

BFT Consensus

From Academic Paper to Mainnet

Alberto Sonnino

Alberto Sonnino

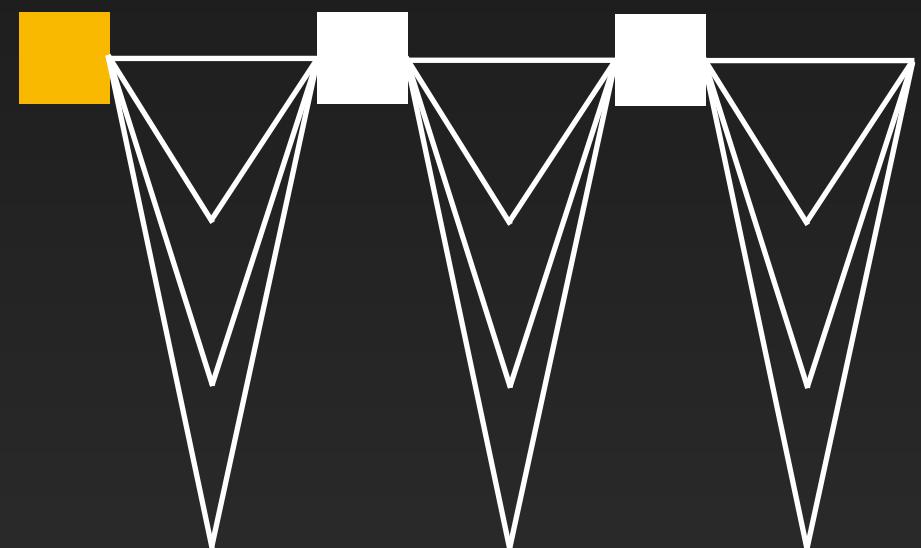
Research Scientist

- PhD from UCL (George Danezis & Jens Groth)
- Co-founded Chainspace
- At Libra / Diem from day 1
- Now building the Sui blockchain

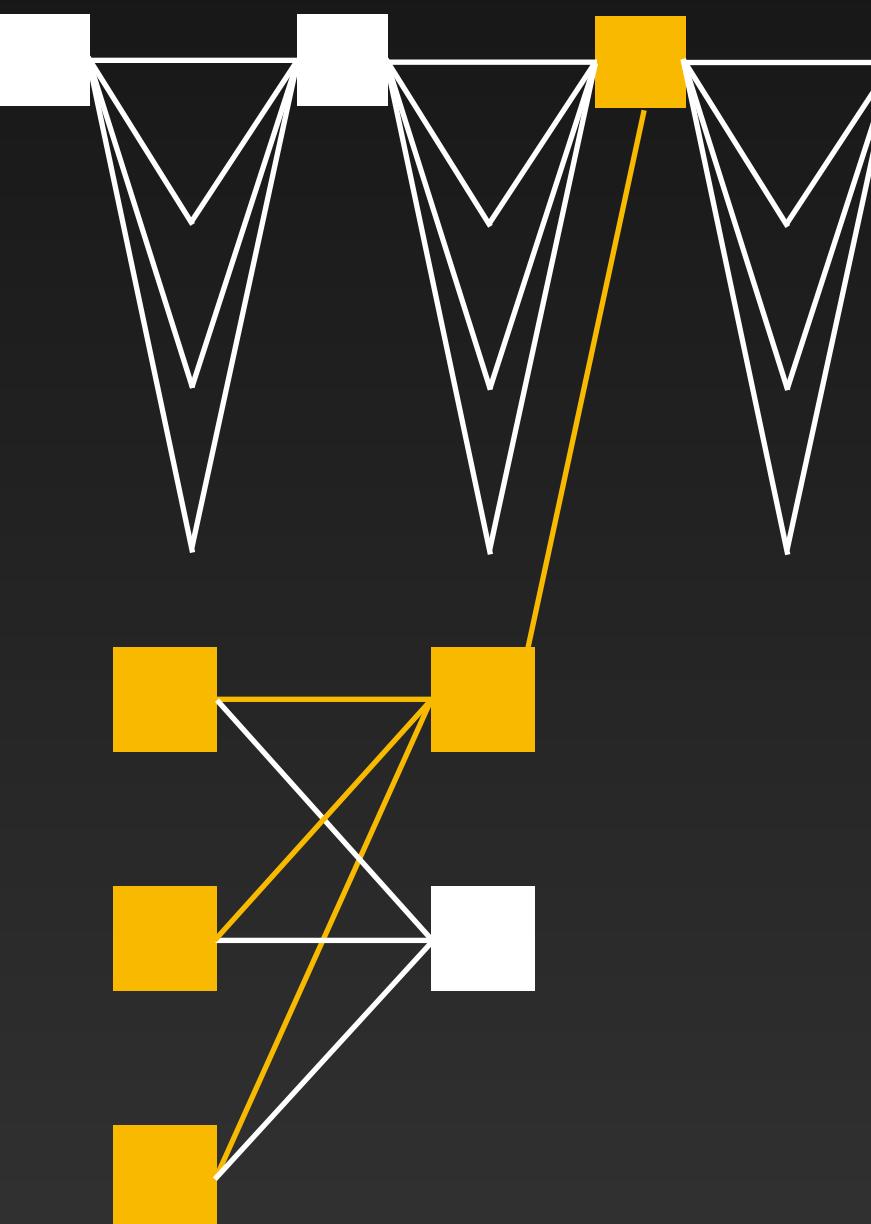
2019



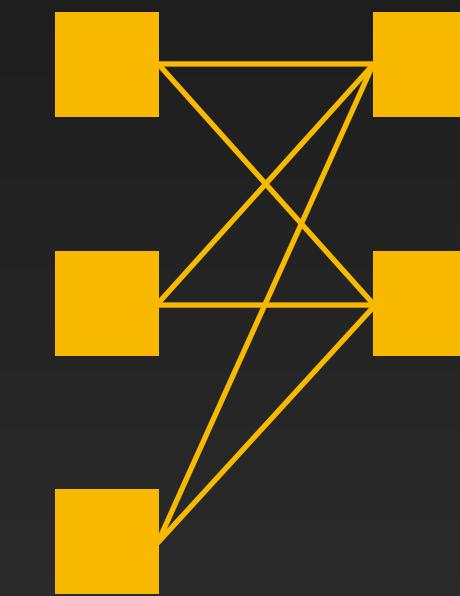
2024



HotStuff



HotStuff + Mempool

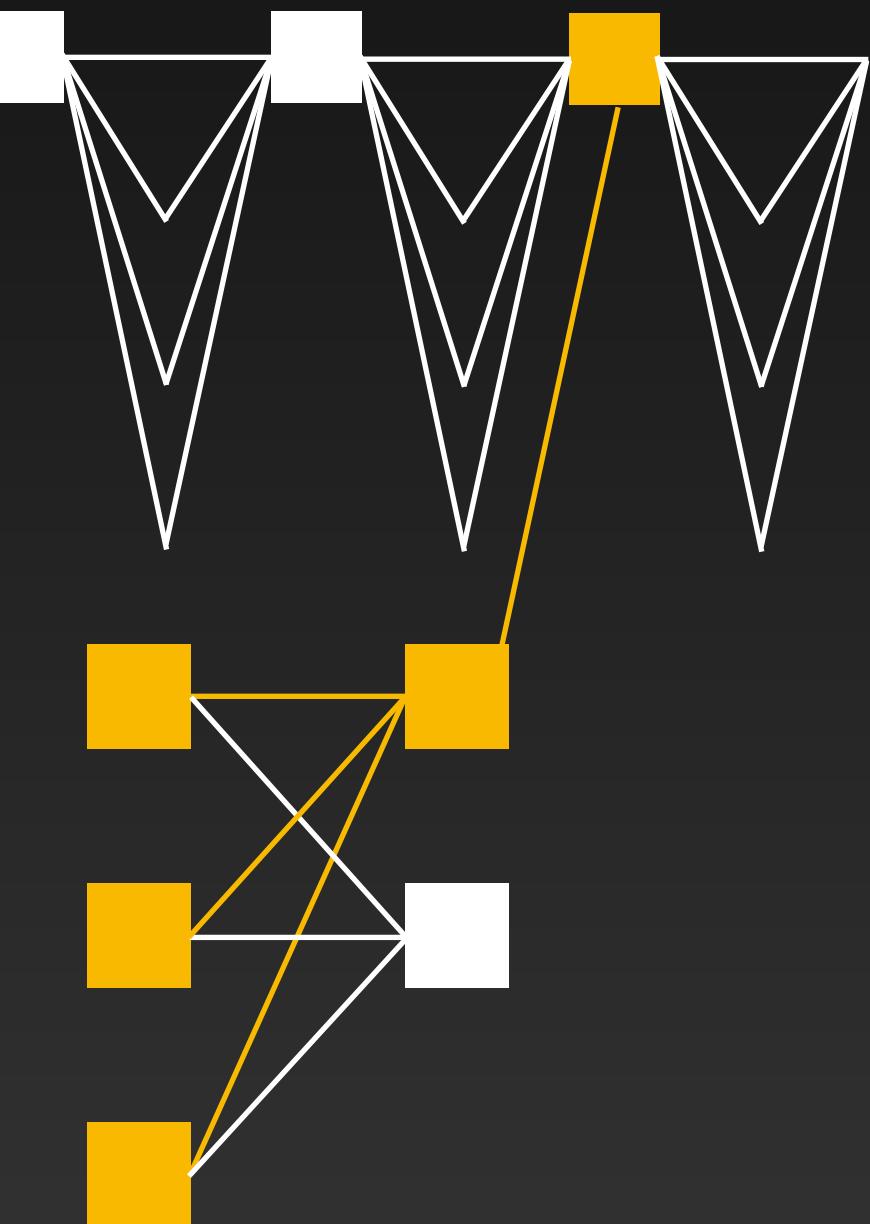
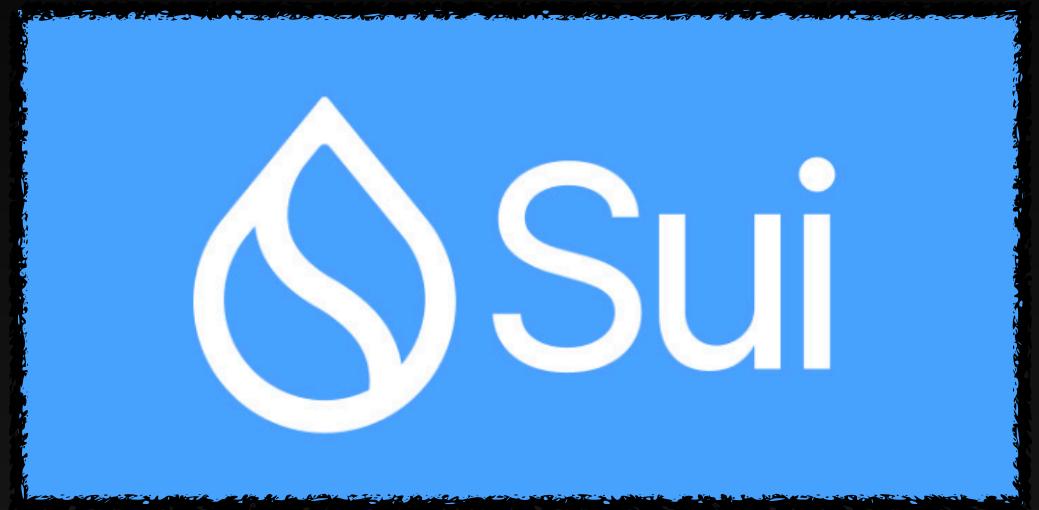


**Bullshark,
Mysticeti**

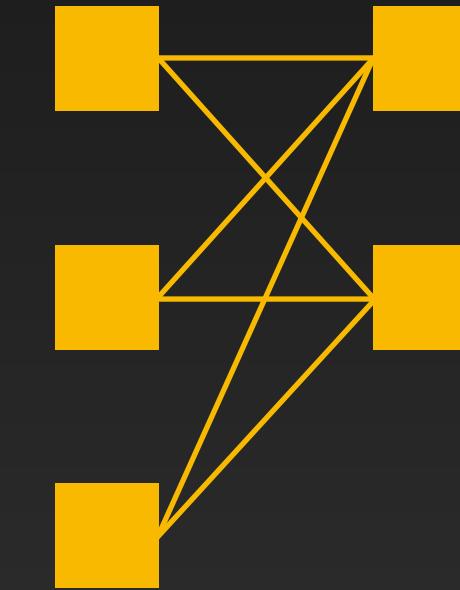
2019



2024



- Lessons learned
- Open research challenges



Research Gifts



(please keep it short)

Byzantine Fault Tolerance



Byzantine Fault Tolerance



$> 2/3$



Partial Synchrony

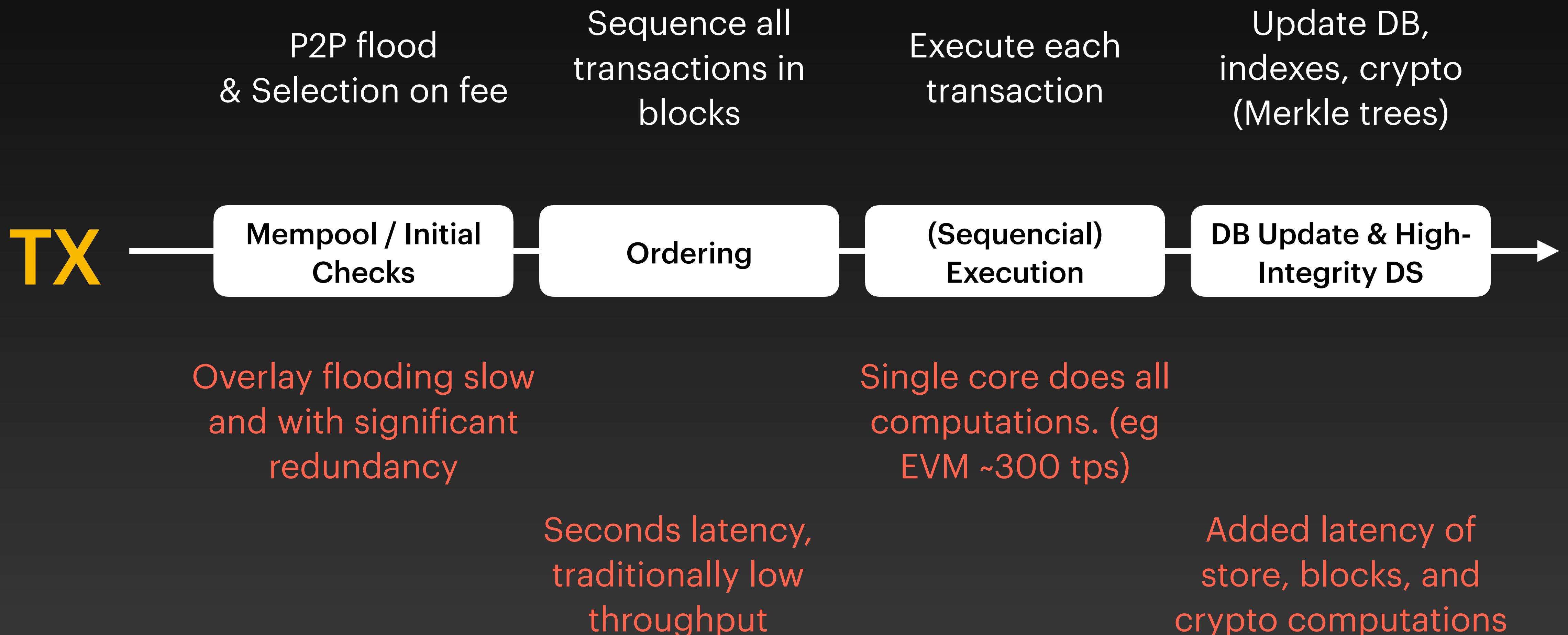


Research Questions

1. Network model?

Lessons Learned

Typical Blockchain



Typical Blockchain

**P2P flood
& Selection on fee**

**Sequence all
transactions in
blocks**

**Mempool / Initial
Checks**

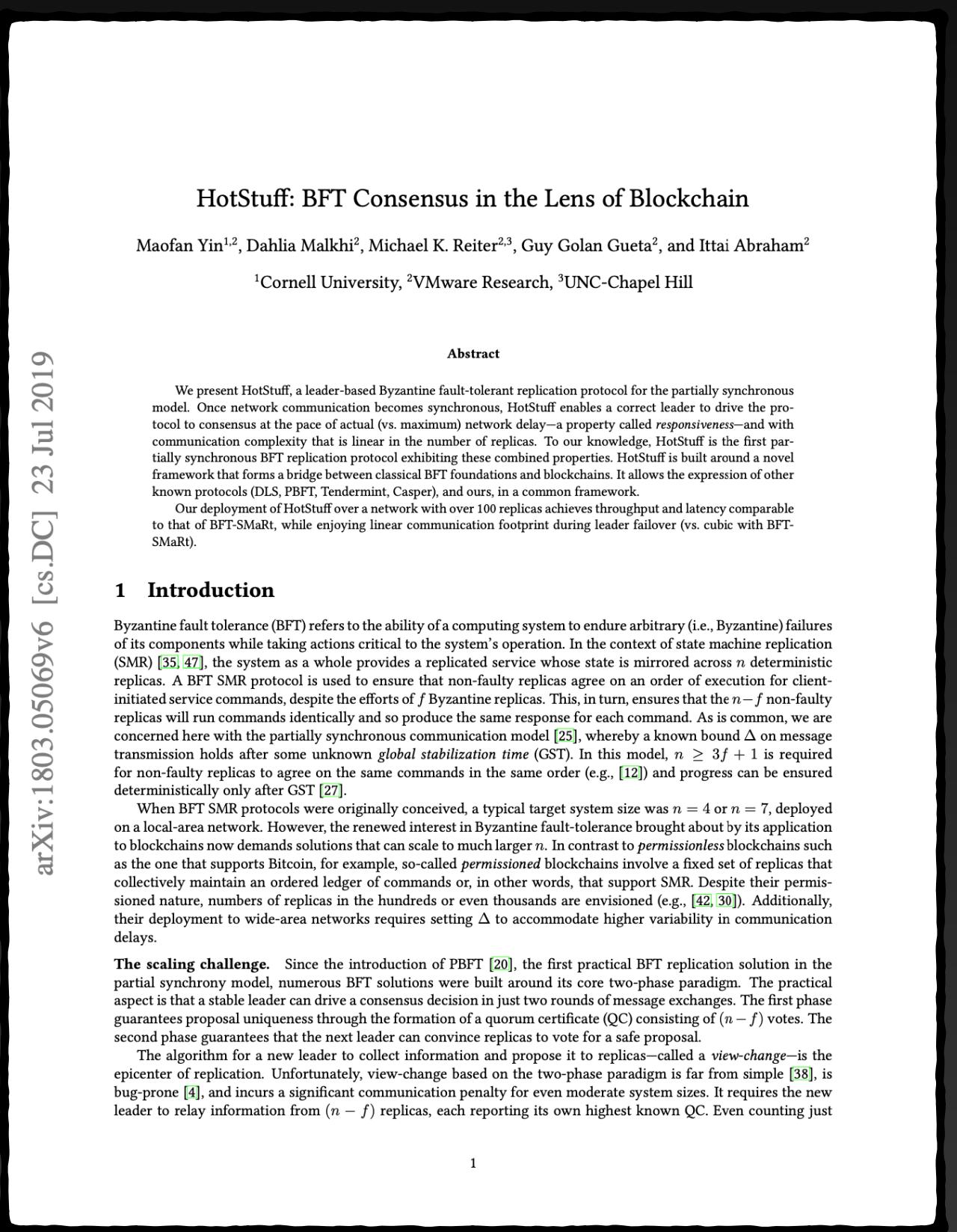
Ordering

**Overlay flooding
slow and with
significant
redundancy**

**Seconds latency,
traditionally low
throughput**

Libra, 2019

HotStuff



arXiv:1803.05069v6 [cs.DC] 23 Jul 2019

HotStuff: BFT Consensus in the Lens of Blockchain

Maofan Yin^{1,2}, Dahlia Malkhi², Michael K. Reiter^{2,3}, Guy Golan Gueta², and Ittai Abraham²

¹Cornell University, ²VMware Research, ³UNC-Chapel Hill

Abstract

We present HotStuff, a leader-based Byzantine fault-tolerant replication protocol for the partially synchronous model. Once network communication becomes synchronous, HotStuff enables a correct leader to drive the protocol to consensus at the pace of actual (vs. maximum) network delay—a property called *responsiveness*—and with communication complexity that is linear in the number of replicas. To our knowledge, HotStuff is the first partially synchronous BFT replication protocol exhibiting these combined properties. HotStuff is built around a novel framework that forms a bridge between classical BFT foundations and blockchains. It allows the expression of other known protocols (DLS, PBFT, Tendermint, Casper), and ours, in a common framework.

Our deployment of HotStuff over a network with over 100 replicas achieves throughput and latency comparable to that of BFT-SMaRt, while enjoying linear communication footprint during leader failover (vs. cubic with BFT-SMaRt).

1 Introduction

Byzantine fault tolerance (BFT) refers to the ability of a computing system to endure arbitrary (i.e., Byzantine) failures of its components while taking actions critical to the system's operation. In the context of state machine replication (SMR) [35, 47], the system as a whole provides a replicated service whose state is mirrored across n deterministic replicas. A BFT SMR protocol is used to ensure that non-faulty replicas agree on an order of execution for client-initiated service commands, despite the efforts of f Byzantine replicas. This, in turn, ensures that the $n - f$ non-faulty replicas will run commands identically and so produce the same response for each command. As is common, we are concerned here with the partially synchronous communication model [25], whereby a known bound Δ on message transmission holds after some unknown *global stabilization time* (GST). In this model, $n \geq 3f + 1$ is required for non-faulty replicas to agree on the same commands in the same order (e.g., [12]) and progress can be ensured deterministically only after GST [27].

When BFT SMR protocols were originally conceived, a typical target system size was $n = 4$ or $n = 7$, deployed on a local-area network. However, the renewed interest in Byzantine fault-tolerance brought about by its application to blockchains now demands solutions that can scale to much larger n . In contrast to *permissionless* blockchains such as the one that supports Bitcoin, for example, so-called *permissioned* blockchains involve a fixed set of replicas that collectively maintain an ordered ledger of commands or, in other words, that support SMR. Despite their permissioned nature, numbers of replicas in the hundreds or even thousands are envisioned (e.g., [42, 30]). Additionally, their deployment to wide-area networks requires setting Δ to accommodate higher variability in communication delays.

The scaling challenge. Since the introduction of PBFT [20], the first practical BFT replication solution in the partial synchrony model, numerous BFT solutions were built around its core two-phase paradigm. The practical aspect is that a stable leader can drive a consensus decision in just two rounds of message exchanges. The first phase guarantees proposal uniqueness through the formation of a quorum certificate (QC) consisting of $(n - f)$ votes. The second phase guarantees that the new leader can convince replicas to vote for a safe proposal.

The algorithm for a new leader to collect information and propose it to replicas—called a *view-change*—is the epicenter of replication. Unfortunately, view-change based on the two-phase paradigm is far from simple [38], is bug-prone [4], and incurs a significant communication penalty for even moderate system sizes. It requires the new leader to relay information from $(n - f)$ replicas, each reporting its own highest known QC. Even counting just

1



HashGraph



arXiv:2102.01167v1 [cs.LO] 1 Feb 2021

Verifying the Hashgraph Consensus Algorithm

Karl Crary
Carnegie Mellon University

Abstract

The Hashgraph consensus algorithm is an algorithm for asynchronous Byzantine fault tolerance intended for distributed shared ledgers. Its main distinguishing characteristic is it achieves consensus without exchanging any extra messages; each participant's votes can be determined from public information, so votes need not be transmitted.

In this paper, we discuss our experience formalizing the Hashgraph algorithm and its correctness proof using the Coq proof assistant. The paper is self-contained; it includes a complete discussion of the algorithm and its correctness argument in English.

1 Introduction

Byzantine fault-tolerance is the problem of coordinating a distributed system while some participants may maliciously break the rules. Often other challenges are also present, such as a lack of communication. The challenge is the construction of a variety of new applications, such as cryptocurrencies. Such applications rely on *distributed shared ledgers*, a form of Byzantine fault-tolerance in which a set of transactions are placed in a globally-agreed total order that is *immutable*. The latter means that once a transaction enters the order, no new transaction can enter at an earlier position.

A distributed shared ledger makes it possible for all participants to agree, at any point in the order, on the current owner of a digital commodity such as a unit of cryptocurrency. A transaction transferring ownership is valid if the commodity's current owner authorizes the transaction. (The authorization mechanism—presumably using a digital signature—is beyond the scope of the ledger itself.) Because the order is total, one transaction out of any pair has priority. Thus we can show that a commodity's chain of ownership is uniquely determined. Finally, because the order is immutable, the chain of ownership cannot change except by adding new transactions at the end.

Algorithmic Byzantine consensus (under various assumptions) have existed for some time, indeed longer than the protocol has been named [12, 9]. Practical algorithms are more recent; in 1999, Castro and Liskov [6] gave an algorithm that when installed into the NFS file system slowed it only 3%. As Byzantine consensus algorithms have become more practical, they have been tailored to specific applications. Castro and Liskov's algorithm was designed for fault-tolerant state machine replication [13] and probably would not perform well under the workload of a distributed shared ledger.

However, in the last few years there have arisen Byzantine fault-tolerance algorithms suitable for distributed shared ledgers, notably HoneyBadgerBFT [10], BEAT [7], and—the subject of this paper—Hashgraph [2]. Moreover, the former two each claim to be the first practical *asynchronous* BFT algorithm (with different standards of practicality). Hashgraph does not claim to be first, but is also practical and asynchronous.

In parallel with that line of work has been the development of distributed shared ledgers based on *proof of work*, beginning with Bitcoin [11]. The idea behind proof of work is to maintain agreement on the ledger by maintaining a list of blocks of transactions, and to ensure that the list does not become a tree. To ensure this, the rules state that (1) the longest branch defines the list, and (2) to create a new block, one must solve a difficult computational problem that takes the last old header as part of its input. The problem's solution is much easier to verify than to obtain, so when one learns of a new block, one's incentive is to restart work from the new head rather than continue work from the old head.

Bitcoin and some of its cousins are widely used, so in a certain sense they are indisputably practical. They are truly permissionless, in a way that the BFT algorithms, including Hashgraph, cannot quite claim. Nevertheless, they offer severely limited throughput. Bitcoin is limited to seven transactions per second and has a latency of one hour, while its BFT competitors all do several orders of magnitude better. Proof-of-work systems are also criticized for being wasteful: an enormous amount of electricity is expended on block-creation efforts that nearly always fail. Finally—more to the point of this paper—the theoretical properties of proof of work are not well understood.

The Hashgraph consensus algorithm is designed to support high-performance applications of a distributed shared ledger. Like the other BFT systems, it is several orders of magnitude faster than proof of work. Actual performance depends very much on configuration choices (e.g., how many peers, geographic distribution, tradeoff between latency and throughput, etc.), but in all configurations published in Miller, et. al [10] (for HoneyBadgerBFT) and Duan, et al. [7] (for BEAT), the Hashgraph algorithm equals or exceeds the published performance figures [4]. A frequently cited throughput goal is to equal the Visa credit-card network. According to Visa's published figures, Hashgraph can

Too much impact?

- Patent your work
- Send the patent around
- Ask companies to cite your patented work (ideally in public)

Libra, 2019

arXiv:1803.05069v6 [cs.DC] 23 Jul 2019

HotStuff: BFT Consensus in the Lens of Blockchain

Maofan Yin^{1,2}, Dahlia Malkhi², Michael K. Reiter^{2,3}, Guy Golan Gueta², and Ittai Abraham²

¹Cornell University, ²VMware Research, ³UNC-Chapel Hill

Abstract

We present HotStuff, a leader-based Byzantine fault-tolerant replication protocol for the partially synchronous model. Once network communication becomes synchronous, HotStuff enables a correct leader to drive the protocol to consensus at the pace of actual (vs. maximum) network delay—a property called *responsiveness*—and with communication complexity that is linear in the number of replicas. To our knowledge, HotStuff is the first partially synchronous BFT replication protocol exhibiting these combined properties. HotStuff is built around a novel framework that forms a bridge between classical BFT foundations and blockchains. It allows the expression of other known protocols (DLS, PBFT, Tendermint, Casper), and ours, in a common framework.

Our deployment of HotStuff over a network with over 100 replicas achieves throughput and latency comparable to that of BFT-SMaRt, while enjoying linear communication footprint during leader failover (vs. cubic with BFT-SMaRt).

1 Introduction

Byzantine fault tolerance (BFT) refers to the ability of a computing system to endure arbitrary (i.e., Byzantine) failures of its components while taking actions critical to the system's operation. In the context of state machine replication (SMR) [35, 47], the system as a whole provides a replicated service whose state is mirrored across n deterministic replicas. A BFT SMR protocol is used to ensure that non-faulty replicas agree on an order of execution for client-initiated service commands, despite the efforts of f Byzantine replicas. This, in turn, ensures that the $n - f$ non-faulty replicas will run commands identically and so produce the same response for each command. As is common, we are concerned here with the partially synchronous communication model [25], whereby a known bound Δ on message transmission holds after some unknown *global stabilization time* (GST). In this model, $n \geq 3f + 1$ is required for non-faulty replicas to agree on the same commands in the same order (e.g., [12]) and progress can be ensured deterministically only after GST [27].

When BFT SMR protocols were originally conceived, a typical target system size was $n = 4$ or $n = 7$, deployed on a local-area network. However, the renewed interest in Byzantine fault-tolerance brought about by its application to blockchains now demands solutions that can scale to much larger n . In contrast to *permissionless* blockchains such as the one that supports Bitcoin, for example, so-called *permissioned* blockchains involve a fixed set of replicas that collectively maintain an ordered ledger of commands or, in other words, that support SMR. Despite their permissioned nature, numbers of replicas in the hundreds or even thousands are envisioned (e.g., [42, 30]). Additionally, their deployment to wide-area networks requires setting Δ to accommodate higher variability in communication delays.

The scaling challenge. Since the introduction of PBFT [20], the first practical BFT replication solution in the partial synchrony model, numerous BFT solutions were built around its core two-phase paradigm. The practical aspect is that a stable leader can drive a consensus decision in just two rounds of message exchanges. The first phase guarantees proposal uniqueness through the formation of a quorum certificate (QC) consisting of $(n - f)$ votes. The second phase guarantees that the next leader can convince replicas to vote for a safe proposal.

The algorithm for a new leader to collect information and propose it to replicas—called a *view-change*—is the epicenter of replication. Unfortunately, view-change based on the two-phase paradigm is far from simple [38], is bug-prone [4], and incurs a significant communication penalty for even moderate system sizes. It requires the new leader to relay information from $(n - f)$ replicas, each reporting its own highest known QC. Even counting just

1

HotStuff

✓ Linear

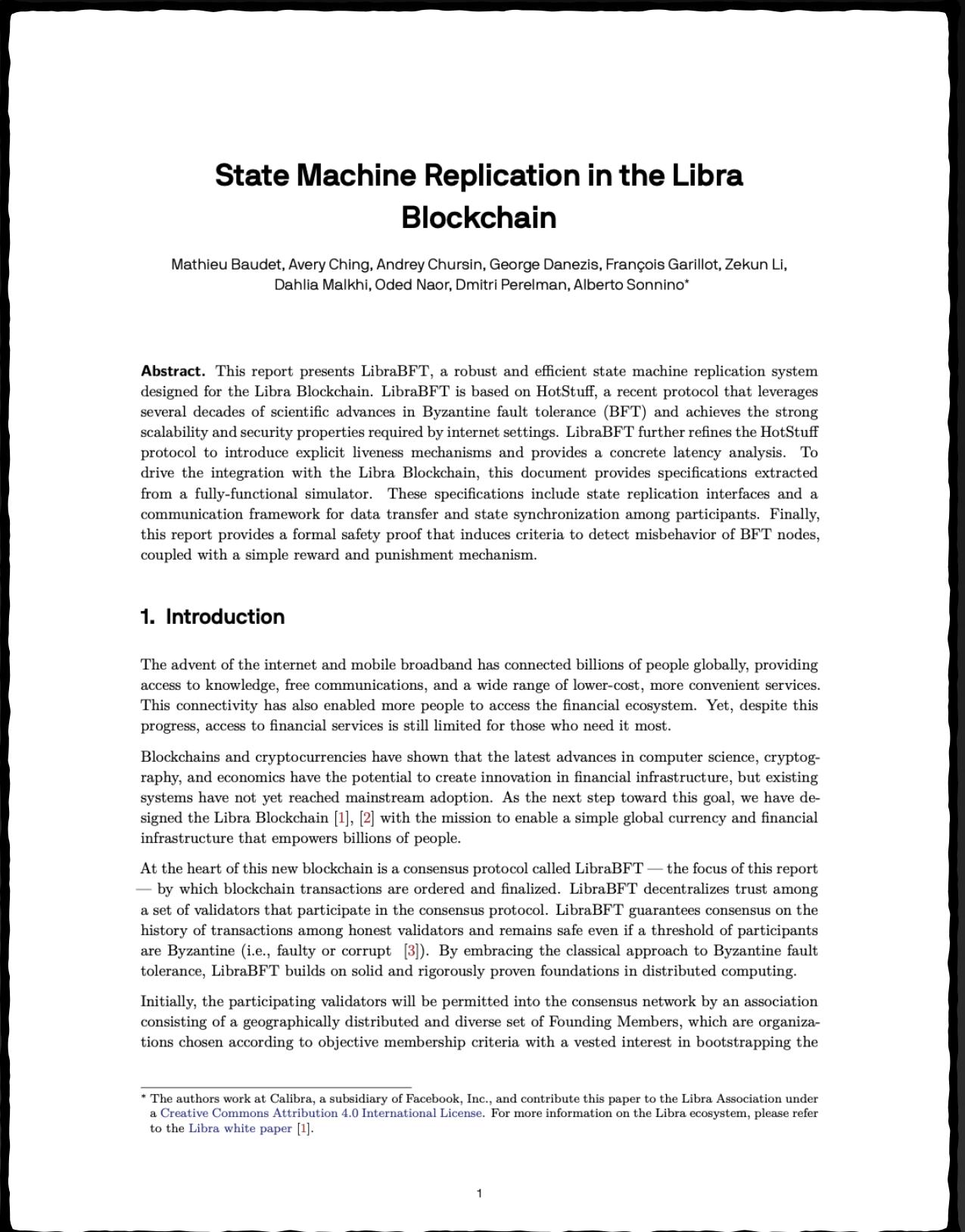
✓ Clearly isolated components

HashGraph

✗ Impossible to garbage collect

✗ Unclear block synchroniser

The first 6 months...



State Machine Replication in the Libra Blockchain

Mathieu Baudet, Avery Ching, Andrey Chursin, George Danezis, François Garillot, Zekun Li, Dahlia Malkhi, Oded Naor, Dmitri Perelman, Alberto Sonnino*

Abstract. This report presents LibraBFT, a robust and efficient state machine replication system designed for the Libra Blockchain. LibraBFT is based on HotStuff, a recent protocol that leverages several decades of scientific advances in Byzantine fault tolerance (BFT) and achieves the strong scalability and security properties required by internet settings. LibraBFT further refines the HotStuff protocol to introduce explicit liveness mechanisms and provides a concrete latency analysis. To drive the integration with the Libra Blockchain, this document provides specifications extracted from a fully-functional simulator. These specifications include state replication interfaces and a communication framework for data transfer and state synchronization among participants. Finally, this report provides a formal safety proof that induces criteria to detect misbehavior of BFT nodes, coupled with a simple reward and punishment mechanism.

1. Introduction

The advent of the internet and mobile broadband has connected billions of people globally, providing access to knowledge, free communications, and a wide range of lower-cost, more convenient services. This connectivity has also enabled more people to access the financial ecosystem. Yet, despite this progress, access to financial services is still limited for those who need it most.

Blockchains and cryptocurrencies have shown that the latest advances in computer science, cryptography, and economics have the potential to create innovation in financial infrastructure, but existing systems have not yet reached mainstream adoption. As the next step toward this goal, we have designed the Libra Blockchain [1], [2] with the mission to enable a simple global currency and financial infrastructure that empowers billions of people.

At the heart of this new blockchain is a consensus protocol called LibraBFT — the focus of this report — by which blockchain transactions are ordered and finalized. LibraBFT decentralizes trust among a set of validators that participate in the consensus protocol. LibraBFT guarantees consensus on the history of transactions among honest validators and remains safe even if a threshold of participants are Byzantine (i.e., faulty or corrupt [3]). By embracing the classical approach to Byzantine fault tolerance, LibraBFT builds on solid and rigorously proven foundations in distributed computing.

Initially, the participating validators will be permitted into the consensus network by an association consisting of a geographically distributed and diverse set of Founding Members, which are organizations chosen according to objective membership criteria with a vested interest in bootstrapping the

* The authors work at Calibra, a subsidiary of Facebook, Inc., and contribute this paper to the Libra Association under a Creative Commons Attribution 4.0 International License. For more information on the Libra ecosystem, please refer to the Libra white paper [1].

SMR in the Libra Blockchain

- The LibraBFT/DiemBFT pacemaker
- Codesign the pacemaker with the rest

Research Questions

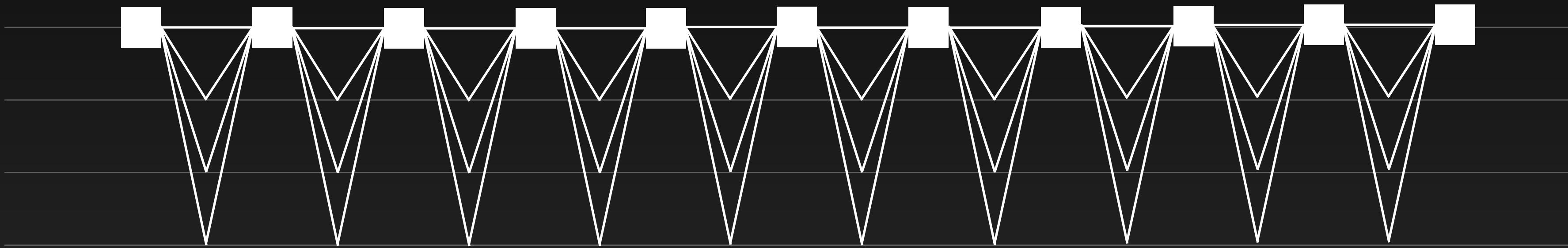
1. Network model?

Lessons Learned

1. Modularisation is a design strategy

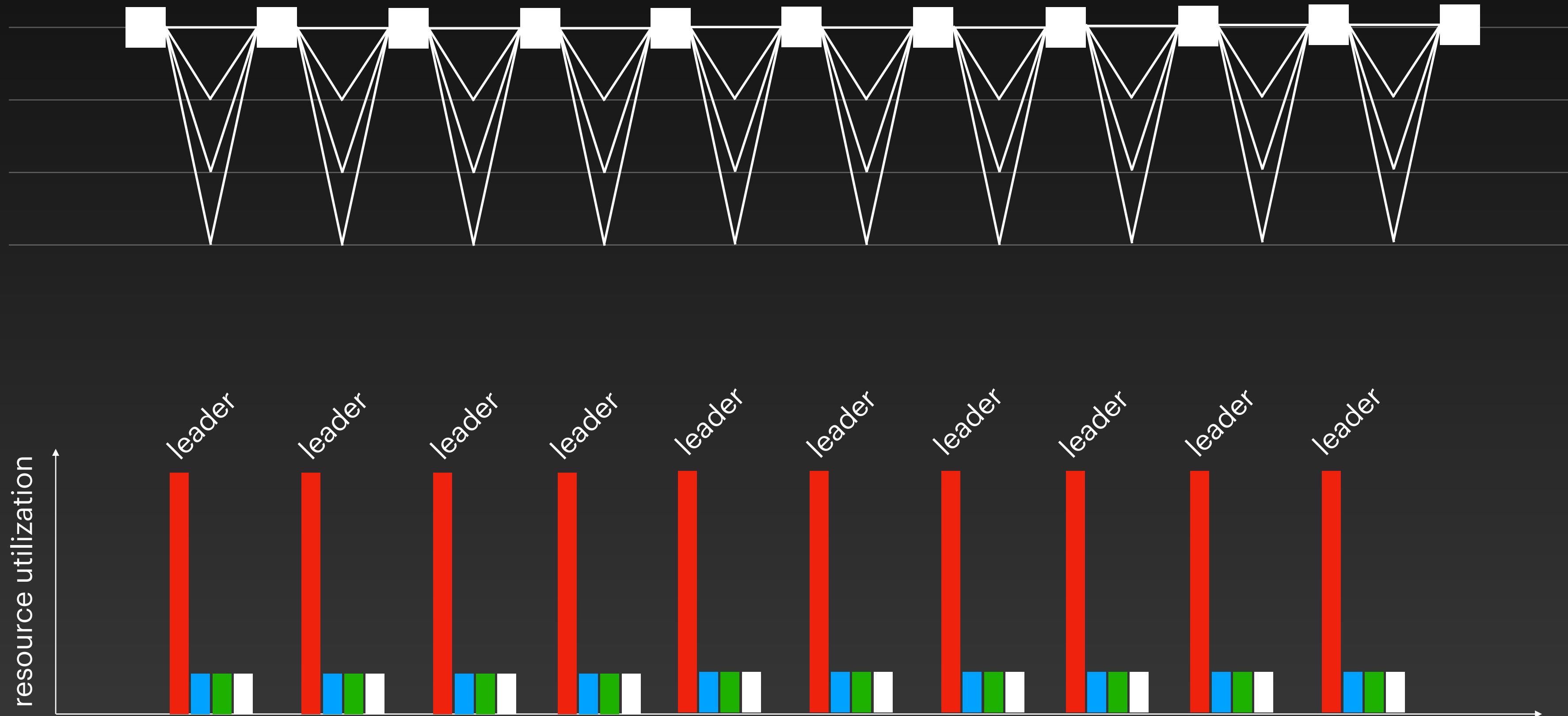
HotStuff

Typical leader-based protocols



HotStuff

Uneven resource utilisation



Research Questions

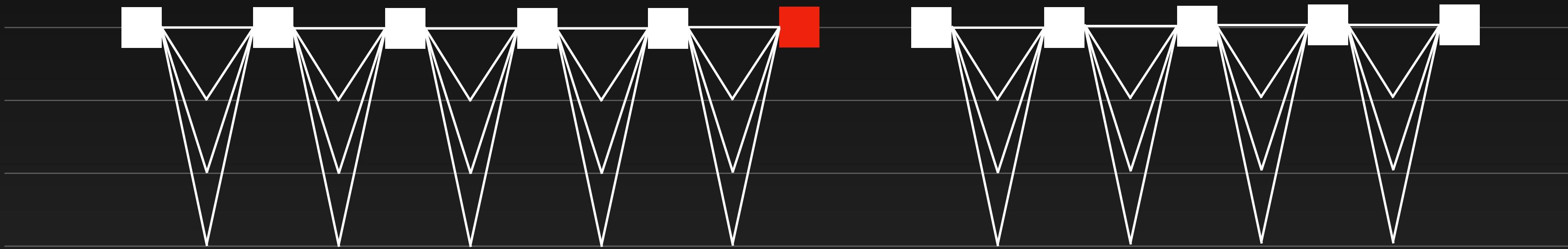
1. Network model?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads allocation

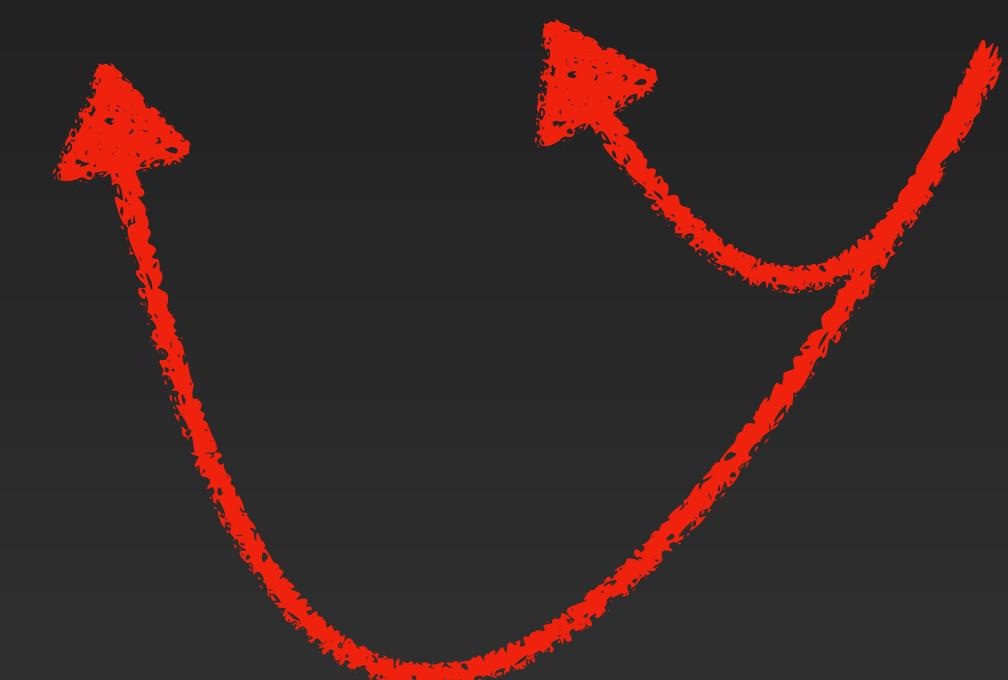
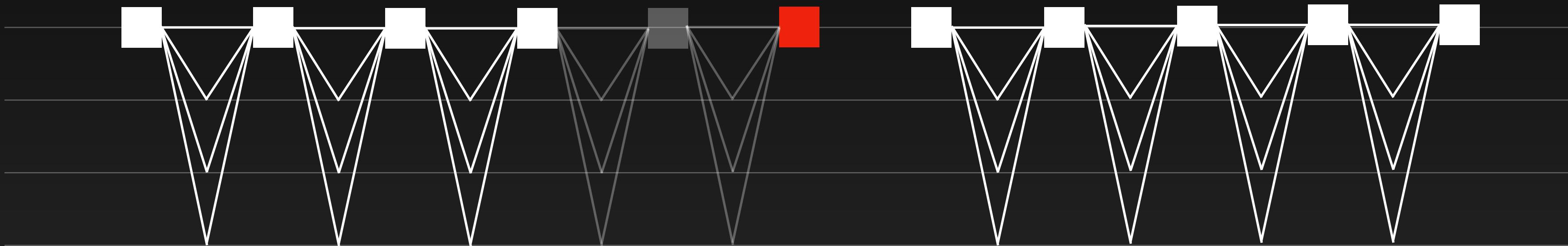
Leader-Driven Consensus

Fragility to faults and asynchrony

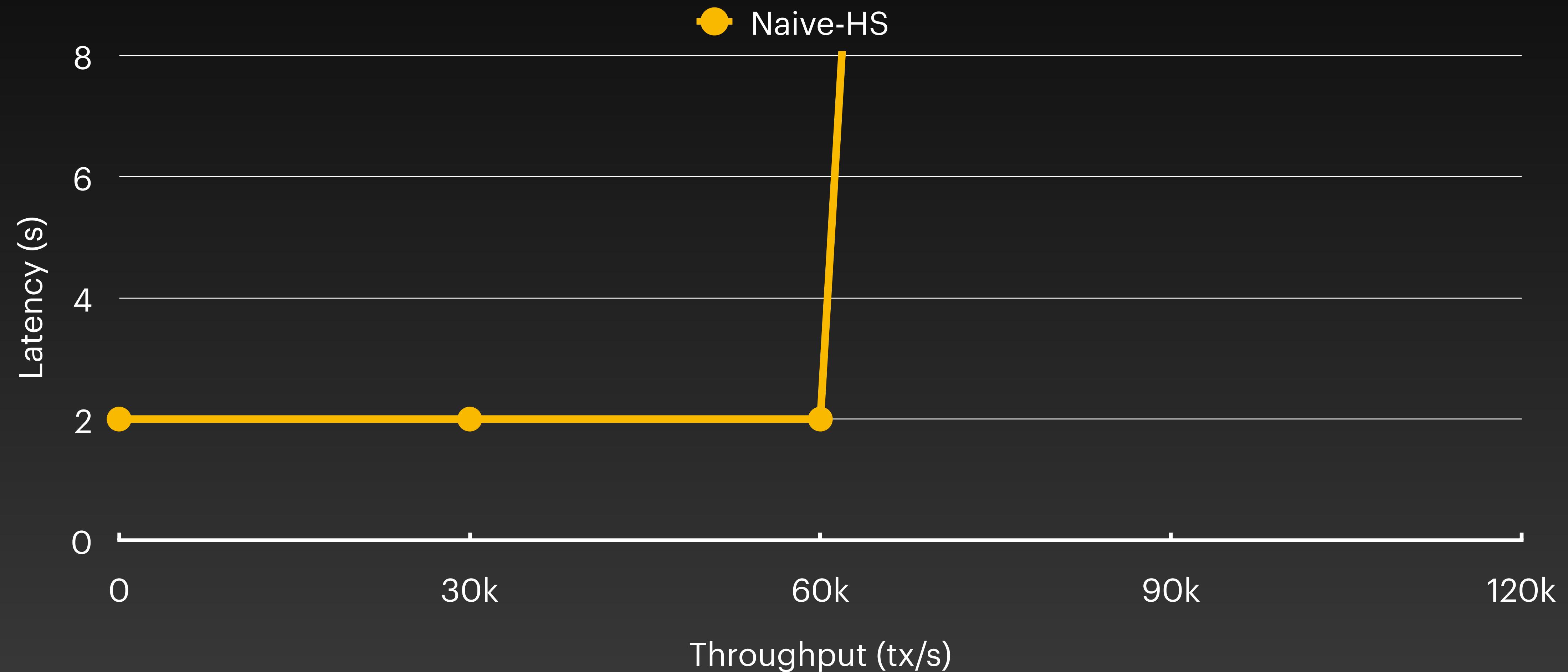


Leader-Driven Consensus

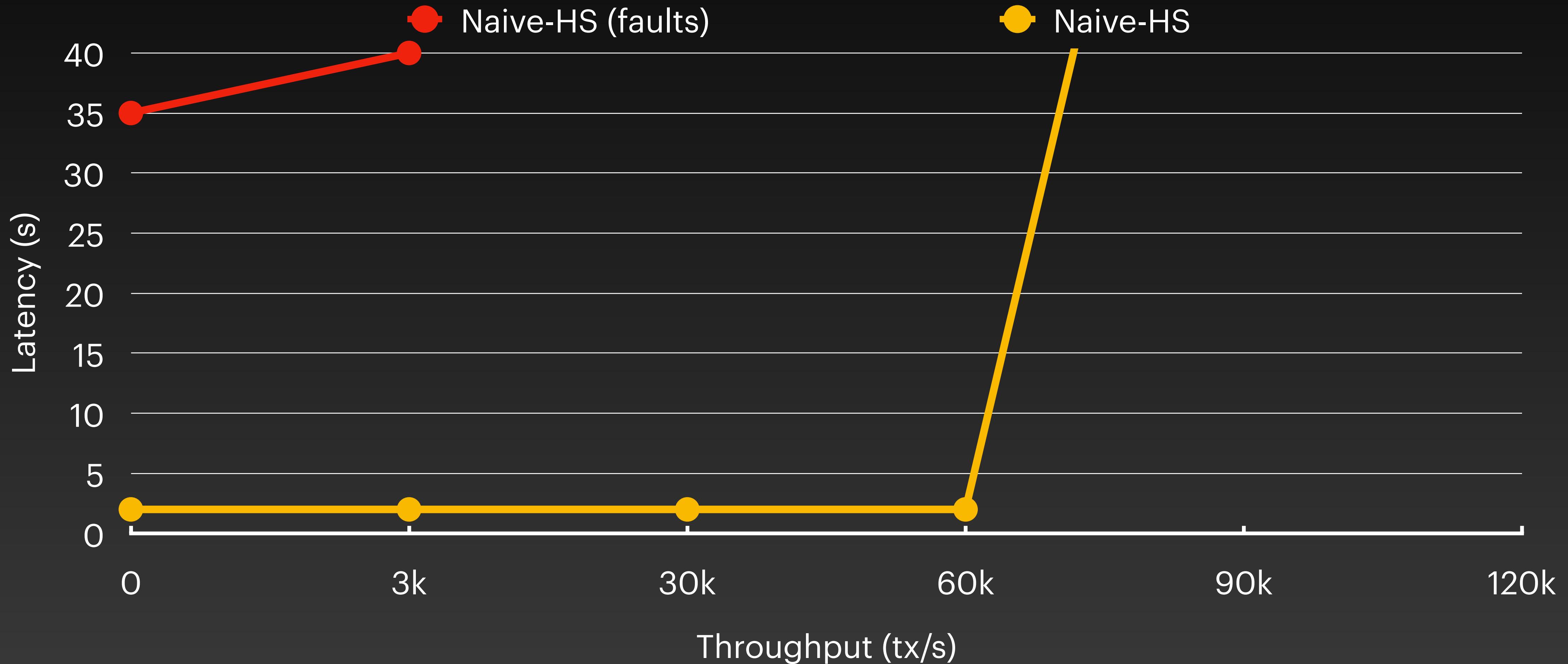
Fragility to faults and asynchrony



Performance



Performance



Research Questions

1. Network model?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early

Libra, 2019

HotStuff: BFT Consensus in the Lens of Blockchain

Maofan Yin^{1,2}, Dahlia Malkhi², Michael K. Reiter^{2,3}, Guy Golan Gueta², and Ittai Abraham²

¹Cornell University, ²VMware Research, ³UNC-Chapel Hill

Abstract

We present HotStuff, a leader-based Byzantine fault-tolerant replication protocol for the partially synchronous model. Once network communication becomes synchronous, HotStuff enables a correct leader to drive the protocol to consensus at the pace of actual (vs. maximum) network delay—a property called *responsiveness*—and with communication complexity that is linear in the number of replicas. To our knowledge, HotStuff is the first partially synchronous BFT replication protocol exhibiting these combined properties. HotStuff is built around a novel framework that forms a bridge between classical BFT foundations and blockchains. It allows the expression of other known protocols (DLS, PBFT, Tendermint, Casper), and ours, in a common framework.

Our deployment of HotStuff over a network with over 100 replicas achieves throughput and latency comparable to that of BFT-SMaRt, while enjoying linear communication footprint during leader failure (vs. cubic with BFT-SMaRt).

1 Introduction

Byzantine fault tolerance (BFT) refers to the ability of a computing system to endure arbitrary (i.e., Byzantine) failures of its components while taking actions critical to the system's operation. In the context of state machine replication (SMR) [35, 47], the system as a whole provides a replicated service whose state is mirrored across n deterministic replicas. A BFT SMR protocol is used to ensure that non-faulty replicas agree on an order of execution for client-initiated service commands, despite the efforts of f Byzantine replicas. This, in turn, ensures that the $n - f$ non-faulty replicas will run commands identically and so produce the same response for each command. As is common, we are concerned here with the partially synchronous communication model [25], whereby a known bound Δ on message transmission holds after some unknown *global stabilization time* (GST). In this model, $n \geq 3f + 1$ is required for non-faulty replicas to agree on the same commands in the same order (e.g., [12]) and progress can be ensured deterministically only after GST [27].

When BFT SMR protocols were originally conceived, a typical target system size was $n = 4$ or $n = 7$, deployed on a local-area network. However, the renewed interest in Byzantine fault-tolerance brought about by its application to blockchains now demands solutions that can scale to much larger n . In contrast to *permissionless* blockchains such as the one that supports Bitcoin, for example, so-called *permissioned* blockchains involve a fixed set of replicas that collectively maintain an ordered ledger of commands or, in other words, that support SMR. Despite their permissioned nature, numbers of replicas in the hundreds or even thousands are envisioned (e.g., [42, 30]). Additionally, their deployment to wide-area networks requires setting Δ to accommodate higher variability in communication delays.

The scaling challenge. Since the introduction of PBFT [20], the first practical BFT replication solution in the partial synchrony model, numerous BFT solutions were built around its core two-phase paradigm. The practical aspect is that a stable leader can drive a consensus decision in just two rounds of message exchanges. The first phase guarantees proposal uniqueness through the formation of a quorum certificate (QC) consisting of $(n - f)$ votes. The second phase guarantees that the next leader can convince replicas to vote for a safe proposal.

The algorithm for a new leader to collect information and propose it to replicas—called a *view-change*—is the epicenter of replication. Unfortunately, view-change based on the two-phase paradigm is far from simple [38], is bug-prone [4], and incurs a significant communication penalty for even moderate system sizes. It requires the new leader to relay information from $(n - f)$ replicas, each reporting its own highest known QC. Even counting just

arXiv:1803.05069v6 [cs.DC] 23 Jul 2019

1

HotStuff (naive mempool)

- Linear
- Clearly isolated components
- Uneven resource utilisation
- Fragile to faults and asynchrony
- Unspecified components (pacemaker)

Libra, 2021

Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus

George Danezis
Mysten Labs & UCL

Alberto Sonnino
Mysten Labs

Abstract
We propose separating the task of reliable transaction dissemination from transaction ordering, to enable high-performance Byzantine fault-tolerant quorum-based consensus. We design and evaluate a mempool protocol, Narwhal, specializing in high-throughput reliable dissemination and storage of causal histories of transactions. Narwhal tolerates an asynchronous network and maintains high performance despite failures. Narwhal is designed to easily scale-out using multiple workers at each validator, and we demonstrate that there is no foreseeable limit to the throughput we can achieve.

Composing Narwhal with a partially synchronous consensus protocol (Narwhal-HotStuff) yields significantly better throughput even in the presence of faults or intermittent loss of liveness due to asynchrony. However, loss of liveness can result in higher latency. To achieve overall good performance when faults occur we design Tusk, a zero-message overhead asynchronous consensus protocol, to work with Narwhal. We demonstrate its high performance under a variety of configurations and faults.

As a summary of results, on a WAN, Narwhal-HotStuff achieves over 130,000 tx/sec at less than 2-sec latency compared with 1,800 tx/sec at 1-sec latency for HotStuff. Additional workers increase throughput linearly to 600,000 tx/sec without any latency increase. Tusk achieves 160,000 tx/sec with about 3 seconds latency. Under faults, both protocols maintain high throughput, but Narwhal-HotStuff suffers from increased latency.

CCS Concepts: Security and privacy → Distributed systems security.

Keywords: Consensus protocol, Byzantine Fault Tolerant

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
EuroSys '22, April 5–8, 2022, Rennes, France
© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9162-7/22/04... \$15.00
<https://doi.org/10.1145/3492321.3519594>

ACM Reference Format:
George Danezis, Lefteris Kokoris-Kogias, Alberto Sonnino, and Alexander Spiegelman. 2022. Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus . In *Seventeenth European Conference on Computer Systems (EuroSys '22), April 5–8, 2022, Rennes, France*. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3492321.3519594>

1 Introduction
Byzantine consensus protocols [15, 19, 21] and the state machine replication paradigm [13] for building reliable distributed systems have been studied for over 40 years. However, with the rise in popularity of blockchains there has been a renewed interest in engineering high-performance consensus protocols. Specifically, to improve on Bitcoin's [33] throughput of only 4 tx/sec early works [29] suggested committee based consensus protocols. For higher throughput and lower latency committee-based protocols are required, and are now becoming the norm in proof-of-stake designs.

Existing approaches to increasing the performance of distributed ledgers focus on creating lower-cost consensus algorithms culminating with HotStuff [38], which achieves linear message complexity in the partially synchronous setting. To achieve this, HotStuff leverages a leader who collects, aggregates, and broadcasts the messages of other validators. However, theoretical message complexity should not be the only optimization target. More specifically:

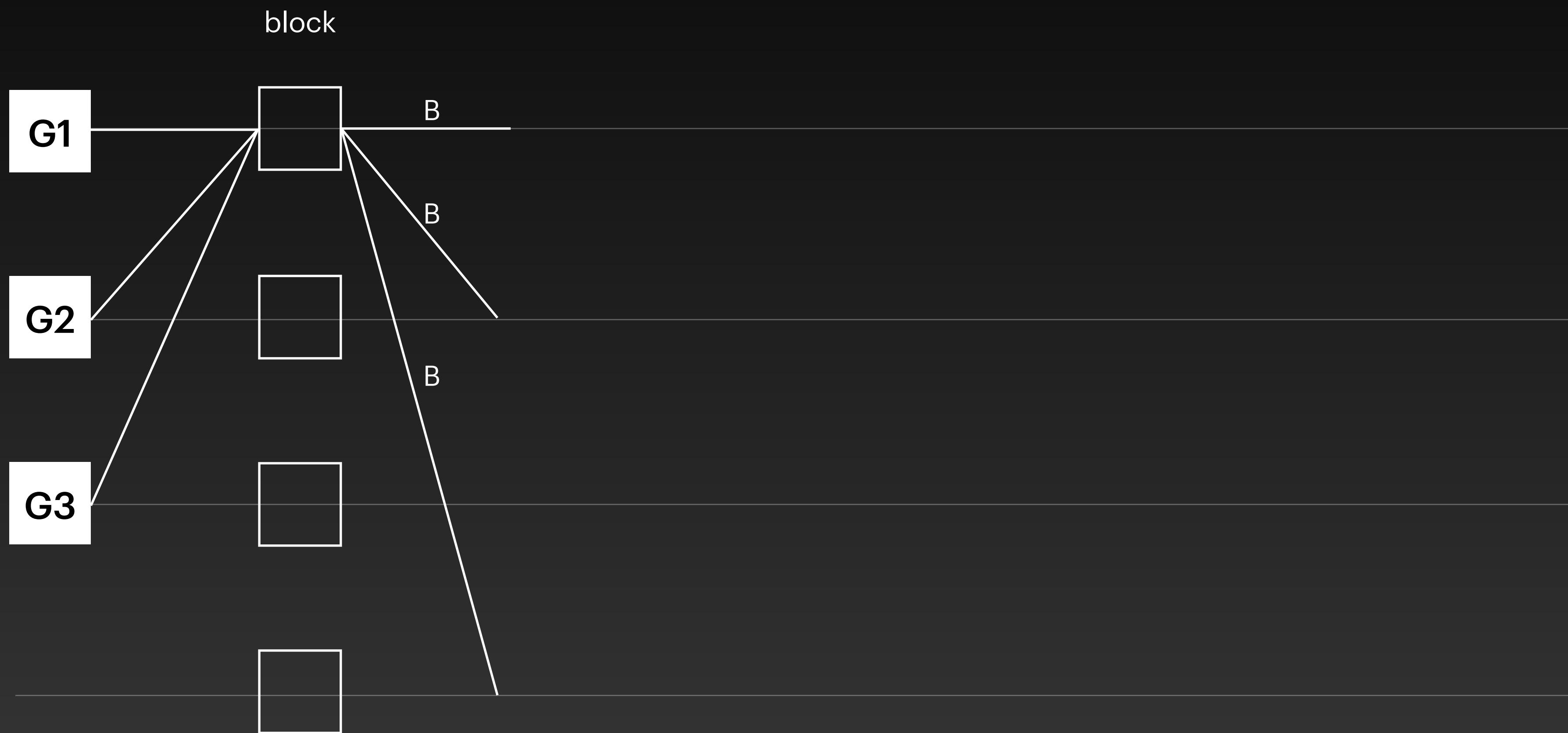
- Any (partially-synchronous) protocol that minimizes overall message number, but relies on a leader to produce proposals and coordinate consensus, fails to capture the high load this imposes on the leader who inevitably becomes a bottleneck.
- Message complexity counts the number of *metadata* messages (e.g., votes, signatures, hashes) which take minimal bandwidth compared to the dissemination of bulk transaction data (blocks). Since blocks are orders of magnitude larger (10MB) than a typical consensus message (100B), the asymptotic message complexity is practically amortized for fixed mid-size committees (up to ~ 50 nodes).

Additionally, consensus protocols have grouped a lot of functions into a monolithic protocol. In a typical distributed

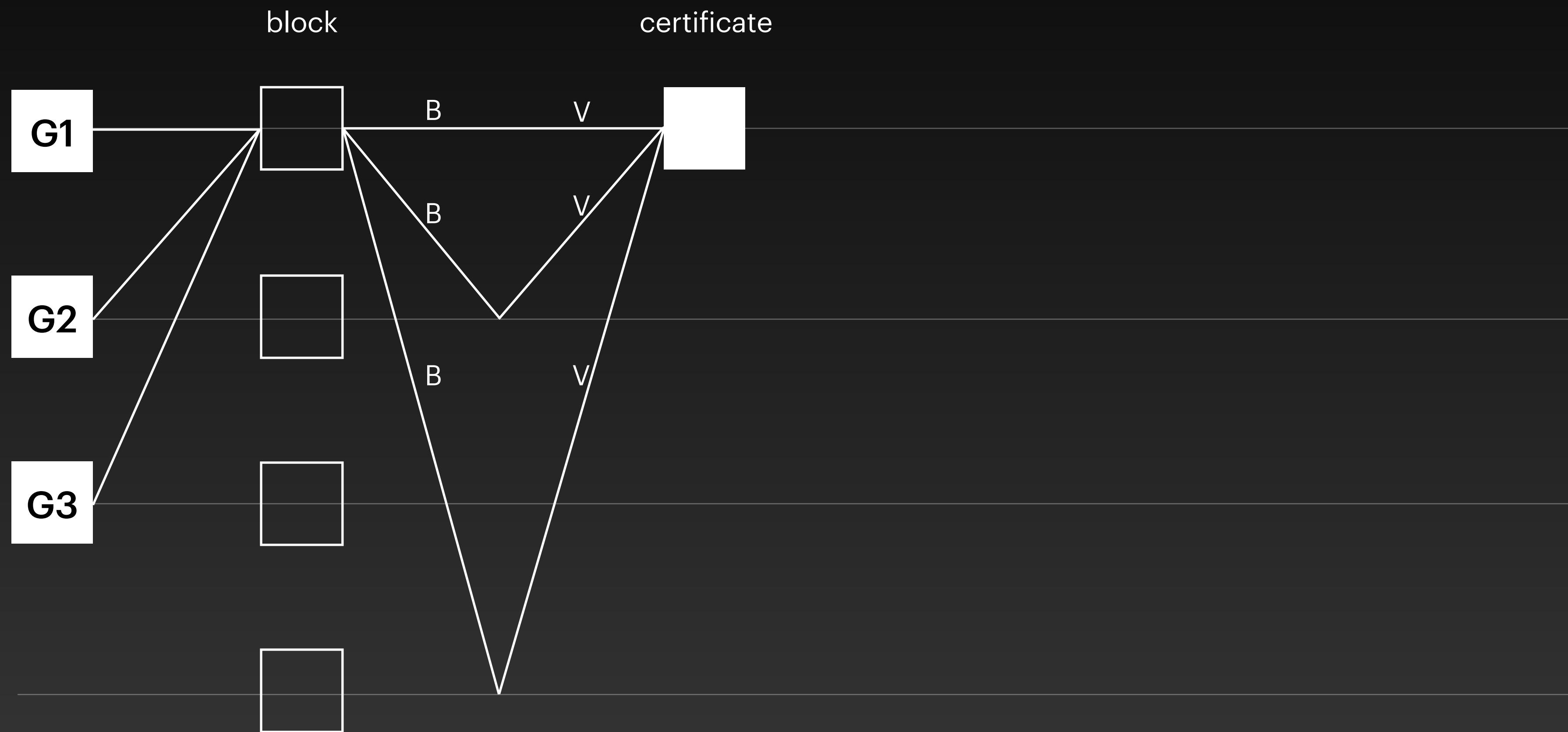
Narwhal

- Quadratic but even resource utilisation
- Separation between consensus and data dissemination

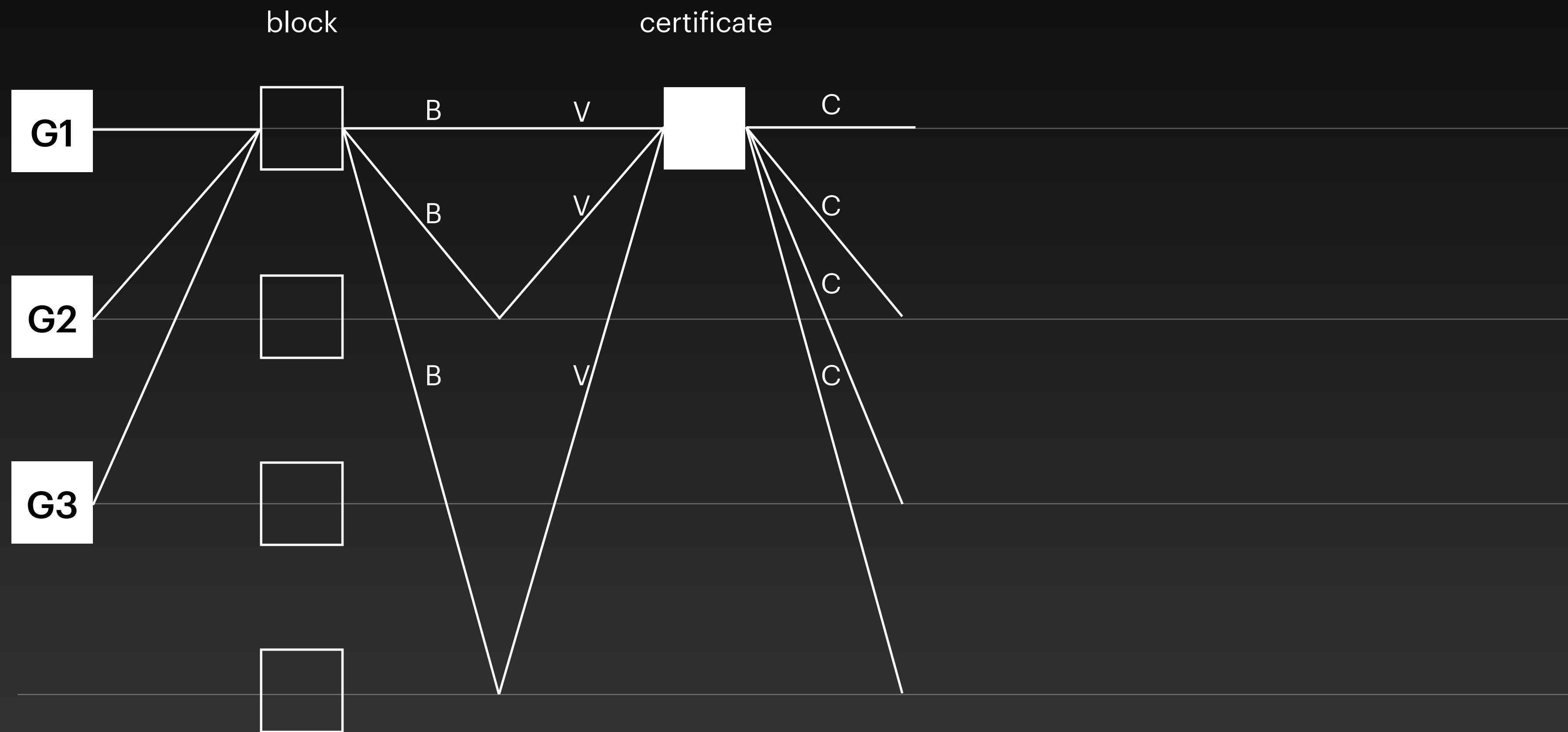
Narwhal



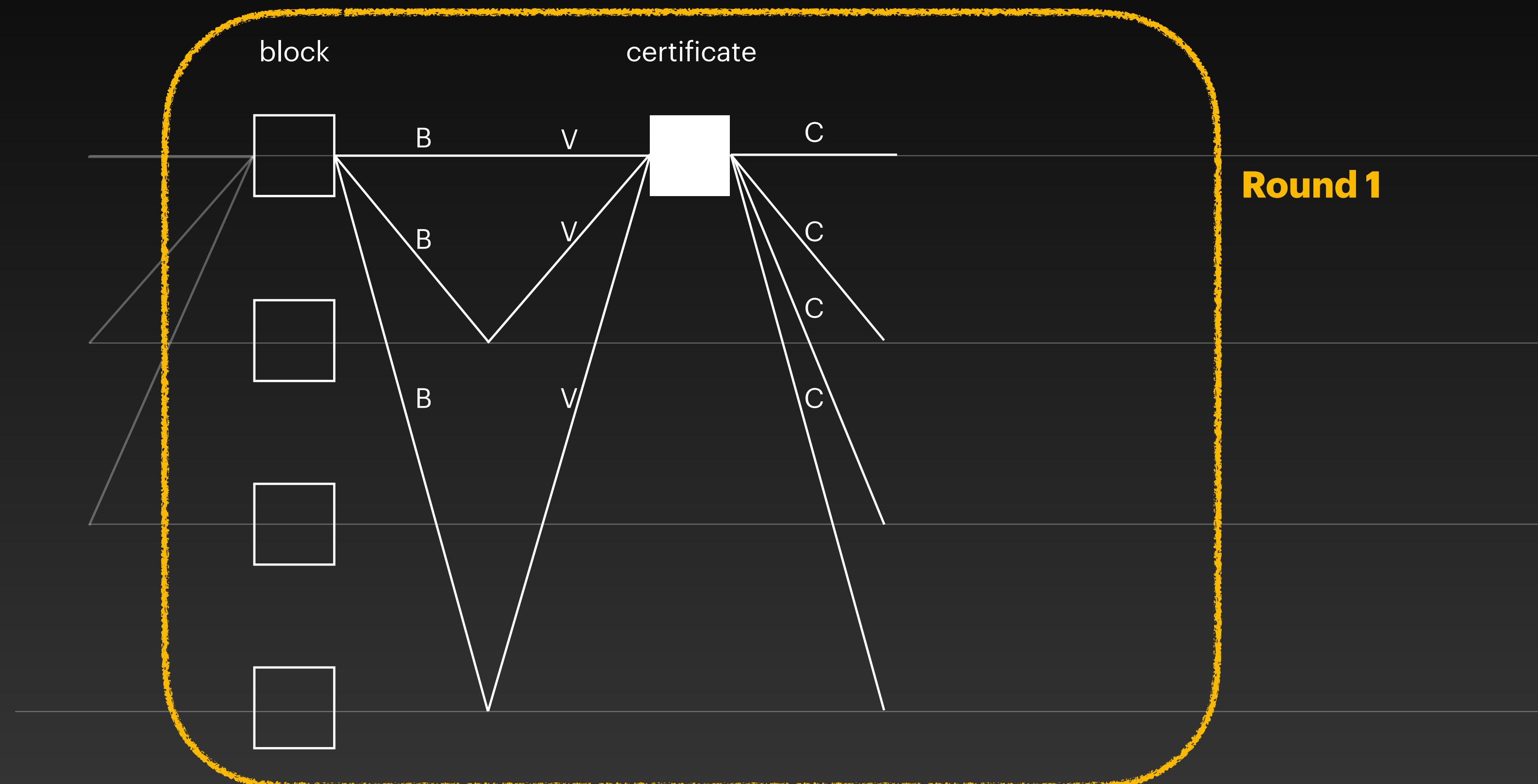
Narwhal



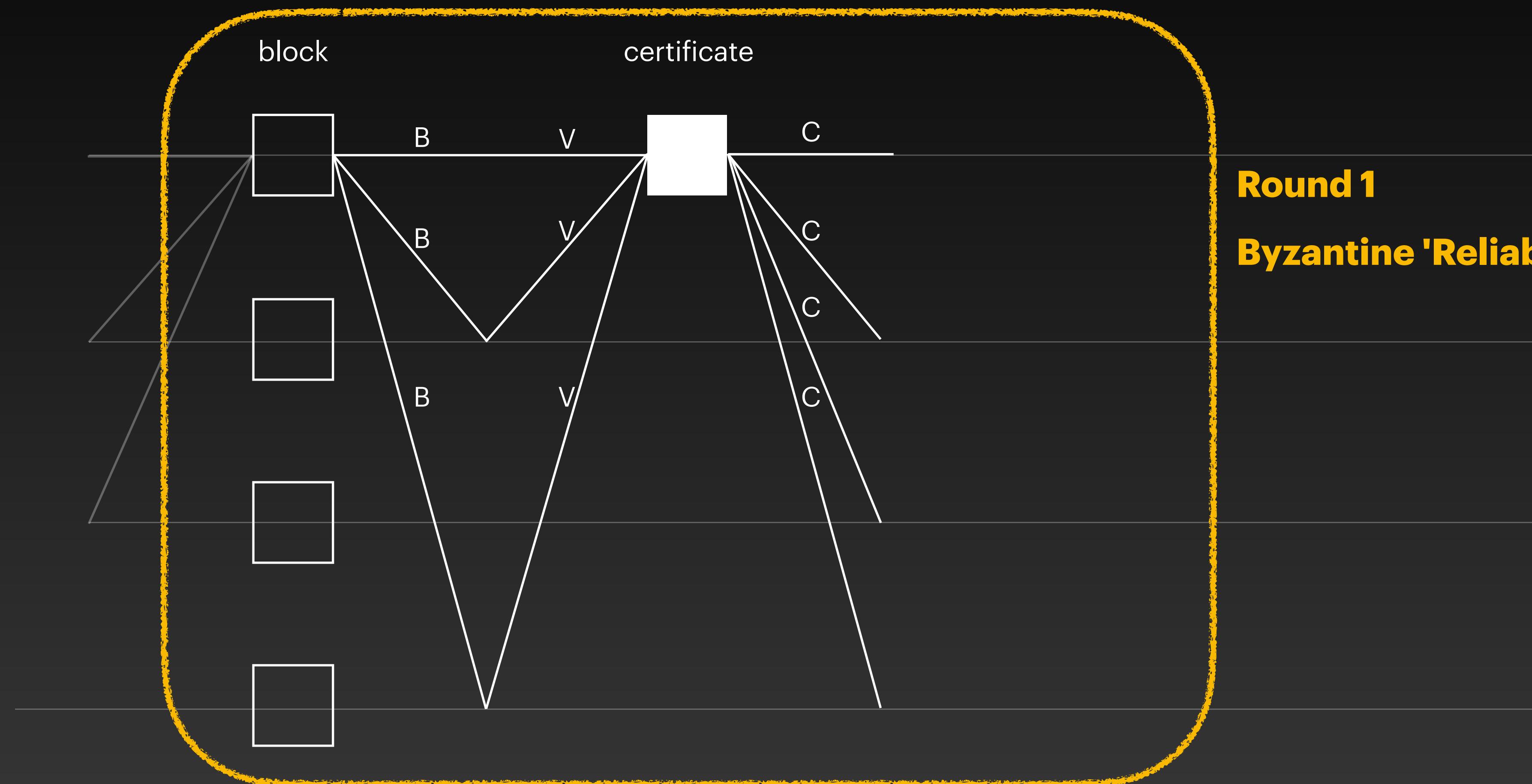
Narwhal



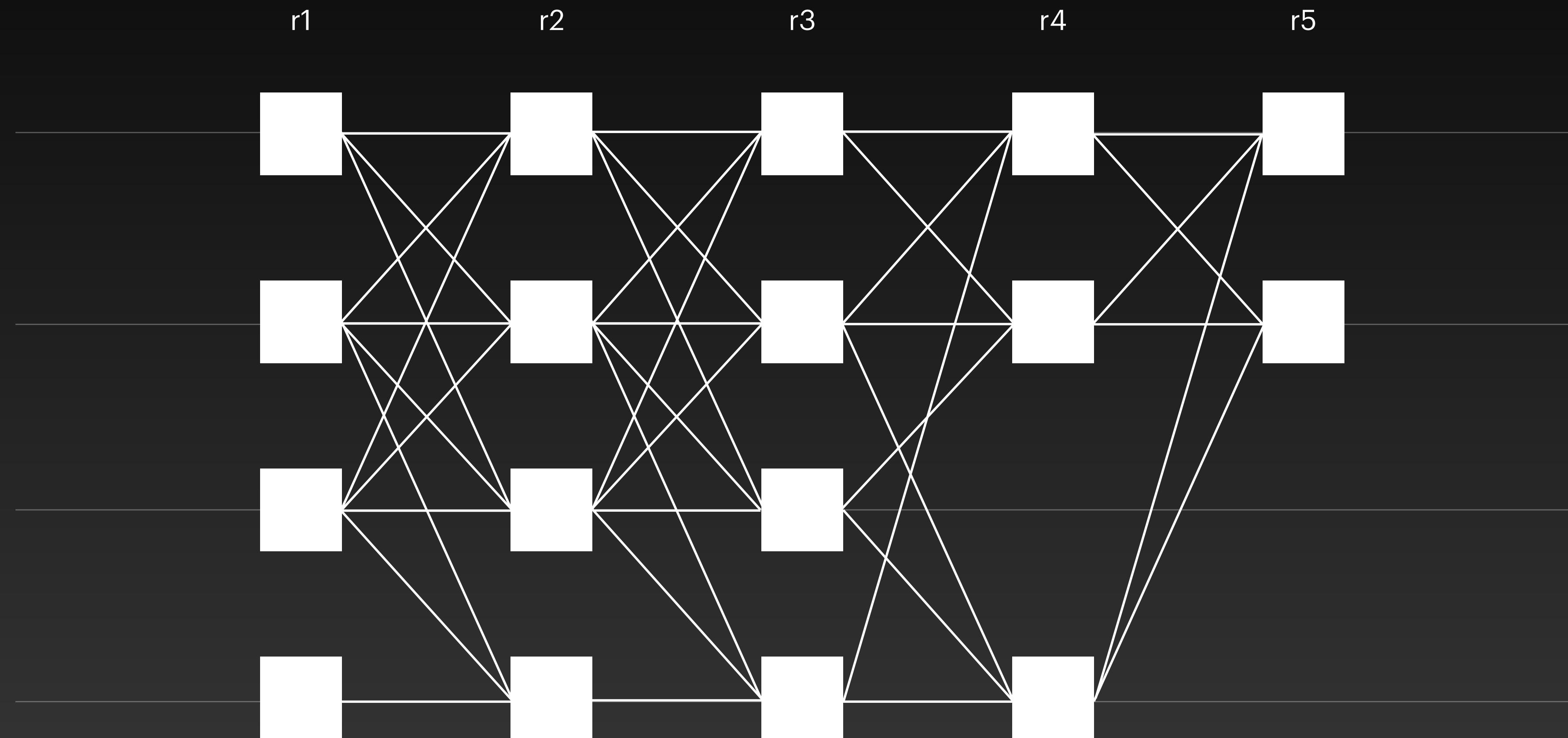
Narwhal



Narwhal



Narwhal



Research Questions

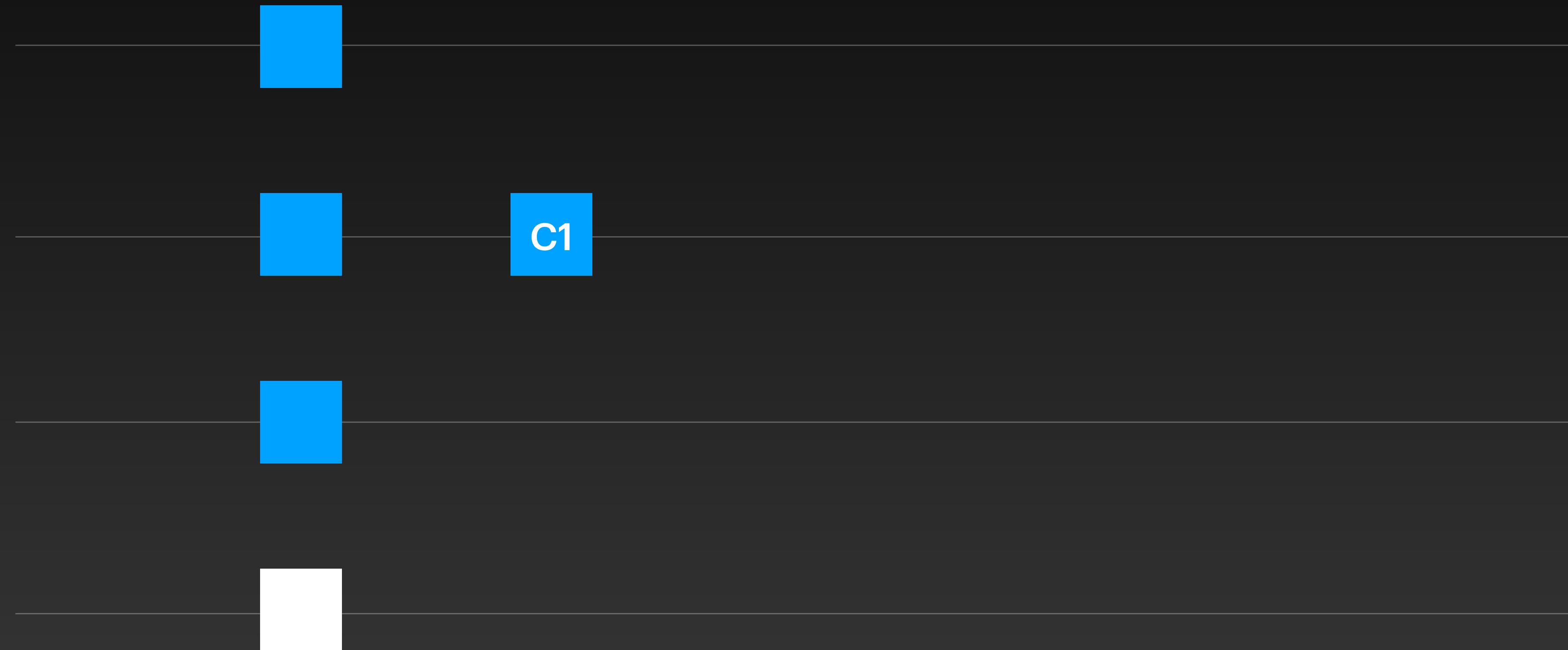
1. Network model?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage

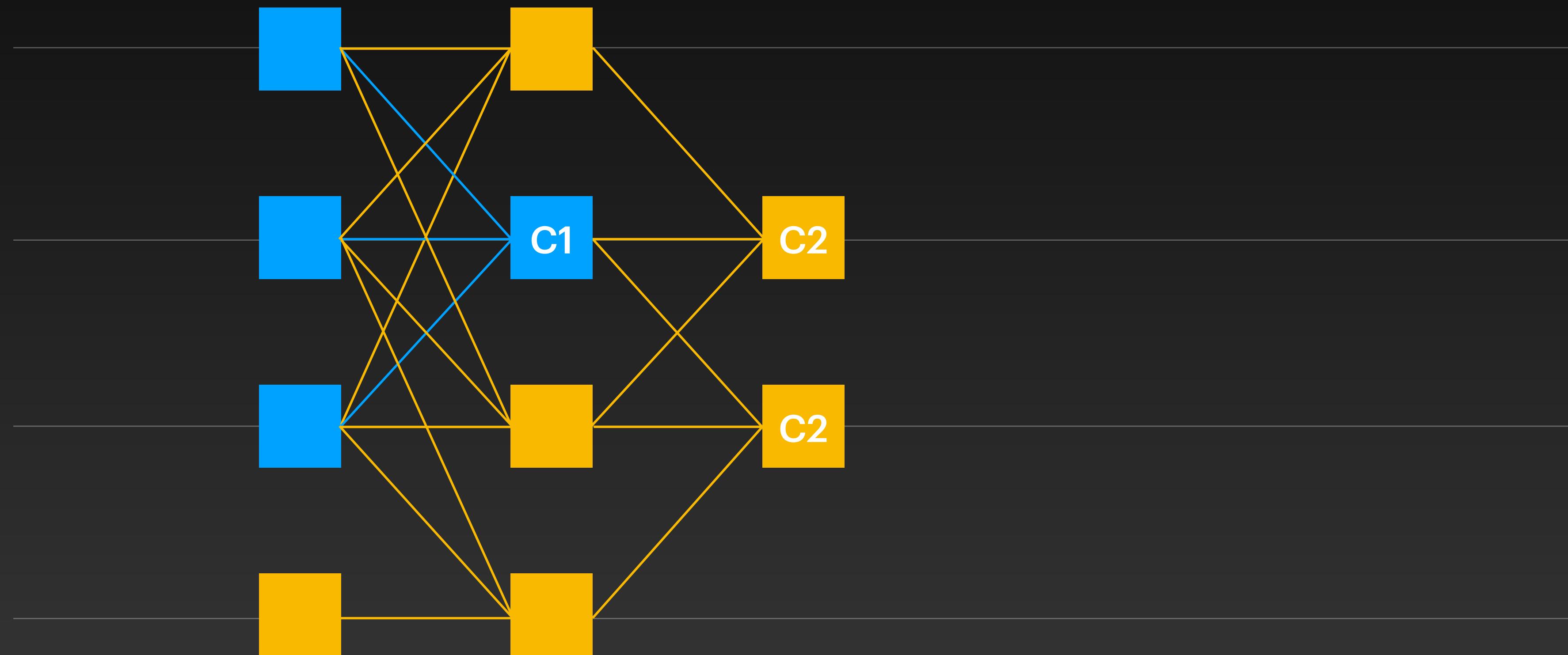
HotStuff on Narwhal

Enhanced commit rule



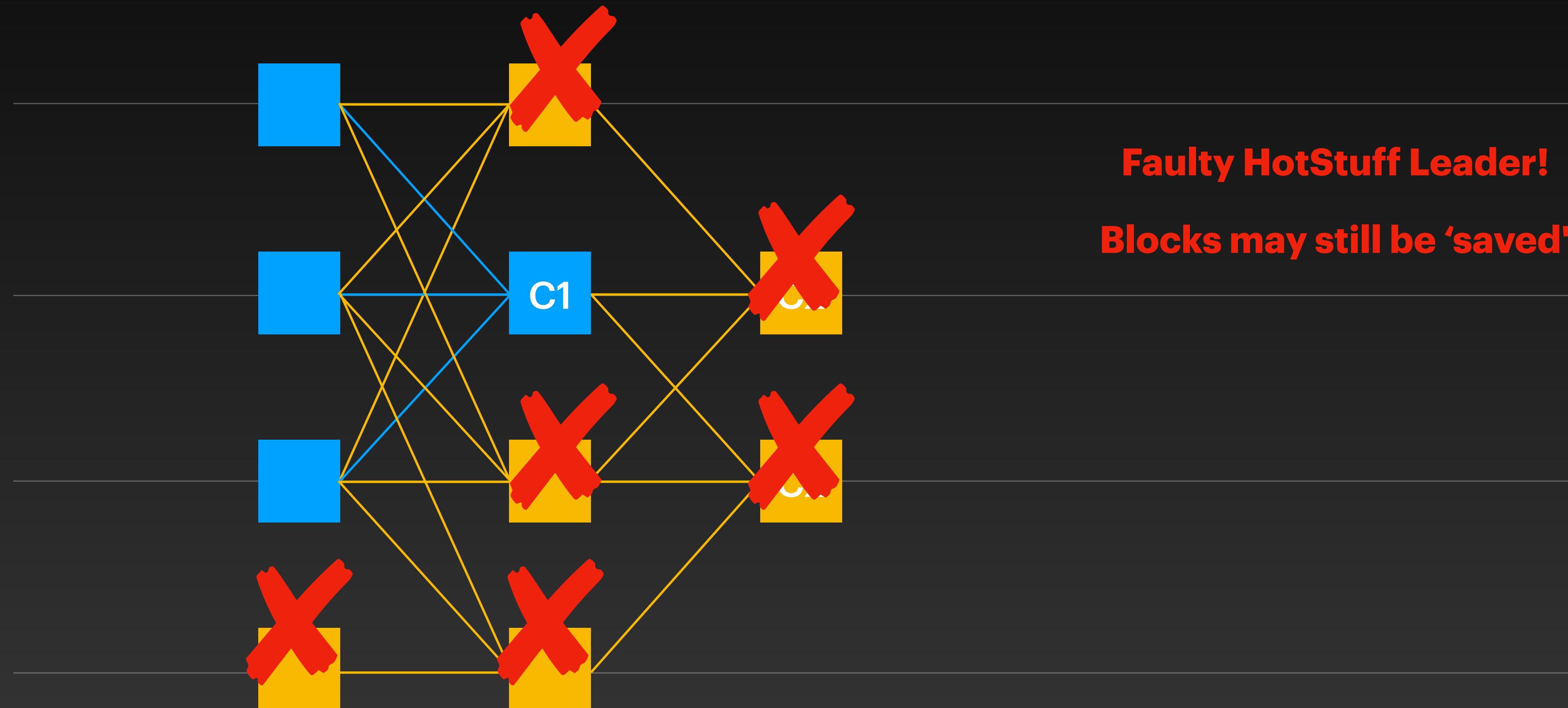
HotStuff on Narwhal

Enhanced commit rule



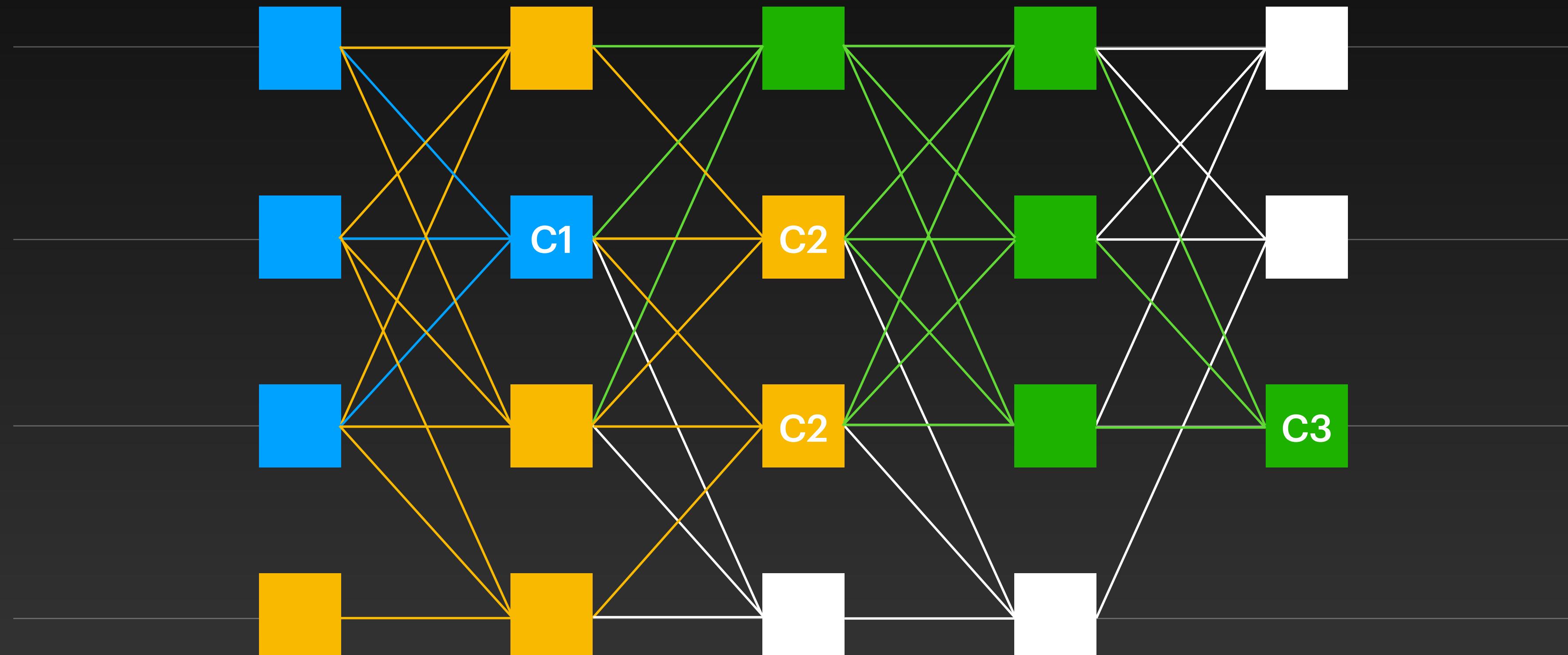
HotStuff on Narwhal

Enhanced commit rule



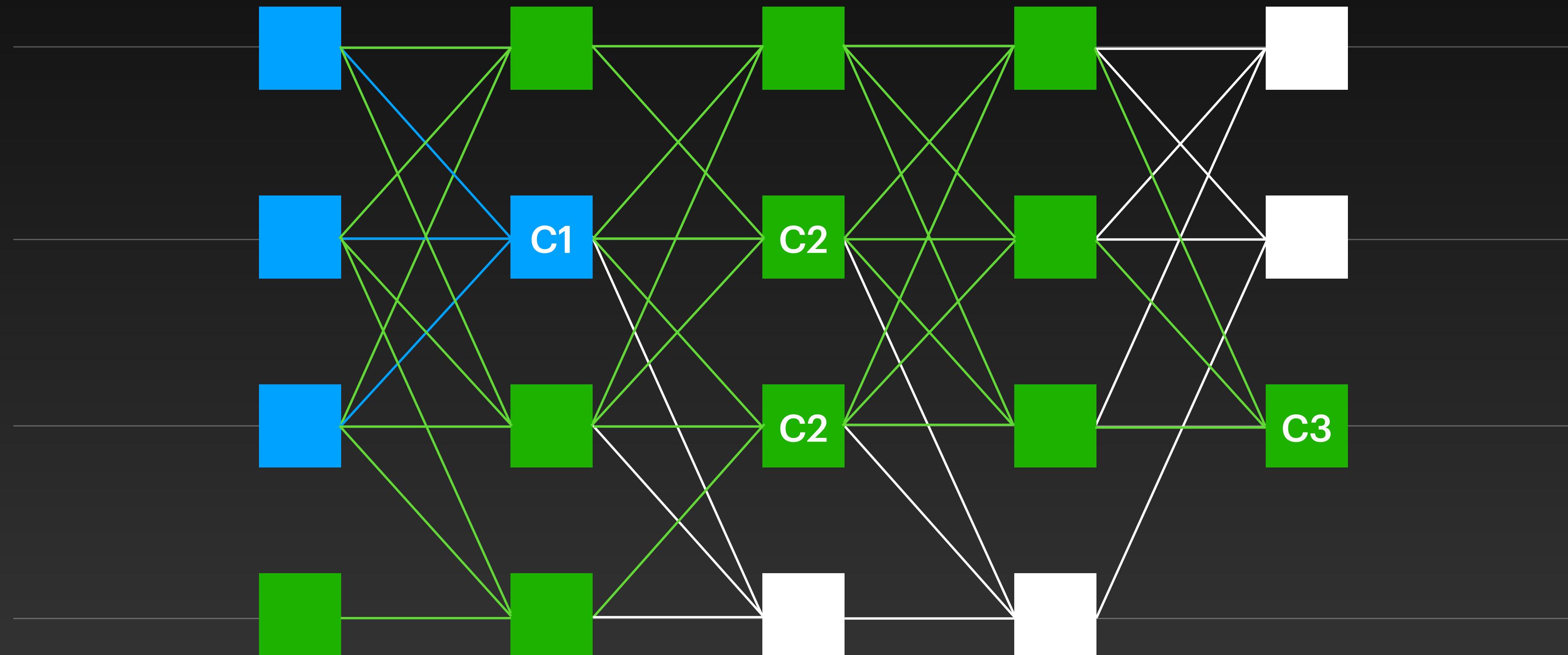
HotStuff on Narwhal

Enhanced commit rule

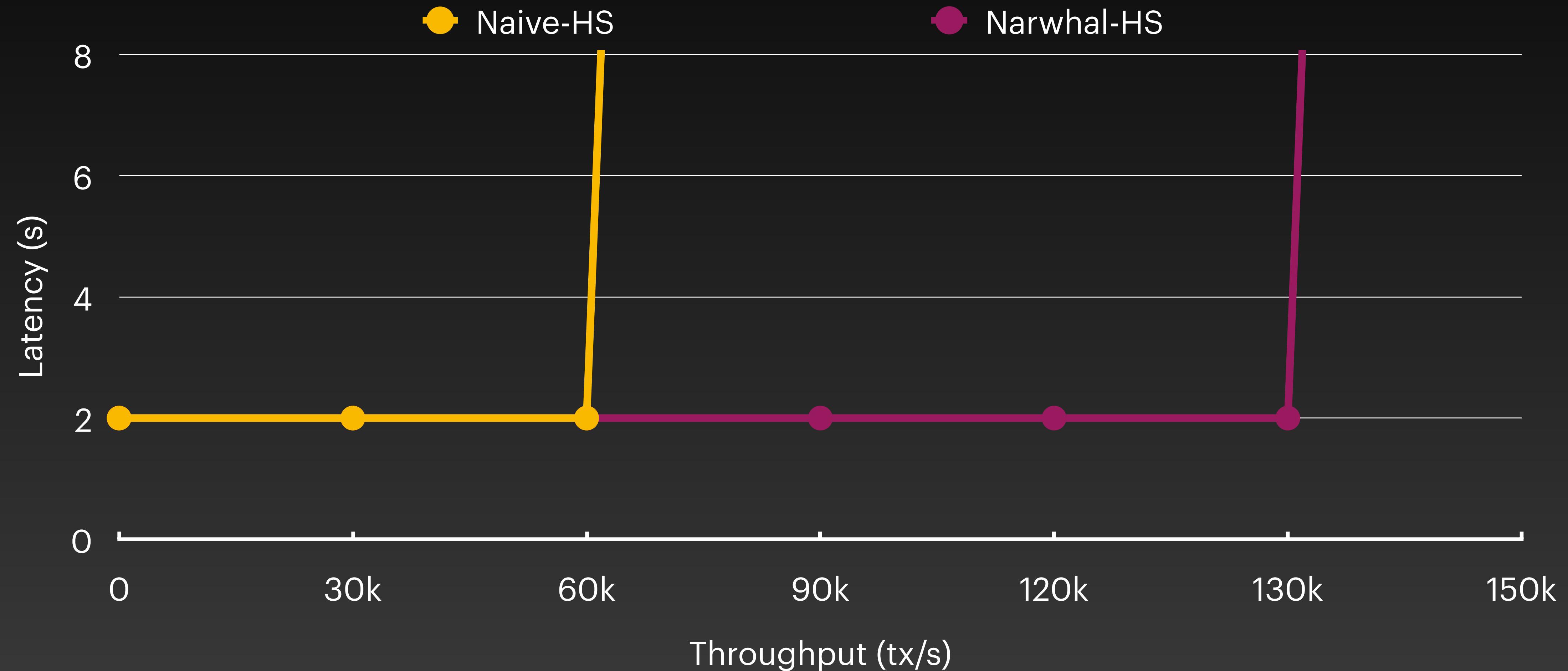


HotStuff on Narwhal

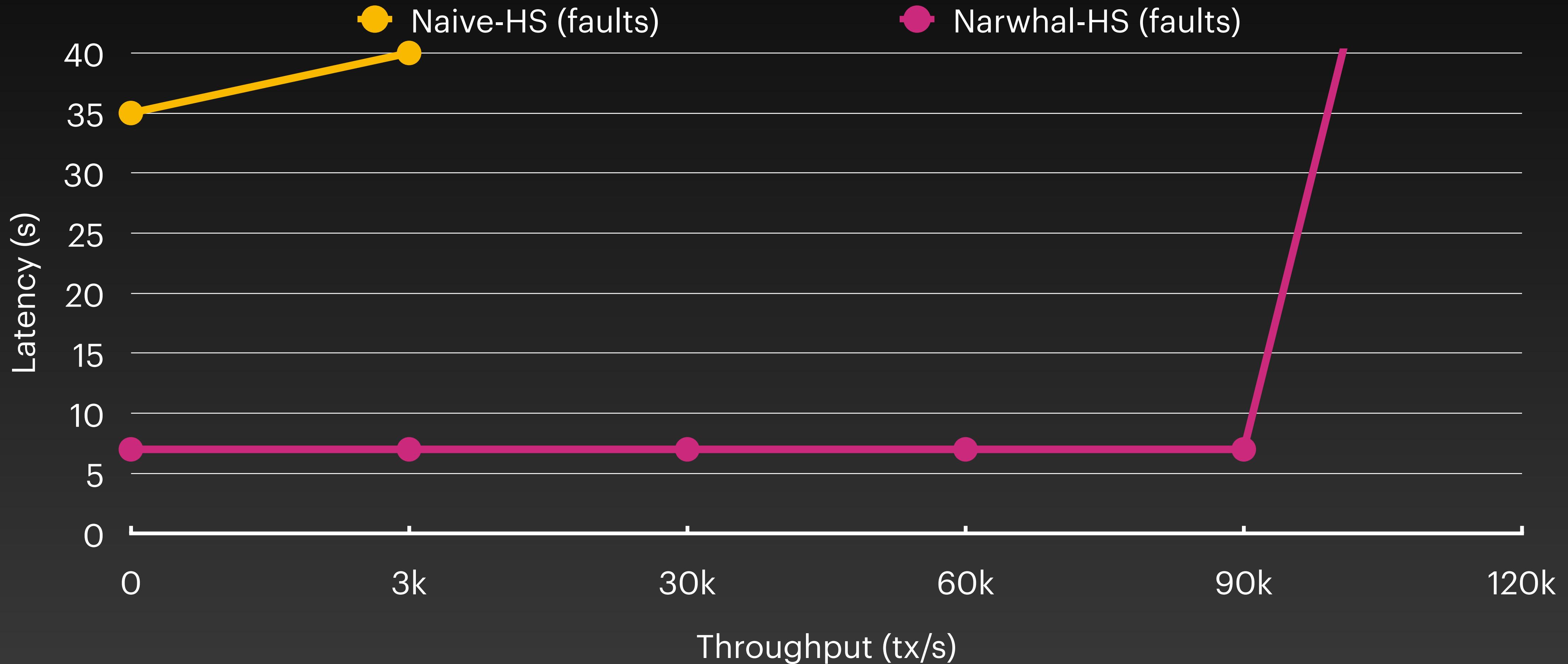
Enhanced commit rule



Performance



Performance



Libra, 2021

Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus

George Danezis
Mysten Labs & UCL

Alberto Sonnino
Mysten Labs

Abstract
We propose separating the task of reliable transaction dissemination from transaction ordering, to enable high-performance Byzantine fault-tolerant quorum-based consensus. We design and evaluate a mempool protocol, Narwhal, specializing in high-throughput reliable dissemination and storage of causal histories of transactions. Narwhal tolerates an asynchronous network and maintains high performance despite failures. Narwhal is designed to easily scale-out using multiple workers at each validator, and we demonstrate that there is no foreseeable limit to the throughput we can achieve.

Composing Narwhal with a partially synchronous consensus protocol (Narwhal-HotStuff) yields significantly better throughput even in the presence of faults or intermittent loss of liveness due to asynchrony. However, loss of liveness can result in higher latency. To achieve overall good performance when faults occur we design Tusk, a zero-message overhead asynchronous consensus protocol, to work with Narwhal. We demonstrate its high performance under a variety of configurations and faults.

As a summary of results, on a WAN, Narwhal-HotStuff achieves over 130,000 tx/sec at less than 2-sec latency compared with 1,800 tx/sec at 1-sec latency for HotStuff. Additional workers increase throughput linearly to 600,000 tx/sec without any latency increase. Tusk achieves 160,000 tx/sec with about 3 seconds latency. Under faults, both protocols maintain high throughput, but Narwhal-HotStuff suffers from increased latency.

CCS Concepts: Security and privacy → Distributed systems security.

Keywords: Consensus protocol, Byzantine Fault Tolerant

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
EuroSys '22, April 5–8, 2022, RENNES, France
© 2022 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-1-4503-9162-7/22/04... \$15.00
<https://doi.org/10.1145/3492321.3519594>

ACM Reference Format:
George Danezis, Lefteris Kokoris-Kogias, Alberto Sonnino, and Alexander Spiegelman. 2022. Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus . In *Seventeenth European Conference on Computer Systems (EuroSys '22), April 5–8, 2022, RENNES, France*. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3492321.3519594>

1 Introduction
Byzantine consensus protocols [15, 19, 21] and the state machine replication paradigm [13] for building reliable distributed systems have been studied for over 40 years. However, with the rise in popularity of blockchains there has been a renewed interest in engineering high-performance consensus protocols. Specifically, to improve on Bitcoin's [33] throughput of only 4 tx/sec early works [29] suggested committee based consensus protocols. For higher throughput and lower latency committee-based protocols are required, and are now becoming the norm in proof-of-stake designs.

Existing approaches to increasing the performance of distributed ledgers focus on creating lower-cost consensus algorithms culminating with HotStuff [38], which achieves linear message complexity in the partially synchronous setting. To achieve this, HotStuff leverages a leader who collects, aggregates, and broadcasts the messages of other validators. However, theoretical message complexity should not be the only optimization target. More specifically:

- Any (partially-synchronous) protocol that minimizes overall message number, but relies on a leader to produce proposals and coordinate consensus, fails to capture the high load this imposes on the leader who inevitably becomes a bottleneck.
- Message complexity counts the number of *metadata* messages (e.g., votes, signatures, hashes) which take minimal bandwidth compared to the dissemination of bulk transaction data (blocks). Since blocks are orders of magnitude larger (10MB) than a typical consensus message (100B), the asymptotic message complexity is practically amortized for fixed mid-size committees (up to ~ 50 nodes).

Additionally, consensus protocols have grouped a lot of functions into a monolithic protocol. In a typical distributed

Narwhal

- Quadratic but even resource utilisation
- Separation between consensus and data dissemination
- High engineering complexity

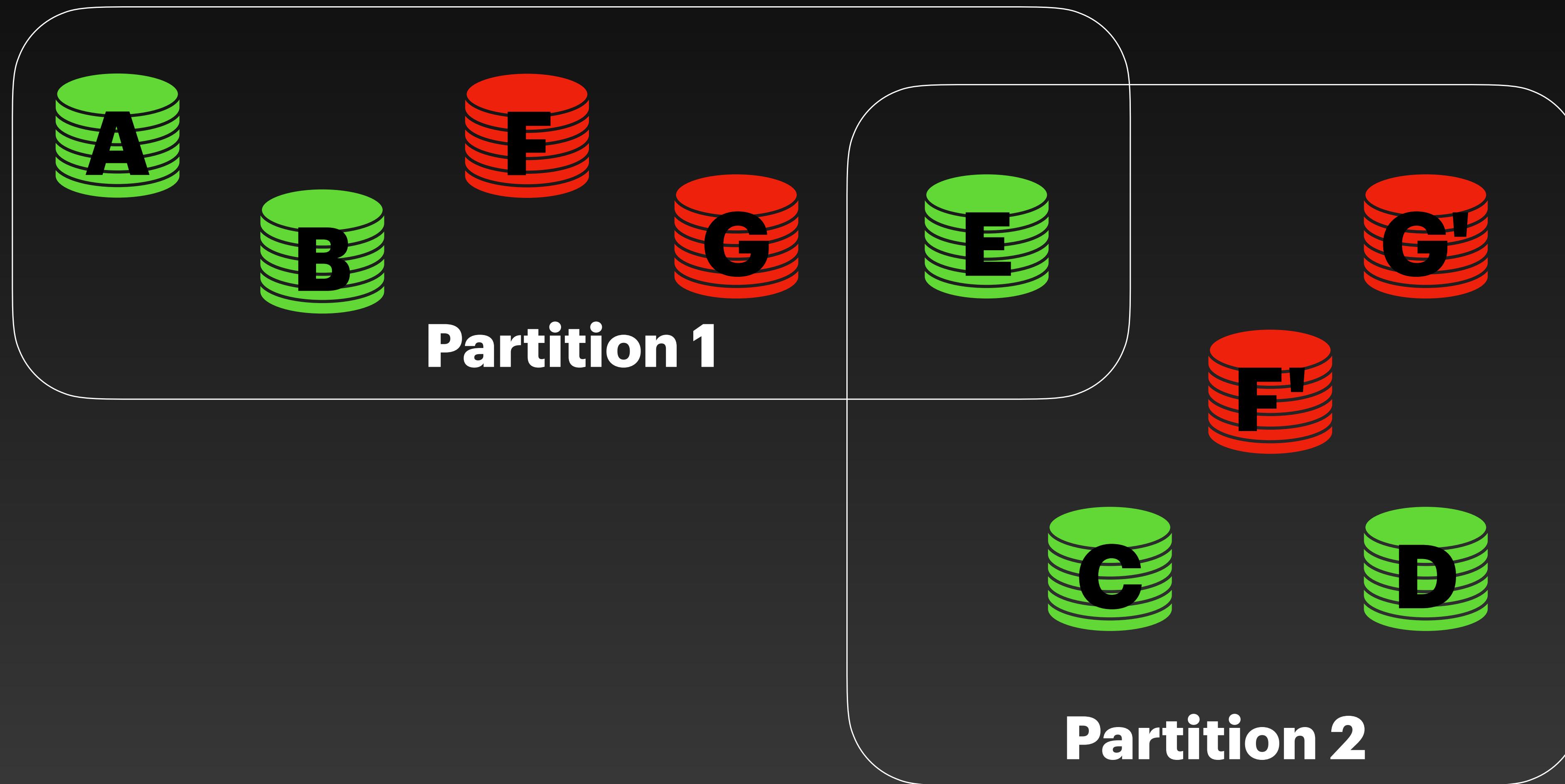
Research Questions

1. Network model?
2. BFT testing?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage

Twins



DagRider

All You Need is DAG

Idit Keidar
Technion

Oded Naor*
Technion

ABSTRACT
We present *DagRider*, the first asynchronous Byzantine Atomic Broadcast protocol that optimizes round-trip communication complexity, and optimal time complexity. DagRider is correct processes eventually get delivered. We construct two layers. In the inner layer, processes broadcast their messages and do not rely on asymmetric cryptographic assumption. Therefore, when using a deterministic threshold-based coin implementation, the safety properties of our DAG protocol are post-quantum secure.

ACM Reference Format:
Idit Keidar, Eleftherios Kokoris-Kogias, Oded Naor, and Alexander Spiegelman. 2021. DagRider: A DAG-based Asynchronous Byzantine Fault-Tolerant Protocol. In *Proceedings of the 2021 ACM SIGSAC Conference on Principles of Distributed Computing (PODC '21)*, July 26–30, 2021, Virtual Event, Italy. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3468044.3492211>

1 INTRODUCTION
The need for replicated systems in scalable geo-replicated Byzantine fault-tolerant consensus systems has motivated an enormous amount of study on the Byzantine State-of-the-Art. Bellal et al. [1] have shown that the state-of-the-art consensus protocols in replicated systems have been studied for over 40 years. However, with the rise in popularity of blockchains there has been a renewed interest in replicated consensus protocols in blockchains. Specifically, to implement Bitcoiin's [13] throughput of only twice early works [29] suggested consensus protocols. For higher throughput and lower latency committee-based consensus protocols are required, and we are interested in the proof of concept. In this paper, we demonstrate that all proposals by correct processes are eventually included.

Asynchronous consensus protocols for the Byzantine consensus problem [12, 16, 26] have been considered too costly or complicated to be used in practical SMR solutions. In fact, two recent papers [1, 29] proposed consensus protocols that can be easily implemented by a single node to commit a single value. To compare DagRider to the state-of-the-art asynchronous Byzantine consensus protocols, we consider SMR implementations that are an isolated sequence of the VABA and Dumbo protocols to independently agree on every slot. To achieve this, DagRider uses the same consensus logic as VABA and Dumbo [33], presented asynchronous solutions with (1) optimal resilience, (2) constant time complexity, and (3) optimal quasideterminism and low amortized communication complexity for the leader. In this paper, we follow this recent line of work.

• DagRider is the first asynchronous Byzantine consensus protocol that minimizes overall message number, but relies on a leader to propose proposals and coordinate consensus fails to capture the high load this imposes on the leader who inevitably becomes a bottleneck.

• Message complexity counts the number of *metadata* messages (e.g., voter, signatures, hashes) which take minimal bandwidth compared to the size of the full transaction data (blocks). Since blocks are orders of magnitude larger (10MB) than a typical consensus message (100B), the asymptotic message complexity is proportional for fixed mid-size committees (up to ~50 nodes).

• Additionally, consensus protocols have grouped a lot of functions into a monolithic protocol. In a typical distributed

• DagRider uses a reliable broadcast abstraction as a basic building block, different instantiations yield different complexities. For example, if we use the classic BABA broadcast [11] to propose a

*Oded Naor is grateful to the Hiroaki Fujisawa (Cyber Security Research Center for providing a research grant. Part of Oded's work was done while at NTT Research.

Permissions to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are made available to the author and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with permission is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2021, copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-8548-0/21/07...\$15.00.
<https://doi.org/10.1145/3468044.3492211>

Tusk

Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus

George Danezis
Mythen Labs & UCL

Alexander Spiegelman
Mythen Labs

Abstract
We propose separating the task of reliable transaction dissemination from transaction ordering, to enable high-performance Byzantine fault-tolerant quorum-based consensus. We design two layers. In the inner layer, Narwhal is a mempool, Narwhal tolerates up to two layers of faults in the layer above it, and it can cast their previous and future actions on a Directed Acyclic Graph (DAG) of the communication among them. In the second layer, Tusk preserves their DAGs and totally order all proposals with no extra communication.

ACM Reference Format:
George Danezis, Eleftherios Kokoris-Kogias, Alberto Sonnino, and Alexander Spiegelman. 2021. Narwhal and Tusk: A DAG-based Mempool and Efficient BFT Consensus. In *Seventeenth European Conference on Computer Systems (EuroSys '21)*, April 5–6, 2021, Rennes, France. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3479221.3479994>

1 Introduction
Recent work on consensus protocols [15, 19, 21] and the state-of-the-art distributed replicated consensus protocols [1, 20] have been studied for over 40 years. However, with the rise in popularity of blockchains there has been a renewed interest in replicated consensus protocols with high performance. Specifically, to implement Bitcoiin's [13] throughput of only twice early works [29] suggested consensus protocols. For higher throughput and lower latency committee-based consensus protocols are required, and we are interested in the proof of concept. In this paper, we demonstrate that all proposals by correct processes are eventually included.

As a summary of results, on a WAN, Narwhal-HotStuff achieves over 100K transaction per second (TPS) with a proposal rate of 1,800 TPS at 1 sec latency for HotStuff. Additional workers increase throughput linearly to 600,000 TPS without any latency increase. Tusk achieves 160,000 TPS with 1 sec seconds latency. Under faults, both protocols maintain high throughput, but Narwhal-HotStuff suffers from increased latency.

CCS Concepts: • Security and privacy → Distributed systems security.

Keywords: Consensus protocol, Byzantine Fault Tolerant

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are made available to the author and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with permission is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2021, copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-8548-0/21/07...\$15.00.
<https://doi.org/10.1145/3468044.3492203>

Bullshark

Bullshark: DAG BFT Protocols Made Practical

George Danezis
Mythen Labs & UCL

Alexander Spiegelman
Aptos

Abstract
We present Bullshark, the first directed acyclic graph (DAG) based asynchronous Byzantine fault-tolerant (BFT) consensus protocol that is optimized for the common synchronous case. Like previous DAG-based BFT protocols [19, 25], Bullshark requires no extra communication to achieve consensus at the top of the DAG. That is, parties can simply receive the view of the top-most node in the DAG and cast the DAG by locally interpreting their view of it without sending any extra messages. This is, once we build the DAG, implementing consensus is as simple as running a local update function on each node.

The pioneering work of Hashgraph [4] constructed an unstructured DAG, where each message refers to previous ones, and used batches of messages as local coin flips to safely order the DAG in each node. Bullshark is a structured DAG, where each node is a round-based DAG and encodes a shared randomness in each round via a threshold signature scheme to achieve constant latency in expectation. Every node in the DAG has at most n vertices (one for each party), each of which contains a block of transactions as well as a timestamped log of the previous $n-1$ rounds of consensus. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably receives $n-1$ blocks from other nodes. Note that building the DAG requires honest parties to broadcast vertices even if they have no transactions to propose. However, the edges of the DAG are not necessarily valid. In this sense it is not different from other BFT protocols in which parties send explicit vote messages, which contain no transactions as well. Remarkably, by using the DAG, Bullshark achieves a much simpler and more efficient edge interpretation logic of DagRider to totally order the DAG spans over less than 30 lines of pseudocode.

DAG-based consensus protocols like the original atomic broadcast (BAB), which achieves optimal amortized communication complexity ($O(n)$ per transaction), post quantum safety, and some notion of fairness (such as liveness) are well known. However, the challenge is how to achieve this while maintaining a low overhead. In this paper, we show that by using an unstructured DAG, we can achieve the same performance as BAB while maintaining a low overhead. Our main contribution is to show that the DAG-based consensus is as simple as BAB, yet it is more efficient and has a lower overhead. We also introduce a simple and partially synchronous version of Bullshark which we evaluate against the state-of-the-art. The implemented protocol is embarrassingly simple [20] and achieves multi-valued validated Byzantine agreement (MVBA) to match the throughput of VABA [1]. It is highly efficient, achieving for example 125,000 transaction per second with a 2 seconds latency for a deployment of 50 parties. In the same setting, VABA pays a steep 50% latency increase as it optimizes for asynchrony.

ACM Reference Format:
George Danezis, Alberto Sonnino, and Alexander Spiegelman. 2021. Bullshark: DAG BFT Protocols Made Practical. In *Proceedings of ACM Conference on Computer Systems, Los Angeles, CA, USA, November 2021* (CompoNS '21), pp. 17–31. <https://doi.org/10.1145/3468044.3492222>

1 Introduction
Ordering transactions in a distributed Byzantine environment via a consensus mechanism has become one of the most timely research areas in recent years due to the blooming Blockchain use-case. A recent line of work [8, 19, 21, 25, 33, 40] proposed an elegant way to separate the consensus from the ordering.

For example, in the context of the VABA and Dumbo protocols, the consensus part is responsible for the safety properties, while the ordering part is responsible for the liveness properties. In this paper, we propose a new ordering protocol that is designed to be as simple as possible. We call it Bullshark.

CCS Concepts: • Security and privacy → Byzantine Fault Tolerant

Keywords: Consensus protocol, Byzantine Fault Tolerant

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are made available to the author and the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with permission is permitted. To copy otherwise, or to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
© 2021, copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-8548-0/21/07...\$15.00.
<https://doi.org/10.1145/3468044.3492204>

Dumbo-NG

Dumbo-NG: Fast Asynchronous BFT Consensus with Throughput-Oblivious Latency

Yiying Guo*
ISCA & UCAS

Yuan Liu*
ISCA

Zherui Lu*
USID

Qiang Tang*
UCAS

Jing Xu*
ISCA

Zhenfeng Zhang*
ISCA

Abstract
We present Dumbo-NG, a novel asynchronous BFT consensus algorithm that is designed to bring down the infrastructure of distributed ledger for mission-critical applications. Such decentralized business is envisioned as critical global infrastructure maintained by a set of multi-party consensus protocols. Dumbo-NG is throughput-oblivious, that is, only ready that the honest nodes incur no communication blow-up by letting the honest nodes to be redundant for deployment over the Internet.

We present Dumbo-NG, a novel asynchronous BFT consensus (throughput-oblivious) algorithm. The technical core is a non-trivial direct broadcast from asynchronous BFT consensus to a multi-valued validated Byzantine agreement (MVBA) to match the throughput of VABA [1]. It is highly efficient, achieving for example 125,000 transaction per second with a 2 seconds latency for a deployment of 50 parties. In the same setting, VABA pays a steep 50% latency increase as it optimizes for asynchrony.

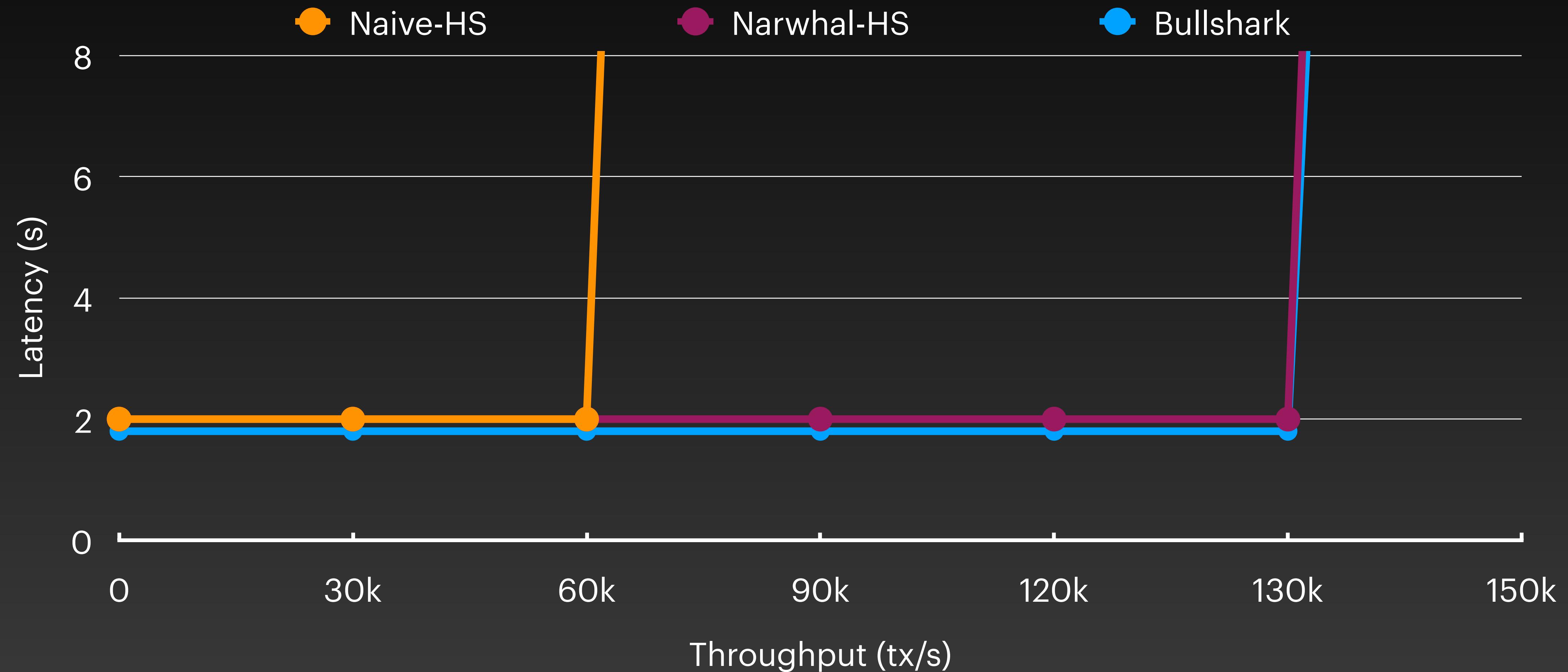
ACM Reference Format:
Yiying Guo*, Yuan Liu*, Zherui Lu*, Qiang Tang*, Jing Xu*, and Zhenfeng Zhang*. 2021. Dumbo-NG: Fast Asynchronous BFT Consensus with Throughput-Oblivious Latency. In *Proceedings of ACM Conference on Computer Systems, Los Angeles, CA, USA, November 2021* (CompoNS '21), pp. 17–31. <https://doi.org/10.1145/3468044.3492222>

1 INTRODUCTION
The large success of Bitcoin [6] and Blockchain [19, 24] leads to an increasing tendency to lay down the infrastructures of distributed ledger for mission-critical applications. Such decentralized business is envisioned as critical global infrastructure maintained by a set of multi-party consensus protocols. Dumbo-NG is throughput-oblivious, that is, only ready that the honest nodes incur no communication blow-up by letting the honest nodes to be redundant for deployment over the Internet.

Asynchronous BFT for indispensible robustness. The consensus of distributed ledger structure has to thrive in a highly adversarial environment. When the applications atop it are critical financial and banking services, some nodes can be well controlled by malicious parties. In this case, the unreliable Internet might become part of the attack surface due to network flooding, misconfiguration and even network attacks. To cope with the adversarial environment, many consensus mechanisms are proposed. Specifically, Byzantine fault-tolerant (BFT) consensus [4, 20, 25, 47, 58, 59, 60] are arguably the most suitable candidates. They can realize high security guarantee to ensure liveness (as well as safety) despite the network delay and the delay may be very long. In contrast, many (partial) synchronous consensus protocols [5, 6, 8, 15, 27, 44, 45, 64, 73] such as PBFT [26] and HotStuff [75] might cause a significant performance degradation due to the communications without making any progress [36, 60] when unluckily encountering an asynchronous network adversary.

1.1 Practical obstacles of adopting asynchronous BFT consensus
Undoubtedly, it is far from trivial to change current practical asynchronous consensus algorithms. One of such practical issues was widely adopted due to serious efficiency concerns. The seminal FLP “impossibility” [36] proves that no deterministic consensus can be achieved in an asynchronous environment. Since the 1980s, many attempts to break this barrier have been made. In particular, the “impossibility” by randomized protocols, but most of them focused on theoretical feasibility, and unfortunately, several attempts of implementations [23, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 628, 629, 630, 631, 632, 633, 634, 635, 636, 636, 637, 638, 639, 639, 640, 641, 642, 643, 643, 644, 645, 646, 646, 647, 648, 649, 649, 650, 65

Performance



Research Questions

1. Network model?
2. BFT testing?
3. Consensus-exec interface?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy

By that time...



← Post

Reply

Pinned



David Marcus

@davidmarcus



...

How Libra Was Killed.

I never shared this publicly before, but since [@pmarca](#) opened the floodgates on [@joerogan](#)'s pod, it feels appropriate to shed more light on this.

As a reminder, Libra (then Diem) was an advanced, high-performance, payments-centric blockchain paired with a stablecoin that we built with my team at [@Meta](#). It would've solved global payments at scale. Prior to announcing the project, we spent months briefing key regulators in DC and abroad. We then announced the project in June 2019 alongside 28 companies. Two weeks later, I was called to testify in front of both the Senate Banking Committee and the House Financial Services Committee, which was the starting point of two years of nonstop work and changes to appease lawmakers and regulators.

By spring of 2021 (yes they slow played us at every step), we had addressed every last possible regulatory concern across financial crime, money laundering, consumer protection, reserve management, buffers,

By that time...



Sui

Aptos

Linera

...

Fundraising with papers
seems to work

Sui, 2022

Over a year for mainnet

- Lack of checkpoints
- Lack of epoch-change
- Lack of crash-recovery

Research Questions

1. Network model?
2. BFT testing?
3. Consensus-exec interface?
4. Storage architecture?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy
6. Epoch change is not an add-on

Sui, 2023

- Latency was too high
- Crash faults were predominant
- Building Bullshark was still too complex

Shoal

Shoal: Improving DAG-BFT Latency And Robustness

Alexander Spiegelman
Aptos
Rati Gelashvili
Aptos

Abstract

The Narwhal system is a state-of-the-art Byzantine fault-tolerant scalable architecture that involves constructing a directed acyclic graph (DAG) of messages among a set of validators (or Blockleaders). Recently, it was proposed as a consensus protocol on top of the Narwhal's DAG that can support over 100k transactions per second. Unfortunately, the high throughput of Bullshark comes with a latency price due to the DAG-based consensus mechanism, which compares to the state-of-the-art leader-based BFT consensus protocols.

We introduce Shoal, a protocol-agnostic framework for enhancing Narwhal-based consensus. By incorporating leader replacement, Shoal achieves a round time that is significantly reduced. Moreover, the combination of properties of the DAG construction and the leader replacement mechanism enables the protocol to scale in all but extremely adversarial scenarios in practice, a property we name "prevalent responsiveness". It strictly subsumes the established and often desired "optimistic responsiveness" property for BFT consensus.

We evaluate Shoal instantiated with Bullshark, the fastest existing Narwhal-based consensus protocol, in an open-source Blockchain project and provide experimental evaluations demonstrating up to 40% latency reduction in the failure-free execution. We also evaluate the execution with failures against the vanilla Bullshark implementation.

CCS Concepts - Security and privacy → Distributed systems security

Keywords: Consensus Protocol, Byzantine Fault Tolerance, ACM Reference Format

Alexander Spiegelman, Rati Gelashvili, and Zekun Li
2023. Shoal: Improving DAG-BFT Latency And Robustness

1 Introduction

Byzantine fault tolerant (BFT) systems, including consensus protocols [13, 23, 24, 29] and state machine replication [7, 10, 22–24], have been designed to tolerate up to t faulty nodes as a means of constructing reliable distributed systems. Recently, the advent of Blockchains has underscored the significance of high performance. While Bitcoin handles approximately 10 transactions per second (TPS), the state-of-the-art consensus-based blockchains [30, 31, 43, 44] are now engaged in a race to deliver a scalable BFT system with the utmost throughput and minimal latency.

Sailfish

Sailfish: Towards Improving the Latency of DAG-based BFT

Nibesh Shrestha
n.shrestha@supraclouds.com
Supra Research
Balaji Arun
Aptos
Rati Gelashvili
Aptos
Zekun Li
Aptos

Historically, the prevailing belief has been that reducing communication complexity was the key to unlocking high performance, low-latency consensus. However, this did not result in dramatic improvements in the throughput. For example, the state-of-the-art Hotstuff [46] protocol in this line of work supports over 100k transactions per second. Unfortunately, the high throughput of Bullshark comes with a latency price due to the DAG-based consensus mechanism, which compares to the state-of-the-art leader-based BFT consensus protocols.

We introduce Shoal, a protocol-agnostic framework for enhancing Narwhal-based consensus. By incorporating leader replacement, Shoal achieves a round time that is significantly reduced. Moreover, the combination of properties of the DAG construction and the leader replacement mechanism enables the protocol to scale in all but extremely adversarial scenarios in practice, a property we name "prevalent responsiveness". It strictly subsumes the established and often desired "optimistic responsiveness" property for BFT consensus.

We evaluate Shoal instantiated with Bullshark, the fastest existing Narwhal-based consensus protocol, in an open-source Blockchain project and provide experimental evaluations demonstrating up to 40% latency reduction in the failure-free execution. We also evaluate the execution with failures against the vanilla Bullshark implementation.

CCS Concepts - Security and privacy → Distributed systems security

Keywords: Consensus Protocol, Byzantine Fault Tolerance, ACM Reference Format

Alexander Spiegelman, Balaji Arun, Rati Gelashvili, and Zekun Li
2023. Shoal: Improving DAG-BFT Latency And Robustness

1. Introduction

Byzantine fault tolerant (BFT) systems, including consensus protocols, form the core underpinning for blockchains. At a high level, a BFT-SMR enables a group of n parties to agree on a sequence of values, even if a bound of t up to $\frac{n}{2}$ of these parties are Byzantine (arbitrarily malicious). Over the last decade, significant progress has been made in improving the efficiency of these consensus protocols. In particular, non-equivocable round-based directed acyclic graph (DAG), a concept initially introduced by Aleph [21]. In this model, each validator runs its own DAG locally, and each vertex is assigned to $n-1$ vertices in the preceding round. Each vertex is disseminated via an efficient reliable broadcast implementation, ensuring that malicious validators cannot distinguish different vertices to different validators within the same round. This allows the DAG to be constructed without the need for consensus between validators. In their words, "Shoal is designed to be a DAG-based consensus protocol that can commit with a latency overhead of 3δ (where δ represents the round time)" [12]. Thus, Shoal can achieve linear communication complexity [17, 18] under optimistic conditions (such as honest leader).

Due to periods of network asynchrony, each validator may observe a slightly different portion of the DAG at any

CM

Cordial Miners: Fast and Efficient Consensus for Every Eventuality

Idit Keidar
Technion
Oded Naor
Technion and StarkWare
Ouri Poupko
Ben-Gurion University
Ehud Shapiro
Weizmann Institute of Science

Agree on the proposed values and ensure that the leader keeps miners busy. Finally, determine the next round. This approach results in two drawbacks. First, there is an uneven scheduling of work among the parties. While the leader is sending a proposal, the other parties' processors and their miners are not used. Second, there is no incentive across parties. Second, in typical leader-based protocols progress stops if the leader fails and until it is replaced. Several techniques proposed in the literature can mitigate these problems. These include the use of erasure coding techniques [2], [42] or the data availability scheme [26], [27], [29] to disseminate the data more efficiently.

Recently, a novel approach known as DAG-based BFT has emerged [5], [18], [28], [33], [34], [40], [46]. These protocols enable a partial ordering of parties and fast path protocols to demonstrate their low latency and resource efficiency. Our paper proposes to extend the DAG-based consensus to multiple leaders within a single round and commits all proposals in parallel. This is done by maintaining the DAG even if a party fails during a round. Consequently, these protocols have demonstrated improved throughput compared to their leader-based counterparts under the same framework [19], [20]. However, existing DAG-based protocols incur a high latency compared to their "leader-heavy" counterparts [22], [23], [30], [37], [51]. Is high latency inherent in DAG-based consensus? Addressing this question is the key goal of this paper.

All existing DAG-based consensus protocols [12, 13] have a problem of ordering vertices. In each round, every party can create a potential DAG vertex containing transactions and edges pointing to vertices from previous rounds. These protocols rely on committing a designated "leader vertex" and order other non-leader vertices based on the leader vertex. The question is which leaders are designated and how fast the leader vertices directly influences the commit latency.

Supporting a leader vertex in each round, State-of-the-art protocols [12, 13] require many message or memory operations per round. In fact, depending on a leader vertex in each round particularly difficult. In their words, "Shoal is designed to be a DAG-based consensus protocol that can commit with a latency overhead of 3δ (where δ represents the round time)" [12]. Thus, Shoal can achieve linear communication complexity [17, 18] under optimistic conditions (such as honest leader).

Most of these protocols design a DAG where each designated leader vertex is mainly responsible for proposing transactions and driving the protocol forward while other parties

Mysticeti

MYSTICETI: Reaching the Latency Limits with Uncertified DAGs

Kushal Babel*, Andrej Chursin*, George Danezis†, Anastasis Kichidis†, Lefteris Kokoris-Kogias*, Arun Koshy*, Alberto Sonnino†, Mingwei Tian†

*Cornell Tech, †C3, *Mythen Labs, †University College London (UCL), †IST Austria

Abstract—We introduce MYSTICETI-C, the first DAG-based Byzantine consensus protocol to achieve the lower bounds of latency. It reaches the limits of consensus in terms of latency and resource efficiency. MYSTICETI-C achieves high throughput improvement by avoiding consensus overhead. It uses DAGs to achieve consensus with a novel commit rule such that every block can be committed without delay, resulting in a latency limit of 4 ms in a crash-recovered and under crash failure. We further extend MYSTICETI-C to MYSTICETI-FPC, which incorporates a fast commit path to achieve high throughput. MYSTICETI-FPC is the first DAG-based consensus protocol to reach the limits of latency. MYSTICETI-FPC minimizes the number of signatures and messages by weaving the fast path into the consensus logic. MYSTICETI-FPC also minimizes the latency subsequently resulting in better performance. We prove the safety and liveness of MYSTICETI-C. We evaluate both MYSTICETI-C and protocols with them with state-of-the-art consensus and fast path protocols to demonstrate their low latency and resource efficiency. MYSTICETI-C is the first Byzantine consensus protocol to achieve WAN latency of 4.6 ms for consensus commit while simultaneously achieving a throughput of 200K TPS. Finally, we report on integrating MYSTICETI-C as the consensus layer in the Sul blockchain [67], resulting in over 4x latency reduction.

I. INTRODUCTION
Several recent blockchains, such as Sul [67], [12], have adopted DAG-based consensus [15], [56], [34], [30], [17], [12], [58], [44]. By design, these consensus protocols scale well in terms of throughput, with a performance of 100K tps of transaction processing per second. They are also highly asynchronous [33], [25]. This, however, comes at a high latency of around 2-3 seconds, which can hinder user experience and prevent low-latency applications.

MYSTICETI-C is the first of uncertified DAG-based consensus protocols [12, 13], where each vertex is delivered through consistent broadcast [14], have high latency for three main reasons: (1) the consensus process requires multiple rounds to reach a consensus between validators; (2) the signatures and re-broadcast certificates. This leads to higher latency than traditional consensus protocols [31], [64], [15]; (2) blocks commit on a "per block" basis, which requires a large amount of resources and results in low throughput. Additionally, they are fragile to faults and implementation mistakes due to their complex nature.

This work presents MYSTICETI, a family of DAG-based protocols allowing to safely commit distributed transactions in a Byzantine setting that focuses on low-latency and low-CPU overhead. MYSTICETI consists of two variants: (1) a consensus protocol based on a threshold logical clock [29] of 2 δ . MYSTICETI, that commits every block as early as it can be decided. MYSTICETI-C solves all of the above challenges as (1) it is the first safe DAG-based consensus protocol that does not require explicit certificates, committing blocks within the

Techniques

- Many leaders per round
- Leaders every round
- Uncertified DAG

Discussion



Shoal/shoal++

- Low latency
- Easier synchroniser
- Leverage existing code

Sailfish/BBCA

- Lower latency
- Easy synchroniser
- Flexible

CM/Mysticeti

- Even Lower latency
- Graceful crash faults
- Simpler, less CPU

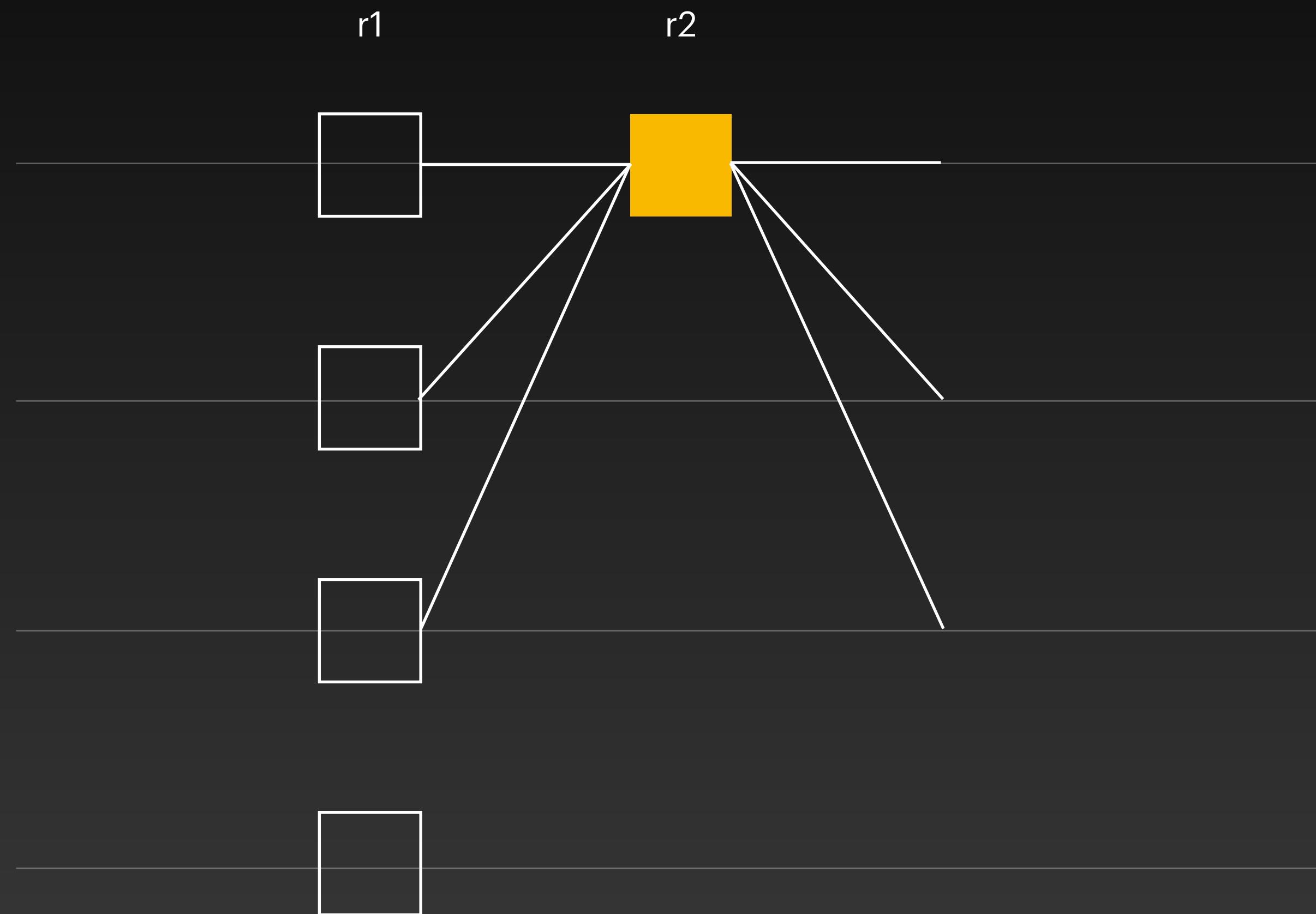
Research Questions

1. Network model?
2. BFT testing?
3. Consensus-exec interface?
4. Storage architecture?
5. Block synchroniser?

Lessons Learned

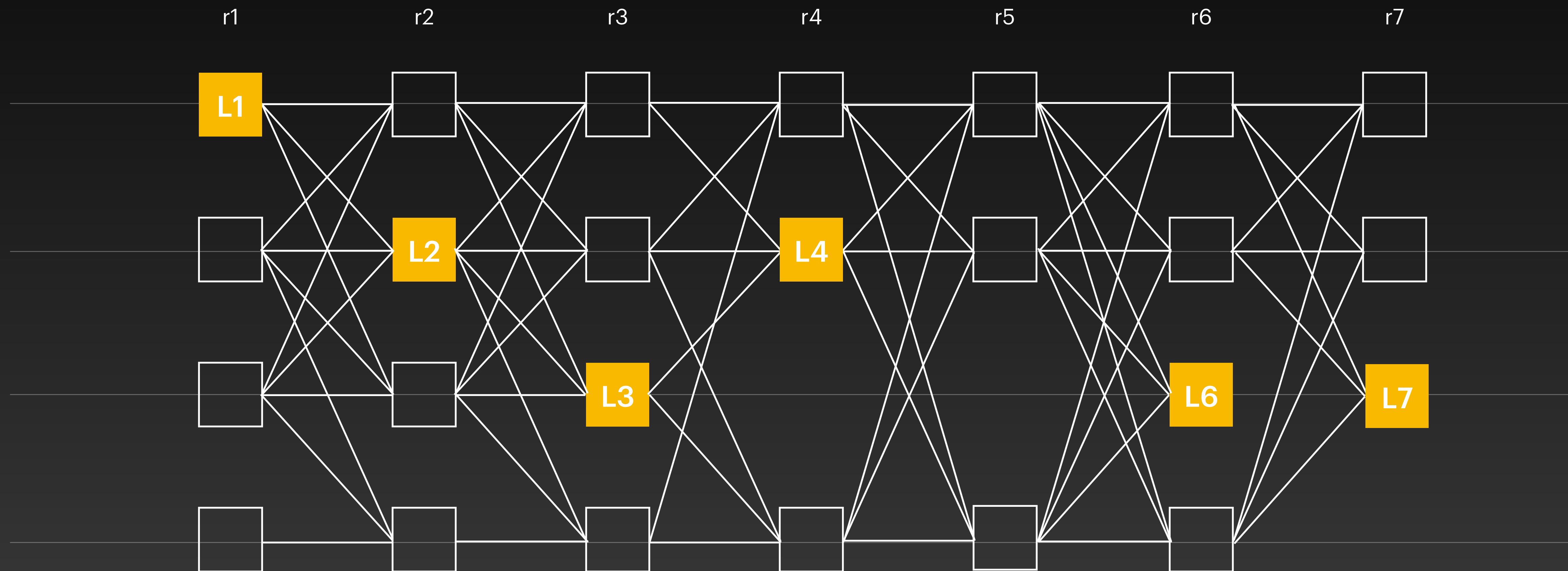
1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy
6. Epoch change is not an add-on

Uncertified DAG

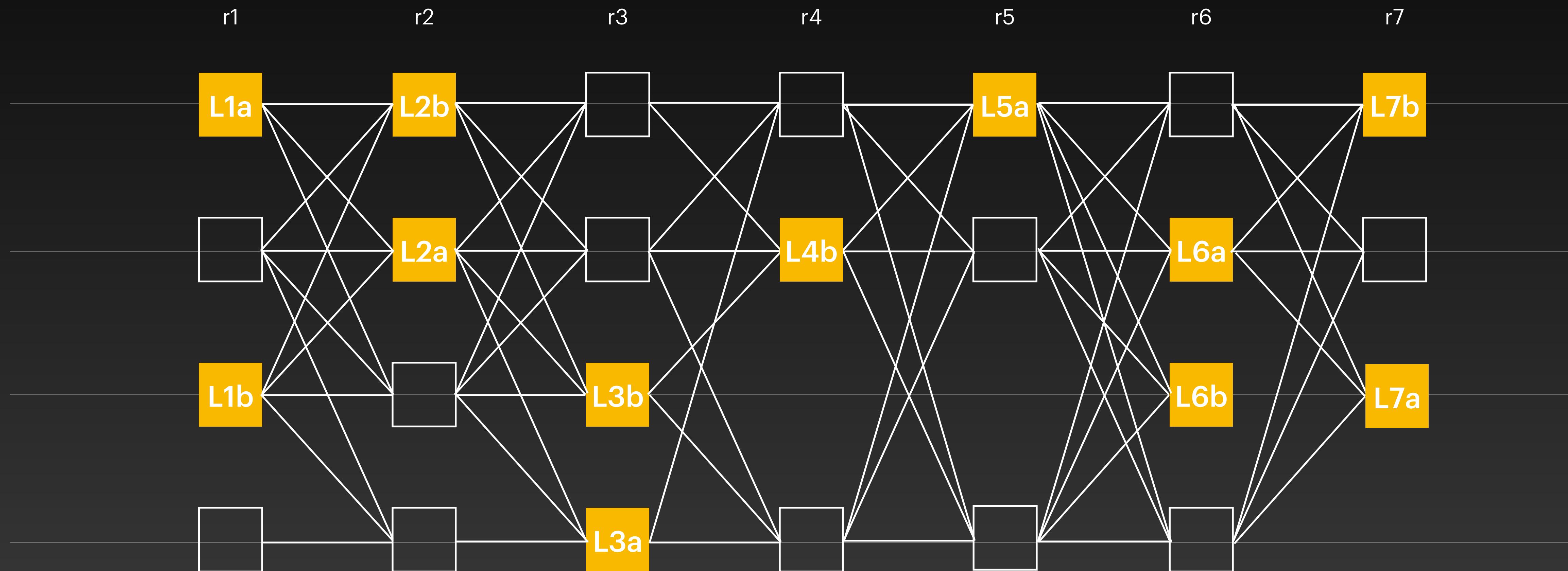


- Round number
- Author
- Payload (transactions)
- Signature

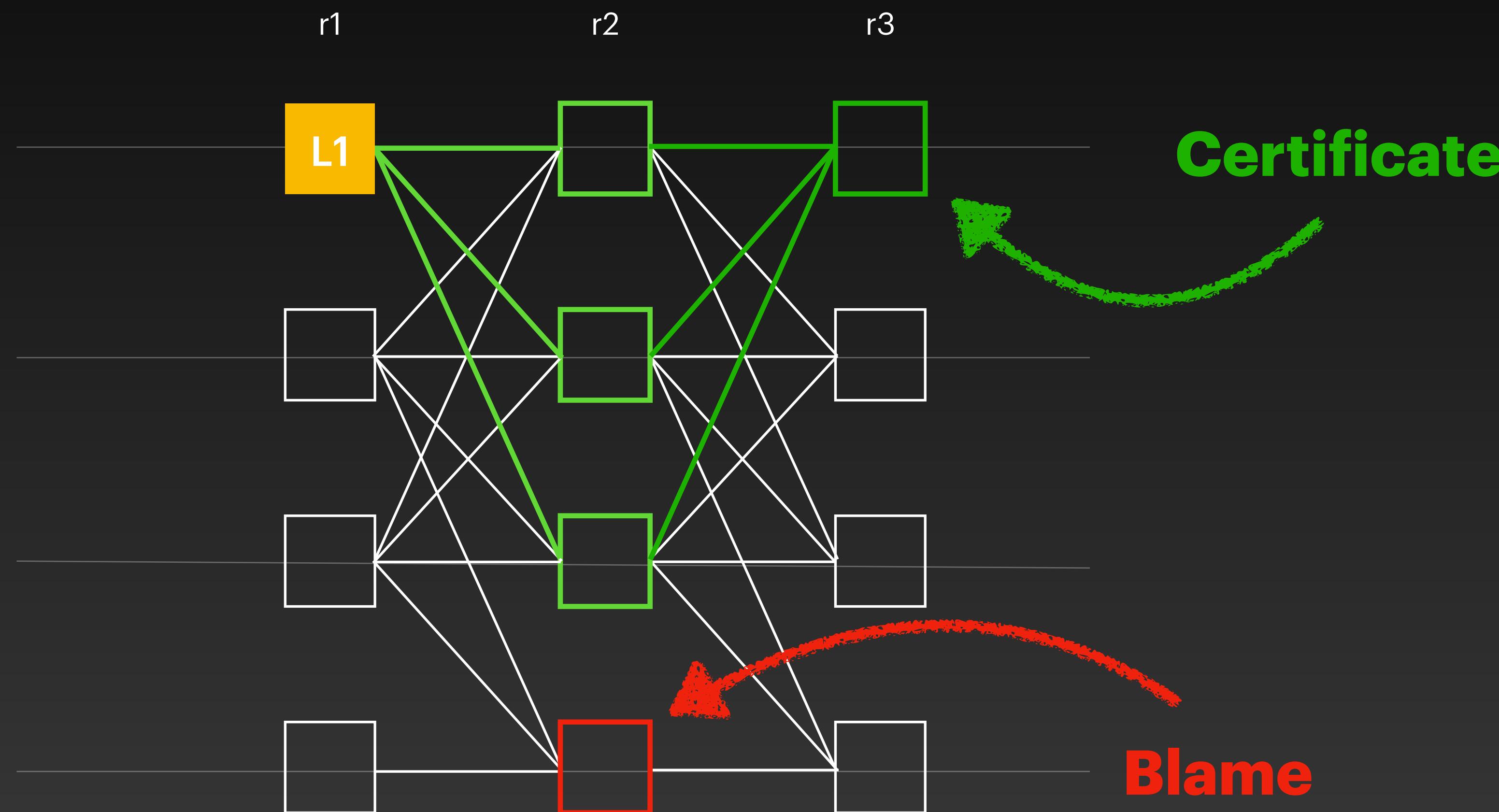
Uncertified DAG



Uncertified DAG



Interpreting DAG Patterns



Direct Decision Rule

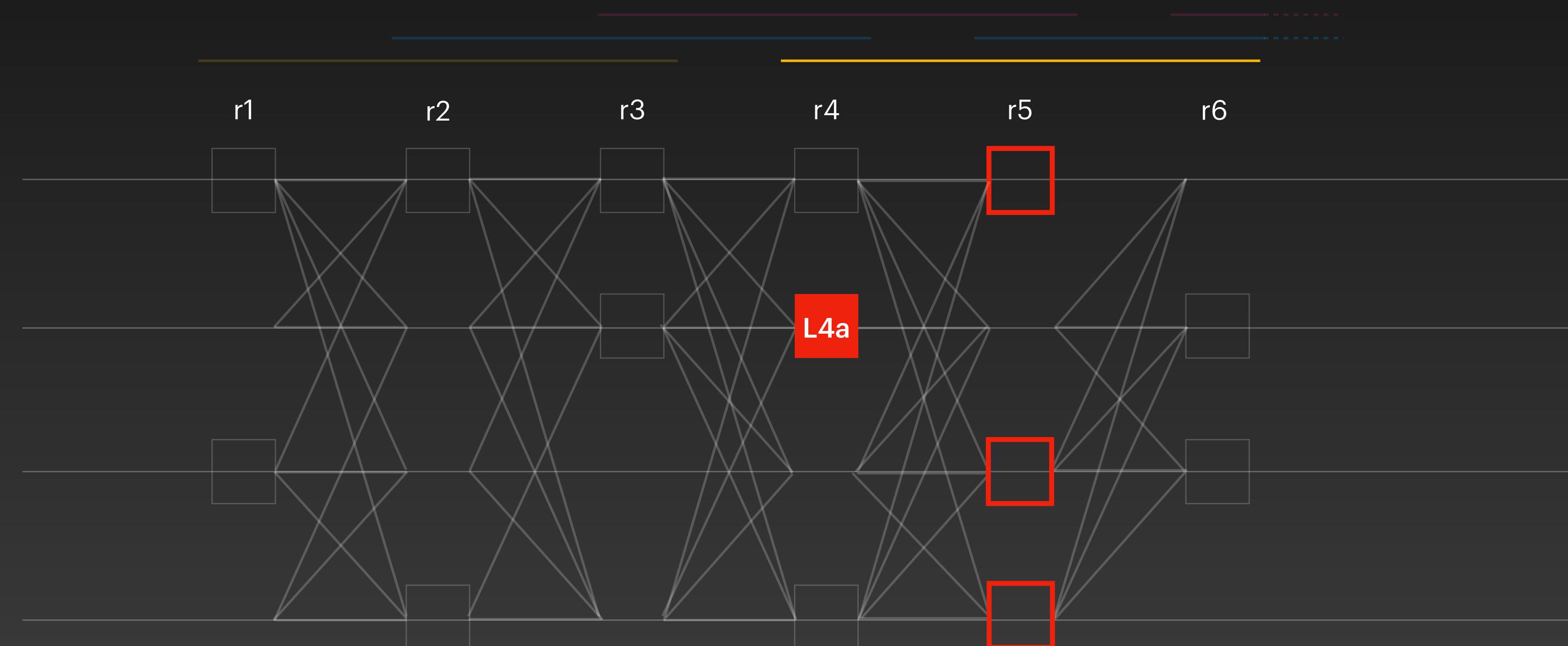
On each leader starting from highest round:

- **Skip** if $2f+1$ blames
- **Commit** if $2f+1$ certificates
- **Undecided** otherwise

Direct Decision Rule

On each leader starting from highest round:

- **Skip** if $2f+1$ blames
- **Commit** if $2f+1$ certificates
- **Undecided** otherwise

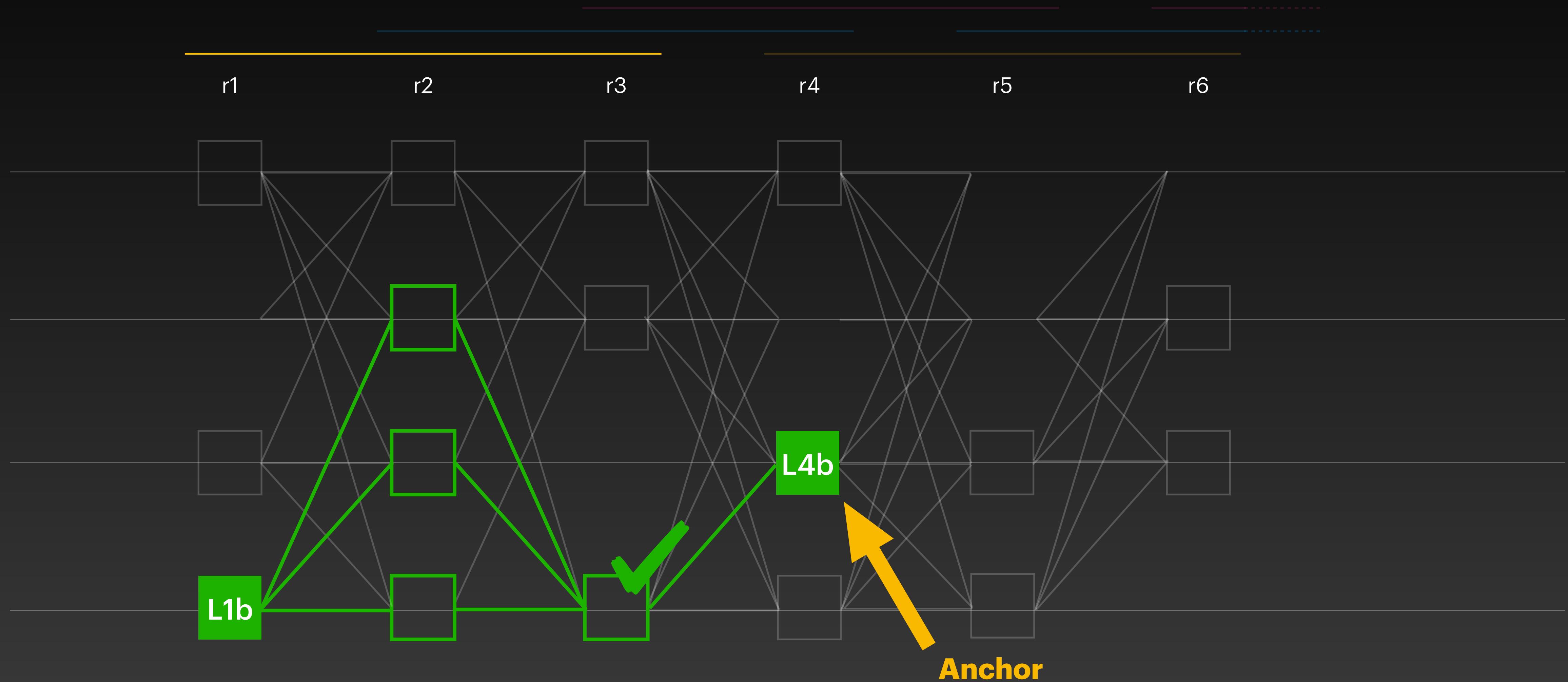


Direct Decision Rule

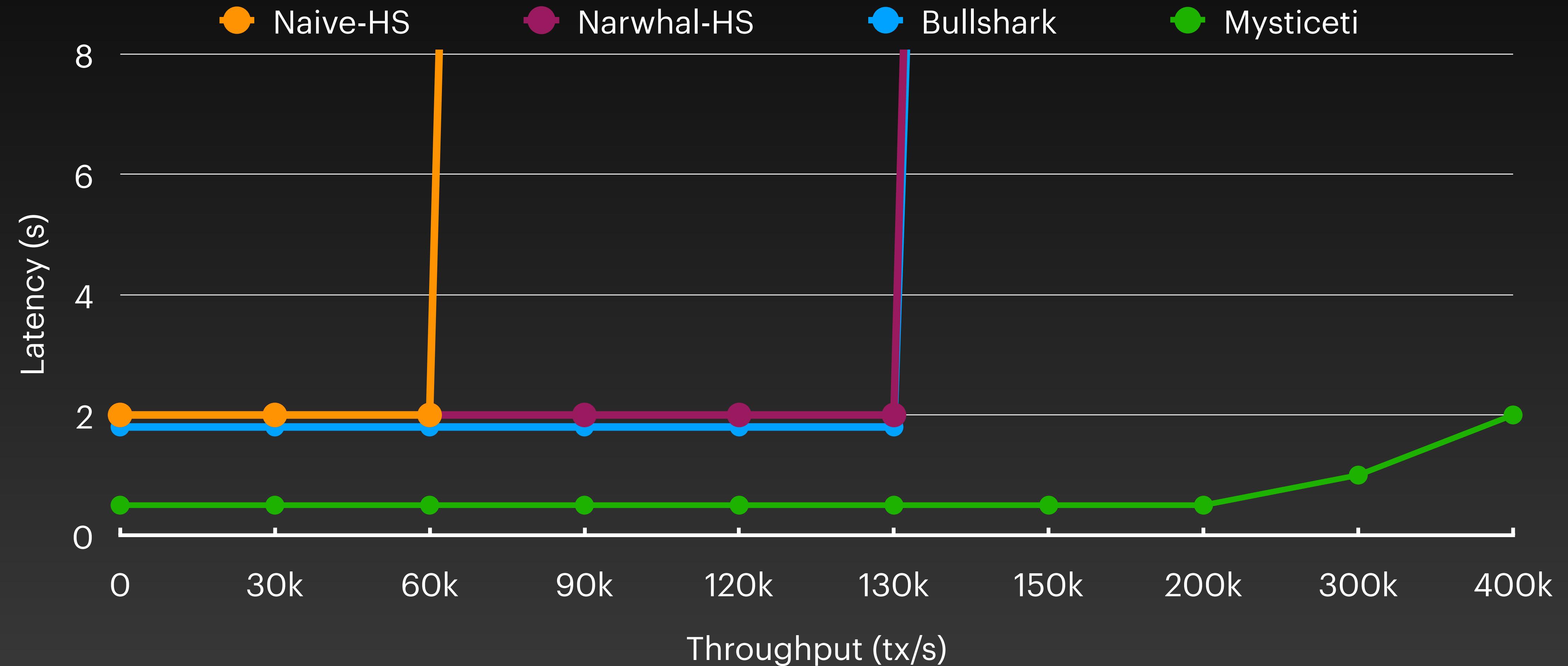
On each leader starting from highest round:

- **Skip** if $2f+1$ blames
 - **Commit** if $2f+1$ certifies
 - **Undecided** otherwise

Apply Indirect Rule



Performance



Engineering Benchmarks

Protocol	Committee	Load/TPS	P50	P95
Bullshark	137	5k	2.89 s	4.60 s
Mysticeti	137	5k	397 ms	690 ms

We ran it for 24h and it looks good 👍

Research Questions

1. Network model?
2. BFT testing?
3. Consensus-exec interface?
4. Storage architecture?
5. Block synchroniser?
6. Realistic benchmarks?
7. Efficient reads?

Lessons Learned

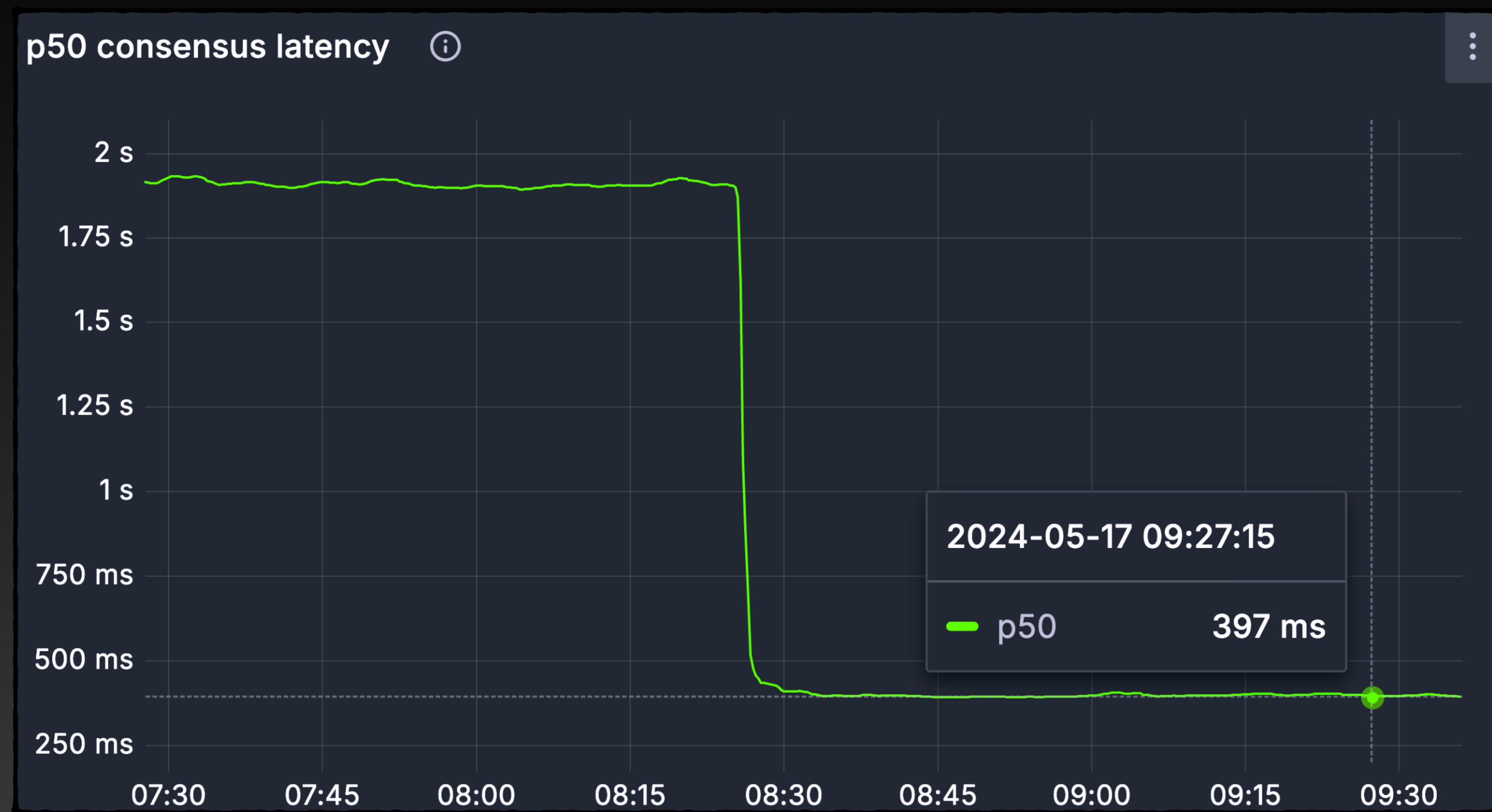
1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy
6. Epoch change is not an add-on

Testing Strategy



- Compare performance & robustness
- Test mainnet change bullshark -> mysticeti
- Prepare for the worst mysticeti -> bullshark

The Sui Mainnet



Conclusion

1. Consensus + Naive mempool
2. Consensus + Better mempool
3. Flatten consensus into mempool
4. Remove flattening overhead

alberto@mystenlabs.com

EXTRA:

Research in Industry

Projects Roadmap



Dmitri Perelman Oct 18th at 5:55 AM

In tomorrow's Research <> Core Eng syncup, [@Mark Logan](#) is going to share top of mind of Core Eng pain points and current struggles. See you 



2



Projects Roadmap

 **Dmitri Perelman** Oct 18th at 5
In tomorrow's Research <> C
going to share top of mind of
struggles. See you 

 2 

< **Thread** # sui-core-internal

 **Dmitri Perelman** Oct 18th at 5:55 AM
In tomorrow's Research <> Core Eng syncup, [@Mark Logan](#) is
going to share top of mind of Core Eng pain points and current
struggles. See you 

 2 

2 replies

 **John Martin** Oct 18th at 6:16 AM
Can I get an invite to this 

 2 

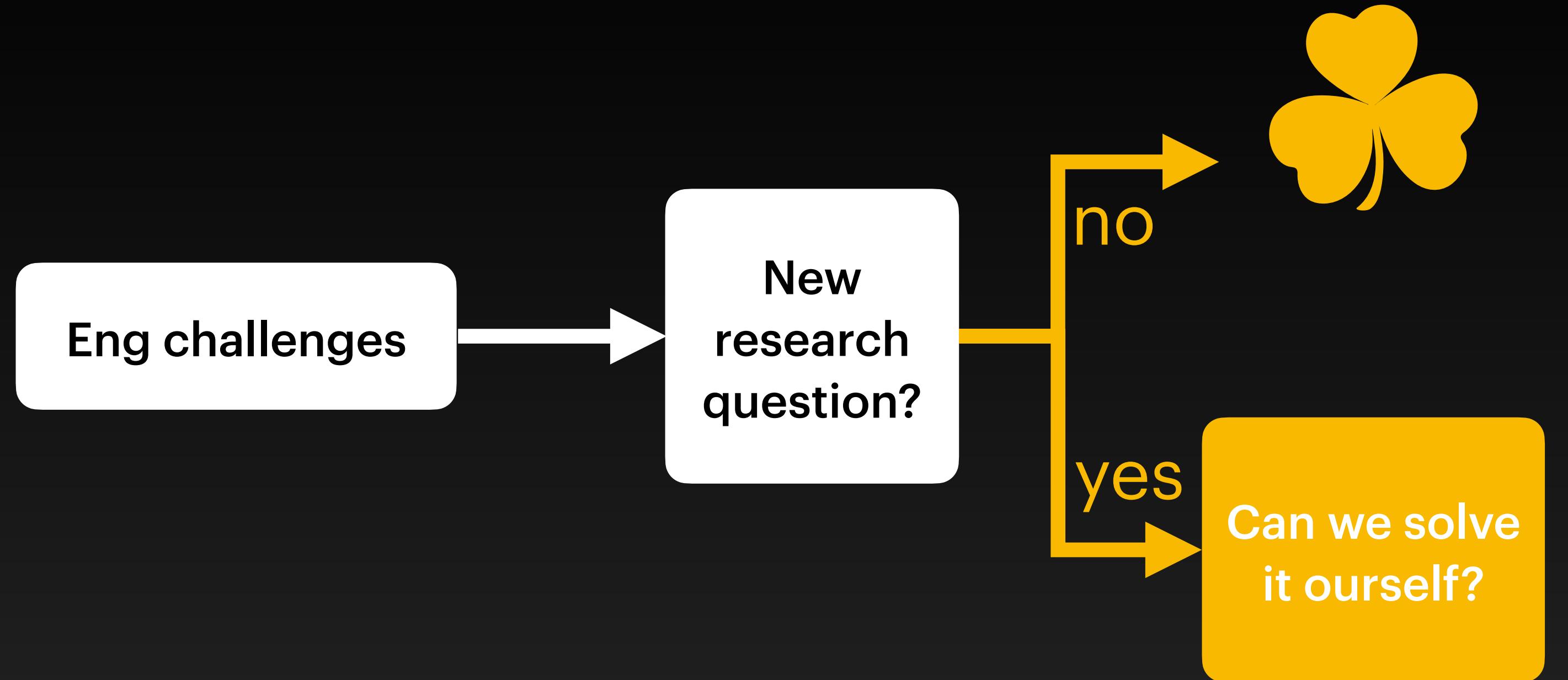
 **Dmitri Perelman** Oct 18th at 7:36 AM
You're in the invite list!

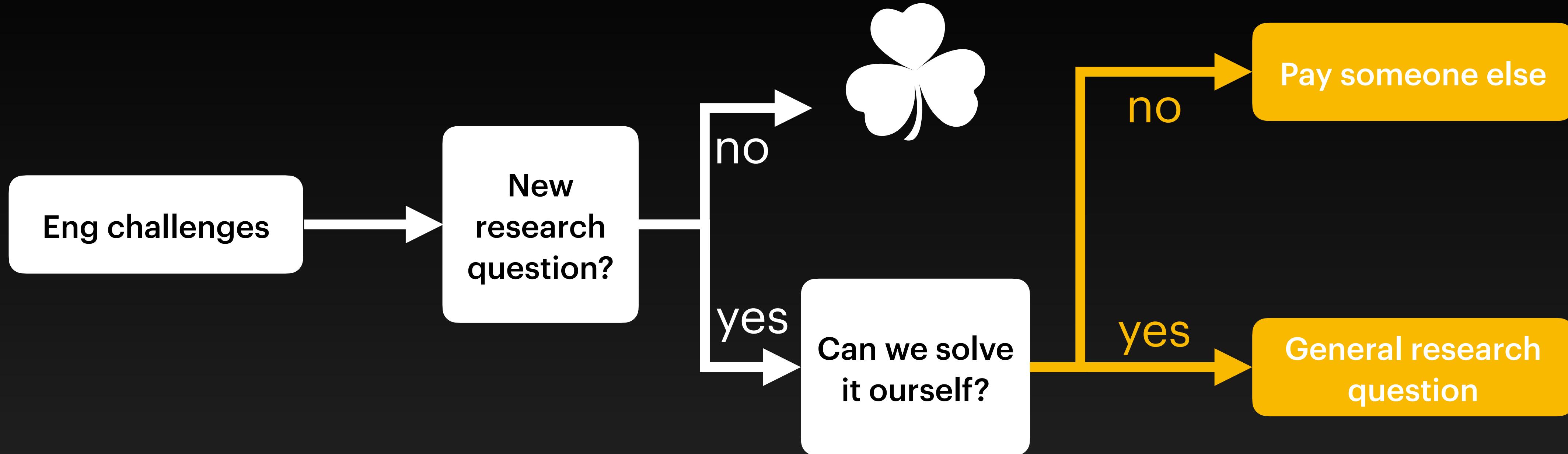
 1 

Eng challenges



New
research
question?

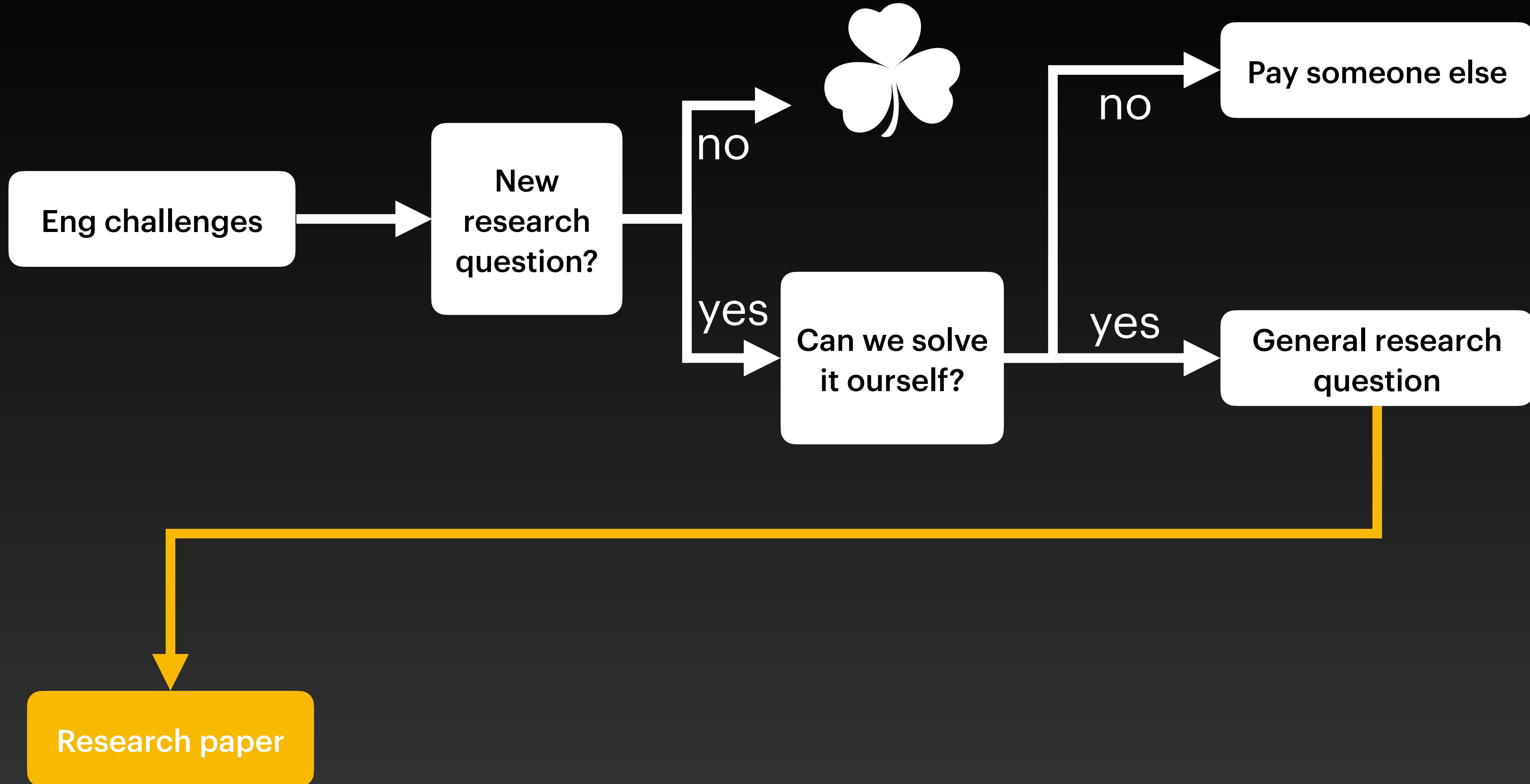


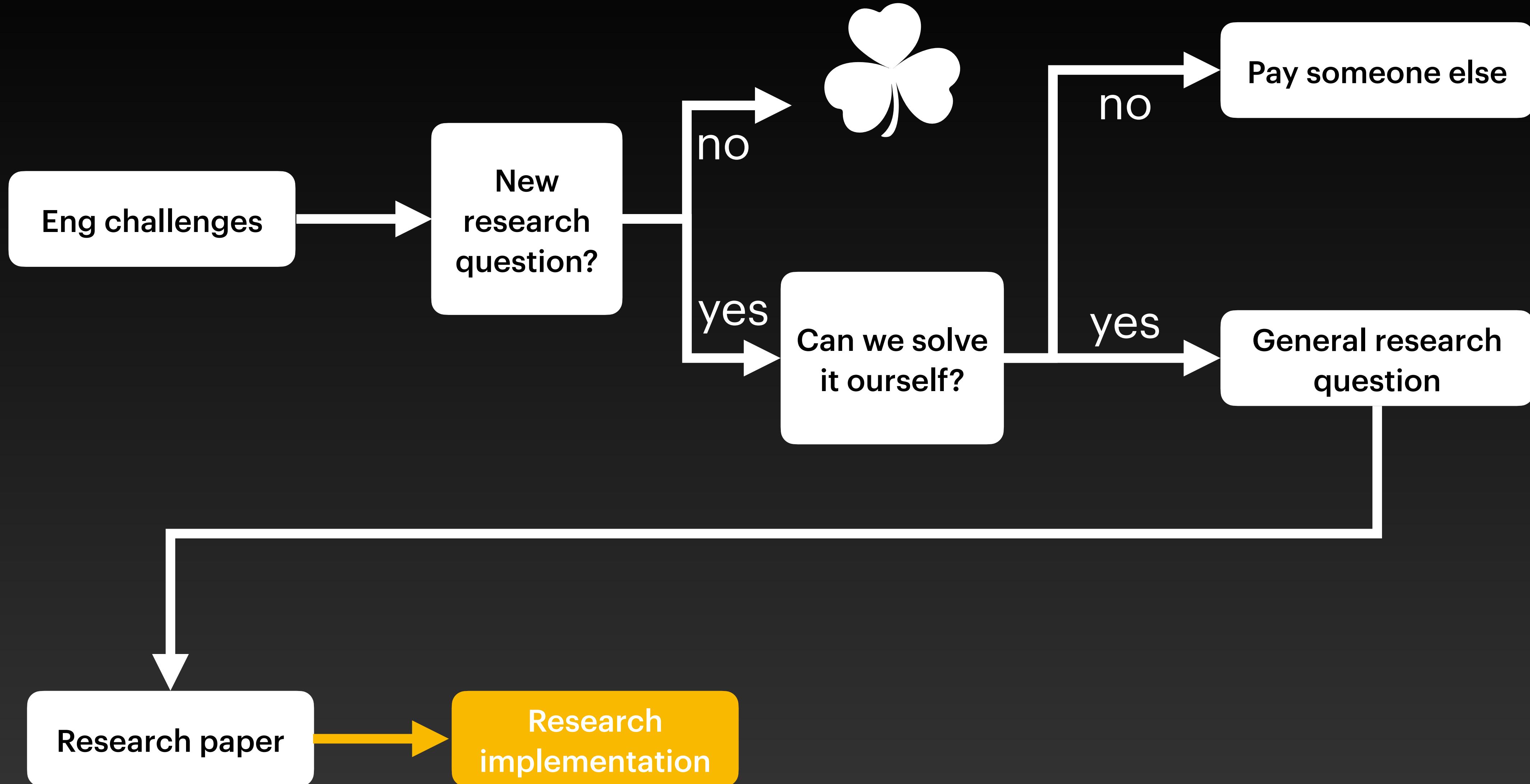


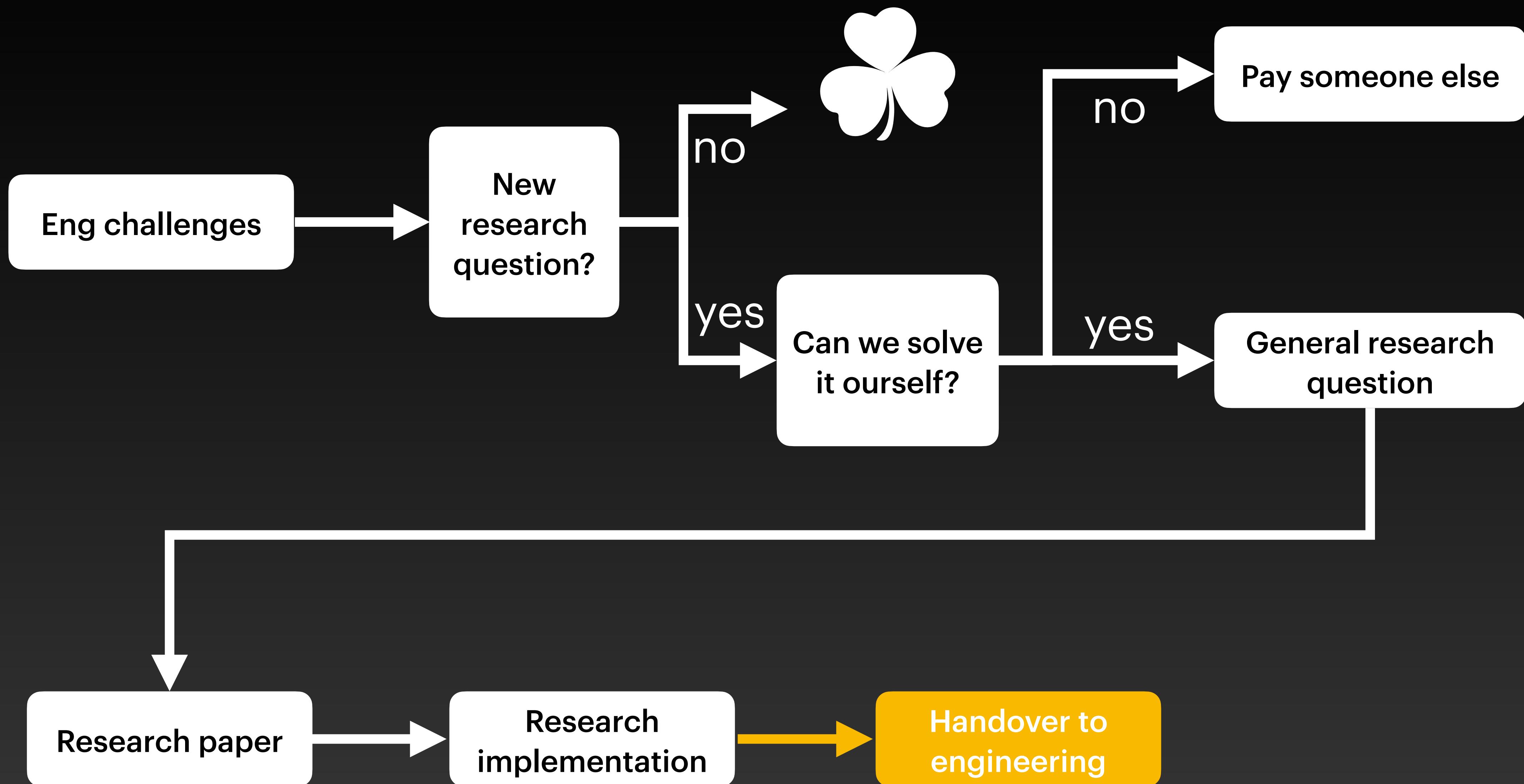
Research Gifts

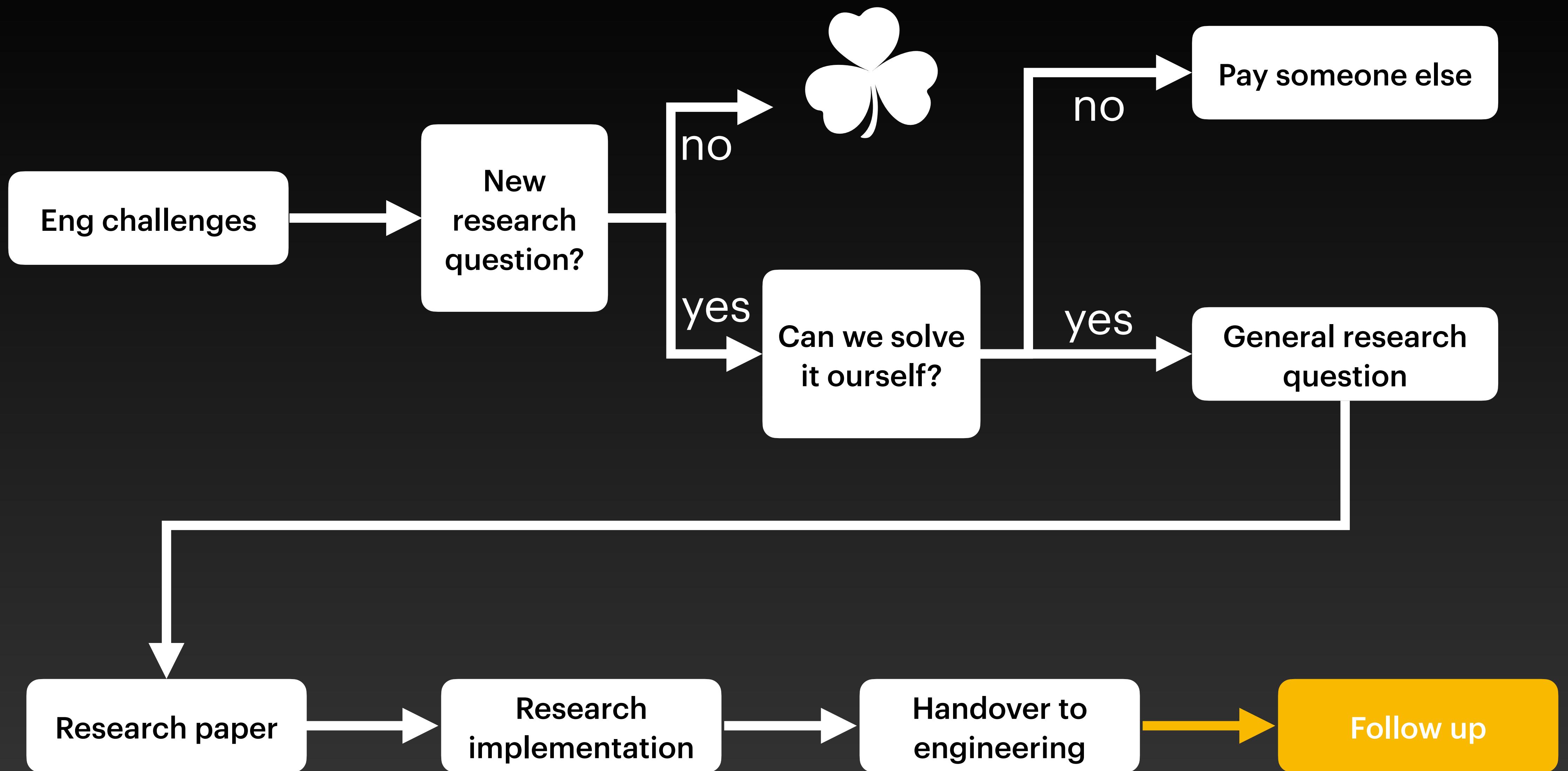


(please keep it short)









[Chainspace: A Sharded Smart Contracts Platform](#)

NDSS • Adopted by chainspace.io

[Coconut: Threshold Issuance Selective Disclosure ...](#)

NDSS • Adopted by chainspace.io, Ketl, Nym, ...

[Replay Attacks and Defenses against Cross-shard ...](#)

EuroS&P • Adopted by chainspace.io

[FastPay: High-Performance Byzantine Fault Tolerant ...](#)

AFT • Adopted by Sui, Linera

[Twins: BFT Systems Made Robust](#)

OPODIS • Adopted by Diem, Aptos, Chainlink

[Fraud Proofs: Maximising Light Client Security and ...](#)

FC • Adopted by Ethereum, Celestia

[Jolteon and Ditto: Network-Adaptive Efficient Consensus ...](#)

FC • Adopted by Flow, Diem, Aptos, Monad

[Be Aware of Your Leaders](#)

FC • Adopted by Diem, Aptos

[Subset-optimized BLS Multi-signature with Key Aggregation](#)

FC • Adopted by Fastcrypto

[Narwhal and Tusk: A DAG-based Mempool and Efficient ...](#)

EuroSys • Adopted by Sui, Aptos, Fleek, Aleo

Best paper award

[Bullshark: DAG BFT Protocols Made Practical](#)

CCS • Adopted by Sui, Aleo, Fleek

[Zef: Low-latency, Scalable, Private Payments](#)

WPES • Adopted by Linera

[Parakeet: Practical Key Transparency for End-to-End ...](#)

NDSS • Adopted by WhatsApp

IETF Applied Networking Research Prize

[HammerHead: Leader Reputation for Dynamic Scheduling](#)

ICDCS • Adopted by Sui

[Fastcrypto: Pioneering Cryptography via Continuous ...](#)

LTB • Adopted by Fastcrypto

[Sui Lutris: A Blockchain Combining Broadcast and ...](#)

CCS • Adopted by Sui

Distinguished paper award

[Mysticeti: Reaching the Limits of Latency with Uncertified ...](#)

NDSS • Adopted by Sui

Research Questions

1. Network model?
2. BFT testing?
3. Consensus-exec interface?
4. Storage architecture?
5. Block synchroniser?
6. Realistic benchmarks?
7. Efficient reads?

Lessons Learned

1. Modularisation is a design strategy
2. Tasks-threads relationship
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy
6. Epoch change is not an add-on
7. Writing papers to explore designs

EXTRA:

Benchmarks

Implementation

- Written in Rust
- Networking: Tokio (TCP)
- Storage: custom WAL
- Cryptography: ed25519-consensus

<https://github.com/mystenlabs/mysticeti>

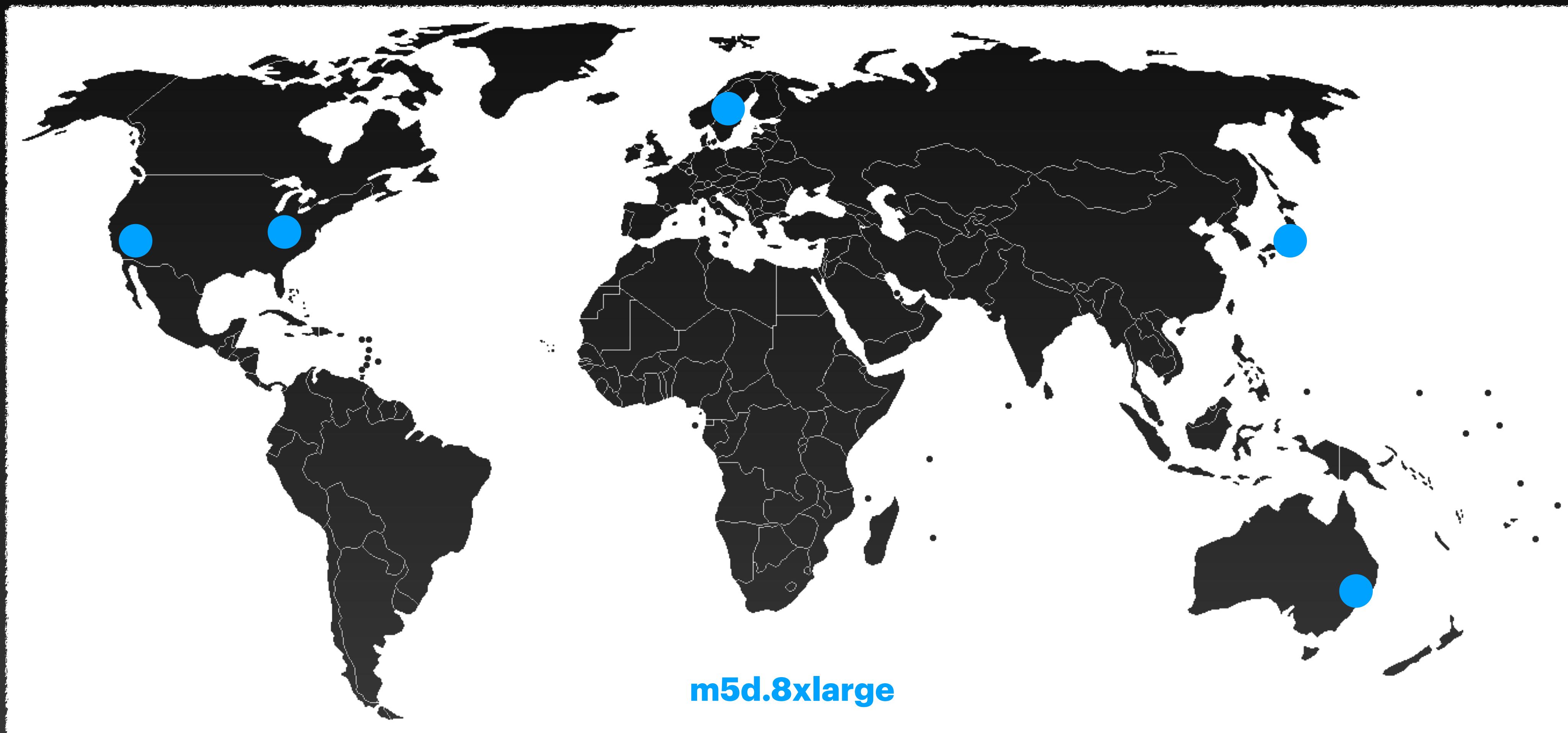
Implementation

- Synchronous core
- One Tokio task per peer (limiting resource usage)
- DTE simulator

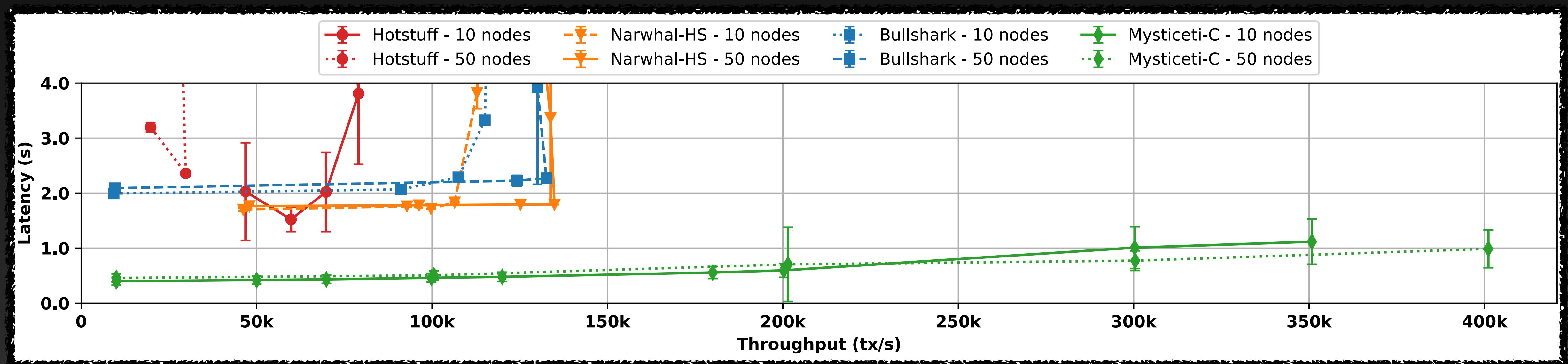
<https://github.com/mystenlabs/mysticeti>

Evaluation

Experimental setup on AWS



Prototype Benchmarks



Prototype Benchmarks

