in the system, bottleneck easily, and only a threshold of them can be compromised before the system loses safety. Brokers, instead, can

---

<sup>1</sup> Payloads as small as 12 bytes can have real-world applications (see Section 8.2.1). A 5 Gbit/s link can receive 52 millions such payloads per second.

<sup>2</sup> Up to some minor differences, line-rate can be considered a real-world measure of oracularity, the property we introduced in Part II to describe the ability of a distributed system to match the complexity of its centralized counterpart.

be spun up by anyone, outside of Chop Chop’s security perimeter, to meet the load produced by clients.

**Evaluation.** We evaluate Chop Chop in a cross-cloud, geo-distributed environment including 320 medium-sized AWS EC2 machines and 64 OVH machines. We simulate up to 257 million clients and consider 12 experimental environments. Setting up each environment requires the installation of 13 TB of synthetic workload. A naive installation using `scp` from a single machine would take 68 hours. We designed `silk`, a one-to-many peer-to-peer file transfer tool optimized for high latency connections, to install the files in 30 minutes instead.

We compare Chop Chop’s throughput and end-to-end latency against its baselines in multiple real-world scenarios including server failures, adverse network conditions, and applications running. In all scenarios, Chop Chop’s throughput outperforms its closest competitor by up to two orders of magnitude, with no penalty in terms of latency. When put under stress, Chop Chop orders, authenticates and deduplicates upwards of 43,600,000 messages per second with a mean latency of 3.6 seconds. Except under the most adverse network conditions and proportions of faulty clients, Chop Chop still achieves millions of operations per second.

**Applications.** Unlike most Atomic Broadcast implementations [22, 55, 140, 156], Chop Chop does not offload authentication and deduplication to the application. This allows Chop Chop-based applications to focus entirely on their core logic without ever engaging in expensive, and easy to get wrong, cryptography. To showcase this, we implement three simple applications to evaluate on top of Chop Chop: a Payment system, an Auction house and an instance of the game “Pixel war”. These three simple applications (300 lines of logic) work effectively with messages as small as 8 bytes, further underlying the communication overhead represented by public keys, signatures and sequence numbers in non-distilled systems. Both Payments and Pixel war inherit Chop Chop’s throughput, respectively processing over 32 and 35 million operations per second. Even the Auction house, which is single-threaded, achieves 2.3 million operations per second. (These applications are meant as examples, and further optimization is beyond the scope of this Part.)

**Contributions.** We identify authentication and deduplication as the main bottlenecks of batched Atomic Broadcast; we introduce distilled batches to extend the amortizing properties of batching to authentication and deduplication; we present distillation, an interactive protocol to produce distilled batches, and identify the opportunity to offload it to an untrusted set of brokers; we implement Chop Chop, an Atomic Broadcast system taking advantage of distillation; we thoroughly evaluate Chop Chop, improving state-of-the-art Atomic Broadcast throughput by two orders of magnitude, maintaining near line-rate performance up to 40 million requests per second; we showcase Chop Chop through a Payment system, an Auction house and an instance of the “Pixel war” game, respectively achieving 32, 2.3 and 35 million operations per second.

**Roadmap.** Section 8.2 introduces Atomic Broadcast, discusses classic batching mechanisms and highlights the cost of authenticating and deduplicating messages in the resulting batches. Section 8.3 presents distilled batches and, for pedagogical purposes, introduces a simplified failure-free version of Chop Chop’s protocol. Section 8.4 describes Chop Chop’s fault-tolerant protocol in further detail. Section 8.5 discusses Chop Chop’s implementation. Section 8.6 discusses Chop Chop’s empirical evaluation, highlighting the challenges of such large scale experiments. We summarize related work in Section 8.7.

## 8.2 Atomic Broadcast

In an Atomic Broadcast system [45], *clients* broadcast messages that are delivered by *servers*.

**Properties [34].** Correct servers deliver the same messages in the same order (*agreement*). Messages from correct clients are eventually delivered (*validity*). Spurious messages cannot be attributed to correct clients (*integrity*). No message is delivered more than once (*no duplication*).

### 8.2.1 Cost of Atomic Broadcast

Informally, Atomic Broadcast’s most distinctive property, agreement, is also the most challenging to satisfy. Correct servers must coordinate to *order* messages without compromising liveness. A great deal of research effort has been put in developing ordering techniques, optimizing for latency [99, 114] or communication complexity [113, 124].

Integrity and no duplication, instead, allow for simple solutions. Clients can ensure integrity by *authenticating* their messages using digital signatures: servers simply ignore incorrectly authenticated messages. For no duplication, clients can tag each message with a strictly increasing *sequence number*: after ordering, servers discard old messages as replays.

Both techniques—we call them *classic authentication* and *classic sequencing*—are non-interactive, easy to implement, and agnostic of the protocol employed to order messages. Arguably due to the simplicity and effectiveness of classic authentication and sequencing, most Atomic Broadcast implementations overlook integrity and no duplication entirely: they offload authentication and sequencing to the application, focusing on the more challenging task of ordering.

**Batching for ordering.** Lacking an efficient technique to minimize its complexity, ordering could be Atomic Broadcast’s main bottleneck.<sup>3</sup> The well-known strategy of *batching*, however, is both general and effective at amortizing the agreement cost of an Atomic Broadcast

---

<sup>3</sup>Byzantine Atomic Broadcast among  $n$  participants cannot be achieved with a bit complexity smaller than  $\Theta(n^2)$  [59].

implementation [43, 138].

Broadly speaking, batching is orchestrated by a *broker* as follows [55].<sup>4</sup> Over a small window of time, the broker collects multiple client-issued messages in a batch, which it disseminates to the servers; the broker then submits to an underlying instance of Atomic Broadcast a cryptographic hash of the batch it collected; upon delivering the hash of a batch from Atomic Broadcast, a server retrieves the batch, and delivers to the application all the messages it contains. Because the size of a hash is constant, the cost of ordering a batch does not depend on its size: as batches become larger, the cost of ordering each message goes to zero. In practice, batching can effectively eliminate the cost of ordering in any real-world implementation of Atomic Broadcast.

**Cost of integrity and no duplication.** Batching does not efficiently uphold integrity and no duplication. Regardless of how many messages are batched together, the cost of classic authentication and sequencing stays constant: one public key, one signature and one sequence number for each message.

In practice, these costs dominate the computation and communication budget of a batched Atomic Broadcast system (see Section 8.3.2). On the one hand, signatures are among the most CPU-intensive items in the standard cryptographic toolbox, dwarfing in particular symmetric primitives such as hashes and ciphers. On the other, public keys, signatures and sequence numbers can easily account for the majority of a batch's size.

To illustrate these costs, consider the example of a payment system. A payment operation requires three fields: sender, recipient, and amount. Sender and recipient fit in 4 B each if the system serves less than 4 billion users. Amount needs 4 B for payments between 1 cent and 40 millions. Hence, a payment can be encoded in just 12 B. Using public keys to identify sender and recipient ( $2 \times 32$  B using Ed25519 [21, 90]) and attaching a signature (64 B) and a sequence number (8 B) to each message inflates payloads to 140 B. For payments, *91% of the bandwidth is spent on integrity and no duplication.*

### 8.2.2 Existing Mitigations

Chop Chop integrates the two following techniques to reduce the bandwidth and CPU cost of authentication.

**Short identifiers.** Repeated public keys consume a significant slice of a server's communication budget. A workaround is to have servers store public keys in an indexed *directory* [4]. Upon first joining the system, a client announces its public key via Atomic Broadcast to *sign up*. Upon delivering a sign-up message, a server appends the new public key to its directory.

---

<sup>4</sup>In the literature, servers usually play the role of brokers. As we discuss in Section 8.4, however, Chop Chop minimizes its load on the servers by offloading brokerage to a separate, trustless set of processes.

## Chapter 8. Chop Chop

---

The same public key appears at the same position in the directory of all correct servers thanks to Atomic Broadcast's agreement. Having signed up, a client uses its position in the directory as identifier instead of its public key.

In the previous example of a payment system, using such identifiers reduces a payment size by 40%, from 140 B to 84 B. However, a signature per payment must still be transmitted.

**Pooled signature verification.** Authenticating a batch by verifying its signatures is a computationally intensive task for a server [43, 142]. However, Red Belly[53] and Mir [142] showed that not all servers need to authenticate all batches. Indeed, assuming at most  $f$  faulty servers, a broker optimistically asks only  $f + 1$  servers to authenticate a batch to be certain to reach at least one correct server. If  $f + 1$  servers do not reply by a timeout, the broker extends its request to  $f$  additional servers, thus reaching at least  $f + 1$  correct servers.

A correct server that authenticates a batch sends back to the broker a *witness shard*, i.e., a signed statement that the batch is correctly signed. The broker aggregates  $f + 1$  identical shards into a *witness*, which it sends to the other  $2f$  servers. Because every witness contains at least one correct shard, the servers can trust the witness instead of verifying the batch.

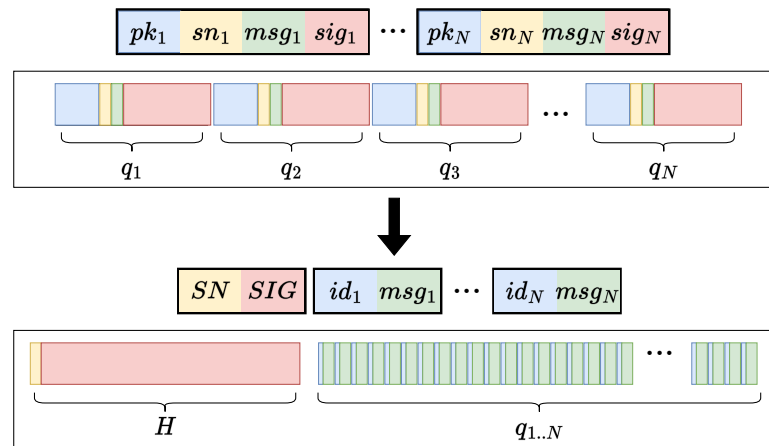
Assuming  $3f + 1$  servers, this technique shaves up to two-thirds off the system's authentication complexity.

### 8.3 Distilled Batches

Chop Chop's main contribution is *distillation*, a set of techniques aimed at extending the amortizing properties of batches to authentication and sequencing.

**Background: multi-signatures.** Chop Chop makes use of multi-signature schemes [88] to authenticate batches. Secret keys produce signatures that can be verified against the corresponding public keys. Public keys and signatures, however, can be *aggregated*. Let  $(p_1, r_1), \dots, (p_n, r_n)$  be distinct key pairs, and  $s_1, \dots, s_n$  be signatures produced by  $r_1, \dots, r_n$  on the *same* message  $m$ :  $p_1, \dots, p_n$  (resp.,  $s_1, \dots, s_n$ ) can be aggregated into a constant-sized aggregate public key  $p$  (resp., aggregate signature  $s$ ).

Remarkably,  $s$  can be verified in constant time against  $p$  and  $m$  [25, 115]. Chop Chop uses BLS multi-signatures [25] which can be aggregated cheaply and non-interactively: even a non-signing process can compute  $p$  (resp.,  $s$ ) once provided with  $p_1, \dots, p_n$  (resp.,  $s_1, \dots, s_n$ ) by computing a single multiplication over an elliptic curve.



**Figure 8.2: Full distillation in action.** With classic authentication and sequencing, each payload  $q_i$  contains a public key  $pk_i$ , a sequence number  $sn_i$ , a message  $msg_i$  and a signature  $sig_i$ . In the fully distilled case, each  $q_i$  reduces to just  $id_i$  and  $msg_i$ : one header  $H$ , composed of one aggregate sequence number  $SN$  and one aggregate signature  $SIG$ , is sufficient for the entire batch. Bars are to scale if small messages are broadcast using Ed25519 for signatures and BLS12-381 for uncompressed multi-signatures:  $sn_i$  and  $SN$  are 8 B,  $msg_i$  is 8 B,  $pk_i$  is 32 B,  $sig_i$  is 64 B,  $SIG$  is 192 B.

### 8.3.1 Distillation at a Glance

In brief, distillation aims to produce *distilled batches*. A distilled batch has some of its signatures (resp., sequence numbers) replaced by an *aggregate signature* (resp., *aggregate sequence number*). When maximally successful, distillation produces a *fully distilled batch*, where all signatures (resp., sequence numbers) have been replaced by a *single* aggregate signature (resp., sequence number). As we discuss below, distilled batches are vastly cheaper for servers to receive and process. Figure 8.2 depicts the effect of distillation on a batch.

**Full distillation (failure-free).** For pedagogical purposes, we introduce distillation under the assumption that all processes are correct. We detail Chop Chop’s fault-tolerant distillation techniques in Section 8.4.2, optimized and adapted to the Byzantine setting. As in the classic batching case, a set  $\chi_1, \dots, \chi_b$  of clients submit their messages  $m_1, \dots, m_b$  to a broker  $\beta$ . Each  $\chi_i$  selects for its message  $m_i$  a sequence number  $k_i$  (greater than any sequence number it previously used), then sends  $(k_i, m_i)$  to  $\beta$ . Upon receiving all  $(k_i, m_i)$ -s,  $\beta$  computes the aggregate sequence number

$$k = \max_i k_i$$

then builds the *batch proposal*

$$B = [(x_1, k, m_1), \dots, (x_b, k, m_b)]$$

## Chapter 8. Chop Chop

---

where  $x_i$  is  $\chi_i$ 's numerical identifier in the system (see Section 8.2.2).  $\beta$  then sends  $B$  back to every  $\chi_i$ . Upon receiving  $B$ ,  $\chi_i$  produces a multi-signature  $s_i$  for the hash  $H(B)$  of  $B$ , which it sends back to  $\beta$ . Having collected all multi-signatures,  $\beta$  computes the aggregate signature

$$s = \prod_i s_i$$

In doing so,  $\beta$  obtains the fully distilled batch

$$\tilde{B} = [s, k, ((x_1, m_1), \dots, (x_b, m_b))]$$

Upon receiving  $\tilde{B}$ , any server now can: compute  $B$  by inserting  $k$  between each  $(x_i, m_i)$ ; compute  $H(B)$ ; use each  $x_i$  to retrieve  $\chi_i$ 's public key  $p_i$  from its directory; compute the aggregate public key

$$p = \prod_i p_i$$

and finally verify  $s$  against  $p$  and  $H(B)$ .

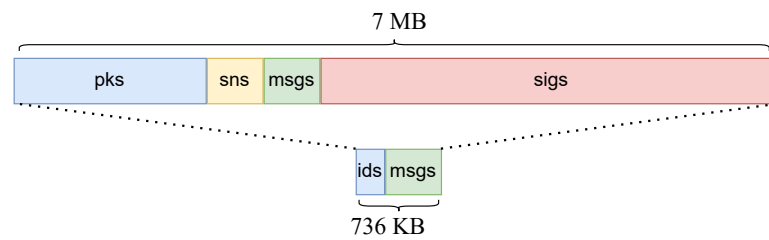
**Distillation outcome.** Having engaged with  $\beta$  to distill the batch, every  $\chi_i$  multi-signs the *same* message  $H(B)$  and updates its sequence number to the *same*  $k$ . This allows  $\beta$  to authenticate and sequence all of  $\tilde{B}$  using  $s$  and  $k$  only.

**Distillation safety.** The proposed distillation protocol has no safety drawback. First, because  $(x_i, k, m_i)$  appears in  $B$ ,  $\chi_i$  still gets to authenticate  $m_i$ . Intuitively,  $\chi_i$ 's multi-signature on  $H(B)$  publicly authenticates *whatever message in  $B$  is attributed to  $\chi_i$ ,  $m_i$*  in this case. Second, because  $k \geq k_i$ ,  $k$  is still a valid sequence number for  $m_i$ . Sequence number distillation might cause  $\chi_i$  to skip some sequence numbers whenever any  $\chi_j$  issues some  $k_j > k_i$ . Contiguity of sequence numbers, however, is not a requirement for deduplication. As with classic sequencing,  $\chi_i$  produces—and servers deliver—messages with strictly increasing sequence numbers; servers disregard all other messages as replays.

### 8.3.2 Distillation Microbenchmark

Having discussed how distilled batches are produced, we now estimate the significance of their effect by means of a back-of-the-envelope calculation and a simple microbenchmark on AWS. Consider a setting where 100 million clients broadcast 8-byte messages, e.g., to issue payments (see Section 8.2.1). We compare *classic authentication and sequencing*, where clients are identified by their public keys, messages are individually signed and sequenced, against *fully distilled batches* where clients are identified by a numerical identifier and each batch contains only one aggregated signature and sequence number. We use Ed25519 [90] for signatures (32 B public keys, 64 B signatures) and BLS12-381 [27] for multi-signatures (192 B uncompressed signatures). We use uncompressed BLS multi-signatures to save computation





**Figure 8.3: Full distillation of a batch of 65,536 payloads (sizes to scale).** The aggregate signature and aggregate sequence number do not appear as a result of their small size.

time at the cost of storage space (96 B compressed vs. 192 B uncompressed).

**Communication complexity.** Payloads are 112 B per message in the classic case (32 B of public key, 8 B of sequence number, 8 B of message, 64 B of signature) vs. 11.5 B in the fully distilled case (28 bits = 3.5 B of identifier to represent 257M clients, 8 B of message). Assuming batches of 65,536 messages (Figure 8.3), classic batches are exactly 7 MB long, while fully distilled batches are 736 KB long including aggregate signature and sequence number.

**Computation complexity.** Running at maximum load, an Amazon EC2 c6i.8xlarge instance authenticates  $16.2 \pm 0.4$  classic batches per second using Ed25519's batch verification for 65,536 signatures. The same machine authenticates  $457.1 \pm 0.3$  fully distilled batches per second: each authentication requires the aggregation of 65,536 BLS12-381 public keys and the verification of one BLS12-381 multi-signature.

**Summary.** By the order-of-magnitude calculations above, fully distilled batches hold the promise to reduce the costs of authentication and sequencing by a factor 9.7 for network bandwidth, and 28.2 for CPU. Chop Chop aims to deliver on that promise for a real-world fault-tolerant system.

## 8.4 Chop Chop

This section overviews Chop Chop's architecture and protocol, and provides arguments for its correctness.

**Overview.** Chop Chop involves three types of processes (Figure 8.4): broadcasting clients, delivering servers and a layer of broadcast-facilitating brokers between them. Servers run an Atomic Broadcast instance among themselves, to which brokers submit messages. Chop Chop is agnostic to the implementation of Atomic Broadcast used by the servers. On top of the provided broker-to-server Atomic Broadcast, Chop Chop implements a much faster client-to-server Atomic Broadcast: clients submit messages to the servers, aided by brokers.

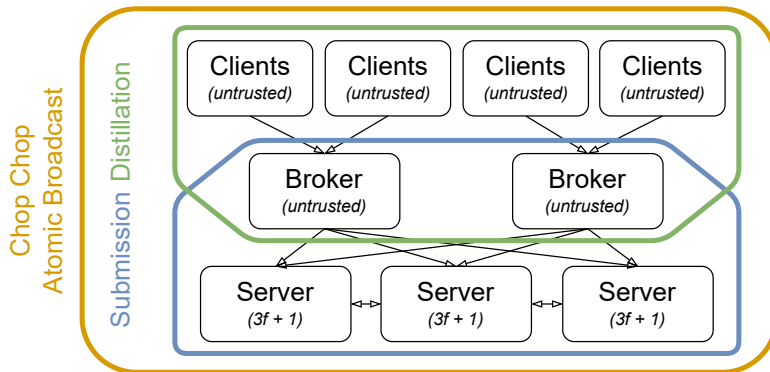


Figure 8.4: Chop Chop architecture.

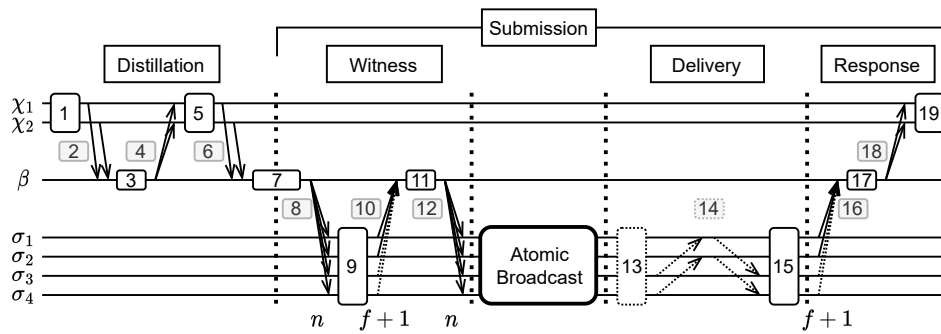
Chop Chop’s protocol unfolds in two phases: *distillation* (Section 8.4.2) and *submission* (Section 8.4.3). In the distillation phase, clients interact with a broker to gather their messages in a distilled batch (see Section 8.3). In the submission phase, the broker disseminates the distilled batch to the servers and submits the batch’s hash to the server-run instance of Atomic Broadcast. Upon delivering its hash from Atomic Broadcast, servers retrieve the batch and deliver its messages. Chop Chop’s contributions mainly focus on the distillation phase. Chop Chop’s submission strategy closely resembles prior batch-based Atomic Broadcast implementations [55, 70, 140].

### 8.4.1 Architecture and Model

Chop Chop augments the architecture of a classic Atomic Broadcast, as described in Section 8.2, with novel brokers.

**Clients and servers.** *Clients* broadcast messages to a (distinct) set of *servers*. We assume that less than one third of servers can be faulty and behave in an arbitrary manner, i.e., be Byzantine [104], while all clients can be faulty. For simplicity, servers form a fixed set that is known by all correct processes at system startup. Chop Chop can be extended for reconfiguration thanks to its modular use of Atomic Broadcast [22, 103] (Figure 8.4). Clients issue messages after broadcasting their public keys to the system (see Section 8.2.2).

**Brokers.** We discussed in Section 8.3 how both classic and distilled batches are assembled by a broker. The role of brokers is traditionally taken by servers. Given the additional strain put on brokers by Chop Chop’s interactive distillation protocol, however, having servers be brokers would result in a waste of scarce, trusted resources. Importantly, however, distillation is *trustless*. On the one hand, agreement rests entirely on Chop Chop’s underlying Atomic Broadcast instance, for which brokers are only clients. On the other hand, as we argue in Sections 8.4.2 and 8.4.4, a faulty broker cannot compromise integrity or no duplication:



**Figure 8.5: Overview of the Chop Chop protocol between two clients ( $\chi_1, \chi_2$ ), a broker ( $\beta$ ) and four servers ( $\sigma_1$ – $\sigma_4$ ).** The protocol is comprised of 19 steps (#1–#19) and of an underlying instance of Atomic Broadcast such as BFT-SMaRt or HotStuff.

distilled batches are publicly authenticated, and correct clients cannot be tricked into using stale sequence numbers. Hence, *brokers need no trust*: a broker either does its job correctly or produces distilled batches that are visibly malformed, and easily discarded by all correct servers.

This observation is of paramount importance to the performance of Chop Chop: *because distillation is heavy but trustless, brokers should be distinct from servers*. Along with clients and servers, we thus assume a third, *independent set of brokers*, sitting between clients and servers, to accelerate Atomic Broadcast by assembling client messages in distilled batches. We assume that at least one broker is correct; the system loses liveness but not safety if all brokers are faulty.

**Network.** Chop Chop guarantees that the batches collected and submitted to servers by correct brokers are well-formed even in asynchrony, but achieves full distillation when the network is synchronous (see Section 8.4.2). Chop Chop inherits the network requirements of its underlying Atomic Broadcast.

### 8.4.2 Distillation Phase

We introduced in Section 8.3 a simplified, failure-free distillation protocol. This section describes how Chop Chop renders distillation tolerant to arbitrary failures and improves its performance via a sequence of improvements, each addressing a shortcoming of the simplified protocol. The complete fault-tolerant protocol of Chop Chop is depicted in Figure 8.5.

In the failure-free distillation protocol: clients  $\chi_1, \dots, \chi_b$  send their messages  $m_1, \dots, m_b$ , with sequence numbers  $k_1, \dots, k_b$  (#2) to a broker  $\beta$  (Figure 8.5, #1);  $\beta$  identifies the maximum submitted sequence number  $k$  and builds a batch proposal  $B = [(x_1, k, m_1), \dots, (x_b, k, m_b)]$  (#3);  $\beta$  disseminates  $B$  to  $\chi_1, \dots, \chi_b$  (#4); each  $\chi_i$  produces a multi-signature  $s_i$  on  $H(B)$  (#5), which it sends to back  $\beta$  (#6);  $\beta$  aggregates  $s_1, \dots, s_n$  into an aggregate  $s$ , thus producing a fully

## Chapter 8. Chop Chop

---

distilled batch  $\tilde{B} = [s, k, ((x_1, m_1), \dots, (x_b, m_b))]$  (#7).

**Background: Merkle trees.** Chop Chop uses Merkle trees [117] to hash batches. An  $l$ -element vector  $z_1, \dots, z_l$  is hashed into a *root*  $r$ , used as commitment. For each  $i$ ,  $z_i$ 's value can be proved by means of a proof of inclusion  $p_i$ , verifiable against  $r$  and  $z_i$ . Proofs of inclusions are  $O(\log l)$  in size and are verified in  $O(\log l)$  time.

**What if a broker forges messages?** A faulty  $\beta$  could try to falsely attribute to some  $\chi_i$  a message  $m'_i \neq m_i$ .  $\beta$  could do so by replacing  $m_i$  with  $m'_i$  in  $B$ , then having  $\chi_i$  sign  $H(B)$ , thus implicitly authenticating  $m'_i$ . This is easily fixed by having  $\chi_i$  check that  $m_i$  correctly appears in  $B$  before signing  $H(B)$ .

**Can a broker avoid sending the entire batch?** A clear inefficiency of the simplified protocol is that  $\beta$  has to convey all of  $B$  back to each  $\chi_i$ . This is fixed using Merkle trees. Upon assembling  $B$ ,  $\beta$  computes the Merkle root  $r$  of  $B$ , along with the Merkle proof  $p_i$  for each  $(x_i, k, m_i)$  in  $B$ . Instead of sending  $B$  to all clients,  $\beta$  just sends  $r$ ,  $k$  and  $p_i$  to each  $\chi_i$ . Upon receiving  $r$ ,  $k$  and  $p_i$ ,  $\chi_i$  checks  $p_i$  against  $r$  and  $(x_i, k, m_i)$ , producing  $s_i$  on  $r$  only if the check succeeds. If  $\chi_i$  signs  $r$ , then  $(x_i, k, m_i)$  is necessarily an element of  $B$ . Importantly, however,  $\beta$  could inject  $(x_i, k, m'_i \neq m_i)$  somewhere else in  $B$ , while still providing  $\chi_i$  only with the proof for  $(x_i, k, m_i)$ . This is solved by having servers ignore every distilled batch where two or more messages are attributed to the same client. This way, if  $\chi_i$  signs  $r$ , then either  $m_i$  is the only message in  $B$  attributed to  $\chi_i$ , or  $\tilde{B}$  is rejected by all servers as malformed: either way, integrity is upheld.

**What if a client does not multi-sign?** Under the assumption that  $\chi_1, \dots, \chi_b$  are correct,  $\beta$  can safely wait until it collects all  $s_1, \dots, s_b$ . This policy is clearly flawed in the Byzantine setting: a single crashed client can prevent  $\beta$  from ever aggregating  $s$ . Furthermore, lacking an assumption of synchrony,  $\beta$  cannot exclude from  $\tilde{B}$  those clients that do not sign  $r$  by some timeout: consistently slow clients would always be excluded, and validity would be lost. This issue is fixed by the fallback mechanism introduced in the following.

**Fault-tolerant distillation.** Upon first sending  $(k_i, m_i)$  to  $\beta$  (#2),  $\chi_i$  also sends an individual, non-aggregable signature  $t_i$  for  $(x_i, k_i, m_i)$ , which  $\beta$  stores.  $\beta$  then waits for  $s_i$ -s on  $r$  until either all  $s_i$ -s are collected, or a timeout expires. For every  $s_i$  that ends up missing, due to  $\chi_i$  being crashed or delayed,  $\beta$  attaches  $(k_i, t_i)$  to  $\tilde{B}$ . Upon receiving  $\tilde{B}$ , a server first checks each individual signature  $t_i$  against the corresponding  $(x_i, k_i, m_i)$ . The server then checks  $s$  against the public keys of the clients for which an individual signature  $t_i$  was not given, i.e., the public keys of all clients that signed  $r$  in time.

In summary: fast, correct clients who successfully produce their  $s_i$ -s in time authenticate their message by multi-signing  $r$ ; slow or crashed clients still get their messages through,

individually authenticated by the  $t_i$ -s that they originally produced. Full distillation is achieved whenever the network is synchronous and all clients are correct, which we argue is the case in practice for the majority of a system's lifetime. When the network is asynchronous, however, a fraction of clients might fail to produce their  $s_i$  in time, resulting in a *partially distilled batch*. At the limit where all clients fail to sign  $r$  in time,  $\bar{B}$  reduces to a classic batch, degrading server-side performance to pre-distillation levels. We underline that safety and liveness are preserved regardless of synchrony.

**What if a broker replays messages?** A problem introduced by the last fix is that  $\chi_i$  authenticates both  $k_i$  and  $k$  as sequence numbers for  $m_i$ , allowing a faulty  $\beta$  to play  $m_i$  twice, hence breaking Atomic Broadcast's no duplication. This is fixed by having each client engage in the broadcast of only one message at a time. This way, while  $\beta$  can replay  $m_i$ , it can only do so consecutively: all sequence numbers  $\chi_i$  authenticates for  $m_i$  belong to a range that does not contain sequence numbers for any other message  $m_{i' \neq i}$  issued by  $\chi_i$ .

This observation is key to the following fix: along with the last sequence number  $\bar{k}_\chi$  each client  $\chi$  used, a correct server  $\sigma$  stores the last message  $\bar{m}_\chi$  that  $\chi$  broadcast; upon ordering a message  $m$  with sequence number  $k$  from  $\chi$ ,  $\sigma$  delivers  $m$  if and only if  $k > \bar{k}_\chi$  and  $m \neq \bar{m}_\chi$ . In doing so,  $\sigma$  discards all consecutive replays of  $\bar{m}_\chi$ , thus preventing replays in general.

**What if a client broadcasts too frequently?** The last fix relies on clients broadcasting one message at a time. Depending on latency, a client broadcasting too frequently might accrue an ever-growing queue of pending messages. This issue is fixed by flushing application messages in bursts, akin to Nagle's buffering algorithm for TCP.

**What if a client submits the largest possible sequence number?** Assuming that a finite number of bits (e.g., 64) are allocated to representing sequence numbers, a faulty client  $\chi_m$  could set its  $k_m$  to the largest possible sequence number  $k_{max}$  (e.g.,  $2^{64} - 1$ ). In doing so,  $\chi_m$  would force all other  $\chi_i$ -s to update their sequence number to  $k_{max}$ . Since correct clients only use strictly increasing sequence numbers, no  $\chi_i$  could ever broadcast again: sequence numbers would run out. Proving the *legitimacy* of sequence numbers fixes this issue.

**Legitimate sequence numbers.** By the rule that we established, no more than one message from the same client can appear in the same batch. Moreover, correct clients always tag their messages with the smallest sequence number they have not yet used, i.e., the largest they have used plus one. By induction, we then have that unless some client misbehaves, no client ever needs to use a sequence number larger than the number of batches ever delivered by the servers: the largest sequence number any client submits to the very first batch is 0, therefore no client submits a sequence number larger than 1 to the second batch, and so on. This observation allows us to define as *legitimate* any sequence number smaller than the number

## Chapter 8. Chop Chop

---

of batches servers have delivered at any given time.

**Legitimacy proofs.** This definition of legitimacy allows for the generation of *legitimacy proofs*: upon delivering the  $n$ -th batch, a server publicly states so with a signature. By collecting  $f + 1$  server signatures stating that the  $n$ -th batch was delivered into a certificate  $l_n$ , any process can publicly prove that any sequence number smaller than  $n$  is legitimate.

Upon initially submitting  $k_i$  (#2),  $\chi_i$  also sends to  $\beta$  a certificate  $l_n$ , for any  $n > k_i$ ;  $\beta$  ignores client submissions that lack such certificate, except when  $k_i = 0$  since no certificate is needed. Upon sending  $k$  back to all  $\chi_i$ -s (#4),  $\beta$  attaches the highest  $l_{\hat{n}}$  it collected:  $l_{\hat{n}}$  proves that  $k$  is legitimate since  $\hat{n} > k$ .  $\chi_i$  signs  $r$  (#5) only if  $k$  is proved legitimate by  $l_{\hat{n}}$ .

This technique ensures correct clients always use legitimate sequence numbers. Because legitimate sequence numbers grow only with the number of batches delivered by the servers, no correct client is forced to skip too far ahead, compromising its own liveness.

**What if a broker crashes?** If  $\beta$  fails to engage in the protocol, each  $\chi_i$  can submit its message to any other broker.

### 8.4.3 Submission Phase

The submission phase ensures that all servers efficiently deliver a distilled batch, and that all broadcasting clients receive a proof that their messages were delivered.

**Witness.** Having gathered a distilled batch  $\tilde{B}$  (#7),  $\beta$  moves on to have  $f + 1$  servers sign a *witness shard* for  $\tilde{B}$ . In signing a witness shard for  $\tilde{B}$ , a server  $\sigma$  simultaneously makes two statements. First,  $\tilde{B}$  is *well-formed*:  $\sigma$  successfully verified  $\tilde{B}$ 's signatures and found all messages in  $\tilde{B}$  to have a different sender. Second,  $\tilde{B}$  is *retrievable*:  $\sigma$  stores  $\tilde{B}$  and makes it available for retrieval, should any other server need it. We call a *witness* for  $\tilde{B}$  the aggregation of  $f + 1$  witness shards for  $\tilde{B}$ . Because any set of  $f + 1$  processes includes a correct process, when presented with a witness for  $\tilde{B}$  any server can trust  $\tilde{B}$  to be well-formed and retrievable.

As discussed in Section 8.2.2, witnesses optimize server-side computation. Only  $f + 1$  servers need to engage in the expensive checks required to safely witness  $\tilde{B}$ . All other servers can trust  $\tilde{B}$ 's witness, saving trusted CPU resources.

In order to collect a witness for  $\tilde{B}$ ,  $\beta$  sends  $\tilde{B}$  to all servers (#8). Optimistically,  $\beta$  asks only  $f + 1$  servers to sign a witness shard for  $\tilde{B}$ , progressively extending its request to  $2f + 1$  servers upon expiration of suitable timeouts. Upon receiving  $\tilde{B}$  (#9), a correct server  $\sigma$  stores  $\tilde{B}$ . If asked to witness  $\tilde{B}$ ,  $\sigma$  checks that  $\tilde{B}$  is well-formed and sends back to  $\beta$  its witness shard for  $\tilde{B}$  (#10).  $\beta$  collects and aggregates  $f + 1$  shards into a witness for  $\tilde{B}$  (#11), then submits  $\tilde{B}$ 's hash and witness to the server-run Atomic Broadcast (#12).

**Delivery.** Upon delivering  $\tilde{B}$ 's hash and witness from Atomic Broadcast (#13), a correct server  $\sigma$  retrieves  $\tilde{B}$ , either from its local storage (if it directly received  $\tilde{B}$  from  $\beta$  at #8) or from another server (#14). Because  $\tilde{B}$  is retrievable,  $\sigma$  is guaranteed to eventually find a server to pull  $\tilde{B}$  from. Having retrieved  $\tilde{B}$  (#15),  $\sigma$  delivers all non-duplicate messages in  $\tilde{B}$  (see Section 8.4.2 for how  $\sigma$  detects duplicates).

**Response.** Finally,  $\sigma$  signs a *delivery certificate*, listing the messages in  $\tilde{B}$  that  $\sigma$  delivered.  $\sigma$  sends its signature back to  $\beta$  (#16). By agreement of Atomic Broadcast, all correct servers deliver the same subset of messages in  $\tilde{B}$ . As such,  $\beta$  is guaranteed to eventually collect  $f + 1$  signatures on the same delivery certificate (#17). Upon doing so,  $\beta$  distributes a copy of  $\tilde{B}$ 's delivery certificate to  $\chi_1, \dots, \chi_b$  (#18). Armed with  $\tilde{B}$ 's delivery certificate, a correct  $\chi_i$  can publicly prove the delivery of  $m_i$  (#19) and safely broadcast its next message.

#### 8.4.4 Correctness

This section summarizes Chop Chop's correctness analysis.

##### Safety

The safety of Chop Chop is given by its agreement, integrity and no duplication properties (see Section 8.2).

**Agreement.** Chop Chop inherits agreement from its underlying, server-run instance of Atomic Broadcast. A correct server delivers messages only upon delivering the hash of a batch from the server-run Atomic Broadcast. Upon doing so, a correct server retrieves the full batch, checks its hash, and delivers all its messages in order of appearance. All correct servers deliver the same messages in the same order assuming cryptographic hashes are collision-resistant.

**Integrity.** A correct server only delivers messages included in a batch witnessed by  $f + 1$  servers, i.e., by at least one correct server. A correct server witnesses a batch only if: no more than one message in the batch is attributed to the same client; every client in the batch authenticates its message with a signature or the root of the batch's Merkle tree with a multi-signature. A correct client multi-signs the root of a batch's Merkle tree only upon receiving a proof of the inclusion of its message in the batch. As such, if a correct client multi-signs the root of a batch's Merkle tree, either the batch contains only the client's intended message or it is not witnessed. In summary, a correct server delivers a message  $m$  from a correct client  $\chi$  only if  $\chi$  broadcast  $m$ .

## Chapter 8. Chop Chop

---

**No duplication.** A correct client only broadcasts one message at a time. As such, while the client might attach multiple sequence numbers to the same message (different brokers may propose different aggregate sequence numbers for the client to authenticate) the sequence numbers the client attaches to each message belong to distinct ranges. A correct server delivers client messages only in increasing order of sequence number, and ignores repeated messages. This means that a correct server delivers at most one message from each sequence number range. In summary, no server delivers a correct client's message more than once.

### Liveness

The liveness of Chop Chop is given by its validity property.

**Validity.** If a correct client submits its message to a correct broker, the message is guaranteed to eventually be delivered by all correct servers: even if the client fails to engage in distillation in a timely manner, its message is still included in a batch which gets disseminated, witnessed and delivered by all correct servers. Faulty brokers can clearly refuse to service (specific) clients. Upon expiration of a suitable timeout, however, a correct client submits its message to a different broker. As we assume that at least one broker is correct, all correct clients are eventually guaranteed to find a correct broker and get their messages delivered by all correct servers.

### Other Attacks

As we outlined in Section 8.4.4, Chop Chop satisfies all properties of Atomic Broadcast. In this section, we consider other attacks an adversary might deal to impair Atomic Broadcast's performance and fairness [93] in Chop Chop.

**Denial of service.** A faulty broker may refuse to service clients, thus forcing them to fall back on other brokers, increasing latency. A faulty broker may also submit deliberately non-distilled batches to servers to force them to waste trusted resources to receive and verify individual signatures. While handling DoS is beyond the scope of this Part, Chop Chop is amenable to accountability mechanisms [83]. Brokers could be asked to stake resource to join the system. Correct, high-performance brokers could be rewarded, akin to gas fees in Ethereum [153]. Brokers that accrue a reputation of misbehavior or slowness could be banned and lose their initial stake.

**Front-running.** A faulty broker might impact fairness by front-running messages of interest [54, 160]. While front-running resistance is beyond the scope of this Part, Chop Chop is compatible as-is with existing mechanisms to mitigate or prevent front-running, most notably schemes that have clients submit encrypted messages whose content is revealed only after



delivery [119, 158]. Importantly, these encrypt-order-reveal schemes could be selectively employed only for those messages that are vulnerable to front-runs, e.g., messages used for stock trading [128]. Maintaining Chop Chop's throughput while providing quorum-enforced fairness for every message [159] opens a valuable future avenue of research.

### 8.5 Implementation Details

A straightforward implementation of the protocol we presented in Section 8.4 would not achieve the throughput and latency we observe in Section 8.6. In this section, we discuss some of the techniques and optimizations required on the way to practically achieving Chop Chop's full potential. (Many optimizations are however left out due to space constraints).

**Code.** Chop Chop is implemented in Rust, totaling 8,900 lines of code. The main libraries Chop Chop depends on are: `tokio` 1.12 for an asynchronous, event-based runtime; `rayon` 1.5 for worker-based parallel computation; `serde` 1.0 for serialization and deserialization; `blake3` 1.0 for cryptographic hashes; `ed25519-dalek` 1.2 for EdDSA signatures on Curve25519 [90]; `blst` 0.3.5 for multi-signatures on the BLS12-381 curve [27]. Chop Chop also depends on in-house libraries: `talk` (9,800 lines of code) for basic distributed computing and high-level networking and cryptography; `zebra` (7,100 lines of code) for Merkle-tree based data structures.

#### 8.5.1 Broker

The goal of a Chop Chop broker is to produce batches as distilled as possible (to minimize server load), as large as possible (to amortize ordering), and as quickly as possible (to minimize latency). Our target is for a broker to assemble one fully distilled batch of 65,536 messages (736 KB, see Figure 8.3) per second, with a 1 second distillation timeout.

**Reliable UDP.** Short-lived TCP connections between broker and clients are easier to work with, but unfeasible for the broker to handle. Assuming an end-to-end broadcast time of up to 10 seconds, the broker would need to maintain upwards of 600,000 simultaneous TCP connections, which preliminary tests immediately proved unfeasible on the hardware we have access to. This makes UDP the only option for client-broker communication. However, UDP lacks the reliability properties of TCP, and tests showed non-negligible packet loss even within the same AWS EC2 availability zone. As we discussed in Section 8.4.2, message loss immediately translates to partial distillation. We address this issue by means of an in-house, ACK-based, message retransmission protocol based on UDP that also smoothens the rate of outgoing packets.

## Chapter 8. Chop Chop

---

**EdDSA batch verification.** To avoid spoofing, all client messages are authenticated with signatures. At the target rate, however, individually verifying each signature is unfeasible for a broker. Luckily, `ed25519-dalek` allows for more efficient batched verification. A broker buffers the client messages it receives and authenticates them in batches.

**Tree-search invalid multi-signatures.** Clients contributing to the same batch produce matching multi-signatures for the batch's root. At the target rate the broker cannot independently verify each multi-signature. We tackle this problem by gathering multiple matching multi-signatures on the leaves of a binary tree: internal nodes aggregate their children. For each tree, the broker verifies the root multi-signature, recurring only on the children of an invalid parent. This allows to identify invalid multi-signatures in logarithmic time while enabling batched verification in the good case.

**Caching legitimacy proofs.** Clients justify their sequence numbers with legitimacy proofs. Again, the broker cannot verify each proof in time. We address this problem by having the broker verify a legitimacy proof only if higher than the highest it previously observed. As a result, a faulty client might get away with submitting an invalid legitimacy proof but, importantly, not an illegitimate sequence number.

### 8.5.2 Server

The goal of a Chop Chop server is to process distilled batches as quickly as possible without overflowing its memory.

**Batch garbage collection.** Servers update each other on which batches they delivered. A server garbage-collects a batch, both messages and metadata, as soon as it is delivered by all other servers. We underline that, even if a single server fails to deliver a batch, the others cannot garbage-collect it as the slow server might be correct. This is an inherent limitation of Atomic Broadcast: agreement without synchrony can be ensured only in the infinite-memory model.

**Identifier-sorted batching.** No two messages from the same client must appear in the same batch. To simplify processing, brokers sort the messages in a batch by client identifier. Servers reject batches whose identifiers are not strictly increasing, thus verifying that all identifiers are distinct in constant size and in linear time. Sorting messages by identifier also enables parallel deduplication: messages are split by identifier range, chunks are deduplicated independently.

## 8.6 Evaluation

We evaluate Chop Chop focusing on the following research questions (RQs): What workload can Chop Chop sustain (Section 8.6.3)? What are the benefits of Chop Chop’s distillation (Section 8.6.4)? How does Chop Chop scale to different numbers of servers (Section 8.6.5)? How efficiently does Chop Chop use resources overall (Section 8.6.6)? How does Chop Chop perform under adverse conditions, such as server failures (Section 8.6.7)? What performance can applications achieve using Chop Chop (Section 8.6.8)?

### 8.6.1 Baselines

We compare Chop Chop against four baselines:

- HotStuff [156]: an Atomic Broadcast protocol designed for high-throughput (written in C++);
- BFT-SMaRt [22]: an Atomic Broadcast protocol, similar to PBFT [43], designed for low-latency (written in Java);
- Narwhal-Bullshark: the DAG-based Atomic Broadcast protocol Bullshark [140] with the state-of-the-art high-throughput mempool Narwhal [55] (written in Rust);
- Narwhal-Bullshark-sig: akin to Narwhal-Bullshark but with Narwhal modified to authenticate messages, thus matching Chop Chop’s guarantees.

We deploy Chop Chop with two distinct underlying Atomic Broadcast protocols (Figure 8.5): HotStuff and BFT-SMaRt.

**HotStuff and BFT-SMaRt.** Evaluating HotStuff and BFT-SMaRt allows us to assess the base performance of an Atomic Broadcast protocol and determine how much acceleration Chop Chop provides. We evaluate Chop Chop on top of the same implementations of HotStuff<sup>5</sup> and BFT-SMaRt<sup>6</sup> we benchmark against. These implementations are production-ready and do not use state-of-the-art mempool protocols, only some basic form of batching. When evaluated stand-alone, each message in these systems includes 80 B of header composed of a client identifier (8 B), a sequence number (8 B), and a signature (64 B) verified by the servers. Both systems use batches of 400 messages, i.e., of 34.4 KB.

**Narwhal-Bullshark.** As a state-of-the-art mempool, Narwhal<sup>7</sup> is a close point of comparison for Chop Chop. Servers in Narwhal scale out following a primary-workers model: each server is paired with one or several workers into a server group. Similarly to Chop Chop, Narwhal greatly accelerates its underlying Atomic Broadcast (here, Bullshark). Unlike Chop Chop,

<sup>5</sup>Authors’ implementation of HotStuff in C++: <https://github.com/hot-stuff/libhotstuff>

<sup>6</sup>Authors’ implementation of BFT-SMaRt in Java: <https://github.com/bft-smart/library>

<sup>7</sup>Authors’ implementation of Bullshark and Narwhal in Rust: <https://github.com/asonnino/narwhal/tree/bullshark>

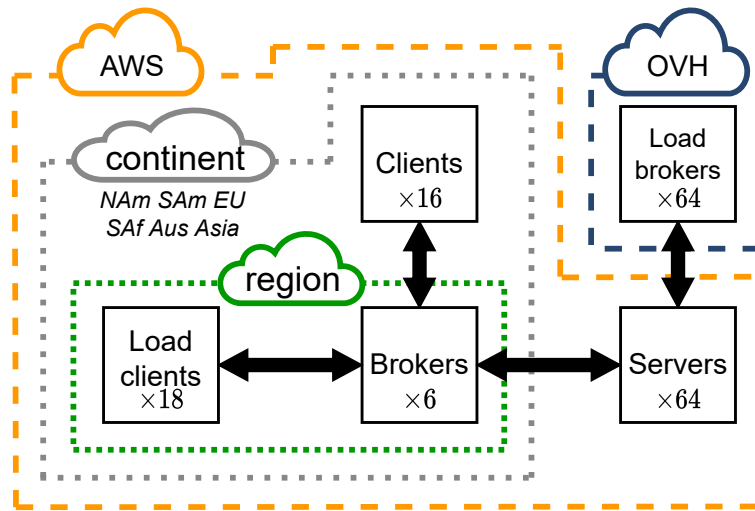


Figure 8.6: Cross-cloud deployment summary.

however, Narwhal leaves the responsibility of authenticating and deduplicating messages to the application.

**Narwhal-Bullshark-sig.** For a better comparison, we also benchmark Narwhal-Bullshark-sig: Narwhal-Bullshark where messages are authenticated by Narwhal in a state-of-the-art way, i.e., using batched, multi-core Ed25519 signature verification. Each message includes an 80 B header as for HotStuff and BFT-SMaRt. As for Narwhal-Bullshark, the remaining parameters are the default ones, e.g., 500 KB batches.

### 8.6.2 Setup

Unless otherwise specified—in Sections 8.6.5 and 8.6.6—the Chop Chop benchmarks involve 64 c6i.8xlarge AWS servers, of 32 Intel vCPUs each, geo-distributed across 14 regions. Brokers assemble, and servers process, batches of 65,536 messages. Each message is 8 B in length, resulting in 736 KB batches (Figure 8.3). Baselines always use the same set of server machines as their Chop Chop counterpart. All experiments run with maximum resilience, e.g., the system survives 21 faulty servers out of 64. Figure 8.6 overviews the used deployment.

**Matching trusted vs. total resources.** Unlike its baselines, Chop Chop leverages *untrusted* resources, brokers, to boost its performance. Lacking a well-defined conversion between trusted and untrusted resources, two extremes can be taken to compare Chop Chop with its baselines: we can either match trusted resources, e.g., same number of Chop Chop servers as Narwhal workers, or match total resources, e.g., same number of servers and brokers in Chop Chop as workers in Narwhal.

Intuitively, the first approach considers untrusted resources to be free while the second considers untrusted resources to be as costly as trusted resources. We use the first approach in Sections 8.6.3 to 8.6.5, 8.6.7 and 8.6.8 to stress Chop Chop, provisioning the system with enough brokers to bottleneck servers. We use the second approach in Section 8.6.6 to assess how efficiently Chop Chop uses its hardware resources, trusted or not.

**Load clients and load brokers.** We show in Section 8.6.3 that Chop Chop servers handle up to 43.6 million operations per second with an average latency of 3.6 seconds. To produce this level of workload, a real-world deployment would require over 700 brokers, each handling around 200,000 clients broadcasting back-to-back thus totaling hundreds of millions of machines. As we cannot experiment at such a scale, we introduce two new actors: *load clients* and *load brokers*.<sup>8</sup>

Load clients connect to brokers and simulate thousands of concurrent client requests. Most system evaluation typically use this approach to stress the system and measure latency. However, we explicitly separate clients from load clients in this evaluation. Clients run on very small machines—less powerful than most smartphones—to provide more accurate end-to-end latency measurements. We similarly split clients from load clients in all baseline runs.

Load brokers are unique to Chop Chop. Even using load clients, we could not deploy enough brokers to bottleneck Chop Chop’s servers. Load brokers work around this limitation, submitting batches of pre-generated messages directly to the servers. Free from interactions with clients and expensive cryptography, a load broker puts on the servers a load equivalent to that of tens of brokers working at full capacity.

Using load clients and load brokers, we manage to show that brokers can quickly generate large batches of messages, and servers can process large numbers of batches.

**Cross-cloud deployment.** All servers are deployed on AWS, balanced across 14 regions: Cape Town, São Paulo, Bahrain, Canada, Frankfurt, Northern Virginia, Northern California, Stockholm, Ohio, Milan, Oregon, Ireland, London, and Paris. For system sizes of 8 in Section 8.6.5, we distribute servers across the first 8 regions from the list, which constitute the most adversarial setup with the highest pairwise latency.

Load brokers are placed in a separate cloud provider, OVH, for two purposes. First, it provides a better representation of Internet load than a single-cloud deployment. AWS operates under its own AS so any AS peering bottlenecks would be bypassed by an AWS-only deployment. Second, OVH is one of the few cloud providers with enough peering with AWS to stress Chop Chop without charging for egress bandwidth, saving us from using AWS’ costly bandwidth. The final cost amounted to 25,000 USD in AWS credits. Using OVH saved us more than 70,000

---

<sup>8</sup>Throughout the remainder of this section, “brokers” and “clients” denote real brokers and real clients. The term “load” is always used explicitly.

## Chapter 8. Chop Chop

---

USD since each of Chop Chop’s data point on a figure would have cost 1,700 USD in AWS egress bandwidth—21 TB at 0.08 USD per GB  $\approx$  1,700 USD.

For all experiments, we deploy one broker in each continent (Cape Town, São Paulo, Tokyo, Sydney, Frankfurt, and Northern Virginia) and one client in each of the 14 regions above, plus Tokyo and Sydney. Clients connect to their nearest broker. We configure the network for geo-distribution and high load, e.g., TCP buffer sizes<sup>9</sup> and UDP parameters.

All baselines run on the same parameters. For Narwhal-Bullshark, we collocate each server with one of the workers in its server group. We reproduced Narwhal-Bullshark’s original experiments [140] and matched the results.

**Hardware.** All servers, brokers and load clients run on c6i.8xlarge machines with an Intel Xeon Platinum 8375C (32 virtual CPUs, 16 physical cores, 2.9 GHz baseline, 3.5 GHz turbo), 64 GB of memory and 12.5 Gb/s of bandwidth. We selected these machines since they provide good performance and are in the same “commodity” price range as those chosen initially for Chop Chop’s main baseline: Narwhal-Bullshark. Clients run on t3.small machines: 2 vCPUs, 1 physical core, 2 GB of memory, and up to 5 Gb/s bandwidth—of which they use less than 1 KB/s. All machines run Linux Ubuntu 20.04 LTS on the AWS patched version of the Linux kernel 5.15.0, except for the load brokers on OVH which run on Linux kernel 5.4.0—the same kernel was not available.

**Challenges.** The most significant evaluation challenges arose from the scale of the targeted deployment. The setup and orchestration alone required simultaneous handling of up to 320 machines across two different cloud providers and 25 regions, as well as transferring 13TB of files—mostly public keys and pre-generated batches—for each of the 12 setups. To handle this, we developed a new command-line tool to efficiently deploy distributed systems: `silk`. Among other things, we use `silk` for peer-to-peer-style file transfer over aggregated TCP connections, and for grouped process control. With `silk`, transferring all files from a single machine takes around 30 minutes, compared to 68 hours with `scp`.

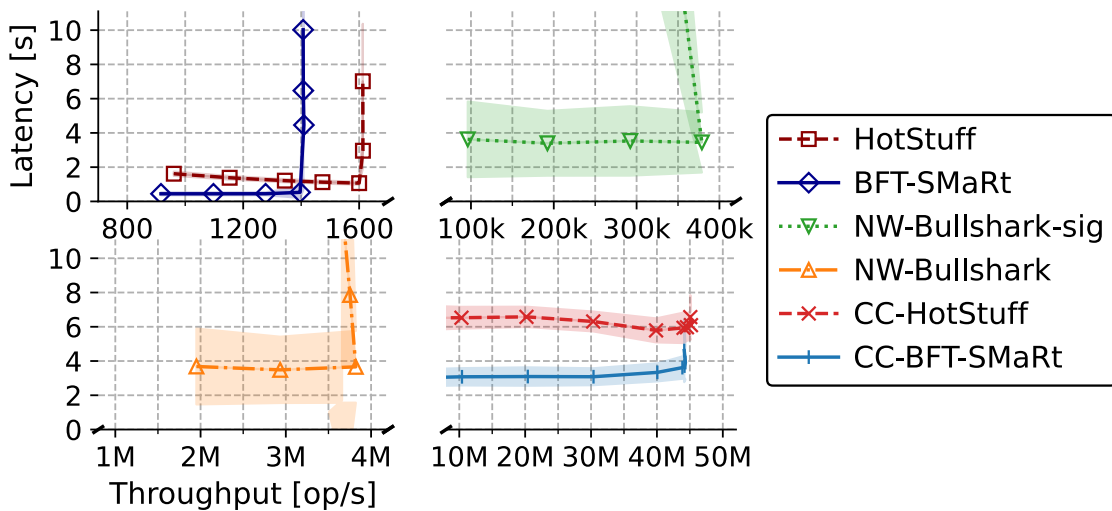
Additional challenges came from the real-world nature of the targeted deployment. First, the connection between OVH and AWS’s Asia and Pacific regions was particularly unstable at certain times of day especially when close to saturation. For example, Tokyo’s connection was frequently degraded between 3pm and 5pm UTC. Second, the performance of some machines sometimes deviated from their specifications. As an example, in a setup size of 64, we observed around 2 machines operating with a 10% lower CPU turbo clock rate than specified. Considering these variations, we increased the number of servers a broker initially asks for witness shards (see Section 8.4.3) by a margin, e.g.,  $f + 5$  instead of  $f + 1$ . This improves system stability—i.e., lower latency variability—while slightly reducing maximum throughput.

---

<sup>9</sup>IBM’s TCP Tuning Guide: <https://www.ibm.com/docs/en/linux-on-systems?topic=recommendations-network-performance-tuning>

Unless otherwise specified, we set the margin to 4 in all experiments, i.e.,  $f + 5$ .

**Plots.** Every data point is the mean of 5 runs of 2 minutes each (after excluding warmup and cooldown, the relevant cross-section is at least 1 minute). All plots further depict one standard deviation from the mean using either colored shaded areas or black error bars (which may be too small to notice).



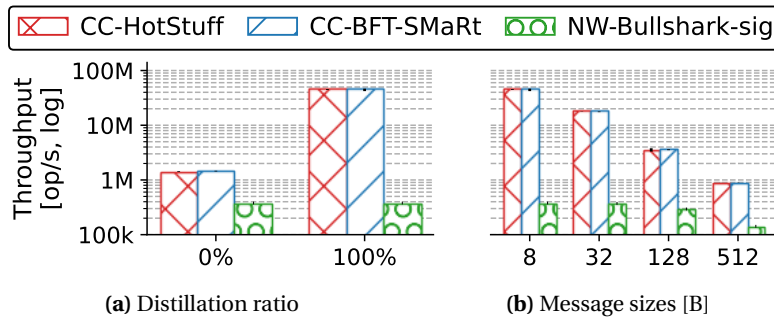
**Figure 8.7: Throughput-latency of Chop Chop and of notable Atomic Broadcast systems under various input rates.**

### 8.6.3 RQ1 – Load Handling

Figure 8.7 shows the latency and throughput of Chop Chop and all its baselines for various input rates of 8 B messages. The variability is represented using shaded areas.

**Baselines.** Both BFT-SMaRt and HotStuff showcase stable performances under low loads, respectively achieving around 1,400 and 1,600 operations per second. BFT-SMaRt’s latency is consistently better than HotStuff’s up to its inflection point (0.45–0.53 s vs. 1.2–1.6 s). We measure up to 3.8M op/s for Narwhal-Bullshark and up to 382k op/s for Narwhal-Bullshark-sig. The difference in respective throughput highlights the cost of authentication for servers: verifying signatures reduces the throughput of Narwhal-Bullshark by one order of magnitude. We observe a latency of around 3.6 s for both Narwhal-Bullshark and Narwhal-Bullshark-sig.

**Chop Chop.** Chop Chop achieves close to 44M op/s while running on top of both HotStuff and BFT-SMaRt. Chop Chop’s latency range is 3.0–3.6 s with BFT-SMaRt and 5.8–6.5 s with HotStuff. Notably, the latency of Chop Chop-HotStuff decreases under high load. This is due to the internal batching mechanism of the HotStuff implementation: buffers fill faster under



**Figure 8.8: Throughput of Chop Chop and authenticated Narwhal with Bullshark (log scale) when (a) Chop Chop has no distillation and with (b) varying message size.**

higher load, thus avoiding timeouts. This has an immediate impact on Chop Chop, which feeds HotStuff at a low rate: HotStuff alone accounts for over 60% of Chop Chop-HotStuff’s overall latency. BFT-SMaRt makes a better fit for Chop Chop, as its throughput is sufficient for Chop Chop’s needs, and its latency is lower than HotStuff’s.

**Mem pools’ trade-off.** In comparison to BFT-SMaRt and HotStuff, Chop Chop trades latency in favor of throughput. This trade-off is mostly explained by batching and distillation. When assembling a batch, a broker has to wait twice: once to collect enough messages to fill a batch, and once to collect all multi-signatures from clients engaging in distillation. We set both waits’ timeout to 1 second. Notably, Narwhal-Bullshark seems to incur a similar latency cost, as Chop Chop’s latency approximately matches that of Narwhal-Bullshark, even though Chop Chop needs an extra round trip between clients and broker (Figure 8.5, #4–#6).

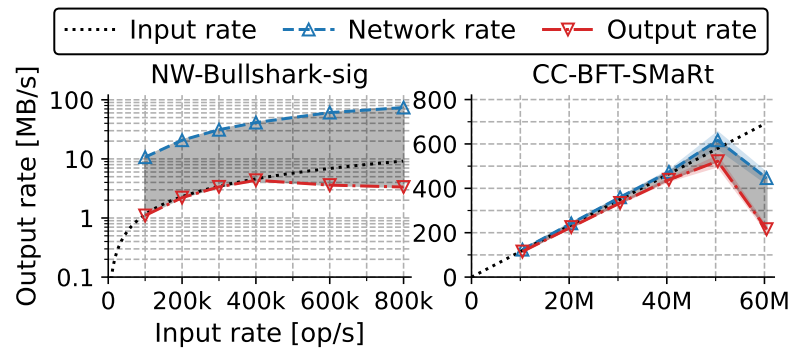
#### 8.6.4 RQ2 – Distillation Benefits

We showcase the benefits of distillation by: evaluating throughput with and without distillation, evaluating distillation for messages of different sizes, and observing the impact of distillation on network bandwidth to achieve line rate.

**Distillation vs. mitigations.** Along with distillation, Chop Chop makes use of two techniques available in the literature to mitigate the cost of Atomic Broadcast’s authentication: short identifiers and pooled signature verification (see Section 8.2.2).

Figure 8.8a breaks down Chop Chop’s throughput, measuring how significantly distillation alone contributes to Chop Chop’s performance. When no message is distilled, Chop Chop’s servers bottleneck at 1.5M op/s, 3.9× higher than Narwhal-Bullshark-sig. This result is in line with both systems bottlenecking on server CPU, as the technique employed by Chop Chop to mitigate authentication complexity has only one third of the servers verify each client signature. (We conjecture that the additional 1.3× factor may be owed to engineering





**Figure 8.9: Throughput efficiency of authenticated Narwhal with Bullshark (left, log scale) and Chop Chop with BFT-SMaRt (right, linear scale) with various input rates.**

differences.) When batches are fully distilled, Chop Chop’s throughput grows to 44M op/s, accounting for the additional 29-fold boost to Chop Chop’s performance.

**Distillation for larger messages.** Figure 8.8b illustrates Chop Chop’s maximum throughput for message sizes of 8 B to 512 B which may be relevant to applications that cannot work around smaller message sizes, e.g., many smart contracts. Chop Chop’s throughput is similar with BFT-SMaRt and HotStuff, decreasing at an approximately 1-to-1 ratio as the message size increases: 44.3M op/s for 8 B, 17.6M op/s for 32 B, 3.5M op/s for 128 B and 890k op/s for 512 B.

This is in line with expectations. As we discuss in Section 8.3.2, a server should receive  $\sim b$  bytes in order to deliver a  $b$ -bytes message in a large, fully distilled batch, as full distillation amortizes to zero the communication cost of authenticating and sequencing each message. For 8 B messages, servers encounter a CPU bottleneck slightly before the link between load brokers and servers is saturated. This explains why the throughput decreases only  $2.52\times$  when messages grow to 32 B: all remaining server-bound bandwidth is used to convey messages (as messages are larger) while the load on server CPUs is reduced (as less messages are delivered overall). The system remains communication-bottlenecked as the size of the messages increases, and throughput starts decreasing linearly with message size, e.g., Chop Chop’s throughput for 512 B messages is  $4.00\times$  smaller than for 128 B.

By contrast, Narwhal-Bullshark-sig bottlenecks on server CPUs longer, due to signature verification, maintaining a stable throughput until 512 B messages finally fill server links. Overall, Narwhal-Bullshark-sig’s throughput only decreases from 382k op/s for 8 B messages to 142k op/s for 512 B messages, which matches their non-authenticated evaluation with 512 B messages. The gap between Chop Chop and Narwhal-Bullshark-sig at 512 B messages can be mostly attributed to Chop Chop’s more efficient use of server bandwidth: unlike Narwhal, Chop Chop offloads the dissemination of batches to external brokers. Narwhal’s use of worker-to-worker communication in its common path also makes it more prone to be affected by AWS’s various upload limitations, e.g., AWS upload bandwidth is half the stated download bandwidth, and there are network credit limits for “burst” uploading.

## Chapter 8. Chop Chop

---

**Line rate.** Figure 8.9 illustrates Chop Chop’s near line-rate network use by depicting its input, network and output rates:

- Input rate measures the total bytes of useful information—i.e., client identifiers and messages—that clients, load clients and load brokers all broadcast per time unit;
- Network rate measures the ingress bandwidth of servers at their network interface, i.e., useful information captured by the input rate as well as the Atomic Broadcast’s overhead for ordering, authentication and deduplication;
- Output rate, or “goodput”, measures the total bytes of useful information that each server delivers per time unit.

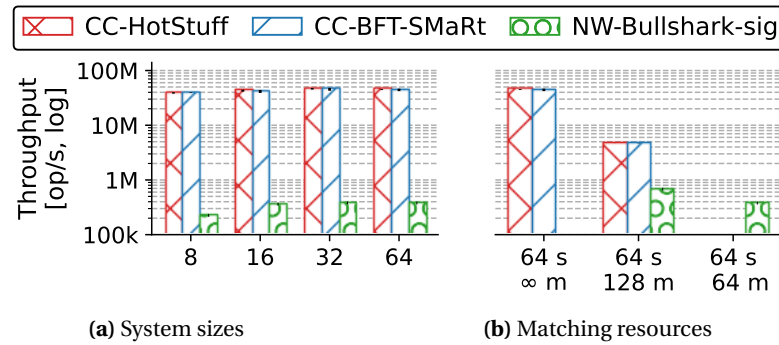
A system with perfect line rate would match all three rates: input rate would match output rate as messages can be delivered in a timely fashion with no backlogging, and output rate would match network rate as a server would only receive useful information, with no overhead due to Atomic Broadcast. The gray-shaded areas in Figure 8.9 highlight this overhead, i.e., the difference between network and input rates. Network and output rates are averaged over all servers.

In this experiment, each of the 257M simulated clients broadcast 8 B messages. This results in 11.5 B of useful information per broadcast as 28 bits = 3.5 B are sufficient to represent every identifier. This conversion is captured by the dotted line which converts the input rate from op/s, represented on the x-axis, to B/s, represented on the y-axis.

For authenticated Narwhal-Bullshark, the output rate closely matches the input rate until signature verification becomes the bottleneck at 378k op/s, shown by the plateauing output rate. The gap between Narwhal-Bullshark-sig’s network and input rates is evident, differing by one order of magnitude (notably in line with our back-of-the-envelope calculation in Section 8.3.2). In contrast, thanks to distillation, Chop Chop practically achieves line-rate up to its maximum throughput. Before its inflection point at 40M op/s, the the overhead of Chop Chop is less than 8%. The drop in output and network rates at 60M op/s is due to servers surpassing their computational capacity: broadcasts stall, server witness verification gets backlogged and brokers, suspecting server faults, ask for more batch witnesses, further stressing servers’ CPUs.

### 8.6.5 RQ3 – Number of Servers

Figure 8.10a illustrates the maximum throughput for systems of 8 ( $f = 2$ ), 16 ( $f = 5$ ), 32 ( $f = 10$ ) and 64 ( $f = 21$ ) servers. For Chop Chop, we adjust the witnessing margin as the system grows by 0, 1, 2, and 4 for 8, 16, 32 and 64 servers respectively (see Section 8.6.2). Both Chop Chop and authenticated Narwhal-Bullshark scale well to 64 servers. Note that, unless trust assumptions are modified, Narwhal-Bullshark-sig only scales vertically: if a Narwhal server or any of its workers are faulty, the entire server group is compromised. Chop Chop, instead, scales horizontally with the number of brokers.



**Figure 8.10: Throughput of Chop Chop and authenticated Narwhal with Bullshark (log scale) when (a) varying system size, and when (b) varying the number of overall machines (“m”) with 64 servers (“s”).** Load brokers in Chop Chop simulate tens of brokers, hence are noted “∞ m”.

### 8.6.6 RQ4 – Overall Efficiency

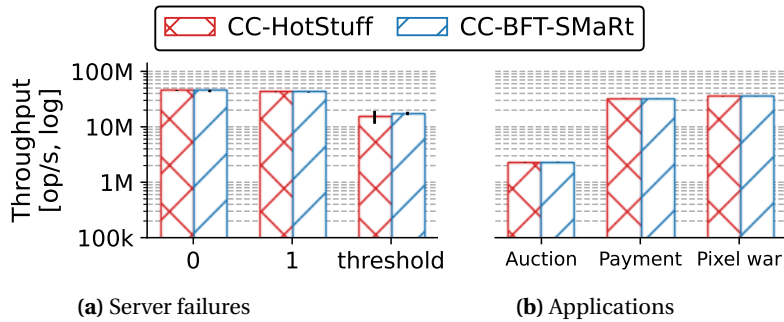
The center cluster of bars in Figure 8.10b compares Chop Chop’s throughput with that of authenticated Narwhal-Bullshark when overall hardware resources are matched. In this setting, both systems have 128 machines at their disposal. Chop Chop is provided with 64 servers, 64 brokers and 0 load brokers. Since a load broker uses pre-generated synthetic data to simulate tens of brokers (see Section 8.6.2), involving load brokers in this experiment would give an unfair advantage to Chop Chop. Narwhal-Bullshark-sig is provided with 128 workers, to match Chop Chop’s total machines, balanced across 64 server groups, to match Chop Chop’s servers. The left and right clusters of bars depict Chop Chop using load brokers and Narwhal-Bullshark-sig with 64 server groups containing 1 worker each, respectively, as in the other experiments.

We observe 4.6M op/s for Chop Chop, with servers reporting around 5% CPU usage. We observe 679k op/s for Narwhal-Bullshark-sig. Chop Chop’s higher throughput is in line with expectations. In Narwhal-Bullshark-sig, workers are trusted, and as such a worker can only contribute to its own server group. Instead, since Chop Chop brokers are untrusted, a broker’s work is useful to all servers.

### 8.6.7 RQ5 – Chop Chop Under Failures

Figure 8.11a depicts Chop Chop’s throughput when some servers crash 30 seconds into the run. Performance drops marginally (from 44M op/s to 43M op/s) with one crash and by 66% (down to 15M op/s) when one-third of the servers crash, resulting in less CPU globally available to witness batches.

Figure 8.8a captures Chop Chop’s performance hit when clients fail to engage in distillation. This could be caused by clients being slow or crashed, or brokers being malicious. Under the most extreme conditions, where no client engages in distillation, the throughput drops from



**Figure 8.11: Throughput of Chop Chop (log scale) with (a) various server failures and for (b) different applications.**

44M op/s to 1.5M op/s.

### 8.6.8 RQ6 – Application Use Cases

Figure 8.11b depicts the maximal stable throughput for various application use cases. In the Auction app, a client can bid an amount on a token it does not own, or take the highest offer it received for an item it owns. The highest amount bid on each token is locked and cannot be used to bid elsewhere. Money bid is transferred when the owner of the token takes the offer, or refunded when the bid is raised by another client. The Auction app is single-threaded and many clients bid on the same token to approximate a real auction. In the Payments app, clients choose a recipient and an amount to transfer. In Pixel war, clients choose a pixel and an RGB color to paint on a 2,048 by 2,048 board. Operations are generated at random.

We observe 2.3M op/s for the Auction, 32M op/s for Payments and 35M op/s for Pixel war. The bottleneck is the application in all cases, thus Chop Chop has sufficient capacity for high, single-application throughput. Chop Chop can also support many separate high-throughput applications simultaneously, making it a fitting Atomic Broadcast candidate to power a universal SMR system (or “Internet computer”).

## 8.7 Related Work

We overview below the state-of-the-art most relevant to Chop Chop, namely Atomic Broadcast systems with high-throughput and efficient signature aggregation schemes.

**High-throughput Atomic Broadcast.** Narwhal [55] is a mempool protocol that separates the reliable distribution of payloads from the communication-expensive ordering in order to accelerate DAG-based Atomic Broadcast [69, 92, 140]. Narwhal utilizes trusted workers to increase throughput while Chop Chop relies on *trustless* brokers, for the same effect, and scales out more efficiently. To circumvent the bottleneck associated with the broadcast leader,

approaches using multiple leaders have been developed—both for crash [64, 120] and arbitrary [7, 12, 141, 142] faults—to scale the broadcast throughput linearly with the number of leaders. Dissemination trees [96, 124] have also been employed to reduce communication cost and maximize network bandwidth utility, while sharded [98, 151] and federated [109] approaches reduce communication cost by promoting local communication in geo-distributed setups. In comparison, Chop Chop shows that an optimal distillation mechanism for batches achieves better performances without adding complexity to the Atomic Broadcast protocol itself.

Other approaches have shown that the underlying hardware of servers can also be exploited for higher throughput, such as FPGA [87, 91] and Intel SGX enclaves [16]. In comparison, Chop Chop uniquely boosts throughput by exploiting *trustless hardware* via brokers. Atomic Broadcast can also be accelerated in data centers by using the topology of the network [105, 127] or even by running within the network itself using P4-programmable switches [56, 97]. In such low latency environments, the processing overhead incurred by the operating system kernel can be bypassed to further increase the throughput of Atomic Broadcast [3, 97, 150].

**Signature aggregation.** Aggregate signatures were first proposed to save space by compacting a large number of signatures into just one [26, 136]. Up until recently, aggregation could also save verification time but only in certain cases: either when the signatures are generated by the same signer [42], or when the signatures are on the same message, i.e., multi-signatures [88]. In the latter case, aggregation mechanisms have been proposed to achieve constant-time verification of aggregated multi-signatures for both BLS [25] and Schnorr [115] signature schemes. In particular, multi-signatures are used in cryptocurrencies to have many servers sign the same batch of payloads [60, 96]. Servers in Chop Chop use rapidly-verifiable BLS multi-signatures [25] for that very purpose. In addition to aggregating server signatures on batches, Chop Chop’s distillation mechanism also aggregates all client signatures in a batch in a way that provides constant-time verification. The theoretical scheme Draft [41] proposed signature aggregation with similar verification performances but is tailored to Reliable Broadcast. It is however unclear how Draft could be implemented as a real-world system without compromising liveness. Indeed, Draft assumes infinite memory to prevent message replay attacks which would rapidly exhaust servers’ memory if Draft were to be deployed to match Chop Chop’s target throughput in our evaluation (see Section 8.6.2). Chop Chop also aggregates client sequence numbers to significantly reduce bandwidth consumption when small messages are broadcast (Figure 8.2). Aggregating sequence numbers is made possible thanks to the ordering of Atomic Broadcast and thanks to novel legitimacy proofs (see Section 8.4.2).



# Conclusions

This thesis focused on the scalability of Byzantine distributed computing, with the goal to develop systems that billions of stakeholders could partake in as servers, or use as clients. On the side of server scalability, we focused on Reliable Broadcast. We introduced Contagion, the first probabilistic Reliable Broadcast protocol to achieve logarithmic per-process computation and communication complexity. At the core of Contagion are samples, a novel alternative to quorums trading intersection guarantees for statistical representativeness. Contagion's analysis is, to our knowledge, the first to study a probabilistic distributed protocol in the Byzantine setting, enabled by a novel proving technique based on decorators. On the side of client scalability, we strove to maximize how efficiently Reliable Broadcast and Consensus servers make use of their resources when subject to a high throughput of requests. In the context of Reliable Broadcast, we introduced the notion of oracularity to describe a system's ability to match the performance of its centralized, trusted counterpart. We then introduced Draft, the first algorithm to achieve oracularity (at least in the good case), effectively achieving zero-cost Byzantine fault tolerance. To do so, we extended the distributed computing model with brokers, an untrusted layer of processes whose goal is to alleviate server complexity by performing heavy but verifiable operations. In the context of Consensus, we extended our broker-based techniques to the Atomic Broadcast abstraction. We introduced Chop Chop, the first near-line-rate system of its kind, processing tens of millions of operations per second on a geo-distributed deployment of medium-sized servers.

Our results, we believe, open interesting new avenues of future research. We outline a few of them below.

**Scaling Consensus servers.** As we discussed in the Introduction to this thesis, a great deal of effort was put towards generalizing our results on scalable Reliable Broadcast to the Consensus class. Doing so was challenging within the limited time frame of our PhD studies. While an algorithm, which we conjecture to be correct, was developed in the partially synchronous model, decorator-based proofs proved unwieldy when tackling the subtleties of probabilistic, adversarial time analysis. At the cost of additional formal effort, we believe a provably secure, logarithmic implementation of Atomic Broadcast should be attainable, completing the stated goal of this thesis.

## Conclusions

---

**Theory of decorators.** We believe decorators, the tool we introduced to prove the correctness of Contagion, could be applied to the broader field of probabilistic Byzantine algorithms. To this end, a more general theory of decorators should be developed, streamlining our ad-hoc, decorator-based results into a usable toolbox of proving techniques. Significant (unpublished) progress was already made towards this goal, framing decorators in the broader context of probabilistic adversarial games.

**Constant-memory Chop Chop.** A global deployment of Chop Chop servers sustains a throughput of tens of millions of operations per second. In the real world, such a high throughput would realistically be generated by billions, maybe tens of billions of clients - even if each client persistently broadcast every two minutes, five billion clients would be necessary to incur a CPU bottleneck. Because each Chop Chop server must store the full list of all client public keys, such a large number of clients might severely strain server memory. More advanced cryptographic techniques, we believe, could enable offloading public key storage to brokers, effectively making Chop Chop a constant-memory algorithm. Doing so while maintaining near-line-rate performance presents a significant theoretical challenge.

**Quantum resistance for reduction and distillation.** The broker-based techniques we presented in this thesis to enhance batch authentication rely on multi-signature aggregation. State-of-the-art multi-signature schemes, such as BLS, generally rest their security on the hardness of the discrete logarithm problem, which is famously vulnerable to quantum attacks. While practical attacks are (to our knowledge) yet to be dealt, the general cryptographic consensus suggests pre-quantum schemes should be progressively phased out in the coming decade. Upgrading our systems with post-quantum authentication while preserving their communication complexity, we believe, will be of great importance in the near future.

**Scaling out resource-intensive services.** All applications showcased in this thesis (such as payments or auctions) have a minimal footprint in terms of computation and bit complexity. This choice was deliberate: our goal was to demonstrate our protocols would scale even in a setting where bottlenecks could not possibly be blamed on application complexity. Real-world applications, however, might be complex enough to severely bottleneck a fully-replicated system. This setting, we believe, calls for partial replication, with a core of secure, fully-replicated servers orchestrating computational resources contributed by the broader Internet. Doing so while maintaining (probabilistic) security would further enable the vision of a secure, planetary-scale computer shared by all.



# Bibliography

- [1] Ittai Abraham, Kartik Nayak, Ling Ren, and Zhuolun Xiang. 2021. Good-Case Latency of Byzantine Broadcast: A Complete Categorization. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [2] Daron Acemoglu and Asu Ozdaglar. 2009. Networks - Lecture 4: Erdős–Rényi Graphs and Phase Transitions. (2009).
- [3] Marcos K. Aguilera, Naama Ben-David, Rachid Guerraoui, Antoine Murat, Athanasios Xygkis, and Igor Zablotchi. 2023. uBFT: Microsecond-Scale BFT Using Disaggregated Memory. In *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*.
- [4] Marcos K. Aguilera, Idit Keidar, Dahlia Malkhi, Jean-Philippe Martin, and Alexander Shraer. 2010. Reconfiguring Replicated Atomic Storage: A Tutorial. *Bulletin of the EATCS* (2010).
- [5] Nicolas Alhaddad, Sisi Duan, Mayank Varia, and Haibin Zhang. 2021. Succinct Erasure Coding Proof Systems. *IACR Cryptology ePrint Archive* (2021).
- [6] Dan Alistarh, Seth Gilbert, Rachid Guerraoui, and Morteza Zadimoghaddam. 2010. How Efficient Can Gossip Be? (On the Cost of Resilient Information Exchange). In *International Colloquium on Automata, Languages and Programming (ICALP)*.
- [7] Salem Alqahtani and Murat Demirbas. 2021. BigBFT: A Multileader Byzantine Fault Tolerance Protocol for High Throughput. In *IEEE International Performance, Computing, and Communications Conference (IPCCC)*.
- [8] Elli Androulaki, Artem Barger, Vita Bortnikov, Christian Cachin, Konstantinos Christidis, Angelo De Caro, David Enyeart, Christopher Ferris, Gennady Laventman, Yacov Manevich, Srinivasan Muralidharan, Chet Murthy, Binh Nguyen, Manish Sethi, Gari Singh, Keith Smith, Alessandro Sorniotti, Chrysoula Stathakopoulou, Marko Vukolić, Sharon Weed Cocco, and Jason Yellick. 2018. Hyperledger Fabric: A Distributed Operating System for Permissioned Blockchains. In *European Conference on Computer Systems (EuroSys)*.
- [9] Hagit Attiya, Amotz Bar-Noy, and Danny Dolev. 1995. Sharing Memory Robustly in Message-Passing Systems. *Journal of the ACM (JACM)* (1995).

## Bibliography

---

- [10] Pierre-Louis Aublin, Rachid Guerraoui, Nikola Knežević, Vivien Quéma, and Marko Vukolić. 2015. The Next 700 BFT Protocols. *ACM Transactions on Computer Systems (TOCS)* (2015).
- [11] Alex Auvolat, Davide Frey, Michel Raynal, and François Taïani. 2021. Byzantine-Tolerant Causal Broadcast. *Theoretical Computer Science* (2021).
- [12] Zeta Avarikioti, Lioba Heimbach, Roland Schmid, Laurent Vanbever, Roger Wattenhofer, and Patrick Wintermeyer. 2020. FnF-BFT: Exploring Performance Limits of BFT Protocols. In *International Colloquium on Structural Information and Communication Complexity (SIROCCO)*.
- [13] Chen Avin, Michael Borokhovich, Keren Censor-Hillel, and Zvi Lotker. 2011. Order Optimal Information Spreading Using Algebraic Gossip. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [14] Baruch Awerbuch and Christian Scheideler. 2009. Towards a Scalable and Robust DHT. *Theory of Computing Systems* (2009).
- [15] Paulo S. L. M. Barreto, Hae Y. Kim, Ben Lynn, and Michael Scott. 2002. Efficient Algorithms for Pairing-Based Cryptosystems. In *Annual International Cryptology Conference (CRYPTO)*.
- [16] Johannes Behl, Tobias Distler, and Rüdiger Kapitza. 2017. Hybrids on Steroids: SGX-Based High Performance BFT. In *European Conference on Computer Systems (EuroSys)*.
- [17] Michael Ben-Or, Ran Canetti, and Oded Goldreich. 1993. Asynchronous Secure Computation. In *Symposium on the Theory of Computing (STOC)*.
- [18] Petra Berenbrink, Robert Elsässer, and Tom Friedetzky. 2008. Efficient Randomised Broadcasting in Random Regular Networks With Applications in Peer-to-Peer Systems. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [19] Petra Berenbrink, Robert Elsässer, and Thomas Sauerwald. 2010. Communication Complexity of Quasirandom Rumor Spreading. In *European Symposium on Algorithms (ESA)*.
- [20] Petra Berenbrink, Robert Elsässer, and Thomas Sauerwald. 2010. Randomised Broadcasting: Memory vs. Randomness. *Theoretical Computer Science* (2010).
- [21] Daniel J. Bernstein. 2006. Curve25519: New Diffie-Hellman Speed Records. In *International Conference on Theory and Practice of Public Key Cryptography (PKC)*.
- [22] Alysso Bessani, Joao Sousa, and Eduardo E.P. Alchieri. 2014. State Machine Replication for the Masses With BFT-SMART. In *Dependable Systems and Networks (DSN)*.
- [23] Kenneth P. Birman, Mark Hayden, Ozgur Ozkasap, Zhen Xiao, Mihai Budiu, and Yaron Minsky. 1999. Bimodal Multicast. *ACM Transactions on Computer Systems (TOCS)* (1999).

- [24] Joseph T.A. Birman K.P. 1987. Reliable Communication in the Presence of Failures. *ACM Transactions on Computer Systems (TOCS)* (1987).
- [25] Dan Boneh, Manu Drijvers, and Gregory Neven. 2018. Compact Multi-Signatures for Smaller Blockchains. In *International Conference on the Theory and Application of Cryptology and Information Security (ASIACRYPT)*.
- [26] Dan Boneh, Craig Gentry, Ben Lynn, and Hovav Shacham. 2003. Aggregate and Verifiably Encrypted Signatures From Bilinear Maps. In *International Conference on the Theory and Application of Cryptographic Techniques (EUROCRYPT)*.
- [27] Dan Boneh, Sergey Gorbunov, Riad S. Wahby, Hoeteck Wee, and Zhenfei Zhang. 2022. BLS Signatures. RCF Draft.
- [28] Edward Bortnikov, Maxim Gurevich, Idit Keidar, Gabriel Kliot, and Alexander Shraer. 2009. Brahms: Byzantine Resilient Random Membership Sampling. *Computer Networks* (2009).
- [29] Elette Boyle, Shafi Goldwasser, and Stefano Tessaro. 2013. Communication Locality in Secure Multi-Party Computation. In *Theory of Cryptography Conference (TCC)*.
- [30] Gabriel Bracha. 1984. An Asynchronous  $[(n-1)/3]$ -Resilient Consensus Protocol. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [31] Gabriel Bracha. 1987. Asynchronous Byzantine Agreement Protocols. *Information and Computation* (1987).
- [32] Gabriel Bracha and Sam Toueg. 1985. Asynchronous Consensus and Broadcast Protocols. *Journal of the ACM (JACM)* (1985).
- [33] Christian Cachin. 2010. State Machine Replication With Byzantine Faults. In *Replication*.
- [34] Christian Cachin, Rachid Guerraoui, and Luís Rodrigues. 2011. *Introduction to Reliable and Secure Distributed Programming*.
- [35] Christian Cachin, Klaus Kursawe, Frank Petzold, and Victor Shoup. 2001. Secure and Efficient Asynchronous Broadcast Protocols. In *Annual International Cryptology Conference (CRYPTO)*.
- [36] Christian Cachin, Klaus Kursawe, and Victor Shoup. 2005. Random Oracles in Constantinople: Practical Asynchronous Byzantine Agreement Using Cryptography. *Journal of Cryptology (JCrypt)* (2005).
- [37] Christian Cachin and Jonathan A Poritz. 2002. Secure Intrusion-Tolerant Replication on the Internet. In *Dependable Systems and Networks (DSN)*.
- [38] Christian Cachin and Stefano Tessaro. 2005. Asynchronous Verifiable Information Dispersal. In *IEEE International Symposium on Reliable Distributed Systems (SRDS)*.

## Bibliography

---

- [39] Martina Camaioni, Rachid Guerraoui, Jovan Komatovic, Matteo Monti, Pierre-Louis Roman, Manuel Vidigueira, and Gauthier Voron. 2023. Carbon: Scaling Trusted Payments with Untrusted Machines. *IEEE Transactions on Dependable and Secure Computing (TDSC)* (2023). (Under review).
- [40] Martina Camaioni, Rachid Guerraoui, Matteo Monti, Pierre-Louis Roman, Manuel Vidigueira, and Gauthier Voron. 2024. Chop Chop: Byzantine Atomic Broadcast to the Network Limit. *USENIX Symposium on Operating Systems Design and Implementation (OSDI)* (2024). (Under minor revisions).
- [41] Martina Camaioni, Rachid Guerraoui, Matteo Monti, and Manuel Vidigueira. 2022. Oracular Byzantine Reliable Broadcast. In *International Symposium on Distributed Computing (DISC)*.
- [42] Jan Camenisch, Susan Hohenberger, and Michael Ostergaard Pedersen. 2012. Batch Verification of Short Signatures. *Journal of Cryptology (JCrypt)* (2012).
- [43] Miguel Castro and Barbara Liskov. 1999. Practical Byzantine Fault Tolerance. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [44] Miguel Castro and Barbara Liskov. 2002. Practical Byzantine Fault Tolerance and Proactive Recovery. *ACM Transactions on Computer Systems (TOCS)* (2002).
- [45] Tushar Deepak Chandra and Sam Toueg. 1996. Unreliable Failure Detectors for Reliable Distributed Systems. *Journal of the ACM (JACM)* (1996).
- [46] Nishanth Chandran, Wutichai Chongchitmate, Juan A. Garay, Shafi Goldwasser, Rafail Ostrovsky, and Vassilis Zikas. 2015. The Hidden Graph Model: Communication Locality and Optimal Resiliency With Adaptive Faults. In *Innovations in Theoretical Computer Science (ITCS)*.
- [47] Fan Chung and Linyuan Lu. 2001. The Diameter of Sparse Random Graphs. *Advances in Applied Mathematics* (2001).
- [48] Pierre Civit, Seth Gilbert, Rachid Guerraoui, Jovan Komatovic, and Manuel Vidigueira. 2023. On the Validity of Consensus. *International Symposium on Distributed Computing (DISC)* (2023).
- [49] Allen Clement, Manos Kapritsos, Sangmin Lee, Yang Wang, Lorenzo Alvisi, Mike Dahlin, and Taylor Riche. 2009. Upright Cluster Services. In *Symposium on Operating Systems Principles (SOSP)*.
- [50] Daniel Collins, Rachid Guerraoui, Jovan Komatovic, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, Yvonne-Anne Pignolet, Dragos-Adrian Seredinschi, Andrei Tonkikh, and Athanasios Xygiakis. 2020. Online Payments by Merely Broadcasting Messages. In *Dependable Systems and Networks (DSN)*.

- [51] Thomas Cover and Joy Thomas. 2005. *Elements of Information Theory, Second Edition*.
- [52] Tyler Crain, Christopher Natoli, and Vincent Gramoli. 2018. Evaluating the Red Belly Blockchain. *Computing Research Repository (CoRR)* (2018).
- [53] Tyler Crain, Christopher Natoli, and Vincent Gramoli. 2021. Red Belly: A Secure, Fair and Scalable Open Blockchain. In *IEEE Symposium on Security and Privacy (S&P)*.
- [54] Philip Daian, Steven Goldfeder, Tyler Kell, Yunqi Li, Xueyuan Zhao, Iddo Bentov, Lorenz Breidenbach, and Ari Juels. 2020. Flash Boys 2.0: Frontrunning in Decentralized Exchanges, Miner Extractable Value, and Consensus Instability. In *IEEE Symposium on Security and Privacy (S&P)*.
- [55] George Danezis, Lefteris Kokoris-Kogias, Alberto Sonnino, and Alexander Spiegelman. 2022. Narwhal and Tusk: A DAG-Based Mempool and Efficient BFT Consensus. In *European Conference on Computer Systems (EuroSys)*.
- [56] Huynh Tu Dang, Pietro Bressana, Han Wang, Ki Suh Lee, Noa Zilberman, Hakim Weatherpoon, Marco Canini, Fernando Pedone, and Robert Soulé. 2020. P4xos: Consensus as a Network Service. *IEEE/ACM Transactions on Networking (TON)* (2020).
- [57] Sourav Das, Zhuolun Xiang, and Ling Ren. 2021. Asynchronous Data Dissemination and Its Applications. In *Annual ACM Conference on Computer and Communications Security (CCS)*.
- [58] Roger Dingledine, Nick Mathewson, and Paul Syverson. 2004. Tor: The Second-Generation Onion Router. In *USENIX Security Symposium (SEC)*.
- [59] Danny Dolev and Rüdiger Reischuk. 1985. Bounds on Information Exchange for Byzantine Agreement. *Journal of the ACM (JACM)* (1985).
- [60] Manu Drijvers, Sergey Gorbunov, Gregory Neven, and Hoeteck Wee. 2020. Pixel: Multi-Signatures for Consensus. In *USENIX Security Symposium (SEC)*.
- [61] Sisi Duan, Michael K. Reiter, and Haibin Zhang. 2018. BEAT: Asynchronous BFT Made Practical. In *Annual ACM Conference on Computer and Communications Security (CCS)*.
- [62] Xavier Défago, André Schiper, and Péter Urbán. 2004. Total Order Broadcast and Multicast Algorithms: Taxonomy and Survey. *ACM Computing Surveys (CSUR)* (2004).
- [63] Robert Elsässer and Dominik Kaaser. 2015. On the Influence of Graph Density on Randomized Gossiping. *IEEE International Parallel and Distributed Processing Symposium (IPDPS)* (2015).
- [64] Vitor Enes, Carlos Baquero, Alexey Gotsman, and Pierre Sutra. 2021. Efficient Replication via Timestamp Stability. In *European Conference on Computer Systems (EuroSys)*.

## Bibliography

---

- [65] Paul Erdős and Alfréd Rényi. 1959. On Random Graphs. *Publicationes Mathematicae* (1959).
- [66] P. Th. Eugster, R. Guerraoui, S. B. Handurukande, P. Kouznetsov, and A.-M. Kermarrec. 2003. Lightweight Probabilistic Broadcast. *ACM Transactions on Computer Systems (TOCS)* (2003).
- [67] Yaacov Fernandess, Antonio Fernández, and Maxime Monod. 2007. A Generic Theoretical Framework for Modeling Gossip-Based Algorithms. *ACM SIGOPS: Operating Systems Review* (2007).
- [68] Michael J. Fischer, Nancy A. Lynch, and Michael S. Paterson. 1985. Impossibility of Distributed Consensus With One Faulty Process. *Journal of the ACM (JACM)* (1985).
- [69] Adam Gagol, Damian Leundefiedniak, Damian Straszak, and Michal Swietek. 2019. Aleph: Efficient Atomic Broadcast in Asynchronous Networks With Byzantine Nodes. In *Conference on Advances in Financial Technologies (AFT)*.
- [70] Fangyu Gai, Jianyu Niu, Ivan Beschastnikh, Chen Feng, and Sheng Wang. 2023. Scaling Blockchain Consensus via a Robust Shared Mempool. In *IEEE International Conference on Data Engineering (ICDE)*.
- [71] Juan Garay, Yuval Ishai, Rafail Ostrovsky, and Vassilis Zikas. 2017. The Price of Low Communication in Secure Multi-Party Computation. In *Annual International Cryptology Conference (CRYPTO)*.
- [72] Juan A Garay, Jonathan Katz, Ranjit Kumaresan, and Hong-Sheng Zhou. 2011. Adaptively Secure Broadcast, Revisited. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [73] Chryssis Georgiou, Seth Gilbert, Rachid Guerraoui, and Dariusz R. Kowalski. 2008. On the Complexity of Asynchronous Gossip. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [74] Chryssis Georgiou, Seth Gilbert, Rachid Guerraoui, and Dariusz R. Kowalski. 2013. Asynchronous Gossip. *Journal of the ACM (JACM)* (2013).
- [75] Chryssis Georgiou, Seth Gilbert, and Dariusz R. Kowalski. 2011. Meeting the Deadline: On the Complexity of Fault-Tolerant Continuous Gossip. *Distributed Computing* (2011).
- [76] Mohsen Ghaffari and Merav Parter. 2016. A Polylogarithmic Gossip Algorithm for Plurality Consensus. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [77] George Giakkoupis, Yasamin Nazari, and Philipp Woelfel. 2016. How Asynchrony Affects Rumor Spreading Time. In *ACM Symposium on Principles of Distributed Computing (PODC)*.

- [78] Yossi Gilad, Rotem Hemo, Silvio Micali, Georgios Vlachos, and Nickolai Zeldovich. 2017. Algorand: Scaling Byzantine Agreements for Cryptocurrencies. In *Symposium on Operating Systems Principles (SOSP)*.
- [79] Rachid Guerraoui, Florian Huc, and Anne-Marie Kermarrec. 2013. Highly Dynamic Distributed Computing With Byzantine Failures. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [80] Rachid Guerraoui, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, and Dragos Seredinschi. 2019. The Consensus Number of a Cryptocurrency. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [81] Rachid Guerraoui, Petr Kuznetsov, Matteo Monti, Matej Pavlovic, and Dragos-Adrian Seredinschi. 2019. Scalable Byzantine Reliable Broadcast. *International Symposium on Distributed Computing (DISC) (2019)*.
- [82] Vassos Hadzilacos and Sam Toueg. 1993. Fault-Tolerant Broadcasts and Related Problems. In *Distributed Systems*.
- [83] Andreas Haeberlen, Petr Kouznetsov, and Peter Druschel. 2007. PeerReview: Practical Accountability for Distributed Systems. In *Symposium on Operating Systems Principles (SOSP)*.
- [84] Bernhard Haeupler, Gopal Pandurangan, David Peleg, Rajmohan Rajaraman, and Zhifeng Sun. 2012. Discovery Through Gossip. In *ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*.
- [85] James Hendricks, Gregory R Ganger, and Michael K Reiter. 2007. Verifying Distributed Erasure-Coded Data. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [86] Maurice Herlihy. 1991. Wait-Free Synchronization. *ACM Transactions on Programming Languages and Systems (TOPLAS) (1991)*.
- [87] Zsolt István, David Sidler, Gustavo Alonso, and Marko Vukolic. 2016. Consensus in a Box: Inexpensive Coordination in Hardware. In *Symposium on Networked Systems Design and Implementation (NSDI)*.
- [88] Kazuharu Itakura and Katsuhiko Nakamura. 1983. A Public-Key Cryptosystem Suitable for Digital Multisignatures. *NEC Research & Development (1983)*.
- [89] Márk Jelasity, Alberto Montresor, and Ozalp Babaoglu. 2009. T-Man: Gossip-Based Fast Overlay Topology Construction. *Computer Networks (2009)*.
- [90] Simon Josefsson and Ilari Liusvaara. 2017. Edwards-Curve Digital Signature Algorithm (EdDSA). RFC 8032.

## Bibliography

---

- [91] Rüdiger Kapitza, Johannes Behl, Christian Cachin, Tobias Distler, Simon Kuhnle, Seyed Vahid Mohammadi, Wolfgang Schröder-Preikschat, and Klaus Stengel. 2012. Cheap-BFT: Resource-Efficient Byzantine Fault Tolerance. In *European Conference on Computer Systems (EuroSys)*.
- [92] Idit Keidar, Eleftherios Kokoris-Kogias, Oded Naor, and Alexander Spiegelman. 2021. All You Need Is DAG. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [93] Mahimna Kelkar, Fan Zhang, Steven Goldfeder, and Ari Juels. 2020. Order-Fairness for Byzantine Consensus. In *Annual International Cryptology Conference (CRYPTO)*.
- [94] Valerie King, Steven Lonargan, Jared Saia, and Amitabh Trehan. 2011. Load Balanced Scalable Byzantine Agreement Through Quorum Building, With Full Information. In *International Conference of Distributed Computing and Networking (ICDCN)*.
- [95] Valerie King, Jared Saia, Vishal Sanwalani, and Erik Vee. 2006. Scalable Leader Election. In *ACM-SIAM Symposium on Discrete Algorithms (SODA)*.
- [96] Eleftherios Kokoris Kogias, Philipp Jovanovic, Nicolas Gailly, Ismail Khoffi, Linus Gasser, and Bryan Ford. 2016. Enhancing Bitcoin Security and Performance With Strong Consistency via Collective Signing. In *USENIX Security Symposium (SEC)*.
- [97] Marios Kogias and Edouard Bugnion. 2020. HovercRaft: Achieving Scalability and Fault-Tolerance for Microsecond-Scale Datacenter Services. In *European Conference on Computer Systems (EuroSys)*.
- [98] Eleftherios Kokoris-Kogias, Philipp Jovanovic, Linus Gasser, Nicolas Gailly, Ewa Syta, and Bryan Ford. 2018. OmniLedger: A Secure, Scale-Out, Decentralized Ledger via Sharding. In *IEEE Symposium on Security and Privacy (S&P)*.
- [99] Ramakrishna Kotla, Lorenzo Alvisi, Mike Dahlin, Allen Clement, and Edmund Wong. 2007. Zyzzyva: Speculative Byzantine Fault Tolerance. In *Symposium on Operating Systems Principles (SOSP)*.
- [100] Klaus Kursawe and Victor Shoup. 2005. Optimistic Asynchronous Atomic Broadcast. In *International Colloquium on Automata, Languages and Programming (ICALP)*.
- [101] Petr Kuznetsov, Yvonne-Anne Pignolet, Pavel Ponomarev, and Andrei Tonkikh. 2021. Permissionless and Asynchronous Asset Transfer. In *International Symposium on Distributed Computing (DISC)*.
- [102] Aptos Labs. 2022. The Aptos Blockchain: Safe, Scalable, and Upgradeable Web3 Infrastructure.
- [103] Leslie Lamport, Dahlia Malkhi, and Lidong Zhou. 2010. Reconfiguring a State Machine. *SIGACT News* (2010).



- [104] Leslie Lamport, Robert Shostak, and Marshall Pease. 1982. The Byzantine Generals Problem. *ACM Transactions on Programming Languages and Systems (TOPLAS)* (1982).
- [105] Jialin Li, Ellis Michael, Naveen Kr. Sharma, Adriana Szekeres, and Dan R. K. Ports. 2016. Just Say NO to Paxos Overhead: Replacing Consensus With Network Ordering. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [106] Meng-Jang Lin, Keith Marzullo, and Stefano Masini. 2000. Gossip Versus Deterministically Constrained Flooding on Small Networks. In *International Symposium on Distributed Computing (DISC)*.
- [107] Joshua Lind, Oded Naor, Ittay Eyal, Florian Kelbert, Emin Gün Sirer, and Peter Pietzuch. 2019. Teechain: A Secure Payment Network With Asynchronous Blockchain Access. In *Symposium on Operating Systems Principles (SOSP)*.
- [108] Shengyun Liu, Paolo Viotti, Christian Cachin, Vivien Quema, and Marko Vukolic. 2016. XFT: Practical Fault Tolerance Beyond Crashes. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [109] Marta Lohava, Giuliano Losa, David Mazières, Graydon Hoare, Nicolas Barry, Eli Gafni, Jonathan Jove, Rafał Malinowsky, and Jed McCaleb. 2019. Fast and Secure Global Payments With Stellar. In *Symposium on Operating Systems Principles (SOSP)*.
- [110] D. Malkhi, M. Merritt, and O. Rodeh. 1997. Secure Reliable Multicast Protocols in a WAN. In *IEEE International Conference on Distributed Computing Systems (ICDCS)*.
- [111] Dahlia Malkhi and Michael Reiter. 1996. A High-Throughput Secure Reliable Multicast Protocol. *IEEE Computer Security Foundations Symposium (CSF)* (1996).
- [112] Dahlia Malkhi and Michael Reiter. 1997. Byzantine Quorum Systems. In *Symposium on the Theory of Computing (STOC)*.
- [113] Dahlia Malkhi, Michael K Reiter, Avishai Wool, and Rebecca N Wright. 2001. Probabilistic Quorum Systems. *Information and Computation* (2001).
- [114] Jean-Philippe Martin and Lorenzo Alvisi. 2005. Fast Byzantine Consensus. In *Dependable Systems and Networks (DSN)*.
- [115] Gregory Maxwell, Andrew Poelstra, Yannick Seurin, and Pieter Wuille. 2019. Simple Schnorr Multi-Signatures With Applications to Bitcoin. *Designs, Codes and Cryptography (DCC)* (2019).
- [116] David Mazieres. 2016. The Stellar Consensus Protocol: A Federated Model for Internet-Level Consensus.
- [117] Ralph C Merkle. 1987. A Digital Signature Based on a Conventional Encryption Function. In *Annual International Cryptology Conference (CRYPTO)*.

## Bibliography

---

- [118] Andrew Miller, Yu Xia, Kyle Croman, Elaine Shi, and Dawn Song. 2016. The Honey Badger of BFT Protocols. In *Annual ACM Conference on Computer and Communications Security (CCS)*.
- [119] Peyman Momeni, Sergey Gorbunov, and Bohan Zhang. 2023. FairBlock: Preventing Blockchain Front-Running With Minimal Overheads. In *Security and Privacy in Communication Networks (SecureComm)*.
- [120] Iulian Moraru, David G. Andersen, and Michael Kaminsky. 2013. There Is More Consensus in Egalitarian Parliaments. In *Symposium on Operating Systems Principles (SOSP)*.
- [121] Achour Mostefaoui, Hamouma Moumen, and Michel Raynal. 2014. Signature-Free Asynchronous Byzantine Consensus With  $t < n/3$  and  $O(n^2)$  Messages. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [122] Satoshi Nakamoto. 2008. Bitcoin: A Peer-to-Peer Electronic Cash System.
- [123] Kartik Nayak, Ling Ren, Elaine Shi, Nitin H Vaidya, and Zhuolun Xiang. 2020. Improved Extension Protocols for Byzantine Broadcast and Agreement. In *International Symposium on Distributed Computing (DISC)*.
- [124] Ray Neiheiser, Miguel Matos, and Luís Rodrigues. 2021. Kauri: Scalable BFT Consensus With Pipelined Tree-Based Dissemination and Aggregation. In *Symposium on Operating Systems Principles (SOSP)*.
- [125] Fernando Pedone and André Schiper. 2002. Handling Message Semantics With Generic Broadcast Protocols. *Distributed Computing* (2002).
- [126] James S Plank and Lihao Xu. 2006. Optimizing Cauchy Reed-Solomon Codes for Fault-Tolerant Network Storage Applications. In *IEEE International Symposium on Network Computing and Applications (NCA)*.
- [127] Dan R. K. Ports, Jialin Li, Vincent Liu, Naveen Kr. Sharma, and Arvind Krishnamurthy. 2015. Designing Distributed Systems Using Approximate Synchrony in Data Center Networks. In *Symposium on Networked Systems Design and Implementation (NSDI)*.
- [128] Kaihua Qin, Liyi Zhou, and Arthur Gervais. 2022. Quantifying Blockchain Extractable Value: How Dark Is the Forest?. In *IEEE Symposium on Security and Privacy (S&P)*.
- [129] HariGovind V Ramasamy and Christian Cachin. 2005. Parsimonious Asynchronous Byzantine-Fault-Tolerant Atomic Broadcast. In *International Conference on Principles of Distributed Systems (OPODIS)*.
- [130] Irving S Reed and Gustave Solomon. 1960. Polynomial Codes Over Certain Finite Fields. *Journal of the Society for Industrial and Applied Mathematics (SIAP)* (1960).

- [131] Michael K. Reiter. 1994. Secure Agreement Protocols: Reliable and Atomic Group Multicast in Rampart. In *Annual ACM Conference on Computer and Communications Security (CCS)*.
- [132] Michael K. Reiter and Kenneth P. Birman. 1994. How to Securely Replicate Services. *ACM Transactions on Programming Languages and Systems (TOPLAS)* (1994).
- [133] Nuno Santos and André Schiper. 2013. Optimizing Paxos With Batching and Pipelining. *Theoretical Computer Science* (2013).
- [134] Christian Scheideler. 2005. How to Spread Adversarial Nodes? Rotate!. In *Symposium on the Theory of Computing (STOC)*.
- [135] Fred B. Schneider. 1990. Implementing Fault-Tolerant Services Using the State Machine Approach: A Tutorial. *ACM Computing Surveys (CSUR)* (1990).
- [136] Claus P. Schnorr. 1991. Efficient Signature Generation by Smart Cards. *Journal of Cryptology (JCrypt)* (1991).
- [137] Victor Shoup. 2000. Practical Threshold Signatures. In *International Conference on the Theory and Application of Cryptographic Techniques (EUROCRYPT)*.
- [138] Atul Singh, Tathagata Das, Petros Maniatis, Peter Druschel, and Timothy Roscoe. 2008. BFT Protocols Under Fire. In *Symposium on Networked Systems Design and Implementation (NSDI)*.
- [139] Suman Sourav, Peter Robinson, and Seth Gilbert. 2018. Slow Links, Fast Links, and the Cost of Gossip. *IEEE International Conference on Distributed Computing Systems (ICDCS)* (2018).
- [140] Alexander Spiegelman, Neil Giridharan, Alberto Sonnino, and Lefteris Kokoris-Kogias. 2022. Bullshark: DAG BFT Protocols Made Practical. In *Annual ACM Conference on Computer and Communications Security (CCS)*.
- [141] Chrysoula Stathakopoulou, Matej Pavlovic, and Marko Vukolić. 2022. State Machine Replication Scalability Made Simple. In *European Conference on Computer Systems (EuroSys)*.
- [142] Chrysoula Stathakopoulou, David Tudor, Matej Pavlovic, and Marko Vukolić. 2022. [Solution] Mir-BFT: Scalable and Robust BFT for Decentralized Networks. *Journal of Systems Research (JSys)* (2022).
- [143] The Diem Team. 2021. DiemBFT V4: State Machine Replication in the Diem Blockchain.
- [144] The DFINITY Team. 2022. The Internet Computer for Geeks. *IACR Cryptology ePrint Archive* (2022).
- [145] The MystenLabs Team. 2022. The Sui Smart Contracts Platform.

## Bibliography

---

- [146] Sam Toueg. 1984. Randomized Byzantine Agreements. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [147] Jelle van den Hooff, David Lazar, Matei Zaharia, and Nickolai Zeldovich. 2015. Vuvuzela: Scalable Private Messaging Resistant to Traffic Analysis. In *Symposium on Operating Systems Principles (SOSP)*.
- [148] Spyros Voulgaris, Márk Jelasity, and Maarten van Steen. 2004. A Robust and Scalable Peer-to-Peer Gossiping Protocol. In *Agents and Peer-to-Peer Computing (AP2PC)*.
- [149] Marko Vukolic. 2010. The Origin of Quorum Systems. *Bulletin of the EATCS* (2010).
- [150] Cheng Wang, Jianyu Jiang, Xusheng Chen, Ning Yi, and Heming Cui. 2017. APUS: Fast and Scalable Paxos on RDMA. In *ACM Symposium on Cloud Computing (SoCC)*.
- [151] Jiaping Wang and Hao Wang. 2019. Monoxide: Scale Out Blockchain With Asynchronous Consensus Zones. In *Symposium on Networked Systems Design and Implementation (NSDI)*.
- [152] Roger Wattenhofer. 2019. *Blockchain Science: Distributed Ledger Technology*.
- [153] Gavin Wood. 2014. Ethereum: A Secure Decentralised Generalised Transaction Ledger.
- [154] Andrew Chi-Chih Yao. 1979. Some Complexity Questions Related to Distributive Computing. In *Symposium on the Theory of Computing (STOC)*.
- [155] Jian Yin, Jean-Philippe Martin, Arun Venkataramani, Lorenzo Alvisi, and Mike Dahlin. 2003. Separating Agreement From Execution for Byzantine Fault Tolerant Services. In *Symposium on Operating Systems Principles (SOSP)*.
- [156] Maofan Yin, Dahlia Malkhi, Michael K. Reiter, Guy Golan Gueta, and Ittai Abraham. 2019. HotStuff: BFT Consensus With Linearity and Responsiveness. In *ACM Symposium on Principles of Distributed Computing (PODC)*.
- [157] B. Zhang, K. Han, B. Ravindran, and E. D. Jensen. 2008. RTQG: Real-Time Quorum-Based Gossip Protocol for Unreliable Networks. In *International Conference on Availability, Reliability and Security (ARES)*.
- [158] Haoqian Zhang, Louis-Henri Merino, Vero Estrada-Galiñanes, and Bryan Ford. 2022. Flash Freezing Flash Boys: Countering Blockchain Front-Running. In *IEEE International Conference on Distributed Computing Systems (ICDCS)*.
- [159] Yunhao Zhang, Srinath Setty, Qi Chen, Lidong Zhou, and Lorenzo Alvisi. 2020. Byzantine Ordered Consensus Without Byzantine Oligarchy. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.
- [160] Liyi Zhou, Kaihua Qin, Christof Ferreira Torres, Duc V Le, and Arthur Gervais. 2021. High-Frequency Trading on Decentralized on-Chain Exchanges. In *IEEE Symposium on Security and Privacy (S&P)*.

# News Readings

- [Cer19] Megan Cerullo. Google Reportedly Mining Millions of Americans Personal Health Data. *CBS News*, 2019.
- [Con18] Nicholas Confessore. Cambridge Analytica and Facebook: The Scandal and the Fallout so Far. *The New York Times*, 2018.
- [cyb23] Significant Cyber Incidents. *Center for Strategic & International Studies*, 2023.
- [ESWV22] P. Ehin, M. Solvak, J. Willemson, and P. Vinkel. Internet Voting in Estonia 2005–2019: Evidence from Eleven Elections. *Government Information Quarterly*, 2022.
- [Far22] Malcomb Farber. Cybercrime Damages to Cost the World \$8 Trillion USD in 2023. *Associated Press*, 2022.
- [Fit22] Afiq Fitri. Big Tech Now Accounts for More than Half of Global Internet Traffic. *Tech Monitor*, 2022.
- [Gay21] Damien Gayle. Facebook Aware of Instagram’s Harmful Effect on Teenage Girls, Leak Reveals. *The Guardian*, 2021.
- [int21] IAB Study Shows Internet Economy is Transforming the U.S. Economy, Creating New Markets and Spurring Job Growth for Large and Small Businesses. *Interactive Advertising Bureau*, 2021.
- [Moz18] Paul Mozur. A Genocide Incited on Facebook, with Posts from Myanmar’s Military. *The New York Times*, 2018.
- [Pri19] Rob Price. The FTC has Approved a Roughly \$5 Billion Settlement with Facebook. *Business Insider*, 2019.
- [Reu20] Reuters. Google Faces \$5 Billion Lawsuit in U.S. for Tracking ‘Private’ Internet Use. *NBC News*, 2020.
- [Run20] Dan Runkevicius. How Amazon Quietly Powers the Internet. *Forbes*, 2020.