

Tendermint: Byzantine Fault Tolerance in the Age of Blockchains

by

Ethan Buchman

A Thesis
presented to
The University of Guelph

In partial fulfilment of requirements
for the degree of
Master of Applied Science
in
Engineering Systems and Computing

Guelph, Ontario, Canada

©Ethan Buchman, May, 2016

ABSTRACT

TENDERMINT: BYZANTINE FAULT TOLERANCE IN THE AGE OF BLOCKCHAINS

Ethan Buchman
University of Guelph, 2016

Advisor:
Professor Graham Taylor

Tendermint is a new protocol for ordering events in a distributed network under adversarial conditions. More commonly known as consensus or atomic broadcast, the problem has attracted significant attention recently due to the widespread success of digital currencies, such as Bitcoin and Ethereum, which successfully solve the problem in public settings without a central authority. Tendermint modernizes classic academic work on the subject to provide a secure consensus protocol with accountability guarantees, as well as an interface for building arbitrary applications above the consensus. Tendermint is high performance, achieving thousands of transactions per second on dozens of nodes distributed around the globe, with latencies of about one second, and performance degrading moderately in the face of adversarial attacks.

Dedicated to Theda.

Preface

The structure and presentation of this thesis was much inspired by Diego Ongaro’s 2014 Doctoral Dissertation, “Consensus: Bridging Theory and Practice”, wherein he specifies and evaluates the Raft consensus algorithm.

Much of the work done in this thesis was done in collaboration with Jae Kwon, who initiated the Tendermint project. Please see the Github repository, at <https://github.com/tendermint/tendermint>, for a more direct account of contributions to the codebase.

Acknowledgments

I learned early in life from Tony Montana that a man has only two things in this world, his word and his balls, and he should break em for nobody. This thesis would not have been completed if I had not given my word to certain people that I would complete it. These include my family, in particular my parents, grandparents, and great uncle Paul, and my primary adviser, Graham, who has, for one reason or another, permitted me a practically abusive amount of flexibility to pursue the topic of my choosing. Thanks Graham.

Were it not for another set of individuals, this thesis would probably have been about machine learning. These include Vlad Zamfir, with whom I have experienced countless moments of discovery and insight; My previous employer and favorite company, Eris Industries, and especially their CEO and COO, Casey Kuhlman and Preston Byrne, for hiring me, mentoring me, and giving me such freedom to research and tinker and ultimately start my own company with technology they helped fund; Jae Kwon, for his direct mentorship in consensus science and programming, for being a great collaborator, and for being the core founder and CEO at Tendermint; Lucius Meredith, for mentoring me in the process calculi; Zach Ramsay, for being, for all intents and purposes, my heterosexual husband; and of course, Satoshi Nakamoto, whomever you are, for sending me down this damned rabbit hole in the first place.

There are of course many other people who have influenced my life during the course of this graduate degree; you know who you are, and I thank you for being that person and for all you've done for me.

Contents

1	Background	5
1.1	Replicated State Machine	5
1.2	Asynchrony	6
1.3	Broadcast and Consensus	8
1.4	Byzantine Fault Tolerance	10
1.5	Cryptography, Trust, and Economics	13
1.6	Blockchain	14
1.7	Process Calculus	14
1.8	The Need For Tendermint	17
2	Tendermint Consensus	18
2.1	Tendermint Overview	18
2.2	Consensus	19
2.2.1	Proposals	22
2.2.2	Votes	22
2.2.3	Locks	23
2.2.4	Formal Specification	25
2.3	Blockchain	28
2.3.1	Why Blocks?	28
2.3.2	Block Structure	29
2.4	Safety	29
2.5	Accountability	31
2.6	Faults and Availability	33
2.7	Conclusion	34
3	Tendermint Subprotocols	35
3.1	P2P-Networking	35
3.2	Consensus Gossip	36

3.2.1	Block Data	36
3.2.2	Votes	37
3.3	Mempool	37
3.4	Syncing the Blockchain	38
4	Building Applications	39
4.1	Background	39
4.2	Tendermint Socket Protocol	40
4.3	Separating Agreement and Execution	43
4.4	Microservice Architecture	44
4.5	Determinism	45
4.6	Termination	45
4.7	Examples	46
4.7.1	Merkleeyes	46
4.7.2	Basecoin	47
4.7.3	Ethereum	47
4.8	Conclusion	48
5	Governance	49
5.1	Government	49
5.2	Validator Set Changes	50
5.3	Punishing Byzantine Validators	51
5.4	Software Upgrades	52
5.5	Crisis Recovery	53
5.6	Conclusion	54
6	Client Considerations	55
6.1	Discovery	55
6.2	Broadcasting Transactions	55
6.3	Mempool	56
6.4	Semantics	57
6.5	Reads	57
6.6	Light Client Proofs	58
7	Implementation	59
7.1	Binary Serialization	59
7.2	Cryptography	60
7.3	Merkle Hash Tree	60

7.4	RPC	61
7.5	P2P Networking	61
7.6	Reactors	61
7.6.1	Mempool	61
7.6.2	Consensus	62
7.6.3	Blockchain	62
7.7	Conclusion	63
8	Performance and Fault Tolerance	64
8.1	Overview	64
8.2	Throughput and Latency	65
8.3	Crash Failures	66
8.4	Random Network Delay	71
8.5	Byzantine Failures	71
8.6	Related Work	73
8.7	Conclusion	73
9	Related Work	75
9.1	Beginnings	75
9.1.1	Faulty Things	76
9.1.2	Clocks	76
9.1.3	FLP	77
9.1.4	Common Coin	78
9.1.5	Transaction Processing	78
9.1.6	Broadcast Protocols	79
9.2	Byzantine	79
9.2.1	Byzantine Generals	79
9.2.2	Randomized Consensus	80
9.2.3	Partial Synchrony	80
9.2.4	PBFT	81
9.2.5	BFT Improvements	82
9.3	Non-Byzantine	83
9.3.1	Paxos	83
9.3.2	Raft	83
9.4	Blockchain	84
9.4.1	Bitcoin	84
9.4.2	Ethereum	84
9.4.3	Proof-of-Stake	85

9.4.4	HyperLedger	85
9.4.5	HoneyBadgerBFT	86
9.5	Conclusion	87
10	Conclusion	88

List of Figures

1.1	Overview of replicated state machine architecture	7
1.2	Byzantine processes tell lies	12
2.1	Overview of Tendermint consensus logic	20
2.2	Formal specification of Tendermint consensus in the π -calculus, part I	26
2.3	Formal specification of Tendermint consensus in the π -calculus, part II	27
4.1	TMSP Message Types	41
4.2	TMSP Architecture	42
8.1	Latency-Throughput trade-off in non-faulty global network . .	67
8.2	Latency-throughput trade-off in non-faulty local network . . .	68
8.3	Latency-Throughput trade-off in non-faulty global network of large machines	69

List of Tables

8.1	Latency statistics under crash faults	70
8.2	Latency statistics under randomized delays	72
8.3	Latency statistics under Byzantine faults	74

chapterIntroduction

The cold, hard truth about computer engineering today is that computers are faulty - they crash, corrupt, slow down, perform voodoo. What's worse, we're typically interested in connecting computers over a network (like the Internet), and networks can be more unpredictable than the computers themselves. These challenges are primarily the concern of "fault tolerant distributed computing", whose aim is to discover principled protocol designs enabling faulty computers communicating over a faulty network to stay in sync while providing a useful service. In essence, to make a reliable system from unreliable parts.

In an increasingly digital and globalized world, however, systems must not only be reliable in the face of unreliable parts, but in the face of malicious or "Byzantine" ones. Over the last decade, major components of critical infrastructure have been ported to networked systems, as have vast components of the world's finances. In response, there has been an explosion of cyber warfare and financial fraud, and a complete distortion of economic and political fundamentals.

In 2009, an anonymous software developer known only as Satoshi Nakamoto introduced an approach to the resolution of these issues that was simultaneously an experiment in computer science, economics, and politics. It was a digital currency called Bitcoin [71]. Bitcoin was the first protocol to solve the problem of fault tolerant distributed computing in the face of malicious adversaries in a public setting. The solution, dubbed a "blockchain", hosts a digital currency, where consent on the order of transactions is negotiated via an economically incentivized cryptographic random lottery based on partial hash collisions. In essence, transactions are ordered in batches (blocks) by those who find partial hash collisions of the transaction data, in such a way that the correct ordering is the one where the collisions have the greatest cumulative difficulty. The solution was dubbed Proof-of-Work (PoW).

Bitcoin's subtle brilliance was to invent a currency, a cryptocurrency, and to issue it to those solving the hash collisions, in exchange for their doing such an expensive thing as solve partial hash collisions. In spirit, it might be assumed that the capacity to solve such problems would be distributed as computing power is, such that anyone with a CPU could participate. Unfortunately, the reality is that the Bitcoin network has grown into the largest supercomputing entity on the planet, greater than all others combined, evaluating only a single function, distributed across a few large data centers running Application Specific integrated circuits (ASICs) produced by

a small number of primarily Chinese companies, and costing on the order of two million USD per day in electricity [7]. Not to mention it takes up to an hour to confirm transactions, is difficult to build on top of, and does not scale in a way which preserves its security guarantees. This is not to mention the internal bout of political struggles resulting from the immaturity of the Bitcoin community’s governance mechanisms.

Despite these troubles, Bitcoin, astonishingly, continues to churn, and its technology, of cryptography and distributed databases and co-operative economics, continues to attract billions in investment capital, both in the form of new companies and new cryptocurrencies, each diverging from Bitcoin in its own unique way.

In 2014, Jae Kwon began the development of Tendermint, which sought to solve the consensus problem, of ordering and executing a set of transactions in an adversarial environment, by modernizing solutions to the problem that have existed for decades, but have lacked the social context to be deployed widely until now.

In early 2015, in an effort led by Eris Industries to bring a practical blockchain solution to industry, the author joined Jae Kwon in the development of the Tendermint software and protocols.

The result of that collaboration is the Tendermint platform, consisting of a consensus protocol, a high-performance implementation in golang, a flexible interface for building arbitrary applications above the consensus, and a suite of tools for deployments and their management. We believe Tendermint achieves a superior design and implementation compared to previous approaches, including that of the classical academic literature [31, 17, 75] as well as Bitcoin [71] and its derivatives [105, 4, 55] by combining the right elements of each to achieve a practical balance of security, performance, and simplicity. We have since started a company to pursue the deployment of the platform at an enterprise level, and have seen tremendous interest.

The primary contributions of this thesis are as follows:

- A complete description of the Tendermint consensus algorithm and platform architecture (Chapters 2-4). Tendermint provides provably optimal Byzantine Fault Tolerant consensus, is easy to understand and reason about, flexible to build on top of, and provides a foundation for advanced and secure socio-political and economic systems.
- Significant contributions to a high performance implementation of the

consensus algorithm in Go, most notably a refactor of the core consensus state machine to be more robust, deterministic, and understandable, as well as countless bug fixes and performance improvements.

- Evaluation of the implementation’s performance and characteristics in normal, faulty, and malicious conditions on large deployments (Chapter 8). Tendermint consensus is mostly asynchronous, with minimal dependence on synchronized clocks, and can tolerate up to a third of its nodes behaving arbitrarily.
- Contributions towards the design and implementation of many other elements necessary for a complete system, such as membership changes, crisis recovery, client design, and security (Chapters 5 and 6).
- A formal specification of Tendermint in the π -calculus and an informal proof of correctness of its safety and accountability (Chapter 2).
- An explanation of blockchains and the Tendermint consensus algorithm in the context of classical academic consensus research (Chapter 9).

All code is available open source at <https://github.com/tendermint/tendermint>, and in associated repositories at <https://github.com/tendermint>. The core is licensed GPLv3 and most of the libraries are Apache 2.0.

The remainder of this thesis introduces the consensus problem, classic solutions, and motivates Tendermint (Chapter 1); presents the Tendermint consensus algorithm, the peer-to-peer network, and the interface for building arbitrary applications on top (Chapters 2-4); mechanisms for governance and membership changes, and how clients interact with Tendermint (Chapters 5 and 6); outlines the implementation in Go and evaluates its fault tolerance and performance (Chapters 7 and 8); and discusses related work (Chapter 9).

Chapter 1

Background

Distributed consensus systems have become a critical component of modern Internet infrastructure, powering every major Internet application at some level or another. This chapter introduces the necessary background material for understanding and discussing these systems. In addition, it introduces the pi-calculus, a formal language for describing concurrent processes, which will be used to specify the Tendermint algorithm in Chapter 2.

1.1 Replicated State Machine

The most common paradigm for studying and implementing distributed consensus is that of the Replicated State Machine, wherein a *deterministic* state machine is replicated across a set of processes, such that it functions as a single state machine despite the failure of some processes [87]. The state machine is driven by a set of inputs, known as *transactions*, where each transaction may or may not, depending on its validity, cause a state transition and return a result. More formally, a transaction is an *atomic* operation on a database, meaning it either completes or doesn't occur at all, and can't be left in an intermediate state [42]. The state transition logic is governed by the state machine's state transition function, which maps a transaction and the current state to a new state and a return value. The state transition function is also sometimes referred to as *application logic*.

It is the responsibility of the consensus protocol to order the transactions so that the resulting *transaction log* is replicated exactly by every process. Using a deterministic state transition function implies that every process will

compute the same state given the same transaction log.

A summary of the replicated state machine architecture is given in Figure 1.1.

Tendermint was motivated from the desire to create a general purpose, high-performance, secure, and robust replicated state machine.

1.2 Asynchrony

The purpose of a fault-tolerant replicated state machine is to co-ordinate a network of computers to stay in sync while providing a useful service, despite the presence of faults.

Staying in sync amounts to replicating the transaction log successfully; providing a useful service amounts to keeping the state machine available for new transactions. These aspects of the system are traditionally known as *safety* and *liveness*, respectively. Colloquially, safety means nothing bad happens; liveness means that something good eventually happens. A violation of safety implies two or more valid, competing transaction logs. Violating liveness implies an unresponsive network.

It is trivial to satisfy liveness by accepting all transactions. And it is trivial to satisfy safety by accepting none. Hence, state machine replication algorithms can be seen to operate on a spectrum defined by these extremes. Typically, processes require some threshold of received information from other processes before they commit a new transaction. In synchronous environments, where we make assumptions about the maximum delay of network messages or the maximum speed of processor clocks, it is easy enough to take turns proposing new transactions, poll for a majority vote, and skip a proposer's turn if they don't propose within the bounds of the synchrony assumptions.

In asynchronous environments, where no such assumptions about network delays or processor speeds are warranted, the trade-off is much more difficult to manage. In fact, the so called FLP impossibility result demonstrates the impossibility of distributed consensus among deterministic asynchronous processes if even a single processes can crash¹[37]. The proof amounts to showing that, because processes can fail, there are valid executions of the protocol in which processes fail at the exact opportune times to prevent consensus. Hence, we have no guarantee of consensus.

¹Prior to FLP, the distinction between sync/async wasn't as prominent

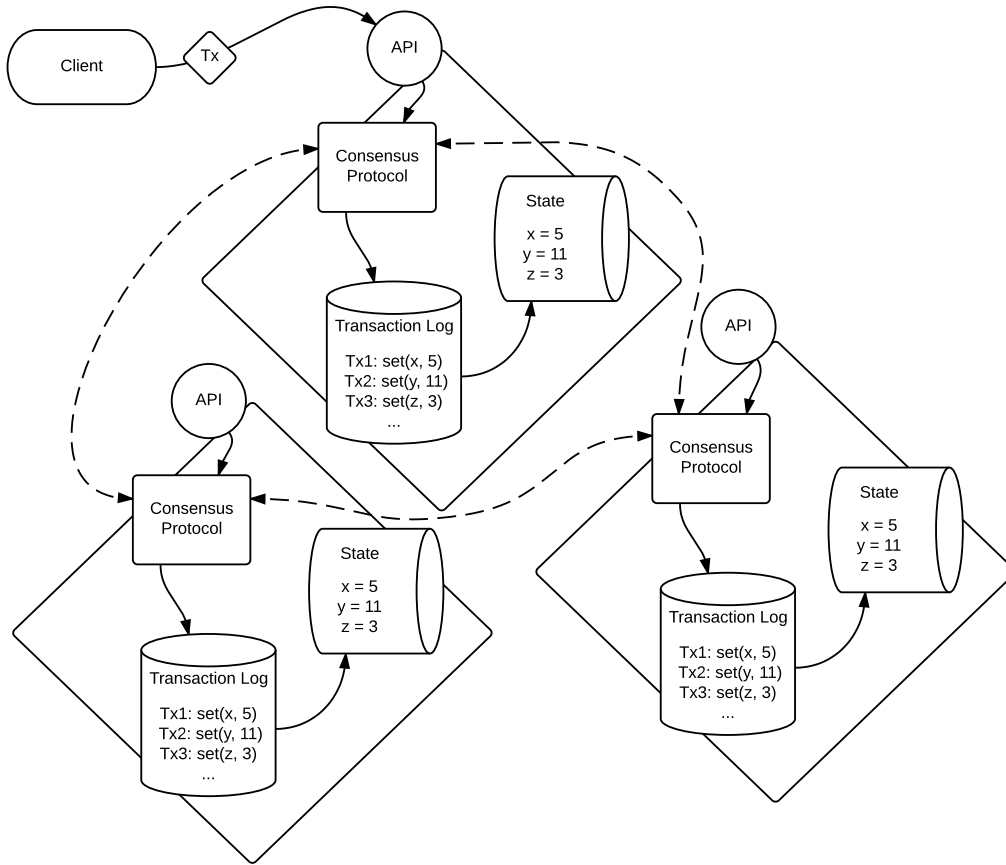


Figure 1.1: A replicated state machine replicates a transaction log and resulting state across multiple machines. Transactions are received from the client, run through the consensus protocol, ordered in the transaction log, and executed against the state. In the figure, each diamond represents a single machine, with dotted lines representing communication between machines to carry out the consensus protocol for ordering transactions.

Typically, synchrony in a protocol is reflected by the use of timeouts to manage certain transitions. In asynchronous environments, where messages can be arbitrarily delayed, relying on synchrony (timeouts) for safety can lead to a fork in the transaction log. Relying on synchrony to ensure liveness can cause the consensus to halt, and the service to become unresponsive. The former case is usually considered more severe, as reconciling conflicting logs can be a daunting or impossible task.

In practice, synchronous solutions are only used where the message latency is under extremely well defined control, for instance between controllers on an airplane [49], or between datacenters that utilizing synchronized atomic clocks [23]. Thus, while many efficient synchronous solutions exist, the general unreliability of computer networks is too great a risk for them to be used in practice without significant additional costs.

There are fundamentally two ways to overcome the FLP impossibility result. The first is to use stronger synchrony assumptions - even rather weak assumptions are sufficient, for instance, that only eventually, crashed processes are suspected of crashing and correct ones are not [19]. Typically, this approach utilizes *leaders*, which play a special co-ordinating role, and which can be skipped if they are suspected of being faulty after some timeout. In practice, such leader-election mechanisms can be difficult to get right.

The second way to overcome FLP is to use non-determinism - include randomization elements such that the probability of coming to consensus tends to 1. While clever, relying on randomization is typically much slower, though certain advanced cryptographic techniques have in recent years achieved tremendous improvements in speed [67]

1.3 Broadcast and Consensus

In order for a process to replicate its state on other processes, it must have access to basic communication primitives which allow it to disseminate, or deliver, information. One of the most useful such primitives is *reliable broadcast*. Reliable broadcast (RBC) is a broadcast primitive satisfying, for message m [19]:

- validity - if a correct process broadcasts m , it eventually delivers m
- agreement - if a correct process delivers m , all correct processes eventually deliver m

- integrity - m is only delivered once, and only if broadcast by its sender

In essence, RBC enables a message to be eventually delivered once on all correct processes.

Another, more useful primitive is *atomic broadcast* (ABC), which satisfies RBC and an additional property [19]:

- total order - if correct processes p and q deliver m and m' , then p delivers m before m' iff q delivers m before m'

Atomic broadcast is thus a reliable broadcast where values are delivered in the same order on each host. Note this is exactly the problem of replicating a transaction log. While colloquially, the problem may be referred to as consensus, the standard definition of the consensus primitive satisfies the following [19]:

- termination - every correct process eventually decides
- integrity - every correct process decides at most once
- agreement - if one correct process decides $v1$ and another decides $v2$, then $v1 = v2$
- validity - if a correct process decides v , at least one process proposed v

Intuitively, consensus and ABC appear remarkably similar, with the critical difference that ABC is a continuous protocol, whereas consensus expects to terminate. That said, it is well known that each can be reduced to the other [19]. Consensus is easily reduced to ABC by deciding the first value to be atomically broadcast. ABC can be reduced to consensus by running many instances of the consensus protocol, in sequence, though certain subtle considerations must be made, especially for handling Byzantine faults. A complete description of the parameter space surrounding the reduction of ABC to consensus remains an open topic of research.

Historically, despite the fact that most use cases actually require ABC, the most widely adopted algorithm has been a consensus algorithm called Paxos, introduced, and proven correct, by Leslie Lamport in the 90s. Paxos simultaneously empowered and confused the discipline of consensus science, on the one hand by providing the first real-world, practical, fault-tolerant consensus algorithm, and on the other by being so difficult to understand

and explain. Each implementation of the algorithm used its own unique bag of ad-hoc techniques to build ABC from Paxos, making the ecosystem difficult to navigate, understand, and utilize. Unfortunately, there was little work on improving the problem framing to make it more understandable, though there were efforts to delineate solutions to the various difficulties [18].

In 2013, Ongaro and Ousterhout published Raft [75], a state machine replication algorithm whose motivating design goal was understandability. Rather than starting from a consensus algorithm, and attempting to build what was needed (ABC), the design of Raft considered first and foremost the transaction log, and sought orthogonal components which could fit together to provide what is ultimately ABC, though it is not described as such.

Paxos has been the staple consensus algorithm for industry, upon which the likes of Amazon [26], Google [10], and others have built out highly available global Internet services. The Paxos consensus sits at the bottom of the application stack, providing a consistent interface to resource management and allocation, operating at much slower time scales than the highly-available applications facing the users.

Since its debut, however, Raft has seen tremendous adoption, especially in the open source community, with implementations in virtually every major language [96], and use as the backbone in major projects, including CoreOs’s distributed Linux distribution [32] and the open source time-series database InfluxDB [51, 45].

Raft’s major divergent design decisions from Paxos was to focus on the transaction-log first, rather than a single value, in particular to allow a leader to persist in committing transactions until he goes down, at which point leadership election can kick in. In some ways, this is similar to the approach taken by blockchains, though the major advantage of blockchains is the ability to tolerate a different kind of fault.

1.4 Byzantine Fault Tolerance

Blockchains have been described as “trust machines” [97] on account of the way they reduce counter party risk through the decentralization of responsibility over a shared database. Bitcoin, in particular, is noted for its ability to withstand attacks and malicious behaviour by any of the participants. Traditionally, consensus protocols tolerant of malicious behaviour were known as

Byzantine Fault Tolerant (BFT) consensus protocols. The term Byzantine was used due to the similarity of the problem to that faced by generals of the Byzantine army attempting to co-ordinate themselves to attack Rome using only messengers, where one of the generals may be a traitor [61].

In a crash fault, a process simply halts. In a Byzantine fault, it can behave arbitrarily. Crash faults are easier to handle, as no process can *lie* to another process. Systems which only tolerate crash faults can operate via simple majority rule, and therefore typically tolerate simultaneous failure of up to half of the system. If the number of failures the system can tolerate is f , such systems must have at least $2f + 1$ processes.

Byzantine failures are more complicated. In a system of $2f + 1$ processes, if f are Byzantine, they can co-ordinate to say arbitrary things to the other $f + 1$ processes. For instance, suppose we are trying to agree on the value of a single bit, and $f = 1$, so we have $N = 3$ processes, A , B , and C , where C is Byzantine, as in Figure 1.2. C can tell A that the value is 0 and tell B that it's 1. If A agrees that its 0, and B agrees that its 1, then they will both think they have a majority and commit, thereby violating the safety condition. Hence, the upper bound on faults tolerated by a Byzantine system is strictly lower than a non-Byzantine one.

In fact, it can be shown that the upper limit on f for Byzantine faults is $f < N/3$ [78]. Thus, to tolerate a single Byzantine process, we require at least $N = 4$. Then the faulty process can't split the vote the way it was able to when $N = 3$.

In 1999, Castro and Liskov published Practical Byzantine Fault Tolerance [17], or *PBFT*, which provided the first optimal Byzantine fault tolerant algorithm in asynchronous networks. It set a new precedent for the practicality of Byzantine fault tolerance in industrial systems by being capable of processing tens of thousands of transactions per second. Despite this success, Byzantine fault tolerance was still considered expensive and largely unnecessary, and the most popular implementation was difficult to build on top of [20]. Hence, despite a resurgence in academic interest, including numerous improved variations [107, 58] not much progress was made in the way of implementations and deployment. Furthermore, PBFT provides no guarantees if a third or more of the network co-ordinates to violate safety.

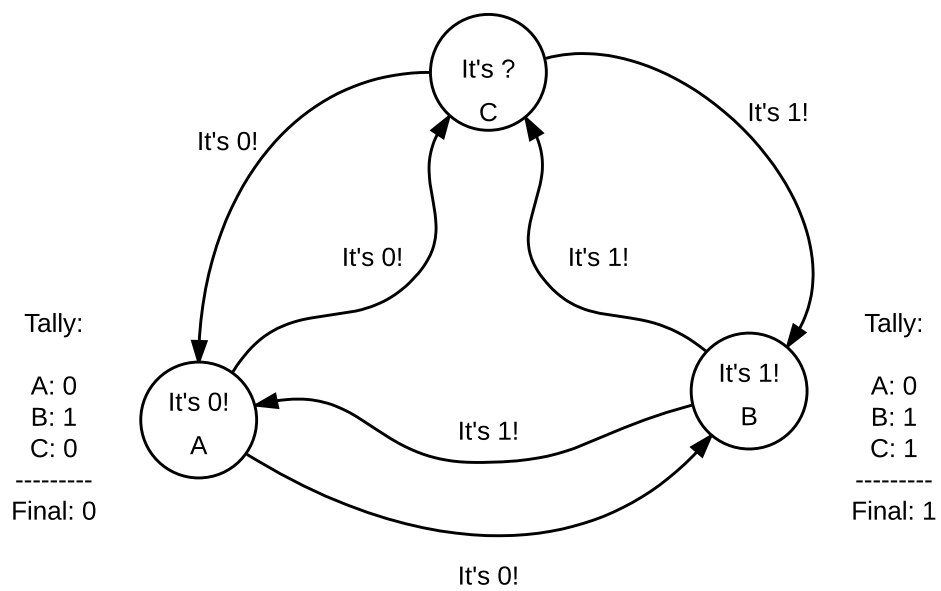


Figure 1.2: A Byzantine process, C, tells A one thing and B another, causing them to come to different conclusions about the network. Here, simple majority vote results in a violation of safety due to only a single Byzantine process.

1.5 Cryptography, Trust, and Economics

Fundamentally, fault tolerance is a problem deriving from a lack of trust - an inability to know how some process will behave. Formally, trust might be defined information theoretically as a means for reducing the entropy of one's model of the world - to trust someone is to optimistically reduce one's uncertainty about the world, enabling more focused attention on higher order forms of organization.

Cryptographic primitives are also fundamentally related to the problem of trust, and may similarly be defined as mechanisms which allow for a massive reduction in entropy - successfully authenticating a cryptographic function collapses a distribution over possible outcomes to a single, or in some cases a small number, of outcomes.

It is well known that civilizations that have greater forms of institutional trust, such as the rule-of-law, have higher productivity and more vibrant economies [108]. The result makes intuitive sense, as being able to trust more about an interaction reduces the space of possible outcomes that need to be actively modelled, making it easier to co-ordinate. Unfortunately, it is becoming increasingly difficult to evaluate the trustworthiness of modern institutions as their complexity has skyrocketed in recent decades, increasing the likelihood that the certainty they allegedly provide is an illusion.

Fortunately, cryptography can form the basis for new institutions of trust in society which may dramatically improve the capacity for human co-ordination at global scale on account of reduced risk of fraudulent and/or unaccountable activity. Of particular interest is the importance of cryptographic primitives in BFT algorithms, both for authentication and for seeding non-determinism.

Most interestingly, economic mechanisms may also serve as means for reducing entropy, in so far as economic agents can be incentivized - which is to say be made more likely to execute a particular behaviour. In fact, Bitcoin's great insight was that cryptographic primitives could be used in conjunction with economic incentives to sufficiently reduce the entropy of a public consensus network to achieve secure replication of state.

A more formal investigation of the information theoretic grounds of trust, cryptography, consensus, and economics, and in particular their inter-relationship, remains for future work.

1.6 Blockchain

A blockchain is, at heart, an integrity-focused approach to Byzantine Fault Tolerant Atomic Broadcast. The Bitcoin blockchain, for instance, uses a combination of economics and cryptographic randomization to provide a strong probabilistic guarantee that safety will not be violated, given a weak synchrony assumption, namely, that blocks are gossipped much more rapidly than they are found via the partial-hash collision lottery. In practice, however, it is well known that Bitcoin’s security guarantees are vulnerable to a number of subtle attacks [24, 33].

The blockchain gets its name from the two key optimizations it employs in solving ABC. The first is that it groups transactions in blocks in order to amortize the high commit latency (on the order of ten minutes) over many transactions. The second is to link blocks via cryptographic hashes into an immutable chain, such that is easy to verify the historical record. Both optimizations are natural improvements to a naive BFT-ABC, the former improving performance, the latter improving tolerance to certain kinds of difficult to model Byzantine faults.

Over the last few years, it has become common to “blockchainize” consensus algorithms, that is, to adopt them to ABC using the blockchain paradigm of hash-linked transaction batches. To the author’s knowledge, Tendermint was the first such proposal, upgrading a well known BFT algorithm from the late 80s [31], though it has since evolved to a consensus algorithm of its own. It has been followed by IBM, which upgraded PBFT to a blockchain [14, 76], and by JP Morgan, which upgraded a BFT version of Raft [9].

1.7 Process Calculus

Distributed systems, where pieces of the system execute concurrently with one another, are notorious for being difficult to design, build, and debug. They are further difficult to formally verify, as most techniques for formal verification, and in fact the very foundations of computer science, have been specifically developed with sequential computation in mind.

Process calculi are a family of models introduced to provide a formal basis for concurrent computation. The most popular calculus, the Communicating Sequential Processes (CSP) [46] forms the theoretical foundation for many modern programming languages, such as Go, which include con-

currency primitives in the language design [89].

In the 80s, Robin Milner introduced the Calculus of Communicating Systems (CCS), designed to be a concurrent analog of the sequential lambda calculus that underlies most functional programming languages. While the lambda calculus has function application as its basic unit of computation, CCS uses communication between two concurrent processes over a shared channel as its basic operational primitive. A more general form of CCS, the π -calculus, enables mobility in the communication graph between processes, such that the channels of communication can themselves be passed along other channels, thereby blurring the distinction between data, variables, and channels. The result is a coherent, minimalistic model of computation more powerful than its sequential predecessors.

The π -calculus has proven to be a highly effective tool for the study of concurrent systems, with applications from business process management [64] to cellular biology [80]. The remarkably simple notation simplifies the description of concurrent protocols. Furthermore, the well known equivalence between computation and logic [2] enables logical systems to be defined complementary to the various process calculi, providing formal means to discuss and verify the properties of systems specified in an appropriate calculus.

Our presentation of the π -calculus is sufficient merely to specify the Tendermint algorithm. For a more complete introduction, see [68].

The grammar of a simple π -calculus, in Backus-Naur form, is as follows:

$$\begin{array}{lll}
 P := & 0 & \text{void} \\
 | & P \mid P & \text{par} \\
 | & \alpha.P & \text{guard} \\
 | & \alpha.P + \alpha.P & \text{guarded-choice} \\
 | & (\nu x)P & \text{fresh} \\
 \\
 & F^s(y) & \text{func} \\
 \\
 \alpha := & \tau & \text{null} \\
 | & x!(y) & \text{send} \\
 | & x?(y) & \text{receive} \\
 | & \text{susp}_i & \text{suspect}
 \end{array}$$

Each grammatical rule is labelled with a reference to its functional meaning. A process may be the empty process, 0 . It may be the parallel composition of two processes, $P \mid P$, denoting two processes running concurrently. A guarded processes, $\alpha.P$, only allows process P to execute after an action, α , has occurred. The action can be a null action, τ , or it can be the sending, $x!(y)$, or receiving, $x?(y)$, of y along x . Guarded choice injects non-determinism into the operation of the calculus, such that the processes $\alpha.P + \beta.Q$ will non-deterministically execute α or β , and then run P or Q , respectively. A new channel, x , can be created via $(\nu x)P$, such that x is only accessible in P . Functional forms $F^s(y)$ allow us to pass variables s and y into the process called F , which may cause it self to execute recursively. Typically, we let s be state-like variables, while y are channels in the calculus. Finally, since we are interested in consensus in asynchronous networks, we employ an abstraction of timeouts known as unreliable failure detectors [19], and model them as a non-deterministic action [72]. The $susp_i$ action is triggered when process i is suspected of having failed - in other words, after some timeout.

Note that we may use $\sum P$ to denote guarded-choice over more than two processes, and $\prod P$ to denote the parallel composition of more than two processes.

An operational semantics defines the actual non-reversible computational steps that a process may execute. Effectively, the only relevant operation is communication, known as the *comm* rule:

$$(x?(y).P \mid x!(z)) \rightarrow P\{z/y\} \quad (1.1)$$

The notation $P\{z/y\}$ means that all occurrences of y in P are replaced with z . In other words, z was sent on x , received as y , and fed to P .

Given a π -calculus process, we can follow its execution by applying the comm rule. For instance,

$$(x?(y).y!(x) \mid x!(z)) \rightarrow z!(x) \quad (1.2)$$

Now, we can use a formal logic to express properties a process might satisfy. For instance, the modal Hennessy–Milner logic can express that a process will satisfy some other logic expression after some or all forms of an action have occurred [69]. By adding more complex operators to the logic, formal systems can be built up which easily describe important properties of distributed systems, such as safety and liveness [92], and localization [15].

Systems written in the π -calculus can then be formally verified to satisfy the relevant properties by model checking software [101].

While we use the π -calculus to specify the Tendermint algorithm, we leave use of an associated formal logic, and the corresponding verification of properties, to future work.

1.8 The Need For Tendermint

The success of Bitcoin and its derivatives, especially Ethereum, and their promise of secure, autonomous, distributed, fault-tolerant execution of arbitrary code has caused virtually every major financial institution on the planet to become interested in the blockchain phenomenon. In particular, there has emerged an understanding of two forms of the technology: On the one hand are the public blockchains, known affectionately as the Big Bad Public Blockchains or BBPBs, whose protocols are dominated by in-built economic incentives bootstrapped by a native currency. On the other are so called private blockchains, which might more accurately be called “consortia blockchains”, and which are effectively improvements on traditional consensus and BFT algorithms through the use of hash trees, digital signatures, peer-to-peer networking, and enhanced accountability.

As the infrastructure of our societies continues to decentralize, and as the nature of business becomes more inter-organizational, there is increasing need for a transparent, accountable, high performance BFT system, which can support applications from finance to domain registration to electronic voting, and which comes equipped with advanced mechanisms for governance and evolution into the future. Tendermint is that solution, optimized for consortia, or inter-organizational logic, but flexible enough to accommodate anyone from private enterprise to global currency, and high-performance enough to compete with the major, non-BFT, consensus solutions available today, such as etcd, consul, and zookeeper, while providing greater resilience, security guarantees, and flexibility to application developers.

A more comprehensive discussion of consensus science and related algorithms is reserved for Chapter 9

Chapter 2

Tendermint Consensus

This chapter presents the Tendermint consensus algorithm and an associated blockchain for atomic broadcast. The BFT consensus problem is described in detail, and a formal specification of Tendermint consensus is given in the π -calculus. The Tendermint blockchain is informally proven to satisfy atomic broadcast. We leave it to future work to capture the full blockchain protocol in a process calculus and to verify its properties.

2.1 Tendermint Overview

Tendermint is a secure state-machine replication algorithm in the blockchain paradigm. It provides a form of BFT-ABC that is furthermore accountable - if safety is violated, it is always possible to verify who acted maliciously.

Tendermint begins with a set of *validators*, identified by their public key, where each validator is responsible for maintaining a full copy of the replicated state, and for proposing new blocks (batches of transactions), and voting on them. Each block is assigned an incrementing index, or *height*, such that a valid blockchain has only one valid block at each height. At each height, validators take turns proposing new blocks in *rounds*, such that for any given round there is at most one valid proposer. It may take multiple rounds to commit a block at a given height due to the asynchrony of the network, and the network may halt altogether if more than one-third of the validators are offline or partitioned. Validators engage in two phases of voting on a proposed block before it is committed, and follow a simple locking mechanism which prevents any coalition of up to one third malicious

validators from compromising safety.

Note that the core round-based voting mechanism is the consensus algorithm, which is strung together into blocks to yield atomic broadcast. Each block contains some metadata, known as its *header*, which includes the hash of the block at the previous height, resulting in a hash chain. The header also includes the block height, local time the block was proposed, and the Merkle root hash of transactions included in the block.

2.2 Consensus

The consensus algorithm can be roughly divided into the following, somewhat orthogonal, components:

- **Proposals:** a new block must be proposed by the correct proposer at each round, and gossiped to the other validators. If a proposal is not received in sufficient time, the proposer should be skipped.
- **Votes:** two phases of voting must occur to ensure optimal Byzantine fault tolerance. They are called *pre-vote* and *pre-commit*. A set of pre-commits from more than two-thirds of the validators for the same block at the same round is a *commit*.
- **Locks:** Tendermint ensures that no two validators commit a different block at the same height, presuming less than one-third of the validators are malicious. This is achieved using a locking mechanism which determines how a validator may pre-vote or pre-commit depending on previous pre-votes and pre-commits at the same height. Note that this locking mechanism must be carefully designed so as to not compromise liveness.

In order to provide tolerance to a single Byzantine fault, a Tendermint network must contain at minimum four validators. Each validator must possess an asymmetric cryptographic key-pair for producing digital signatures. Validators start from a common initial state, which contains the ordered list, \mathcal{L} , of validators. Each validator is identified via their public key, and all proposals and votes must be signed by the respective private key. This ensures that proposals and votes can always be verified by any observer. It is helpful to assume that up to one-third of validators are malicious, co-operating in arbitrary ways to subvert system safety or liveness.

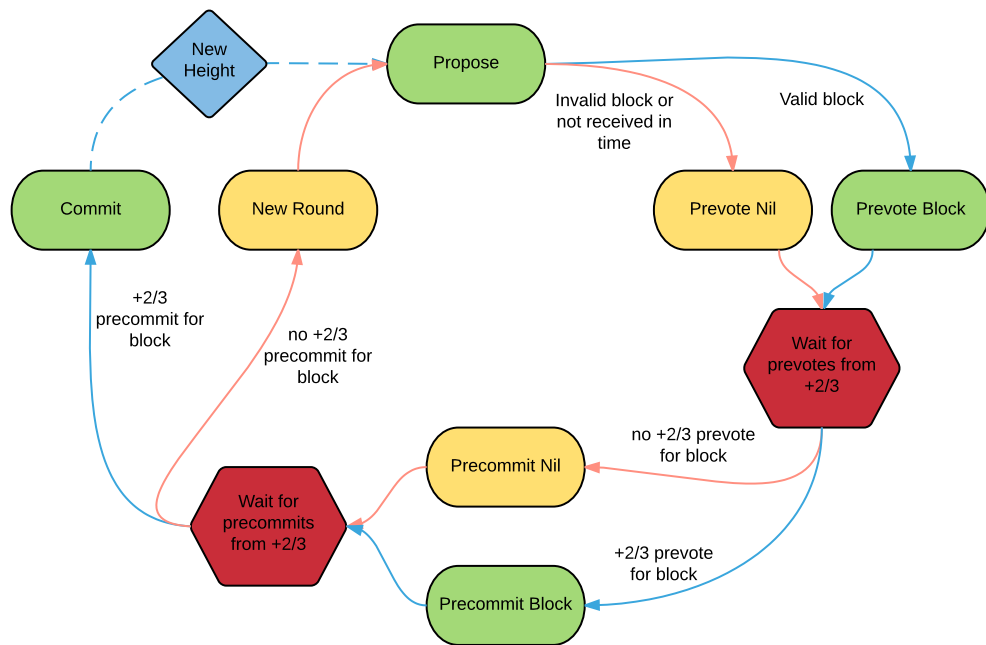


Figure 2.1: After the proposal step, validators only make progress after hearing from two-thirds or more ($+2/3$) of other validators. The dotted arrow extends the consensus into atomic broadcast by moving to the next height.

Consensus begins at round 0; the first proposer is the first validator in \mathcal{L} . The outcome of a round is either a commit, or a decision to move to the next round. With a new round comes the next proposer. Using multiple rounds gives validators multiple opportunities to come to consensus in the event of network asynchrony or validator failures.

In contrast to algorithms which require a form of leader election, Tendermint has a new leader (the proposer) for each round. Validators vote to skip to the next round in the same way they vote to accept the proposal, lending the protocol a uniformity of mechanism that is absent from algorithms with an explicit leader-election program.

The beginning of each round has a weak dependence on synchrony as it utilizes local clocks to determine when to skip a proposer. That is, if a validator does not receive a proposal within a locally measured *TimeoutPropose* of entering a new round, it can vote to skip the proposer. Inherent in this mechanism is the assumption that, at least eventually, the proposal will be delivered within *TimeoutPropose*, which may itself increment with each round. This assumption is discussed more fully in Chapter 9.

After the proposal, rounds proceed in a fully asynchronous manner - a validator makes progress only after hearing from more than two-thirds of the other validators. This relieves any sort of dependence on synchronized clocks or bounded network delays, but implies that the network will halt if one-third or more of the nodes become unresponsive. This circuit of weakly synchronous proposals, followed by asynchronous voting, is depicted in Figure 2.1.

To round-skip safely, a small number of *locking* rules are introduced which force validators to justify their votes. While we don't necessarily require them to broadcast their justifications in real time, we do expect them to keep the data, such that it can be brought forth as evidence in the event that safety is compromised by sufficient Byzantine failures. This accountability mechanism enables Tendermint to provide stronger guarantees in the face of such failure than eg. PBFT, which provides no guarantees if a third or more of the validators are Byzantine.

Validators communicate using a diverse set of messages for managing the blockchain, application state, peer network, and consensus. The core consensus algorithm, however, consists of just two messages:

- *ProposalMsg*: a proposal for a block at a given height and round, signed by the proposer.

- *VoteMsg*: a signed vote for a proposal.

In practice, we use additional messages to optimize the gossiping of block data and votes, as discussed in Chapter 3.

2.2.1 Proposals

Each round begins with a proposal. The proposer for the given round takes a batch of recently received transactions from its local cache (the Mempool, see Chapter 3), composes a block, and broadcasts a signed *ProposalMsg* containing the block. If the proposer is Byzantine, it might broadcast different proposals to different validators.

Proposers are ordered via a simple, deterministic round robin, so only a single proposer is valid for a given round, and every validator knows the correct proposer. If a proposal is received for a lower round, or from an incorrect proposer, it is rejected.

Cycling of proposers is necessary for Byzantine tolerance. For instance, in Raft, if an elected leader is Byzantine and maintains strong network connections to other nodes, it can completely compromise the system, destroying all safety and liveness guarantees. Tendermint preserves safety via the voting and locking mechanisms, and maintains liveness by cycling proposers, so if one won't process any transactions, others can pick up. Perhaps more interestingly, validators can vote through governance modules (see Chapter 5) to remove or replace Byzantine validators.

2.2.2 Votes

Once a complete proposal is received by a validator, it signs a pre-vote for that proposal and broadcasts it to the network. If a validator does not receive a correct proposal within *ProposalTimeout*, it pre-votes for *nil* instead.

In asynchronous environments with Byzantine validators, a single stage of voting, where each validator casts only one vote, is not sufficient to ensure safety. In essence, because validators can act fraudulently, and because there are no guarantees on message delivery time, a rogue validator can co-ordinate some validators to commit a value while others, having not seen the commit, go to a new round, within which they commit a different value.

A single round of voting allows validators to tell each other what they know about the proposal. But to tolerate Byzantine faults (which amounts,

essentially to lies, fraud, deceit, etc.), they must also tell each other what they know about what other validators have professed to know about the proposal. In other words, a second stage ensures that enough validators witnessed the result of the first stage.

A pre-vote for a block is thus a vote to prepare the network to commit the block. A pre-vote for nil is a vote to prepare the network to move to the next round. In an ideal round with an online proposer, more than two-thirds of validators will pre-vote for the proposal. A set of more than two-thirds of pre-votes for a single block at a given round is known as a *polka*¹. A set of more than two-thirds of pre-votes for nil is a *nil-polka*.

When a validator receives a polka (read: more than two-thirds pre-votes for a single block), it has received a signal that the network is prepared to commit the block, and serves as justification for the validator to sign and broadcast a pre-commit vote for that block. Sometimes, due to network asynchrony, a validator may not receive a polka, or there may not have been one. In that case, the validator is not justified in signing a pre-commit for that block, and must therefore sign and publish a pre-commit vote for nil. That is, it is considered malicious behaviour to sign a pre-commit without justification from a polka.

A pre-commit is a vote to actually commit a block. A pre-commit for nil is a vote to actually move to the next round. If a validator receives more than two-thirds pre-commits for a single block, it commits that block, computes the resulting state, and moves on to round 0 at the next height. If a validator receives more than two-thirds pre-commits for nil, it moves on to the next round.

2.2.3 Locks

Ensuring safety across rounds can be tricky, as circumstances must be avoided which would provide justification for two different blocks to be committed at two different rounds at the same height. In Tendermint, this problem is solved via a *locking* mechanism. In essence, once a pre-commit is cast, a validator is *locked* on the associated block, and must follow certain locking rules. There are two rules of locking:

¹The original term used was PoL, or PoLC, for Proof-of-Lock or Proof-of-Lock-Change. The term evolved to polka as it was realized the validators are doing the polka.

- **Prevote-the-Lock:** a validator must pre-vote for the block they are locked on, and propose it if they are the proposer. This prevents validators from pre-committing one block in one round, and then contributing to a polka for a different block in the next round, thereby compromising safety.
- **Unlock-on-Polka:** a validator may only release a lock after seeing a polka (more than two-thirds pre-voting for a single block) at a round greater than that at which it locked. This allows validators to unlock if they pre-committed something the rest of the network doesn't want to commit, thereby protecting liveness, but does it in a way that does not compromise safety, by only allowing unlocking if there has been a polka in a round after that in which the validator became locked.

For simplicity, a validator is considered to have locked on nil at round -1 at each height, so that Unlock-on-Polka implies that a validator cannot pre-commit at a new height until they see a polka.

These rules can be understood more intuitively by way of examples. Consider four validators, A , B , C , D , and suppose there is a proposal for *blockX* at round R . Suppose there is a polka for *blockX*, but A doesn't see it, and pre-commits nil, while the others pre-commit for *blockX*. Now suppose the only one to see all pre-commits is D , while the others, say, don't see D 's pre-commit (they only see their two pre-commits and A 's pre-commit nil). D will now commit the block, while the others go to round $R + 1$. Since any of the validators might be the new proposer, if they can propose and vote for any new block, say *blockY*, then they might commit it and compromise safety, since D already committed *blockX*. Note that there isn't even any Byzantine behaviour here, just asynchrony!

Locking solves the problem by forcing validators to stick with the block they pre-committed, since other validators might have committed based on those pre-commits (as D did in this example). In essence, once more than two-thirds pre-commit a block in a round, the network is locked on that block, which is to say it must be impossible to produce a valid polka for a different block at a higher round. This is direct motivation for Prevote-the-Lock.

Prevote-the-Lock is not sufficient, however. There must be a way to unlock, lest we sacrifice liveness. Consider a round where A and B pre-committed *blockX* while C and D pre-committed nil - a split vote. They all move to the next round, and *blockY* is proposed, which C and D prevote

for. Suppose A is Byzantine, and prevotes for $blockY$ as well (despite being locked on $blockX$), resulting in a polka. Suppose B does not see the polka and pre-commits nil, while A goes off-line and C and D pre-commit $blockY$. They move to the next round, but B is still locked on $blockX$, while C and D are now locked on $blockY$, and since A is offline, they can never get a polka. Hence, we've compromised liveness with less than a third (here, only one) Byzantine validators.

The obvious justification for unlocking is a polka. Once B sees the polka for $blockY$ (which C and D used to justify their pre-commits for $blockY$), it ought to be able to unlock, and hence pre-commit $blockY$. This is the motivation for Unlock-on-Polka, which allows validators to unlock (and pre-commit a new block), if they have seen a polka in a round greater than that in which they locked.

2.2.4 Formal Specification

Now that we have explained the protocol in detail, we provide a formal specification in the π -calculus.

Let $Consensus := \prod_{i=1}^N Y_i$ represent a consensus protocol over a set of N validators, each executing one of a mutually exclusive set of processes, Y_i . Internal state $s = \{r, p, v\}$ consists of a strictly increasing round, r , a proposal p , containing the proposed block for this round; and a set of votes, v , containing all votes at all rounds; We denote by v_r^1 and v_r^2 the set of prevotes and pre-commits, respectively, at round r . We define $proposer(r) = r \bmod N$ to be the index of the proposer at round r . We represent a peer at a particular point in the protocol as $Y_i^{r,p,v}$. Processes Y_i range over PR_i , PV_i , PC_i , respectively abbreviating *propose*, *prevote*, *pre-commit*. We introduce additional sub-functions for PV and PC to capture the recursion, denoted $PV1$, $PV2$, etc.

Peers are connected using broadcast channels for each message type, namely $propose_i$, $prevote_i$, and $pre - commit_i$, as well as a channel for deciding on, or committing, a value, d_i . Via an abuse of notation, a single send on a broadcast channel xxx_i can be received by each process along xxx_i .

We use only two message types: proposals and votes. Each contains a round number, block (hash), and signature, denoted $msg.round$, $msg.block$, $msg.sig$. Note we can absorb the signature into the broadcast channel itself, but we need it for use as evidence in the event of Byzantine behaviour.

The specification is given in two parts, in Figures 2.2 and ??.

$$Consensus := \prod_{i=1}^N PR_i^{0,\emptyset,\emptyset},$$

$$\begin{aligned}
PR_i^{r,p,v} := & \text{if } i = \text{proposer}(r) \text{ then} \\
& \text{propose}_i!(prop) \mid PV_i^{r,prop,v}, \text{ where } prop = \text{chooseProposal}(p) \\
& \text{else if } p \neq \emptyset \text{ then} \\
& \quad PV_i^{r,p,v} \\
& \text{else} \\
& \quad \text{propose}_{\text{proposer}(r)}?(prop).PV_i^{r,prop,v} + \text{susp}_{\text{proposer}(r)}.PV_i^{r,\emptyset,v}
\end{aligned}$$

$$PV_i^{r,p,v} := \text{prevote}_i!(p) \mid (\nu c)(\prod_{j=1}^n \text{prevote}_j?(w).c!(\text{prevote}_j, w) \mid PV_1^{r,p,v}(c))$$

$$\begin{aligned}
PV_1^{r,p,v}(c) := & \text{if } \max_b(|\{w \in v_r^1 : w.\text{block} = b\}|) > \frac{2}{3}N \text{ then} \\
& \quad PC_i^{r,b,v} \\
& \text{else if } |v_r^1| > \frac{2}{3}N \text{ then} \\
& \quad \quad PC_i^{r,\emptyset,v} \\
& \text{else} \\
& \quad c?(pv, vote). \text{ if } vote.\text{round} < r \text{ then} \\
& \quad \quad pv?(w).c!(pv, w) \\
& \quad \text{else if } vote.\text{round} = r \text{ then} \\
& \quad \quad PV_1^{r,p,vote::v}(c) \\
& \text{else} \\
& \quad PR_i^{vote.\text{round},p,vote::v}
\end{aligned}$$

Figure 2.2: Formal specification of Tendermint consensus in the π -calculus, part I. *chooseProposal*(*p*) must return *p* if it is not \emptyset , and otherwise should gather transactions from the mempool as described in Chapter ???. After receiving a proposal or timing out, validators move onto prevote, where they broadcast their prevote and wait to receive prevotes from the others. If a vote is received for a later round, we skip ahead to that round.

$$PC_i^{r,p,v} := precommit_i!(p) \mid (\nu c)(\prod_{j=1}^n precommit_j?(w).c!(w) \mid PC1_i^{r,p,v}(c))$$

$$PC1_i^{r,p,v}(c) := \text{if } max_b(|\{w \in v_r^2 : w.block = b\}|) > \frac{2}{3}N \text{ then}$$

$$d_i!(b)$$

$$\text{else if } |v_r^2| > \frac{2}{3}N \text{ then}$$

$$PR_i^{r+1,\emptyset,v}$$

else

$$c?(pc, vote). \text{ if } vote.round < r \text{ then}$$

$$pc?(w).c!(pc, w)$$

$$\text{else if } vote.round = r \text{ then}$$

$$PC1_i^{r,p,vote::v}(c)$$

else

$$PR_i^{vote.round,p,vote::v}$$

Figure 2.3: Formal specification of Tendermint consensus in the π -calculus, part II. Validators broadcast their pre-commit and wait to receive pre-commits from the others. If a vote is received for a later round, we skip ahead to that round. When more than two-thirds pre-commit for block b , we fire b on channel d_i , signalling the commit, and terminating the protocol.

2.3 Blockchain

Tendermint operates on batches, or blocks, of transactions at a time. Continuity is maintained from one block to the next by explicitly linking each block to the one before it via its cryptographic hash, forming a blockchain. The blockchain contains both the ordered transaction log and evidence that the block was committed by the validators.

2.3.1 Why Blocks?

Consensus algorithms typically commit transactions one at a time by design, and implement batching after the fact. As mentioned in Chapter 1, tackling the problem from the perspective of batched atomic broadcast results in two primary optimizations, which give us more throughput and fault-tolerance:

- Bandwidth optimization: since every commit requires two rounds of communication across all validators, batching transactions in blocks amortizes the cost of a commit over all the transactions in the block.
- Integrity optimization: the hash chain of blocks forms an immutable data structure, much like a Git repository, enabling authenticity checks for sub-states at any point in the history.

Blocks induce another effect as well, which is more subtle but potentially important. They increase the minimum latency of a transaction to that of the whole block, which for Tendermint is on the order of hundreds of milliseconds to seconds. Traditional serializable database systems provide commit latencies on the order of milliseconds to tens of milliseconds. They are able to do this because they are not Byzantine Fault Tolerant, requiring only one round of communication (instead of two) and responses from over half of the replicas (instead of two-thirds). However, unlike the fast commit times interrupted by leader elections in other algorithms, Tendermint provides a more regular pulse that is more responsive to the overall health of the network, in terms of node failures and asynchrony.

What role such pulses might play in the coherence of communicating autonomous systems on the internet is yet to be determined, though purposefully induced latency has shown promise in the financial markets [86].

2.3.2 Block Structure

The purpose of blocks is to contain a batch of transactions, and to link to the previous block. The link comes in two forms: the previous block hash, and the set of pre-commits which caused the previous block to be committed, also known as the *LastCommit*. Thus a block is composed of three parts: the block header, the list of transactions, and the *LastCommit*.

2.4 Safety

Here we sketch a brief proof that Tendermint satisfies atomic broadcast, which is defined as satisfying:

- validity - if a correct process broadcasts m , it eventually delivers m
- agreement - if a correct process delivers m , all correct processes eventually deliver m
- integrity - m is only delivered once, and only if broadcast by its sender
- total order - if correct processes p and q deliver m and m' , then p delivers m before m' iff q delivers m before m'

Note that if we take m to be a block, Tendermint does not satisfy validity, since there is no guarantee that a proposed block is eventually committed, as validators may move to a new round and commit a different block. If we take m to be a batch of transactions in a block, then we can satisfy validity by having validators re-propose the same batch until it is committed. However, to satisfy the first half of integrity we must introduce an additional rule that forbids a correct validator from proposing a block or pre-committing for a block containing a batch of transactions that has already been committed. Fortunately, batches can be indexed by their merkle root, and a lookup performed before proposals and pre-commits.

Alternatively, if we take message m to be a transaction, then we can satisfy validity by asserting a *persistence* property on the mempool, namely, that a transaction persists in the mempool until it is committed. However, to satisfy the first half of integrity we must rely on the application state to enforce some ruleset over transactions such that a given transaction is only valid once. This can be done, for instance, using sequence numbers on

accounts, as is done in ethereum, or by keeping a list of unused resources, each of which can only be used once, as is done in Bitcoin. Since there are multiple approaches, Tendermint does not in itself ensure that a message is only delivered once, but allows the application developer to specify. Note that the second half of integrity is trivially satisfied, since only transactions in blocks proposed by a correct proposer can be committed.

To show that Tendermint satisfies the remaining properties, we introduce a new property, *state machine safety*, and show that a protocol satisfying state machine safety satisfies agreement and total order. State machine safety states that if a correct validator commits a block at some height H , no other correct validator will ever commit a different block at H . Given that all messages are eventually received, this immediately implies agreement, since if a correct validator commits a block B at height H containing a transaction m , all other correct validators will be unable to commit any other block, and hence must eventually commit B , thereby delivering m .

Now, it remains to show that state machine safety satisfies total order, and that Tendermint satisfies state machine safety. To see the former, consider two messages m and m' delivered by validators p and q . State machine safety ensures that p delivers m at height H_m if and only if q delivers m at height H_m , and that p delivers m' at height $H_{m'}$ if and only if q delivers m' at height $H_{m'}$. Without loss of generality, and since height is strictly increasing, let $H_m < H_{m'}$. Then we have that p delivers m before m' if and only if q delivers m before m' , which is exactly the statement of total order.

Finally, to show Tendermint satisfies state machine safety when less than a third of validators are Byzantine, we proceed by way of contradiction. Suppose Tendermint does not satisfy state machine safety, allowing more than one block to be committed at the same height. Then we show that at least one-third of validators must be Byzantine for that to happen, contradicting our assumption.

Consider a correct validator having committed block B at height H and round R . To commit a block means the validator witnessed pre-commits for block B in round R from more than two-thirds of validators. Suppose another block C is committed at height H . We have two options: either it was committed in round R , or round $S > R$.

If it was committed in round R , then more than two-thirds of validators must have pre-committed for it in round R , which means that at least a third of validators pre-committed for both blocks B and C in round R , which is clearly Byzantine. Suppose block C was instead committed in round $S > R$.

Since more than two-thirds pre-committed for B , they are locked on B in round S , and thus must pre-vote for B . To pre-commit for block C , they must witness a polka for C , which requires more than two-thirds to pre-vote for C . However, since more than two-thirds are locked on and required to pre-vote for B , a polka for C would require at least one third of validators to violate Prevote-the-Lock, which is clearly Byzantine. Thus, to violate state machine safety, at least one third of validators must be Byzantine. Therefore, Tendermint satisfies state machine safety when less than a third of validators are Byzantine.

Given the above, then, Tendermint satisfies atomic broadcast.

In future work, we aim to provide a more formal proof of Tendermint's safety property.

2.5 Accountability

An accountable BFT algorithm is one that can identify all Byzantine validators when there is a violation of safety. Traditional BFT algorithms do not have this property, and provide no guarantees in the event safety is compromised. Of course, accountability can only apply when between one-third and two-thirds of validators are Byzantine. If more than two-thirds are Byzantine, they can completely dominate the protocol, and we have no guarantee that a correct validator will receive any evidence of their misdeeds.

Futhermore, accountability can be at best eventual in asynchronous networks - following a violation of safety, the delayed delivery of critical messages may make it impossible to determine which validators were Byzantine until some time after the safety violation is detected. In fact, if correct processes can receive evidence of Byzantine behaviour, but fail irreversibly before they are able to gossip it, there may be cases where accountability is permanently compromised, though in practice such situations should be surmountable with advanced backup solutions.

By enumerating the possible ways in which a violation of safety can occur, and showing that in each case, the Byzantine validators are identifiable, a protocol can be shown to be accountable. Tendermint's simplicity affords it a much simpler analysis than protocols which have to manage leadership elections.

There are only two ways for a violation of safety to occur in Tendermint, and both are accountable. In the first, a Byzantine proposer makes two

conflicting proposals within a round, and Byzantine validators vote for both of them. In the second, Byzantine validators violate locking rules after some validators have already committed, causing other validators to commit a different block in a later round. Note that it is not possible to cause a violation of safety with two-thirds or fewer Byzantine validators using only violations of Unlock-on-Polka - more than a third must violate Prevote-the-Lock for their to be a polka justifying a commit for the remaining honest nodes.

In the case of conflicting proposals and conflicting votes, it is trivial to detect the conflict by receiving both messages, and to identify culprits via their signatures.

In the case of violating locking rules, following a violation of safety, correct validators must broadcast all votes they have seen at that height, so that the evidence can be stitched together. The correct validators, which number something under two-thirds, were collectively privy to all votes which caused the two blocks to be committed. Within those votes, if there are not a third or more validators signing conflicting votes, then there are a third or more violating Prevote-the-Lock.

If a pre-vote or a pre-commit influenced a commit, it must have been seen by a correct validator. Thus, by collecting all votes, violations of Prevote-the-Lock can be detected by matching each pre-vote to the most recent pre-commit by the same validator, unless there isn't one.

Similarly, violations of Unlock-on-Polka can be detected by matching each pre-commit to the polka that justifies it. Note that this means a Byzantine validator can pre-commit before seeing a polka, and escape accountability if the appropriate polka eventually occurs. However, such cases cannot actually contribute to violations of safety if the polka is happening anyways.

The current design provides accountability following a post-crisis broadcast protocol, but it could be improved to allow accountability in real time. That is, a commit could be changed to include not just the pre-commits, but all votes justifying the pre-commits, going all the way back to the beginning of the height. That way, if safety is violated, the unjustified votes can be detected immediately.

2.6 Faults and Availability

As a BFT consensus algorithm, Tendermint can tolerate Byzantine failure in up to (but not including) one-third of validators. This means nodes can crash, send different and contradictory messages to different peers, refuse to relay messages, or otherwise behave arbitrarily, without compromising safety or liveness (with the usual FLP caveat for liveness).

There are two places in the protocol where we can make optimizations for asynchrony by utilizing timeouts based on local clocks: after receiving two-thirds or more pre-votes, but not for a single block or nil, and after receiving two-thirds or more pre-commits, but not for a single block or nil. In each case, we can sleep for some amount of time to give slower or delayed votes a chance to be received, thereby reducing the likelihood of going to a new round without committing a block. Clocks do not need to be synced across validators, as they are reset each time a validator observes votes from two-thirds or more others.

If a third or more of validators crash, the network halts, as no validator is able to make progress without hearing from more than two-thirds of the validator set. The network remains available for reads, but no new commits can be made. As soon as validators come back on-line, they can carry on from where they left in a round. The consensus state-machine should employ a write-ahead log, such that a recovered validator can quickly return to the step it was in when it crashed, ensuring it doesn't accidentally violate a rule.

If a third or more of validators are Byzantine, they can compromise safety a number of ways, for instance, by proposing two blocks for the same round, and voting both of them through to commit, or by pre-committing on two different blocks at the same height but in different rounds by violating the rules on locking. In each case, there is clear, identifiable evidence that certain validators misbehaved. In the first instance, they signed two proposals at the same round, a clear violation of the rules. In the second, they may have pre-voted for a different block in round R than they locked on in $R-1$, a violation of the Prevote-the-Lock rule.

When using economic and governance components to incentivize and manage the consensus (Chapter 5) these additional accountability guarantees become critical.

2.7 Conclusion

Tendermint is a weakly synchronous, Byzantine fault tolerant, state machine replication protocol, with optimal Byzantine fault tolerance and additional accountability guarantees in the event the BFT assumptions are violated. The protocol uses a round-robin approach for proposers, and uses the same mechanism to skip a proposer as to commit a proposed block. Safety is maintained across rounds via a simple locking mechanism.

The presentation of the protocol in this chapter left out many important details, such as the efficient gossiping of blocks, buffering transactions, changes to the validator set, and the interface with application logic. These important topics are taken up in subsequent chapters.

Chapter 3

Tendermint Subprotocols

The presentation of Tendermint consensus in the previous chapter left out a number of details regarding the gossip protocols used to disseminate blocks, votes, transactions, and other peer information. This was done in order to focus in on the consensus protocol itself, without distraction from the hydra of practical software engineering. This chapter describes one particular approach to filling in these details, by implementing components as relatively independent reactors that are multiplexed over each peer connection.

3.1 P2P-Networking

On startup, each Tendermint node receives an initial list of peers to dial. For each peer, a node maintains a persistent TCP connection over which multiple subprotocols are multiplexed in a rate-limited fashion. Messages are serialized into a compact binary representation to be sent on the wire, and connections are encrypted via an authenticated encryption protocol [28].

Each remaining section of this chapter describes a separate reactor that is multiplexed over each peer connection. An additional peer exchange reactor can be run which allows nodes to request other peer addresses from each other and keep track of peers they have connected to before, in order to stay connected to some minimum number of other peers.

3.2 Consensus Gossip

The consensus reactor wraps the consensus state machine, and ensures each node broadcasts to all peers its current state every time it changes. In this way, each node keeps track of the consensus state of all its peers, allowing it to optimize the gossiping of messages to only send peers information they need at the very moment, and which they don't already have. For each peer, a node maintains two routines which continuously check for new information to send the peer, namely, proposals and votes. Information should be gossiped in a “rarest first” manner in order to maximize gossip efficiency and minimize the chance that some information becomes unavailable [62]

3.2.1 Block Data

In Chapter 2, it was assumed that proposal messages include the block. However, since blocks emerge from a single source and can be quite large, this puts undue pressure on the block proposer to upload the data to all other nodes; blocks can be disseminated much more quickly if they are split into parts and gossiped.

A common approach to securely gossiping data, as popularized by various p2p protocols [21, 79], is to use a Merkle tree [65], allowing each piece of the data to be accompanied by a short proof (logarithmic in the size of the data) that the piece is a part of the whole. To use this approach, blocks are serialized and split into chunks of an appropriate size for the expected block size and number of validators, and chunks are hashed into a Merkle tree. The signed proposal, instead of including the entire block, includes just the Merkle root hash, allowing the network to co-operate in gossiping the chunks. A node informs its peers every time it receives a chunk, in order to minimize the bandwidth wasted by transmitting the same chunk to a node more than once.

Once all the chunks are received, the block is deserialized and validated to ensure it refers correctly to the previous block, and that its various checksums, implemented as Merkle trees, are correct. While it was previously assumed that a validator does not pre-vote until the proposal (including the block) is received, some performance benefit may be obtained by allowing validators to pre-vote after receiving a proposal, but before receiving the full block. This would imply that it is okay to pre-vote for what turns out to be an invalid block. However, pre-committing for an invalid block must always

be considered Byzantine.

Peers that are catching up (i.e. are on an earlier height) are sent chunks for the height they are on, and progress one block at a time.

3.2.2 Votes

At each step in the consensus state machine, after the proposal, a node is waiting for votes (or a local timeout) to progress. If a peer has just entered a new height, it is sent pre-commits from the previous block, so it may include them in the next blocks *LastCommit* if it's a proposer. If a peer has pre-voted but has yet to pre-commit, or has pre-committed, but has yet to go to the next round, it is sent pre-votes or pre-commits, respectively. If a peer is catching up, it is sent the pre-commits for the committed block at its current height.

3.3 Mempool

Chapter 2 made little mention of transactions, despite the purpose of a consensus algorithm being to order and execute transactions, as Tendermint operates on blocks on a time, and has no concern for individual transactions, so long as their checksum in the block is correct.

Transactions are managed independently in an in-memory cache, which, following Bitcoin, has come to be known as the *mempool*. Transactions are validated by the application logic when they are received and, if valid, added to the mempool and gossiped using an ordered multicast algorithm. A node maintains a routine for each peer which ensures that transactions in the mempool are sent to the peer in the same order in which they were processed by the node.

Proposers reap transactions from the ordered list in the mempool for new block proposals. Once a block is committed, all transactions included in the block are removed from the mempool, and the remaining transactions are re-validated by the application logic, as their validity may have changed on account of other transactions being committed, which the node may not have had in its mempool.

3.4 Syncing the Blockchain

The consensus reactor provides a relatively slow means of syncing with the latest state of the blockchain, as it was designed for real-time consensus, meaning peers wait to receive all information to commit a single block before worrying about the next block. To accommodate peers that may be more than just a few blocks behind, an additional reactor, the blockchain reactor, allows peers to download many blocks in parallel, enabling a peer to sync hundreds of times faster than via the consensus reactor.

When a node connects to a new peer, the peer sends its current height. The node will request blocks, in order, beginning with its current height, from all peers that self-reported higher heights, and download the blocks concurrently, adding them to the block pool. Another routine continuously attempts to remove blocks from the pool and add them to the blockchain by validating and executing them, two blocks at a time, against the latest state of the blockchain. Blocks must be validated two blocks at a time because the commit for one block is included as the LastCommit data in the next one.

The node continuously queries its peers for their current height, and continues to concurrently request blocks until it has caught up to the highest height among its peers, at which point it stops making requests for peer heights and starts the consensus reactor.

Chapter 4

Building Applications

Tendermint is designed to be a general purpose algorithm for replicating a deterministic state machine. It uses the Tendermint Socket Protocol (TMSP) to standardize communication between the consensus engine and the state machine, enabling application developers to build their state machines in any programming language, and have it automatically replicated via Tendermint's BFT algorithm.

4.1 Background

Applications on the Internet can in general be characterized as containing two fundamental components:

- Engine: handles core security, networking, replication. This is typically a webserver, like Apache or Nginx, when powering a web app, or a consensus algorithm when powering a distributed application.
- State-machine: the actual application logic that processes transactions received from the engine and updates internal state.

This separation of concerns enables application developers to write state-machines in any programming language representing arbitrary applications, on top of an engine which may be specialized for its performance, security, usability, support, and other considerations.

Unlike web-servers and their applications, which often take the form of processes communicating over a socket via the Common Gateway Interface

(CGI) protocol, consensus algorithms have traditionally had much less usable or less general purpose interfaces to build applications on top of. Some, like zookeeper, etcd, consul, and other distributed key-value stores, provide HTTP interfaces to a particular instance of a simple key-value application, with some more interesting features like atomic compare-and-swap operations and push notifications. But they do not give the application developer control of the state-machine code itself.

Demand for such a high-level of control over the state-machine running above a consensus engine has been driven primarily by the success of Bitcoin and the consequent interest in blockchain technology. By building more advanced applications directly into the consensus, users, developers, regulators, etc. can achieve greater security guarantees on arbitrary state-machines, far beyond key-value stores, like currencies, exchanges, supply-chain management, governance, and so on. What has captured the attention of so many is the potential of a system which permits collective enforcement of the execution of code. It is practically a re-invention of many dimensions of the legal system, using distributed consensus algorithms and deterministically executable contracts, rather than policemen, lawyers, judges, juries, and the like. The ramifications for the development of human society are explosive, much as the introduction of the democratic rule of law was in the first place.

Tendermint aims to provide the fundamental interface and consensus engine upon which such applications might be built.

4.2 Tendermint Socket Protocol

The Tendermint Socket Protocol (TMSP) defines the core interface by which the consensus engine communicates with the application state machine. The interface definition consists of a number of message types, specified using Google’s Protocol Buffers [100], that are length-prefixed and transmitted over a socket. A list of message types, their arguments, return values, and purpose is given in Figure 4.1, and an overview of the architecture and message flow is shown in Figure 4.2.

TMSP is implemented as an ordered, asynchronous server, where message types come in pairs of request and response, and where a special message type, Flush, pushes any buffered messages over the connection and awaits all responses.

At the core of the TMSP are two messages: *AppendTx* and *Commit*.

```

type Application interface {
// Return application info
Info() (info string)

// Set application option
SetOption(key string, value string) (log string)

// Append a tx
AppendTx(tx []byte) Result

// Validate a tx for the mempool
CheckTx(tx []byte) Result

// Return the application Merkle root hash
Commit() Result

// Query for state
Query(query []byte) Result

// Signals the beginning of a block
BeginBlock(height uint64)

// Signals the end of a block
// validators: changed validators from app to TendermintCore
EndBlock(height uint64) (validators []*Validator)
}

type CodeType int32

type Result struct {
Code CodeType
Data []byte
Log string // Can be non-deterministic
}

type Validator struct {
PubKey []byte
Power uint64
}

```

Figure 4.1: The TMSP application interface as defined in Go. TMSP messages are defined using Google’s Protocol Buffers, and their serialized form is length prefixed before being sent over the TMSP socket. Return values include a *Code*, similar to an HTTP Status Code, representing any errors, and 0 is used to indicate no error. Messages are buffered client side until a *Flush* message is sent, at which point all messages are transmitted. While the server design is asynchronous, message responses must be correctly ordered and match their request.

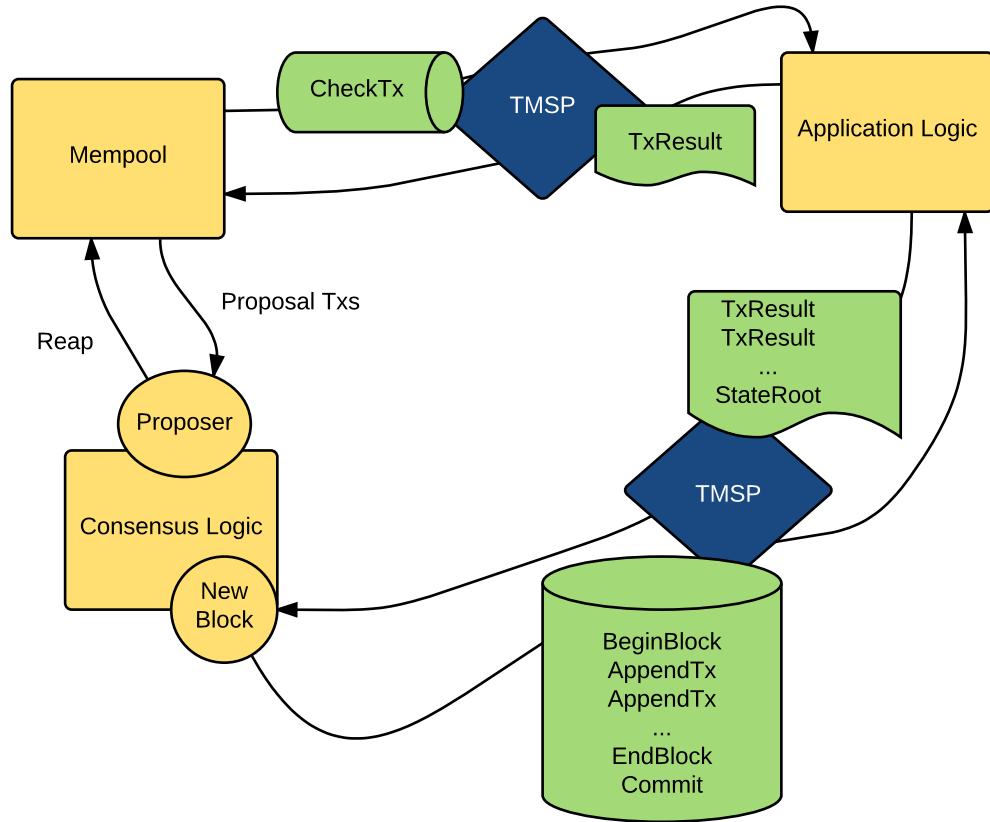


Figure 4.2: The consensus logic communicates with the application logic via TMSP, a socket protocol. Two sockets are maintained, one for the mempool to check the validity of new transactions, and one for the consensus to execute newly committed blocks.

Once a block is decided by the consensus, the engine calls *AppendTx* on each transaction in the block, passing it to the application state-machine to be processed. If the transaction is valid, it will result in a state-transition in the application.

Once all *AppendTx* calls have returned, the consensus engine calls *Commit*, causing the application to commit to the latest state, and persist it to disk.

4.3 Separating Agreement and Execution

Using the TMSP affords us an explicit separation between consensus, or agreement on the order of transactions, and their actual execution in the state-machine. In particular, we achieve consensus on the order first, and then execute the ordered transactions. This separation actually improves the system’s fault tolerance [107]: while $3f + 1$ replicas are still needed for agreement to tolerate f Byzantine failures, only $2f + 1$ replicas are needed for execution. That is, while we still need a two-thirds majority for ordering, we only need a one-half majority for execution.

On the other hand, the fact that transactions are executed after they are ordered results in possibly invalid transactions, which can waste system resources. This is solved using an additional TMSP message, *CheckTx*, which is called by the mempool, allowing it to check whether the transaction would be valid against the latest state. Note, however, that the fact that commits come in blocks at a time introduces complexity in the handling of *CheckTx* messages. In particular, applications are expected to maintain a second state-machine that executes only those rules of the main state-machine pertaining to a transaction’s validity. This second state-machine is updated by *CheckTx* messages and is reset to the latest committed state after every commit. In essence, the second state machine describes the transaction pool’s filter rules.

To some extent, *CheckTx* can be used as an *optimistic execution* returning a result to the transaction sender with the caveat that the result may be wrong if a block is committed with a conflicting transaction before the transaction of interest is committed. This sort of optimistic execution is the focus of an approach to scalable BFT systems that can work quite well for particular applications where conflicts between transactions are rare. At the same time, it adds additional complexity to the client, by virtue of needing to handle possibly invalid results. The approach is discussed further in Chapter

4.4 Microservice Architecture

Adopting separation of concerns as a strategy in application design is generally considered wise practice [50]. In particular, many large scale application deployments today adopt a microservice architecture, wherein each functional component is implemented as a standalone network service, and typically encapsulated in a Linux container (e.g. using Docker) for efficient deployment, scalability, and upgradeability.

Applications running above Tendermint consensus will often be decomposable into microservices. For instance, many applications will utilize a key-value store for storing state. Running the key-value store as an independent service is quite common, in order to take advantage of the data store’s specialized features, such as high-performance data types or Merkle trees.

Another important microservice for applications is a governance module, which manages a certain subset of TMSP messages, enabling the application to control validator set changes. Such a module can become a powerful paradigm for governance in BFT systems.

Some applications may utilize a native currency or account structure for users. It may thus be useful to provide a module which supports basic elements of, for instance, handling digital signatures and managing account dynamics.

The list of possible microservices to compose a complex TMSP application goes on. In fact, one might even build an application which can launch sub-applications using data sent in transactions. For instance, including the hash of a docker image in a transaction, such that the image could be pulled from some file-storage backend and run as a sub-application where future transactions in the consensus could cause it to execute. This is the approach of ethereum, which allows developers to deploy bits of code to the network that can be triggered to run within the Ethereum Virtual Machine by future transactions [105], and of IBM’s recent OpenBlockChain (OBC) project, which allows developers to send full docker contexts in transactions, defining containers that run arbitrary code in response to transactions addressed to them [76].

4.5 Determinism

The most critical caveat about building applications using TMSP is that they must be deterministic. That is, for the replicated state-machine to not compromise safety, every node must obtain the same result when executing the same transaction against the same state.

This is not a unique requirement for Tendermint. Bitcoin, Raft, Ethereum, any other distributed consensus algorithm, and applications like lock-step multi-player gaming must all be strictly deterministic, lest a consensus failure arise.

There are many sources of non-determinism in programming languages, most obviously via random numbers and time, but also, for instance, via the use of floating point precision, and by iteration over hash tables (some languages, such as Go, enforce randomized iteration over hash tables to force programmers to be explicit about when they need ordered data structures). The strict restriction on determinism, and its notable lacking from every major programming language, prompted ethereum to develop its own, Turing-complete, fully deterministic virtual machine, which forms the platform for application developers to build applications above the ethereum blockchain. While deterministic, it has many quirks, such as 32-byte stack words, storage keys, and storage values, and no support for byte-shifting operations - everything is big number arithmetic.

Deterministic programming is well studied in the world of real-time, lock-step, multi-party gaming. Such games constitute another example of replicated state machines, and are quite similar in many ways to consensus algorithms. Application developers building with TMSP are encouraged to study their methods, and to take care when implementing an application. On the one hand, the use of functional programming languages and proof methods can enable the construction of correct programs. On the other, compilers are being built to translate possibly non-deterministic programs to canonically deterministic ones [1].

4.6 Termination

If determinism is critical for preserving safety, termination of transaction execution is critical for preserving liveness. It is, however, not in general possible to determine whether a given program halts for even a single input,

let alone all of them, a problem known as the Halting Problem [98, 25].

Ethereum’s virtual machine solves the problem by *metering*, that is, charging for each operation in the execution. This way, a transaction is guaranteed to terminate when the sender runs out of funds. Such metering may be possible in a more general case, via compilers that compile programs to metered versions of themselves.

It is difficult to solve this problem without significant overhead. In essence, a validator cannot tell if an execution is in an infinite loop or is just slow, but nearly complete. It may be possible to use the Tendermint consensus protocol to decide on transaction timeouts, such that more than two-thirds of validators must agree that a transaction timed out and is thus considered invalid (ie. having no effect on the state). However, we do not pursue the idea further here, leaving it to future work. In the meantime, it is expected that applications will undergo thorough testing before being deployed in any consensus system, and that monitoring and governance mechanisms will be used to resurrect the system in the event of consensus failure.

4.7 Examples

In this section, examples of increasingly more complex TMSP applications are introduced and discussed, with particular focus on *CheckTx* and managing the mempool.

4.7.1 Merkleeyes

A simple example of a TMSP application is a Merkle tree based key-value store. Tendermint provides Merkleeyes, a TMSP application which wraps a self-balancing, Merkle binary search tree. The first byte of a transaction determines if the transaction is a get, set, or remove operation. For get and remove operations, the remaining bytes are the key. For the set operation, the remaining bytes are a serialized list containing the key and value. Merkleeyes may utilize a simple implementation of *CheckTx* that only decodes the transaction, to ensure it is properly formatted. One could also make a more advanced *CheckTx*, where get and remove operations on unknown keys are invalid. Once Commit is called, the latest updates are added into the Merkle tree, all hashes are computed, and the latest state of the tree

is committed to disk.

Note that Merkleeyes was designed to be a module used by other TMSP applications for a Merkle tree based key-value store, rather than a stand alone TMSP application, though the simplicity of the TMSP interface makes it amenable to both.

4.7.2 Basecoin

A more complete example is a simple currency, using an account structure pioneered by Ethereum, where each user has a public key and an account with the balance for that public key. The account also contains a sequence number, which is equal to the number of transactions sent by the account. Transactions can send funds from the account if they include the correct sequence number and are signed by the correct private key. Without the sequence number, the system would be susceptible to replay attacks [93], where a signed transaction debiting an account could be replayed, causing the debit to occur multiple times. Furthermore, to prevent replay attacks in a multi-chain environment, transaction signatures should include a network or blockchain identifier.

An application supporting a currency has naturally more logic than a simple key-value store. In particular, certain transactions are distinctly invalid, such as those with an invalid signature, incorrect sequence number, or sending an amount greater than the sender's account balance. These conditions can be checked in *CheckTx*.

Furthermore, a supplementary application state must be maintained for *CheckTx* in order to update sequence numbers and account balances when there are multiple transactions involving the same accounts in the mempool at once. When commit is called, the supplementary application state is reset to the latest committed state. Any transactions still in the mempool can be replayed via *CheckTx* against the latest state.

4.7.3 Ethereum

Ethereum uses the mechanisms already described to filter transactions out of the mempool, but it also runs some transactions in a virtual machine, which updates state and returns results. The virtual machine execution is not done in *CheckTx*, as it is much more expensive and depends heavily on the ultimate order of transactions as they are included in blocks.

4.8 Conclusion

TMSP provides a simple yet flexible means to build arbitrary applications, in any programming language, that inherit BFT state-machine replication from the Tendermint consensus algorithm. It plays much the same role for a consensus engine and an application that, for instance, CGI plays for Apache and Wordpress. However, application developers must take special care to ensure their applications are deterministic, and that transaction executions terminate.

Chapter 5

Governance

So far, this thesis has reviewed the basic elements of the Tendermint consensus protocol and application environment. Critical elements of operating the system in the real world, such as managing validator set changes and recovering from a crisis, have not yet been discussed.

This chapter proposes an approach to these problems that formalizes the role of governance in a consensus system. As validator sets come to encompass more decentralized sets of agents, competent governance systems for maintaining the network will be increasingly paramount to the network's success.

5.1 Governmint

The basic functionality of governance is to filter proposals for action, typically through a form of voting. The most basic implementation of governance as software is a module that enables users to make proposals, vote on them, and tally the votes. Proposals may be programmatic, in which case they may execute automatically following a successful vote, or they may be non-programmatic, in which case their execution is a manual exercise.

To enable certain actions in Tendermint, such as changing the validator set or upgrading the software, a governance module has been implemented, called Governmint. Governmint is a minimum viable governance application with support for multiple groups of entities, each of which can vote internally on proposals, some of which may result in programmatic execution of actions, like changing the validator set, or upgrading Governmint itself (for instance

to add new proposal types or other voting mechanisms).

The system utilizes digital signatures to authenticate voters, and may use a variety of possible voting schemes. Of particular interest are quadratic voting schemes, where the cost to vote is quadratic in the weight of the vote, which have been shown to have a superior ability to satisfy voter preferences [84].

5.2 Validator Set Changes

Validator set changes are a critical component of real world consensus algorithms that many previous approaches have failed to specify or have been left as a black art. Raft took pains to expound a sound protocol for validator set changes, which required the change pass through consensus, using a new message type. Tendermint takes a similar approach, though it is standardized through the TMSP interface using the *EndBlock* message, which is run after all the *AppendTx* messages, but before *Commit*. If a transaction, or set of transactions, is included in a block with the intended effect of updating the validator set, the application can return a list of validators to update by specifying their public key and new voting power in response to the *EndBlock* message. Validators can be removed by setting their voting power to zero. This provides a generic means for applications to update the validator set without having to specify transaction types.

If the block at height H returns an updated validator set, then the block at height $H + 1$ will reflect the update. Note, however, that the *LastCommit* in block $H + 1$ must utilize the validator set as it was at H , since it may contain signatures from a validator that was removed.

Changes to voting power are applied for $H + 1$ such that the next proposer is affected by the update. In particular, the validator that otherwise should have been the next proposer may be removed. The round robin algorithm should handle this gracefully, simply moving on to the next proposer in line. Since the same block is replicated on at least two-thirds of validators, and the round robin is deterministic, they will all make the same update and expect the same next proposer.

5.3 Punishing Byzantine Validators

One of the salient points of Bitcoin’s design is its incentive structure, in so far as the goal of the protocol was to incentivize validators to behave correctly by rewarding them. While this makes sense in the context of Bitcoin’s consensus protocol, a superior incentive may be to provide strong dis-incentives, such that validators have real *skin-in-the-game* [95], rather than a soft opportunity cost.

Disincentives can be achieved in Tendermint using an approach first proposed by Vitalik Buterin [12] as a so-called Proof-of-Stake protocol. In essence, validators must make a security deposit (“they must bond some stake”) in order to participate in consensus. In the event that they are found to double-sign proposals or votes, other validators can publish evidence of the transgression in the form of a transaction, which the application state can use to change the validator set by removing the transgressor, burning its deposit. This has the effect of associating an explicit economic cost with Byzantine behaviour, and enables one to estimate the cost of violating safety by bribing a third or more of the validators to be Byzantine.

Note that a consensus protocol may specify more behaviours to be punished than just double signing. In particular, we are interested in punishing any strong signalling behaviour which is unjustified - typically, any reported change in state that is not based on the reported state of others. For instance, in a version of Tendermint where all pre-commits must come with the polka that justifies them, validators may be punished for broadcasting unjustified pre-commits. Note, however, that we cannot just punish for any unexpected behaviour - for instance, a validator proposing when it is not their round to propose may be a basis for optimizations which pre-empt asynchrony or crashed nodes.

In fact, a generalization of Tendermint along these two lines, of 1) looser forms of justification and 2) allowing validators to propose before their term, gives rise to a family of protocols similar in nature to that proposed by Vlad Zamfir, under the guise Casper, as the consensus mechanism for a future version of ethereum [109]. A more formal account of the relationship between the protocols, and of the characteristics of anti-Byzantine justifications, remains for future work.

5.4 Software Upgrades

Government can also be used as a natural means for negotiating software upgrades on a possibly decentralized network. Software upgrades on the public Internet are a notoriously challenging operation, requiring careful planning to maintain backwards compatibility for users that don't upgrade right away, and to not upset loyal users of the software by introducing bugs, removing features, adding complexity, or, perhaps worst of all, updating automatically without permission.

The challenge of upgrading a decentralized consensus system is made especially apparent with Bitcoin. While Ethereum has already managed a successful, non-backwards-compatible upgrade, due to its strong leadership and unified community, Bitcoin has been unable to make some needed upgrades, despite a plethora of software engineering ills, on account of a viciously divided community and a lack of strong leadership.

Upgrades to blockchains are typically differentiated as being *soft forks* or *hard forks*, on account of the scope of the changes. Soft forks are meant to be backwards compatible, and to use degrees of freedom in the protocol that may be ignored by users who have not upgraded, but which provide new features to users which do. Hard forks, on the other hand, are non-backwards compatible upgrades that, in Bitcoin's case, may cause violations of safety, and in Tendermint's case, cause the system to halt.

To cope, developers of the Bitcoin software have rolled out a series of soft forks for which validators can vote by signalling in new blocks. Once a certain threshold of validators are signalling for the update, it automatically takes effect across the network, at least for users with a version of the software supporting the update. The utility of the Bitcoin system has grown tremendously on account of these softforks, and is expected to continue to do so on account of upcoming ones. Interestingly, the failure of the community to successfully hard fork the software has on the one hand raised concerns about the long term stability of the system, and on the other triggered excitement and inspiration about the system's resilience to corrupt governance - its ungovernability.

There are many reasons to take the latter stance, given the overwhelming government corruption apparent in the world today. Still, cryptography and distributed consensus provide a new set of tools that enables a degree of transparency and accountability otherwise not imaginable in the paper-pen-handshake world of modern governments, nor even the digital world of the

traditional web, which suffers tremendously from sufficiently robust authentication systems.

In a system using Governmint, developers would be identifiable entities on the blockchain, and may submit proposals for software upgrades. The mechanism is quite similar to that of a Pull Request on Github, only it is integrated into a live running system, and the agreement passes through the consensus protocol. Clients should be written with configurable update parameters, so they can specify whether to update automatically or to require that they are notified first.

Of course, any software upgrade which is not thoroughly vetted could pose a danger to the system, and a conservative approach to upgrades should be taken in general.

5.5 Crisis Recovery

In the event of a crisis, such as a fork in the transaction log, or the system coming to a halt, a traditional consensus system provides little or no guarantees, and typically requires manual intervention.

Tendermint assures that those responsible for violating safety can be identified, such that any client who can access at least one honest validator can discern with cryptographic certainty who the dishonest validators are, and thereby chose to follow the honest validators onto a new chain with a validator set excluding those who were Byzantine.

For instance, suppose a third or more validators violate locking rules, causing two blocks to be committed at height H . The honest validators can determine who double-signed by gossiping all the votes. At this point, they cannot use the consensus protocol, because the basic fault assumptions have been violated. One approach is for each validator to wait until they have all Note that being able to at this point accumulate all votes for H implies strong assumptions about network connectivity and availability during the crisis, which, if it cannot be provided by the p2p network, may require validators use alternative means, such as social media and high availability services, to communicate evidence. A new blockchain can be started by the full set of remaining honest nodes, once at least two-thirds of them have gathered all the evidence.

Alternatively, modifying the Tendermint protocol so that pre-commits require polka would ensure that those responsible for the fork could be pun-

ished immediately, and would not require an additional publishing period. This modification remains for future work.

More complex uses of Governmint are possible for accommodating various particularities of crisis, such as permanent crash failures and the compromise of private keys. However, such approaches must be carefully thought out, as they may undermine the safety guarantees of the underlying protocol. We leave investigation of these methods to future work, but note the importance of the socio-economic context in which a blockchain is embedded, in terms of understanding its ability to recover from crisis.

Regardless of how crisis recovery proceeds, its success depends on integration with clients. If clients do not accept the new blockchain, the service is effectively offline. Thus, clients must be aware of the rules used by the particular blockchain to recover. In the cases of safety violation described above, they must also gather the evidence, determine which validators to remove, and compute the new state with the remaining validators. In the case of the liveness violation, they must keep up with Governmint.

5.6 Conclusion

Governance is a critical element of a distributed consensus system, though competent governance systems remain poorly understood. Tendermint provides governance as a TMSP module called Governmint, which aims to facilitate increased experimentation in software-based governance for distributed systems.

Chapter 6

Client Considerations

This chapter reviews some considerations pertaining to clients that interact with an application hosted on Tendermint.

6.1 Discovery

Network discovery occurs simply by dialing some set of seed nodes over TCP. The p2p network uses authenticated encryption, but the public keys of the validators must be verified somehow out of band, that is, via an alternative medium not within the purview of the protocol. Indeed, in these systems, the genesis state itself must be communicated out of band, and ideally is the only thing that must be communicated, as it should also contain the public keys used by validators for authenticated encryption, which are different than those used for signing votes in consensus.

For validator sets that may change over time, it is useful to register all validators via DNS, and to register new validators before they actually become validators, and remove them after they are removed as validators. Alternatively, validator locations can be registered in another fault-tolerant distributed data store, including possibly another Tendermint cluster itself.

6.2 Broadcasting Transactions

As a generalized application platform, Tendermint provides only a simple interface to clients for broadcasting transactions. The general paradigm is that a client connects to a Tendermint consensus network through a proxy,

which is either run locally on its machine, or hosted by some other provider. The proxy functions as a non-validator node on the network, which means it keeps up with the consensus and processes transactions, but does not sign votes. The proxy enables client transactions to be quickly broadcast to the whole network via the gossip layer.

A node need only connect to one other node on the network to broadcast transactions, but by default will connect to many, minimizing the chances that the transaction will not be received. Transactions are passed into the mempool, and gossiped through the mempool reactor to be cached in the mempool of all nodes, so that eventually one of them will include it in a block.

Note that the transaction does not execute against the state until it gets into a block, so the client does not get a result back right away, other than confirmation that it was accepted into the mempool and broadcast to other peers. Clients should register with the proxy to receive the result as a push notification when it is computed during the commit of a block.

It is not essential that a client connect to the current proposer, as eventually any validator which has the transaction in its mempool may propose it. However, preferential broadcasting to the next proposer in line may lead to lower latency for the transaction in certain cases where the network is under high load. Otherwise, the transaction should be quickly gossiped to every validator.

6.3 Mempool

The mempool is responsible for caching transactions in memory before they are included in blocks. Its behaviour is subtle, and forms a number of challenges for the overall system architecture. First and foremost, caching arbitrary numbers of transactions in the mempool is a direct denial of service attack that could trivially cripple the network. Most blockchains solve this problem using their native currency, and permitting only transactions which spend a certain fee to reside in the mempool.

In a more generalized system, like Tendermint, where there is not necessarily a currency to pay fees with, the system must establish stricter filtering rules and rely on more intelligent clients to resubmit transactions that are dropped. The situation is even more subtle, however, because the rule set for filtering transactions in the mempool must be a function of the application

itself. Hence the *CheckTx* message of TMSP, which the mempool can use to run a transaction against a transient state of the application to determine if it should be kept around or dropped.

Handling the transient state is non-trivial, and is something left to the application developer, though examples are provided in the many example applications. In any case, clients must monitor the state of the mempool (i.e. the unconfirmed transactions) to determine if they need to rebroadcast their transactions, which may occur in highly concurrent settings where the validity of one transaction depends on having processed another.

6.4 Semantics

Tendermint's core consensus algorithm provides only *at-least-once semantics*, which is to say the system is subject to replay attacks, where the same transaction can be committed many times. However, many users and applications expect stronger guarantees from a database system. The flexibility of the Tendermint system leaves the strictness of these semantics up to the application developer. By utilizing the *CheckTx* message, and by adequately managing state in the application, application developers can provide the database semantics that suit them and their users' needs. For instance, as discussed in Chapter 4, using an account based system with sequence numbers mitigates replay attacks, and changes the semantics from *at-least-once* to *exactly-once*.

6.5 Reads

Clients issue read requests to the same proxy node they use for broadcasting transactions (writes). The proxy is always available for reads, even if the network halts. However, in the event of a partition, the proxy may be partitioned from the rest of the network, which continues making blocks. In that case, reads from the proxy might be stale.

To avoid stale reads, the read request can be sent as a transaction, presuming the application permits such queries. By using transactions, reads are guaranteed to return the latest committed state, i.e. when the read transaction is committed in the next block. This is of course much more expensive than simply querying the proxy for the state. It is possible to use heuristics

to determine if a read will be stale, such as if the proxy is well-connected to its peers and is making blocks, or if it's stuck in a round with votes from one-third or more of validators, but there is no substitute for performing an actual transaction.

6.6 Light Client Proofs

One of the major innovations of blockchains over traditional databases is their deliberate use of Merkle hash trees to enable the production of compact proofs of system substates, so called light-client proofs. A light client proof is a path through a Merkle tree that allows a client to verify that some key-value pair is in the Merkle tree with a given root hash. The state's Merkle root hash is included in the block header, such that it is sufficient for a client to have only the latest header to verify any component of the state. Of course, to know that the header itself is valid, they must have either validated the whole chain, or kept up-to-date with validator set changes only and rely on economic guarantees that the state transitions were correct.

Chapter 7

Implementation

The reference implementation of Tendermint is written in Go [81] and hosted at <https://github.com/tendermint/tendermint>. Go is a C-like language with a rich standard library, concurrency primitives for light-weight massively concurrent executions, and a development environment optimized for simplicity and efficiency.

The code uses a number of packages which are modular enough to be isolated as their own libraries. These packages were written for the most part by Jae Kwon, with bug fixes, tests, and the occasional feature contributed by the author. The most important of these packages are described in the following sub-sections.

7.1 Binary Serialization

Tendermint uses a binary serialization algorithm optimized for simplicity and determinism. It supports all integer types (including varints, which are encoded with a one-byte length prefix), strings, byte arrays, and time (unix time with millisecond precision). It also supports arrays of any type and structs (encoded as a list of ordered values, ignoring keys). It is somewhat inspired by Go's type system, especially its use of interface types, which can be implemented as one of many concrete types. Interfaces can be registered and each concrete implementation given a leading type-byte in its encoding.

See <https://github.com/tendermint/go-wire> for more details.

7.2 Cryptography

Consensus algorithms such as tendermint use three primary cryptographic primitives: digital signatures, hash functions, and authenticated encryption. While many implementations of these primitives exist, choosing a cryptography library for enterprise software is no trivial task, given especially the profound insecurity of the world's most used security library, OpenSSL [77].

Contributing to the insecurity of cryptographic systems is the potential deliberate undermining of their security properties by government agencies such as the NSA, who, in collaboration with the NIST, have designed and standardized many of the most popular cryptographic algorithms in use today. Given the apparent unlawfulness of such agencies, as made evident, for instance, by Edward Snowden [43], and a history of trying to compromise public cryptographic standards [63], many in the cryptography community prefer to use algorithms designed in an open, academic environment [.] Tendermint, similarly, uses only such algorithms.

Tendermint uses RIPEMD160 as its cryptographic hash function, which produces 20-byte outputs. It is used in the Merkle trees of transactions and validator signatures, and for computing the block hash. Go provides an implementation in its extended library. RIPEMD160 is also used as one of two hashing functions by Bitcoin in the derivation of addresses from public keys.

As its digital signature scheme, Tendermint uses Schnorr signatures over the ED25519 elliptic curve. ED25519 was designed in the open by Dan Bernstein [6], with the intention of being high performance and easy to implement without introducing vulnerabilities. Bernstein also introduced NaCl, a high level library for doing authenticated encryption that uses the ED25519 curve. Tendermint uses the implementation provided by Go in its extended library.

7.3 Merkle Hash Tree

Merkle trees function much like other tree-based data-structures, with the additional feature that it is possible to produce a proof of membership of a key in the tree that is logarithmic in the size of the tree. This is done by recursively concatenating and hashing keys in pairs until only a single hash is left, the root hash of the tree. For any leaf in the tree, a trail of hashes leading from it to the root serves as proof of its membership. This makes Merkle

trees particularly useful for p2p file-sharing applications, where pieces of a large file can be verified as belonging to the file without having all the pieces. Tendermint uses this mechanism to gossip block parts on the network, where the root hash is included in the block proposal.

Tendermint also provides a self-balancing, Merkle binary tree, modeled after the AVL tree [3], as a TMSP service called Merkleeyes. The IAVL tree can be used for storing state of dynamic size, allowing lookups, inserts, and removals in logarithmic time.

7.4 RPC

Tendermint exposes HTTP APIs for querying the blockchain, network information, and consensus state, and for broadcasting transactions. The same API is available via three methods: GET requests using URI encoded parameters, POST requests using the JSONRPC standard [53], and websockets using the JSONRPC standard. Websockets are the preferred method for high transaction throughput, and are necessary for receiving events.

7.5 P2P Networking

The P2P subprotocols used by Tendermint are described more fully in Chapter 3.

7.6 Reactors

The Tendermint node is composed of multiple concurrent reactors, each managing a state machine sending and receiving messages to peers over the network, as described in Chapter 3. Reactors synchronize by locking shared datastructures, but the points of synchronization are kept to a minimum, so that each reactor runs mostly concurrently with the others.

7.6.1 Mempool

The mempool reactor manages the mempool, which caches transactions before they are packed in blocks and committed. The mempool uses a subset of

the application's state machine to check the validity of transactions. Transactions are kept in a concurrent linked list structure, allowing safe writes and many concurrent reads. New, valid transactions are added to the end of the list. A routine for each peer traverses the list, sending each transaction to the peer, in order, only once. The list is also scanned to collect transactions for a new proposal, and is updated every time a block is committed: committed transactions are removed, uncommitted transactions are re-run through CheckTx, and those that have become invalid are removed.

7.6.2 Consensus

The consensus reactor manages the consensus state machine, which handles proposals, voting, locking, and the actual committing of blocks. The state machine is managed using a few persistent go-routines, which order received messages and enable them to be played back deterministically to debug the state. These go-routines include the readLoop, for reading off the queue of received messages, and the timeoutLoop, for registering and triggering timeout events.

Transitions in the consensus state machine are made either when a complete proposal and block are received, or when more than two-thirds of either pre-votes or pre-commits have been received at a given round. Transitions result in the broadcast of proposals, block data, or votes, which are queued on the internalReqQueue, and processed by the readLoop in serial with messages received from peers. This puts internal messages and peer messages on equal footing as far as being inputs to the consensus state machine, but allows internal messages to be processed faster, as they don't sit in the same queue as those from peers.

7.6.3 Blockchain

The blockchain reactor syncs the blockchain using a much faster technique than the consensus reactor. Namely, validators request blocks of incrementing height until none of their peers have blocks of any higher height. Blocks are collected in a blockpool and synced to the blockchain by a worker routine that periodically takes blocks from the pool and validates them against the current chain.

Once the blockchain reactor finishes syncing up, it turns on the consensus reactor to take over.

7.7 Conclusion

The implementation of Tendermint in Go takes advantage of the language's concurrency primitives, garbage collection, and type safety, to provide a clear, modular, easy to read code base with many reusable components. As will be seen in Chapter 8, the implementation obtains high performance and is robust to many different kinds of fault.

Chapter 8

Performance and Fault Tolerance

Tendermint is designed as a Byzantine fault tolerant state-machine replication algorithm. It guarantees safety so long as less than a third of validators are Byzantine, and guarantees liveness similarly, so long as network messages are eventually delivered, with weak assumptions about network synchrony for gossiping proposals. In this section, we evaluate Tendermint’s fault tolerance empirically by injecting crash faults and Byzantine faults. The goal is to show that the implementation of Tendermint consensus does not compromise safety in the event of such failures, that it suffers minimum performance impact, and that it is quick to recover.

Performance of the Tendermint algorithm can be evaluated in a few key ways. The most obvious measures are the block commit time, which is a measure of finalization latency, and transaction throughput, which measures the network’s capacity. We collect measurements for each on networks with validators distributed over the globe, where the number of validators ranges, in multiples of 2, from 2 to 64.

8.1 Overview

The experiments in this chapter can be reproduced using the repository at https://github.com/tendermint/network_testing. All experiments except one take place in docker containers running on *Amazon EC2* instances of type *t2.medium*. The *t2.medium* has 2 vCPU and 4 GB of RAM. One ex-

periment tests the throughput on larger instances, the *c3.8xlarge*, which has 32 vCPUs and 60 GB of RAM. Instances are distributed across seven datacenters, spanning five continents. A second docker container, responsible for generating transactions, is run on each instance. Transactions are 250 bytes in size (a reasonable size for including a few 32 or 64 byte hashes and signatures), and were constructed to be debuggable, to be quick to generate, and to contain some stochasticity. Thus, the leading bytes are Big-Endian encoded integers representing transaction number and validator index for that instance, the trailing 16 bytes are randomly drawn from the operating system, and the intermediate bytes are just zeros.

A network monitoring tool is used to maintain active websocket connections to each validator’s Tendermint RPC server, and uses its local time when it receives a new committed block for the first time as the official commit time for that block. Experiments were first run without the monitor by copying all data from the validators for analysis and using the local time of the 2/3th validator committing a block as the commit time. Using the monitor is much faster, amenable to online monitoring, and was found to not impact the results so long as only block header information (and not the whole block) was passed over the websockets.

Docker containers on remote machines are easily managed using the *docker-machine* tool, and the `network_testing` repository provides some tools which take advantage of Go’s concurrency features to perform actions on docker containers on many remote machines at once.

Each validator connects directly to each other to avoid confounding effects of network topology.

For experiments involving crash faults or Byzantine behaviour, the number of faulty nodes is given by $N_{fault} = \lfloor (N - 1)/3 \rfloor$, where N is the total number of validators.

8.2 Throughput and Latency

This section describes experiments which measure the raw performance of Tendermint in non-adversarial conditions, where all nodes are online and synced and no accommodations are made for asynchrony. That is, an artificially high `TimeoutPropose` is used (10 seconds), and all other timeout parameters are set to 1 millisecond. Additionally, all mempool activity is disabled (no gossiping of transactions or rechecking them after commits),

and an in-process nil application is used to bypass TMSP. This serves as a control scenario for evaluating the performance drop in the face of faults and/or asynchrony.

Experiments are run on validator set sizes doubling in size from two to 64, and on block sizes doubling from 128 to 32768. Transactions are preloaded on each validator. Each experiment is run for 16 blocks.

As can be seen in Figure 8.1, Tendermint easily handles thousands of transactions per second with around one second block latency, though there appears to be a capacity limit at around ten thousand transactions per second. A block of 16384 transactions is about 4 MB in size, and analysis of network bandwidth shows each connection easily reaching upwards of 20MB/s, though analysis of the logs shows that at high block sizes, validators can spend upwards of two seconds waiting for block parts. Additionally, experiments in single data centers, as shown in Figure 8.2, demonstrate that much higher throughputs are possible, while experiments on much larger machines exhibit more consistent performance, relieving the capacity limit, as shown in Figure 8.3. We leave further investigations of this capacity limit to future work.

In the experiments that follow, various forms of fault are injected and latency statistics presented. Each experiments was run for validator set sizes doubling from 4 to 32, for varying values of TimeoutPropose, and with a block size of 2048 transactions.

8.3 Crash Failures

To evaluate the performance of a network subject to crash failures, every three seconds N_{fault} validators were randomly selected, stopped, and restarted three seconds later.

The results in Table 8.1 demonstrate that performance under this crash failure scenario drops by about 50%, and that larger TimeoutPropose values help mediate latencies. While the average latency increases to about two seconds, the median is closer to one second, and latencies may run as high as ten or twenty seconds, though in one case it was as high as seventy seconds. It is likely that modifying TimeoutPropose to be slightly non-deterministic may ease the probability of such extreme latencies.

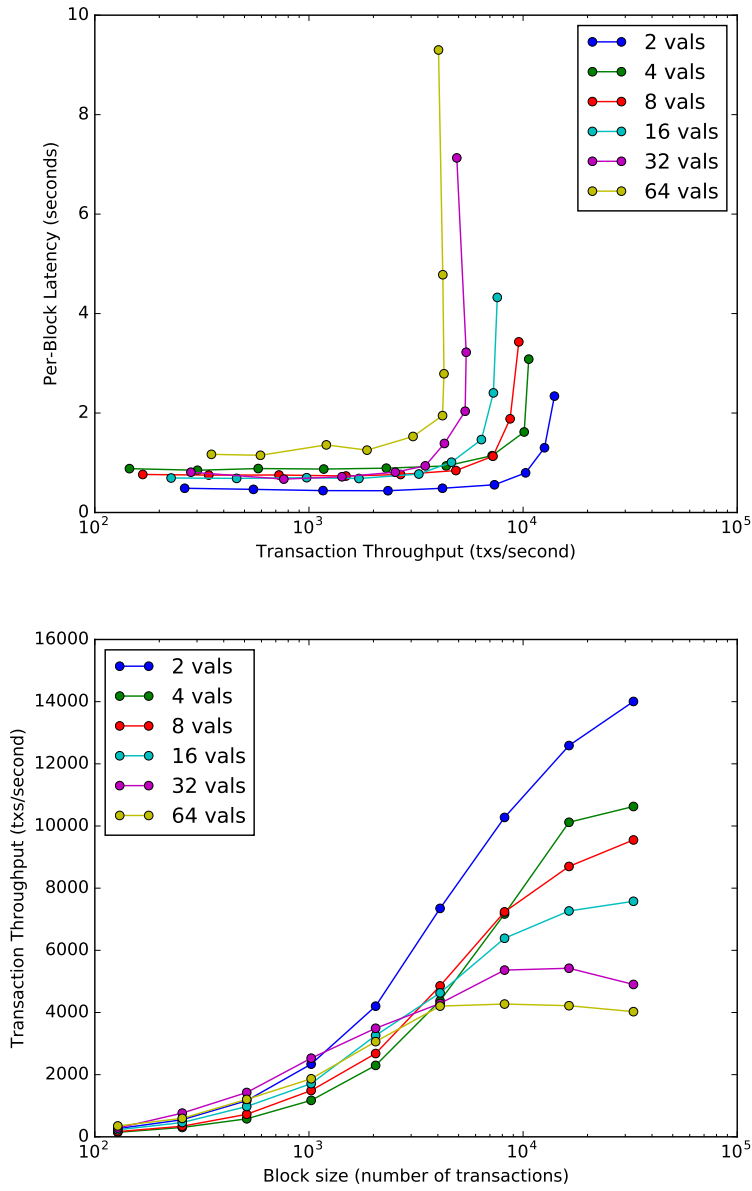


Figure 8.1: Latency-throughput trade-off. Larger blocks incur diminishing returns in transaction throughput, with an ultimate capacity at around 10,000 txs/s

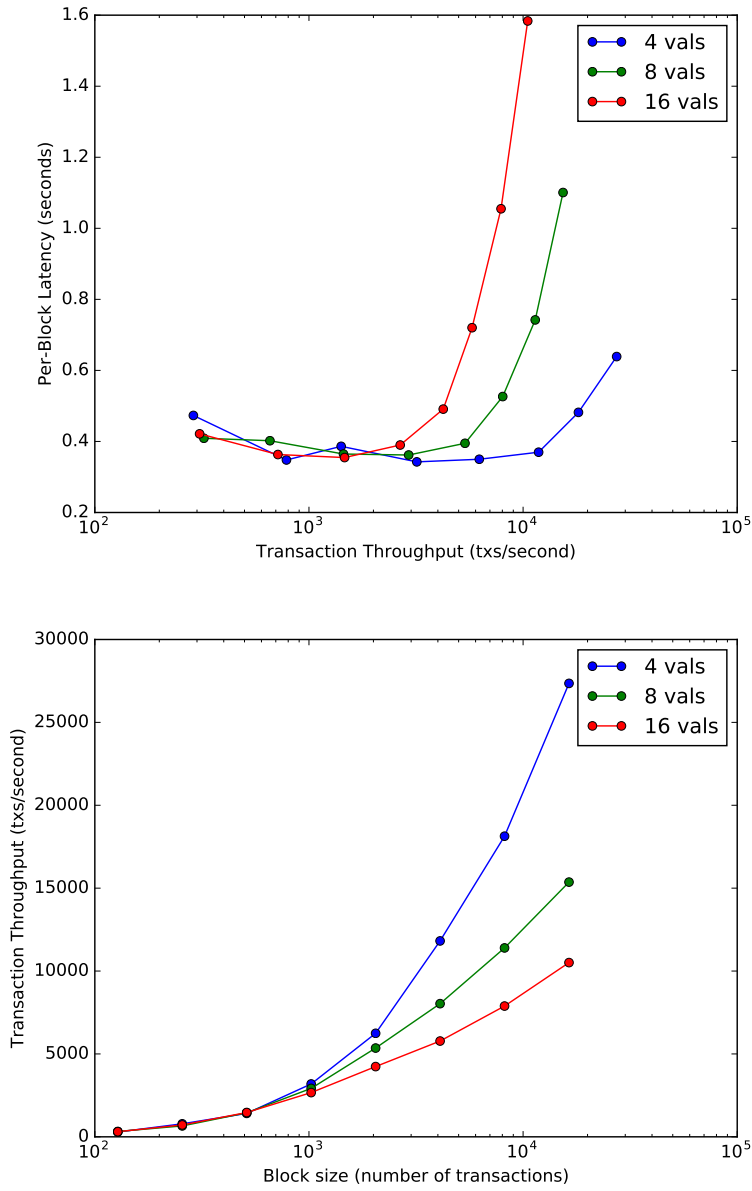


Figure 8.2: Single datacenter. When messages don't need to cross the public Internet, Tendermint is capable of tens of thousands of transactions per second.

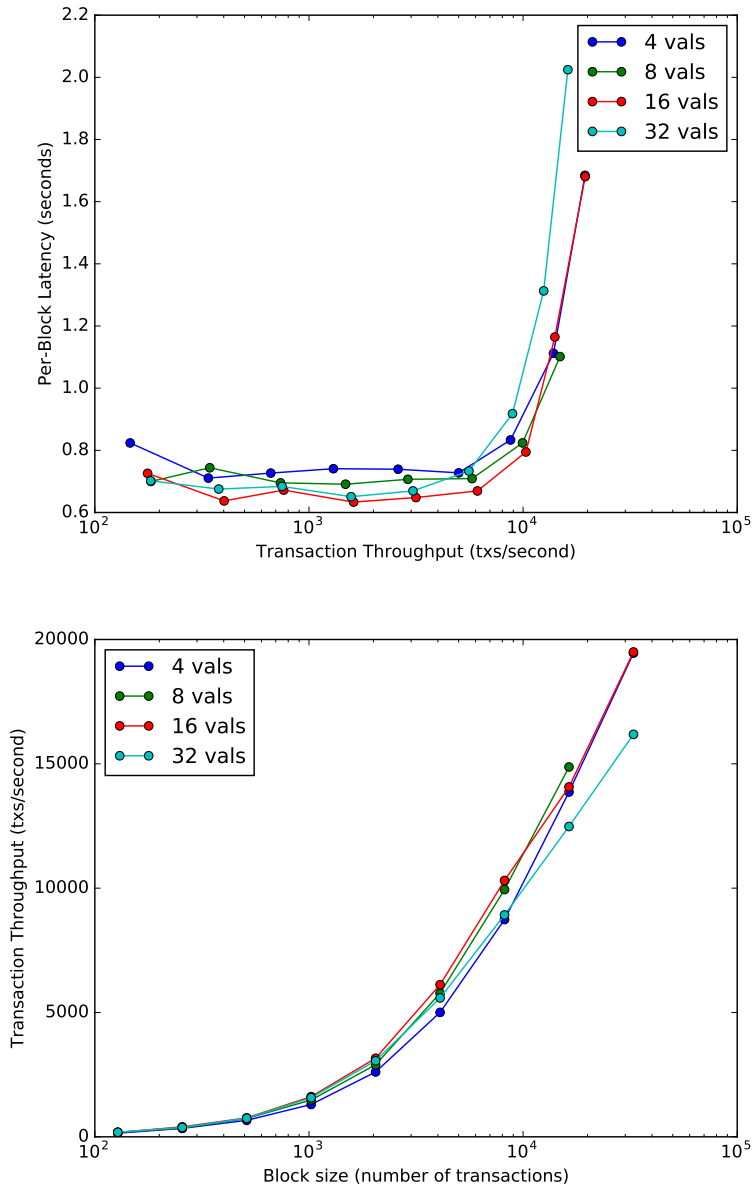


Figure 8.3: Large machines. With 32 vCPU and 60 GB of RAM, transaction throughput increases linearly with block-size, relieving the capacity limits found on smaller machines.

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
500	434	15318	2179	1102	5575
1000	516	18149	2180	1046	5677
2000	473	15067	2044	1049	5479
3000	428	9964	2005	1096	5502

(a) 4 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
500	618	126481	2679	990	5589
1000	570	9832	1763	962	5835
2000	594	8869	1658	968	5481
3000	535	10101	1633	959	5485

(b) 8 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
500	782	21354	1977	1001	5930
1000	758	12659	1761	981	5642
2000	751	21285	2041	1005	6872
3000	719	72406	2395	991	5987

(c) 16 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
500	760	24692	2591	1087	14025
1000	755	19696	2328	1119	9321
2000	852	21044	2178	1141	6514
3000	763	25587	2289	1119	6707

(d) 32 Validators

Table 8.1: Crash-fault latency statistics. Every three seconds, a random selection of N_{fault} validators were crashed, and restarted three seconds later. This crash-restart procedure continued for 200 blocks. Each table reports the minimum, maximum, average, median, and 95th percentile of the block latencies, for varying values of the TimeoutPropose parameter.

8.4 Random Network Delay

Another form of fault, which may be attributed either to Byzantine behaviour or to network asynchrony, is to inject random delays into every read and write to a network connection. In this experiment, before every read and write on every network connection, N_{fault} of the validators slept for X milliseconds, where X was drawn uniformly on $(0, 3000)$. As can be seen in Table 8.2, latencies are similar to the crash failure scenario, though increasing the TimeoutPropose has the opposite effect. Since not all validators were faulty, small values of TimeoutPropose allow faulty validators to be skipped quickly. If all validators were subject to the network delays, larger TimeoutPropose values would be expected to reduce latency since there would be no non-faulty validators to skip to, and more time would be provided to receive delayed messages.

8.5 Byzantine Failures

A more explicit Byzantine failure can be injected through the following modifications to the state machine:

- Conflicting proposals: during its time to propose, a Byzantine validator signs two conflicting proposals and broadcasts each, along with a pre-vote and pre-commit, to separate halves of its connected peers.
- No nil votes: a Byzantine validator never signs a nil-vote.
- Sign every proposal: a Byzantine validator submits a pre-vote and a pre-commit for every proposal it sees, as soon as it sees it.

Taken together, these behaviours explicitly violate the double signing and locking rules. Note, however, that the behaviour is dominated by the broadcast of conflicting proposals, and the eventual committing of one of them. More complex arrangements of Byzantine strategies are left for future work.

Despite the injected Byzantine faults, which would cause many systems to fail completely and immediately, Tendermint maintains respectable latencies, as can be seen from Table 8.3. Since these faults have little to do with asynchrony, there is no real discernible effect from TimeoutPropose. The

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
1000	873	2796	1437	1036	2627
2000	831	4549	1843	1180	4036
3000	921	5782	2273	1251	5491
4000	967	6875	2700	1413	6781

(a) 4 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
1000	870	2840	1449	1040	2786
2000	957	4268	1848	1076	4148
3000	859	5724	2156	1100	5649
4000	897	11859	3055	1093	11805

(b) 8 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
1000	914	5595	1821	1135	5466
2000	950	7782	2490	1165	7650
3000	978	10305	3049	1163	9890
4000	1018	6890	2808	1174	6813

(c) 16 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – ile
1000	1202	8562	2219	1349	5733
2000	1196	7878	2549	1365	7579
3000	1164	10082	3003	1382	9805
4000	1223	17571	3696	1392	12014

(d) 32 Validators

Table 8.2: Random delay latency statistics. N_{fault} validators were set to inject a random delay before every read and write, where the delay time was chosen uniformly on $(0, 3000)$ milliseconds.

performance also falls off with larger validator sets, which may be the result of a naive algorithm for handling Byzantine votes.

8.6 Related Work

The throughput experiments in this chapter were modeled after those in [67], which benchmarks the performance of a PBFT implementation and a new randomized BFT protocol called HoneyBadgerBFT. In their results, PBFT achieves over 15,000 transactions per second on four nodes, but decays exponentially as the number of nodes increases, while HoneyBadgerBFT attains roughly even performance of between 10,000 and 15,000 transactions per second. Block latencies in HoneyBadgerBFT, however, are much higher, closer to 10 seconds for validator sets of size 8, 16, and 32, and even more for larger ones.

A well known tool for studying consensus implementations is Jepsen [52], which is used to test the consistency guarantees of databases by simulating many forms of network partition. Testing Tendermint with Jepsen remains an exciting area for future work.

The author is not aware of any throughput experiments in the face of persistent Byzantine failures, like those presented here.

8.7 Conclusion

The implementation of Tendermint written by the author and Jae Kwon easily achieves thousands of transactions per second on up to 64 nodes on machines distributed around the globe, with latencies mostly in the one to two second range. This is highly competitive with other solutions, and especially with the current state of blockchains, with Bitcoin, for instance, capping out at around 7 transactions per second. Furthermore, our implementation is shown to be robust to both crash faults, message delays, and deliberate Byzantine faults, being able to maintain over a thousand transactions per second in each scenario.

TimeoutPropose	Min	Max	Mean	Median	95 th % – <i>ile</i>
1000	868	3888	1450	1086	3320
2000	929	4375	1786	1272	4166
3000	881	4363	1224	1099	1680
4000	824	8256	1693	1272	2607

(a) 4 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – <i>ile</i>
1000	771	3445	1472	916	3288
2000	731	3661	1426	902	3339
3000	835	6402	1912	962	6155
4000	811	4462	1512	964	3592

(b) 8 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – <i>ile</i>
1000	877	15930	2086	1024	5844
2000	808	5737	1580	1027	4155
3000	919	10533	1801	1110	4174
4000	915	5589	1745	1095	4181

(c) 16 Validators

TimeoutPropose	Min	Max	Mean	Median	95 th % – <i>ile</i>
1000	1594	11730	2680	1854	5016
2000	1496	17801	3430	1874	11730
3000	1504	15963	3280	1736	9569
4000	1490	24836	3940	1773	12866

(d) 32 Validators

Table 8.3: Byzantine-fault latency statistics. Byzantine validators propose conflicting blocks and vote on any proposal as soon as they see it. Each table reports the minimum, maximum, average, median, and 95th percentile of the block latencies, for varying values of the TimeoutPropose parameter.

Chapter 9

Related Work

Byzantine consensus has a rich history that spans cryptography, distributed computing, and economics, but the socio-economic context for its products to be deployed in industry has not existed until recently, at least not outside of traditionally critical real-time systems like aircraft control [47]. On the one hand, the invention of Bitcoin and the coining of the term “blockchain” popularized the notion of a distributed ledger not controlled by a single entity, using cryptography and aligned economic incentives to preserve safety in the face of Byzantine faults. On the other, the continued commoditization of servers, in the form of “The Cloud”, and the invention of Raft, have popularized distributed computing in mainstream developer culture, and brought renewed attention to distributed consensus algorithms as co-ordination hubs in large-scale deployments.

At the intersection are a collection of solutions, typically geared for banking and financial applications, but also for governance, logistics, and other general forms of co-ordination, that draw on classic academic BFT modified and modernized in various ways. This chapter reviews the history and diversity of these ideas, with the goal of providing a rich context within which to understand the blockchain phenomenon.

9.1 Beginnings

Distributed algorithms first emerged in the late 19th century in the telecommunications and railroad industries, in attempts to effectively handle multiple concurrent message streams in a transmission, or multiple trains on the

same set of tracks.

Academic work on the subject appears to have been launched officially by the seminal work of Edsger Dijkstra on the mutual exclusion problem [30], and of Tony Hoare on models for describing communicating processes [46].

A host of concurrency problems with catchy names were popularized around this time, including the cigarette smokers problem [44], where smokers sit around a table, each with a different ingredient, and must successfully roll a full cigarette; the dining philosophers problem [29], where philosophers sitting around a table must take turns eating and thinking, but each can only eat while its neighbours are thinking; and the two-generals or co-ordinated attack problem [38], where two generals must co-ordinate from afar to attack an enemy city at the same time.

These problems served to put the focus on synchronization primitives such as semaphores, mutexes, and communication channels, and would lay the groundwork for a number of advancements over the coming decades.

9.1.1 Faulty Things

Fault tolerant distributed computing effectively emerged in the late seventies out of the effort to utilize microprocessors for aircraft control, resulting in a number of early systems [103, 48]. Today, it is standard for NASA to conduct BFT research [70], and for commercial aircraft to use BFT systems, such as the SAFEbus [49].

Many systems, however, do not require tolerance to Byzantine faults as they are run in controlled environments, where there is presumably no malicious behaviour and the code is written correctly. In these circumstances, which are common in data-centers managed by large companies like Google or Amazon, fault tolerant computing is used to defend against various faults, whether it be a break in a network link, power failure in a server rack, or a dead hard-drive.

9.1.2 Clocks

The problem of distributed consensus, however, did not formally emerge until Leslie Lamport introduced it in his “Time, clocks, and the ordering of events in a distributed system” [60]. In that work, Lamport demonstrated how a partial ordering of events emerges from a definition of causality based on communication [60]. That is, events occurring in concurrent processes,

between communication events, effectively happen at the same time, as they cannot influence one another. Thus, a system of logical clocks can be defined based on the individual sequential processes and the fact that messages are sent before they are received. Events can then be totally ordered by assigning any arbitrary but consistent total ordering above the partial ordering, for instance by assigning each process in the system an index and ordering events which happen at the same logical time by the index of the process in which they happen. The algorithm is quite simple, requiring each process to hear from each other process in order to determine the order of events.

Lamport's work established time as a principle obstacle to designing fault tolerant distributed systems, as synchronizing clocks across geographical locations requires the communication of messages which is ultimately limited by the speed of light. This formulation of the problem has close ties to the relativism of modern physics, wherein frames of reference are relative to an observer and the speed of light imposes a constraint on information propagation.

9.1.3 FLP

As discussed in Chapter 1, one of the primary factors in designing consensus algorithms are assumptions made about network and/or processor synchrony. A synchronous network is one in which messages are delivered within some fixed, known amount of time. Similarly, synchronous processors are one whose clocks stay within some fixed, known number of ticks of each other. In the early days of consensus research, the distinction was not well characterized, though the close relationship between asynchrony and crash failures is apparent even in [60]. Lamport's original consensus algorithm is able to operate in asynchronous environments, so long as all messages are eventually delivered from each process. However, the algorithm is obviously not fault tolerant as the failure of just a single process can halt the algorithm forever.

The intuition behind a single failure thwarting a consensus protocol was given formal ground by Fischer, Lynch, and Patterson, who proved the impossibility of deterministic distributed consensus in asynchronous environments even if a single process fails [37]. The result does not apply to synchronous contexts, as assumptions about network synchrony allow processors to detect failures using timeouts, such that if a process does not respond within some given amount of time it is assumed to have crashed. Furthermore, the result applies to deterministic consensus protocols only, as its proof relies on

the moment when the network goes deterministically from a bivalent state, where not all processes hold the same value, to a univalent one, where they do. Since the point of transition is a deterministic point in time, consensus fails if a single process crashes at that opportune moment.

9.1.4 Common Coin

The FLP result became something of a warning bell to distributed systems scientists, establishing a clear impossibility result at the heart of the emerging field. Later, the approach would be generalized to derive many more impossibility results [36], and significant academic effort would be expended on relaxing either the synchrony or determinism assumptions to derive algorithms which circumvent the result.

In particular, in a short note, Ben Or demonstrated how an algorithm which includes a simple amount of non-determinism can circumvent the FLP result [5]. The algorithm is tolerant to faults of up to half of the processes in asynchronous environments. Essentially, in trying to reach consensus on the value of a single bit, if a process does not receive votes from a majority for the same value, it randomly changes the value it votes for the next round. With everyone changing values, eventually more than half of them will vote the same value. This approach came to be known as a *common coin*, due to the resemblance of the procedure to communally flipping a coin to obtain a shared value.

The problem with Ben Or's common coin is that, in the asynchronous case, the algorithm requires a number of rounds exponential in the number of validators. This was quickly rectified in a follow up by Rabin, who showed how a common coin could be constructed using secret sharing, as pioneered by Shamir [88], to achieve consensus in a fixed number of rounds [85]. The approach is useful for BFT as well, and is discussed more fully in that context in a later section.

9.1.5 Transaction Processing

Parallel to the development of fault tolerant consensus algorithms was the emergence of the first commercial database systems. While they did not at first use the consensus protocols being developed, they built atop the growing body of work in distributed computing and concurrency. In particular is the seminal work of Jim Gray, who introduced the term *transaction* as an atomic

unit of work in a database system [42]. That is, a transaction is either applied in full or not at all.

Gray also introduced other classic features of modern databases, such as the principles of Atomicity, Consistency, Isolation, and Durability, which come part and parcel with the transaction concept [42], and the use of write-ahead-logs, for logging transactions to disk before they are executed in order to recover from faults occurring during transaction execution [41].

In a distributed database setting, this work on transactions, atomicity, and consistency led to a series of approaches for database replication centered around the notion of an *atomic commit*, wherein a transaction is replicated atomically across all machines. These approaches are known as two-phase-commit [41], and its non-blocking alternative, three-phase-commit [90].

Both two-phase and three-phase commit protocols work only in a synchronous setting, where crash failures can be detected, and utilize a co-ordinator process that serves as leader for the protocol.

9.1.6 Broadcast Protocols

The two most important broadcast protocols, RBC and ABC, were introduced in Chapter 1. A taxonomy and survey of solutions to the problem is provided in [27].

9.2 Byzantine

Many fault tolerant protocols focus only on crash failures, as they are the most common, while much less attention has been given to the problem of potentially arbitrary, including malicious, behaviour of software. This more general problem is known as Byzantine Fault Tolerance.

9.2.1 Byzantine Generals

Lamport introduced the problem of Byzantine Fault Tolerance in [78], but gave the problem its name in a later paper by making an analogy with the problem faced by the Byzantine army in co-ordinating to attack an enemy city [61]. The army is composed of multiple divisions, each of which is led by a general. Communication between generals happens only via messenger.

How can the generals agree on a common plan of action if one or some of the generals is a traitor?

The original paper provides the first proof that to tolerate f Byzantine faults, a system must have at least $3f + 1$ nodes. The intuition behind this result was depicted in 1.2 and discussed throughout Chapters 1 and 2. A number of algorithms are provided in both papers as the first solutions to the problem, though they are designed to work only in the synchronous case, where the absence of a message can be detected.

9.2.2 Randomized Consensus

Asynchronous Byzantine consensus saw its first solution in the form of the common coins introduced by Ben Or [5] and Rabin [85]. However, neither solution achieves optimal Byzantine fault tolerance of $3f + 1$ machines for f faults. Ben Or's solution requires $5f + 1$ machines, while Rabin's requires $10f + 1$ machines. The solution was iteratively improved to achieve optimal Byzantine agreement with low overhead [35, 16, 13].

9.2.3 Partial Synchrony

The next major advancement in BFT came in the form of the so called *DLS* consensus algorithms, named after the authors Dwork, Lynch, and Stockmeyer [31]. The innovation of DLS was to define a middle ground between synchrony and asynchrony called *partial synchrony*. Recall from Chapter 1 that a synchrony assumption is one which states that messages are received within some known, finite amount of time, or that processor clocks only deviate from each other by some finite number of ticks. The secret to partial synchrony is to suppose one of the following:

- Messages are guaranteed to be delivered within some fixed but unknown amount of time.
- Messages are guaranteed to be delivered within some known amount of time, beginning an unknown amount of time in the future.

The DLS algorithm proceeds via a series of rounds, each of which is divided into *trying* and *lock-release* phases. Each round has a corresponding proposer, and processes can *lock* on a value at a round if they think the proposer will propose that value. A round begins with processes gossiping the

values they deem acceptable. The proposer will propose a value if it has heard from at least $N - f$ processes that the value is acceptable. Any process which receives the proposed value should lock on it, and send an acknowledgment message that it has done so. If the proposer receives acknowledgment from $f + 1$ processes, it commits the value.

Variations on the basic protocol are discussed for different combinations of assumptions, and many proofs are provided of its soundness. Despite its success, however, DLS algorithms were never widely adopted for BFT. Tendermint’s original design was based on DLS, in particular the version which assumes a partially synchronous network but synchronous processor clocks. In practice, due to the use of protocols like the Network Time Protocol (NTP), synchronized clocks may be a fair assumption. However, NTP is vulnerable to a number of attacks, and protocols which assume synchronous clocks can be slow to recover from crash faults. In the summer of 2015, the core Tendermint consensus protocol was redesigned to be more fully asynchronous, as described in Chapter 2, and has thus come to more closely resemble another BFT algorithm, known as Practical Byzantine Fault Tolerance (PBFT).

9.2.4 PBFT

PBFT was introduced in 1999 [17], and was widely hailed as the first practical BFT algorithm, suitable for use in asynchronous networks, though it does in fact make weak synchrony assumptions which can be violated by a careful adversary [67]. PBFT proceeds through a series of views, where each view has a proposer, known as a primary, that is selected in round-robin order. The primary receives requests from clients, assigns them a sequence number, and broadcasts a signed *pre-prepare* messages to the other processes containing the view and sequence numbers. Replicas accept the *pre-prepare* message if they have not already accepted one for the same view and sequence numbers, assuming the message is for the current view and signed by the correct primary.

Once a *pre-prepare* is accepted, a replica broadcasts a signed *prepare* message. A replica is said to be *prepared* for a given client request when it has received $2f$ *prepare* messages for that request, with the same view and sequence number. The combination of *pre-prepare* and *prepare* ensure a total order on the requests in a single view, according to their sequence number. Once a replica is prepared, it broadcasts a signed *commit* message, which

is accepted so long as its properly signed and the view is correct. When a replica accepts a *commit* message, it runs the client request against the state machine and returns the result to the client.

PBFT employs an additional mechanism to facilitate view changes in the event the primary is faulty. Replicas maintain a timeout, which restarts every time they receive a new client request, and terminates when a *pre-prepare* is received for that request. If no *pre-prepare* is received, the replica times out, and triggers the view change protocol. View change is subtle and somewhat complicated as it requires consensus that the view should be changed, and all client requests since the last commit must be brought into the new view.

Tendermint side-steps these issues through the use of blocks and by changing proposers every block, allowing a proposer to be skipped using the same mechanism used to commit the proposed block. Furthermore, the use of blocks allows Tendermint to include the set of *pre-commit* messages from one block in the next block, removing the need for an explicit *commit* message.

9.2.5 BFT Improvements

Many improvements have been proposed for PBFT since it was published. Some of these focus on so-called *optimistic execution*, where transactions are executed before they are committed in order to provide a low-latency, optimistic reply to clients [58, 39]. The trouble with these approaches is that the responsibility of managing inconsistency is relegated to the client, while presumably the reason they used a consistent consensus protocol in the first place was to avoid that responsibility. Alternatively, this may be a useful approach in low-fault circumstance. The phenomenon is referred to as *zero-conf transactions* in Bitcoin and is widely warned against, given Bitcoin’s weak consistency guarantees.

Others have focused on the possibility of running independent transactions concurrently to achieve higher throughputs [57]. This is the approach that has begun to be researched in the blockchain community, especially by Ethereum, in order to produce a scalable blockchain architecture.

9.3 Non-Byzantine

In parallel to the BFT algorithms, a number of non-BFT algorithms have emerged, and a number of important highly available Internet services have been built on top of them.

9.3.1 Paxos

It is often said in consensus science that there is only one consensus algorithm, and it is Paxos. This is on the one hand a statement of the significance of the Paxos algorithm to the field, and on the other a reflection on the universal foundation of consensus protocols, which is in every case “Paxos-like”.

Lamport introduced Paxos in the early nineties, though the article was not accepted for publication until almost a decade later [59]. Many have pointed out that the algorithm is actually quite similar to Viewstamped Replication, published in the late eighties [73], and that the two represent independent discovery of the same protocol.

The protocols are quite similar to PBFT, which came after them, but require only $2f + 1$ machines to tolerate f faults as they are not BFT. Another similar protocol, the Zookeeper Atomic Broadcast protocol (ZAB) [54] was developed for the Apache Zookeeper distributed key-value store. The similarities and differences of each algorithm are illuminated in [99].

9.3.2 Raft

Non-BFT consensus science received a major improvement with the introduction of Raft [75], which was designed from the ground up to be *understandable*, and which even proved itself to be more understandable than Paxos through a user survey [74].

Raft is similar in spirit to Paxos and Viewstamped Replication, but it emphasizes replicating a transaction log, rather than a single bit, and introduces randomization for more efficient leader elections. Furthermore, Raft’s safety guarantees have been formally proven using the Coq proof assistant [106] and a framework built above Coq, called Verdi, for formally verifying distributed systems [104]. It remains to be seen how Verdi will compare to process calculus based approaches.

9.4 Blockchain

This thesis was motivated by the introduction of blockchain technology, which emerged in the form of Bitcoin, and has since seen many iterations. Few have succeeded in putting the blockchain in context of classical consensus science until recently [102, 14, 67].

9.4.1 Bitcoin

Bitcoin was the first blockchain, introduced in [71]. It solved the atomic broadcast problem in a public, adversarial setting through a clever use of economics. In particular, the order of transactions comes in blocks proposed by those who solve partial hash collisions, where the data being hashed is the block of transactions. Since computing partial hash collisions is expensive, requiring brute force search in a large space, the effort is subsidized by the issuance of a currency, known as bitcoins, with every block. The protocol has been wildly successful, with the currency achieving a market capitalization of roughly five billion dollars (USD), and with many clones of the original that have market capitalizations in the millions.

However, Bitcoin is not without its issues. A number of design flaws make the protocol cumbersome and difficult for application developers to work with it. Furthermore, a number of academic works have shed light on incentive incompatibilities in the protocol, weakening widely held assumptions about the protocol’s security [33, 24].

Numerous approaches have been proposed to improve Bitcoin, including those that change the nature of the partial hash collision function [66], those that change the nature of leadership election in the protocol to improve many features of the economics and underlying performance [34] and those that aim to augment the protocol in an effort to achieve scalability [4, 83].

9.4.2 Ethereum

Ethereum was introduced by Vitalik Buterin as a solution to the proliferation of cryptocurrencies that followed Bitcoin, with different varieties of features [11]. Ethereum sought a more pure mandate: to have no features. Instead, Ethereum provides a Turing complete virtual machine, the Ethereum Virtual Machine (EVM), for transaction execution above the consensus, and provides

a means for users to upload code to the EVM that can execute upon the processing of future transactions. So-called *smart contracts* [94] offer the promise of automatically enforced execution of code in a public setting, using strong cryptography and BFT replication. The Ethereum project was successful in one of the largest crowd-funds to date, over \$18 million USD, and the market capitalization of its native token, ether, which is used to pay for transaction execution and code uploads, has since reached \$1 billion USD.

Ethereum currently uses a modified form of Proof-of-Work called Greedy Heaviest Observed Sub Tree (GHOST) [91], but is planning to move to a more secure economic consensus algorithm modeled around Proof of Stake.

9.4.3 Proof-of-Stake

Proof-of-Stake (PoS) was first proposed as an alternative to Proof-of-Work for use in the PPCoin [56]. Under PoS, proposals are made by, and voted on, those who can prove ownership of some stake of coins in the network. While eliminating the excessive energy costs of PoW, naive implementations of PoS are vulnerable to so called “nothing-at-stake” attacks, wherein validators may propose and vote on multiple blocks at a given height, resulting in a dramatic violation of safety, with no incentive to converge. While the problems with naive PoS are well known [82], many popular cryptocurrencies still use it.

The nothing-at-stake problem can be rectified with a mechanism known as *slasher* [12], whereby validators must place a security deposit in order to be eligible to validate blocks, such that the deposit can be *slashed* if the validator is found to propose or vote for conflicting blocks. Tendermint was the first implementation of such an approach, though other BFT algorithms may work as well.

9.4.4 HyperLedger

The success of Bitcoin, Ethereum, and other cryptocurrencies has inspired an increasingly diverse cross section of society, including regulators, bankers, business executives, auditors, account managers, logisticians, and more. In particular, a recent project under the Linux Foundation, spearheaded by IBM and a new blockchain-based company called Digital Asset Holdings (DAH), seeks to provide a unified blockchain architecture for industrial applications. The project is called HyperLedger, after a company with the same name,

providing a rudimentary implementation of a PBFT-based blockchain, was acquired by DAH.

Two contributions to the HyperLedger initiative are particularly relevant. The first is the combination of Juno and Hopper by the team at JP Morgan. Juno is an implementation of Tangaroa, a BFT version of Raft [22], Hopper is a new virtual machine design, based on linear logic [40] and dependent type systems [8], that aims to provide an execution environment for smart contract systems equipped with a formal logic for making and proving statements about the state of the system, or the behaviour of a contract. Both Juno and Hopper are written in Haskell.

The other project is the OpenBlockchain by IBM, a PBFT-based blockchain written in Go, sporting an application state that supports the deployment of arbitrary docker containers. Since an arbitrary docker container may contain non-determinism, their PBFT implementation was modified with additional steps to preserve safety in the face of possibly non-deterministic execution [14].

Another relevant contribution from IBM is a recent review paper, similar in spirit to this chapter [102].

9.4.5 HoneyBadgerBFT

All Paxos-like consensus protocols, including Raft, PBFT, and Tendermint, despite functioning well in asynchronous environments, are not strictly asynchronous. This is because each one uses a timeout somewhere in the protocol, typically to detect faulty leaders. On the other hand, randomized consensus protocols like the common coin offer solutions that work in a fully asynchronous context, with no timeouts.

All consensus protocols rely one way or another on the eventual delivery of messages. The assumption of asynchrony simply states that there is no upper bound on when a message will be delivered. Most of the time, networks act synchronous, in the sense that most messages are delivered within some bound. The difference between a fully asynchronous protocol and one with timeouts is that an asynchronous protocol can *always make progress* during times when the network is behaving synchronously. This point is illustrated clearly in [67], which introduces HoneyBadgerBFT, the first fully asynchronous blockchain design, based on common coin consensus.

An adversary with arbitrary control over the network, and the ability to crash any one node at a time, can cause PBFT to halt for arbitrarily long.

This can be done by crashing the current primary/proposer/leader during times when the network is synchronous, and bringing it back for periods of asynchrony. The network still eventually delivers messages, with some average synchrony, but with precise timing can stop all system progress. The experiment is carried out on PBFT directly in [67], and would work similarly against Tendermint.

HoneyBadgerBFT utilizes a series of cryptographic techniques, including secret sharing, erasure coding, and threshold signatures to design a high performance asynchronous BFT consensus protocol that over comes this problem, on account of not incurring any synchrony assumptions, as it is fully leaderless. However, it requires a trusted dealer for initial setup and for validator changes, and it relies on relatively new cryptographic assumptions about the hardness of certain problems that have yet to withstand the test of time.

9.5 Conclusion

Tendermint emerges from and complements a rich history of consensus science which spans the gamut of synchrony and fault-tolerance assumptions. The invention of the blockchain and of Raft have rekindled the fire in consensus research and spawned a new generation of protocols and software for co-ordination over the Internet.

Chapter 10

Conclusion

Byzantine Fault Tolerant consensus provides a rich basis upon which to build services that do not depend on centralized, trusted parties, and which may be adopted by society to manage critical components of socioeconomic infrastructure. Tendermint, as presented in this thesis, was designed to meet the needs of such systems, and to do so in a way that is understandably secure and easily high performance, and which allows arbitrary systems to have transactions ordered by the consensus protocol, with minimal fuss.

Careful considerations are necessary when deploying a distributed consensus system, especially one without an agreed upon central authority to mediate potential disputes and reset the system in the event of a crisis. Tendermint seeks to address such problems using explicit governance modules and accountability guarantees, enabling integration of Tendermint deployments into modern legal and economic infrastructure.

There is still considerable work to do. This includes formal verification of the algorithm's guarantees, performance optimizations, and architectural changes to enable the system to increase capacity with the addition of machines. And of course, many, many TMSP applications remain to be built.

We hope that this thesis better illuminates some of the problems in distributed consensus and blockchain architecture, and inspires others to build something better.

Bibliography

- [1] *A Deterministic Version of Javascript*. <https://github.com/NodeGuy/Deterministic.js>.
- [2] Samson Abramsky. “Proofs as processes”. In: *Theoretical Computer Science* 135.1 (1994), pp. 5–9.
- [3] M AdelsonVelskii and Evgenii Mikhailovich Landis. *An algorithm for the organization of information*. Tech. rep. DTIC Document, 1963.
- [4] Adam Back et al. “Enabling blockchain innovations with pegged sidechains”. In: (2014).
- [5] Michael Ben-Or. “Another advantage of free choice (extended abstract): Completely asynchronous agreement protocols”. In: *Proceedings of the second annual ACM symposium on Principles of distributed computing*. ACM. 1983, pp. 27–30.
- [6] Daniel J Bernstein. “Curve25519: new Diffie-Hellman speed records”. In: *Public Key Cryptography-PKC 2006*. Springer, 2006, pp. 207–228.
- [7] *Bitcoin Blockchain Charts*. <https://blockchain.info/charts>.
- [8] Ana Bove and Peter Dybjer. “Dependent types at work”. In: *Language engineering and rigorous software development*. Springer, 2009, pp. 57–99.
- [9] Buckie. *Juno - Smart Contracts Running on a BFT Hardened Raft*. <https://github.com/buckie/juno>. 2016.
- [10] Mike Burrows. “The Chubby lock service for loosely-coupled distributed systems”. In: *Proceedings of the 7th symposium on Operating systems design and implementation*. USENIX Association. 2006, pp. 335–350.
- [11] Vitalik Buterin. *Ethereum white paper: a next generation smart contract & decentralized application platform*. 2013.

- [12] Vitalik Buterin. *Slasher: a punitive proof of stake algorithm*. <https://blog.ethereum.org/2014/a-punitive-proof-of-stake-algorithm/>.
- [13] Christian Cachin, Klaus Kursawe, and Victor Shoup. “Random oracles in constantipole: practical asynchronous Byzantine agreement using cryptography”. In: *Proceedings of the nineteenth annual ACM symposium on Principles of distributed computing*. ACM. 2000, pp. 123–132.
- [14] Christian Cachin, Simon Schubert, and Marko Vukolić. “Non-determinism in Byzantine Fault-Tolerant Replication”. In: *arXiv preprint arXiv:1603.07351* (2016).
- [15] Luis Caires and Luca Cardelli. “A spatial logic for concurrency (part I)”. In: *Information and Computation* 186.2 (2003), pp. 194–235.
- [16] Ran Canetti and Tal Rabin. “Fast asynchronous Byzantine agreement with optimal resilience”. In: *Proceedings of the twenty-fifth annual ACM symposium on Theory of computing*. ACM. 1993, pp. 42–51.
- [17] Miguel Castro, Barbara Liskov, et al. “Practical Byzantine fault tolerance”. In: *Proceedings of the Third Symposium on Operating Systems Design and Implementation*. 1999.
- [18] Tushar D Chandra, Robert Griesemer, and Joshua Redstone. “Paxos made live: an engineering perspective”. In: *Proceedings of the twenty-sixth annual ACM symposium on Principles of distributed computing*. ACM. 2007, pp. 398–407.
- [19] Tushar Deepak Chandra and Sam Toueg. “Unreliable failure detectors for reliable distributed systems”. In: *Journal of the ACM (JACM)* 43.2 (1996), pp. 225–267.
- [20] Nikos Chondros, Konstantinos Kokordelis, and Mema Roussopoulos. “On the practicality of practical Byzantine fault tolerance”. In: *Proceedings of ACM/IFIP/USENIX International Middleware Conference (MIDDLEWARE)*. Springer, 2012, pp. 436–455.
- [21] Bram Cohen. *The BitTorrent protocol specification*. 2008.
- [22] Christopher Copeland and Hongxia Zhong. “Tangaroa: a Byzantine Fault Tolerant Raft”. In: ().

- [23] James C Corbett et al. “Spanner: Google’s globally distributed database”. In: *ACM Transactions on Computer Systems (TOCS)* 31.3 (2013), p. 8.
- [24] Nicolas T Courtois and Lear Bahack. “On subversive miner strategies and block withholding attack in bitcoin digital currency”. In: *arXiv preprint arXiv:1402.1718* (2014).
- [25] Martin Davis. *Computability & unsolvability*. Courier Corporation, 1958.
- [26] Giuseppe DeCandia et al. “Dynamo: amazon’s highly available key-value store”. In: *ACM SIGOPS Operating Systems Review*. Vol. 41. 6. ACM. 2007, pp. 205–220.
- [27] Xavier Défago, André Schiper, and Péter Urbán. “Total order broadcast and multicast algorithms: Taxonomy and survey”. In: *ACM Computing Surveys (CSUR)* 36.4 (2004), pp. 372–421.
- [28] Whitfield Diffie, Paul C Van Oorschot, and Michael J Wiener. “Authentication and authenticated key exchanges”. In: *Designs, Codes and cryptography* 2.2 (1992), pp. 107–125.
- [29] Edsger W. Dijkstra. “Hierarchical ordering of sequential processes”. In: *Acta informatica* 1.2 (1971), pp. 115–138.
- [30] Edsger W Dijkstra. “Solution of a problem in concurrent programming control”. In: *Pioneers and Their Contributions to Software Engineering*. Springer, 2001, pp. 289–294.
- [31] Cynthia Dwork, Nancy Lynch, and Larry Stockmeyer. “Consensus in the presence of partial synchrony”. In: *Journal of the ACM (JACM)* 35.2 (1988), pp. 288–323.
- [32] *ETCD Distributed Key-Value Store Source Code Repository*. <https://github.com/coreos/etcd>
- [33] Ittay Eyal and Emin Gün Sirer. “Majority is not enough: Bitcoin mining is vulnerable”. In: *Financial Cryptography and Data Security*. Springer, 2014, pp. 436–454.
- [34] Ittay Eyal et al. “Bitcoin-ng: A scalable blockchain protocol”. In: *arXiv preprint arXiv:1510.02037* (2015).
- [35] Paul Feldman and Silvio Micali. “Optimal algorithms for Byzantine agreement”. In: *Proceedings of the twentieth annual ACM symposium on Theory of computing*. ACM. 1988, pp. 148–161.

- [36] Michael J Fischer, Nancy A Lynch, and Michael Merritt. “Easy impossibility proofs for distributed consensus problems”. In: *Distributed Computing* 1.1 (1986), pp. 26–39.
- [37] Michael J Fischer, Nancy A Lynch, and Michael S Paterson. “Impossibility of distributed consensus with one faulty process”. In: *Journal of the ACM (JACM)* 32.2 (1985), pp. 374–382.
- [38] Luciano Floridi. “On the logical unsolvability of the Gettier problem”. In: *Synthese* 142.1 (2004), pp. 61–79.
- [39] Rui Garcia, Rodrigo Rodrigues, and Nuno Preguiça. “Efficient middleware for byzantine fault tolerant database replication”. In: *Proceedings of the sixth conference on Computer systems*. ACM. 2011, pp. 107–122.
- [40] Jean-Yves Girard. “Linear logic”. In: *Theoretical computer science* 50.1 (1987), pp. 1–101.
- [41] James N Gray. *Notes on data base operating systems*. Springer, 1978.
- [42] Jim Gray et al. “The transaction concept: Virtues and limitations”. In: *VLDB*. Vol. 81. 1981, pp. 144–154.
- [43] Glenn Greenwald. *No place to hide: Edward Snowden, the NSA, and the US surveillance state*. Macmillan, 2014.
- [44] A Nico Habermann. “On a solution and a generalization of the Cigarette Smokers’ Problem”. In: (1972).
- [45] Hashicorp’s Implementation of Raft in Go. <https://github.com/hashicorp/raft>.
- [46] Charles Antony Richard Hoare. *Communicating sequential processes*. Springer, 1978.
- [47] Albert L Hopkins Jr, Jaynarayan H Lala, and T Basil Smith III. “The evolution of fault tolerant computing at the Charles Stark Draper Laboratory, 1955–85”. In: *The Evolution of fault-tolerant computing*. Springer, 1987, pp. 121–140.
- [48] Albert L Hopkins Jr, T Smith III, and Jaynarayan H Lala. “FTMP—a highly reliable fault-tolerant multiprocess for aircraft”. In: *Proceedings of the IEEE* 66.10 (1978), pp. 1221–1239.
- [49] Kenneth Hoyme and Kevin Driscoll. “SAFEbus (for avionics)”. In: *Aerospace and Electronic Systems Magazine, IEEE* 8.3 (1993), pp. 34–39.

- [50] Walter L Hürsch and Cristina Videira Lopes. “Separation of concerns”. In: (1995).
- [51] *InfluxDB: Scalable datastore for metrics, events, and real-time analytics*. <https://github.com/influxdata/influxdb>.
- [52] *JEPSEN - Distributed Systems Safety Analysis*. <http://jepsen.io>.
- [53] *JSON-RPC*. <http://json-rpc.org/>.
- [54] Flavio P Junqueira, Benjamin C Reed, and Marco Serafini. “Zab: High-performance broadcast for primary-backup systems”. In: *Dependable Systems & Networks (DSN), 2011 IEEE/IFIP 41st International Conference on*. IEEE. 2011, pp. 245–256.
- [55] Sunny King and Scott Nadal. “Ppcoin: Peer-to-peer crypto-currency with proof-of-stake”. In: *self-published paper, August 19* (2012).
- [56] Sunny King and Scott Nadal. “Ppcoin: Peer-to-peer crypto-currency with proof-of-stake”. In: *self-published paper, August 19* (2012).
- [57] Ramakrishna Kotla and Mike Dahlin. “High throughput Byzantine fault tolerance”. In: *Dependable Systems and Networks, 2004 International Conference on*. IEEE. 2004, pp. 575–584.
- [58] Ramakrishna Kotla et al. “Zyzyva: speculative byzantine fault tolerance”. In: *ACM SIGOPS Operating Systems Review*. Vol. 41. 6. ACM. 2007, pp. 45–58.
- [59] Leslie Lamport. “The part-time parliament”. In: *ACM Transactions on Computer Systems (TOCS)* 16.2 (1998), pp. 133–169.
- [60] Leslie Lamport. “Time, clocks, and the ordering of events in a distributed system”. In: *Communications of the ACM* 21.7 (1978), pp. 558–565.
- [61] Leslie Lamport, Robert Shostak, and Marshall Pease. “The Byzantine generals problem”. In: *ACM Transactions on Programming Languages and Systems (TOPLAS)* 4.3 (1982), pp. 382–401.
- [62] Arnaud Legout, Guillaume Urvoy-Keller, and Pietro Michiardi. “Rarest first and choke algorithms are enough”. In: *Proceedings of the 6th ACM SIGCOMM conference on Internet measurement*. ACM. 2006, pp. 203–216.
- [63] Steven Levy. *Crypto: How the Code Rebels Beat the Government—Saving Privacy in the Digital Age*. Penguin, 2001.

- [64] Roberto Lucchi and Manuel Mazzara. “A pi-calculus based semantics for WS-BPEL”. In: *The Journal of Logic and Algebraic Programming* 70.1 (2007), pp. 96–118.
- [65] Ralph C Merkle. “A digital signature based on a conventional encryption function”. In: *Advances in Cryptology—CRYPTO’87*. Springer. 1987, pp. 369–378.
- [66] Andrew Miller et al. “Nonoutsourcable Scratch-Off Puzzles to Discourage Bitcoin Mining Coalitions”. In: *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*. ACM. 2015, pp. 680–691.
- [67] Andrew Miller et al. *The Honey Badger of BFT Protocols*. Tech. rep. Cryptology ePrint Archive 2016/199, 2016.
- [68] Robin Milner, Joachim Parrow, and David Walker. “A calculus of mobile processes, i”. In: *Information and computation* 100.1 (1992), pp. 1–40.
- [69] Robin Milner, Joachim Parrow, and David Walker. “Modal logics for mobile processes”. In: *Theoretical Computer Science* 114.1 (1993), pp. 149–171.
- [70] Paul Miner et al. “A unified fault-tolerance protocol”. In: Springer.
- [71] Satoshi Nakamoto. *Bitcoin: A peer-to-peer electronic cash system*. 2008.
- [72] Uwe Nestmann, Rachele Fuzzati, and Massimo Merro. “Modeling consensus in a process calculus”. In: *CONCUR 2003-Concurrency Theory*. Springer, 2003, pp. 399–414.
- [73] Brian M Oki and Barbara H Liskov. “Viewstamped replication: A new primary copy method to support highly-available distributed systems”. In: *Proceedings of the seventh annual ACM Symposium on Principles of distributed computing*. ACM. 1988, pp. 8–17.
- [74] Diego Ongaro. “Consensus: Bridging theory and practice”. PhD thesis. Stanford University, 2014.
- [75] Diego Ongaro and John Ousterhout. “In search of an understandable consensus algorithm”. In: *2014 USENIX Annual Technical Conference (USENIX ATC 14)*. 2014, pp. 305–319.
- [76] *OpenBlockChain: Blockchain Fabric Code*. <https://github.com/openblockchain/obc-peer>.

- [77] *OpenSSL Vulnerabilities*. <https://www.openssl.org/news/vulnerabilities.html>.
- [78] Marshall Pease, Robert Shostak, and Leslie Lamport. “Reaching agreement in the presence of faults”. In: *Journal of the ACM (JACM)* 27.2 (1980), pp. 228–234.
- [79] Riccardo Petrocco, Johan Pouwelse, and Dick HJ Epema. “Performance analysis of the libswift p2p streaming protocol”. In: *Peer-to-Peer Computing (P2P), 2012 IEEE 12th International Conference on*. IEEE. 2012, pp. 103–114.
- [80] Andrew Phillips and Luca Cardelli. “Efficient, correct simulation of biological processes in the stochastic pi-calculus”. In: *Computational methods in systems biology*. Springer. 2007, pp. 184–199.
- [81] Rob Pike. “The Go Programming Language”. In: *Talk given at Google’s Tech Talks* (2009).
- [82] Andrew Poelstra et al. *Distributed Consensus from Proof of Stake is Impossible*. 2014.
- [83] Joseph Poon and Thaddeus Dryja. *The bitcoin lightning network: Scalable off-chain instant payments*. Tech. rep. Technical Report (draft). <https://lightning.network>, 2015.
- [84] Eric A Posner and E Glen Weyl. “Quadratic voting as efficient corporate governance”. In: *University of Chicago Law Review, Forthcoming* (2013).
- [85] Michael O Rabin. “Randomized byzantine generals”. In: *Foundations of Computer Science, 1983., 24th Annual Symposium on*. IEEE. 1983, pp. 403–409.
- [86] Ronan Ryan. “Beyond Flash Boys: Improving Transparency and Fairness in Financial Markets”. In: *CFA Institute Conference Proceedings Quarterly*. Vol. 32. 4. CFA Institute. 2015, pp. 10–17.
- [87] Fred B Schneider. “Implementing fault-tolerant services using the state machine approach: A tutorial”. In: *ACM Computing Surveys (CSUR)* 22.4 (1990), pp. 299–319.
- [88] Adi Shamir. “How to share a secret”. In: *Communications of the ACM* 22.11 (1979), pp. 612–613.
- [89] *Share Memory By Communicating*. <https://blog.golang.org/share-memory-by-communicating>.

- [90] Dale Skeen and Michael Stonebraker. “A formal model of crash recovery in a distributed system”. In: *Software Engineering, IEEE Transactions on* 3 (1983), pp. 219–228.
- [91] Yonatan Sompolinsky and Aviv Zohar. “Secure high-rate transaction processing in Bitcoin”. In: *Financial Cryptography and Data Security*. Springer, 2015, pp. 507–527.
- [92] Colin Stirling and David Walker. “Local model checking in the modal mu-calculus”. In: *Theoretical Computer Science* 89.1 (1991), pp. 161–177.
- [93] Paul Syverson. “A taxonomy of replay attacks [cryptographic protocols]”. In: *Computer Security Foundations Workshop VII, 1994. CSFW 7. Proceedings*. IEEE. 1994, pp. 187–191.
- [94] Nick Szabo. “Formalizing and securing relationships on public networks”. In: *First Monday* 2.9 (1997).
- [95] Nassim Nicholas Taleb and Constantine Sandis. “The skin in the game heuristic for protection against tail events”. In: *Review of Behavioral Economics* 1 (2014), pp. 1–21.
- [96] *The Raft Consensus Algorithm*. <http://raft.github.io>.
- [97] “The Trust Machine”. In: *The Economist*, 2015.
- [98] Alan Mathison Turing. “On computable numbers, with an application to the Entscheidungsproblem”. In: *J. of Math* 58.345–363 (1936), p. 5.
- [99] Robbert Van Renesse, Nicolas Schiper, and Fred B Schneider. “Vive la différence: Paxos vs. viewstamped replication vs. zab”. In: *Dependable and Secure Computing, IEEE Transactions on* 12.4 (2015), pp. 472–484.
- [100] Kenton Varda. “Protocol buffers: Google’s data interchange format”. In: *Google Open Source Blog, Available at least as early as Jul* (2008).
- [101] Hugo Vieira, Luís Caires, and Ruben Viegas. “The spatial logic model checker user’s manual”. In: (2004).
- [102] Marko Vukolic. “The quest for scalable blockchain fabric: Proof-of-work vs. BFT replication”. In: *Proc. IFIP WG 11.4 Workshop on Open Research Problems in Network Security (iNetSec 2015)*.

- [103] John H Wensley et al. “SIFT: Design and analysis of a fault-tolerant computer for aircraft control”. In: *Proceedings of the IEEE* 66.10 (1978), pp. 1240–1255.
- [104] James R Wilcox et al. “Verdi: A framework for implementing and formally verifying distributed systems”. In: *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*. ACM. 2015, pp. 357–368.
- [105] Gavin Wood. “Ethereum: A secure decentralised generalised transaction ledger”. In: *Ethereum Project Yellow Paper* (2014).
- [106] Doug Woos et al. “Planning for change in a formal verification of the raft consensus protocol”. In: *Proceedings of the 5th ACM SIGPLAN Conference on Certified Programs and Proofs*. ACM. 2016, pp. 154–165.
- [107] Jian Yin et al. “Separating agreement from execution for byzantine fault tolerant services”. In: *ACM SIGOPS Operating Systems Review*. Vol. 37. 5. ACM. 2003, pp. 253–267.
- [108] Paul J Zak and Stephen Knack. “Trust and growth”. In: *The economic journal* 111.470 (2001), pp. 295–321.
- [109] Vlad Zamfir. *Introducing Casper “the Friendly Ghost”*. <https://blog.ethereum.org/2015/08/casper-friendly-ghost/>.