

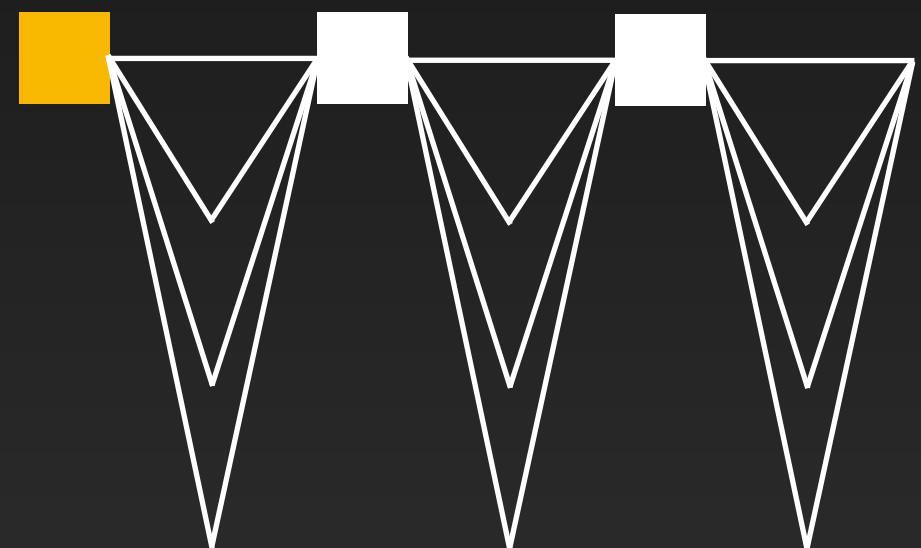
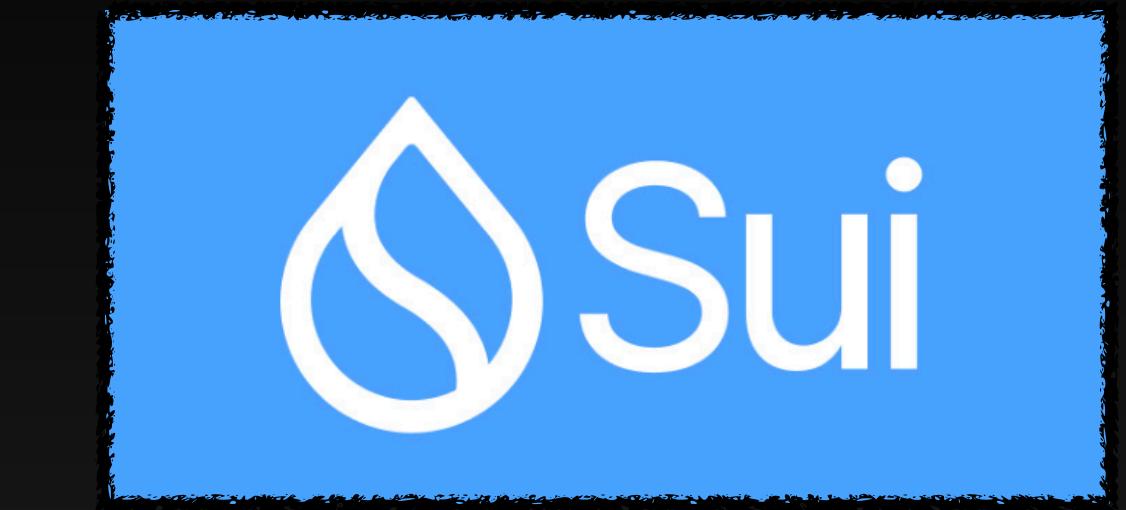
# The Evolution of Sui

## From Academic Paper to Mainnet

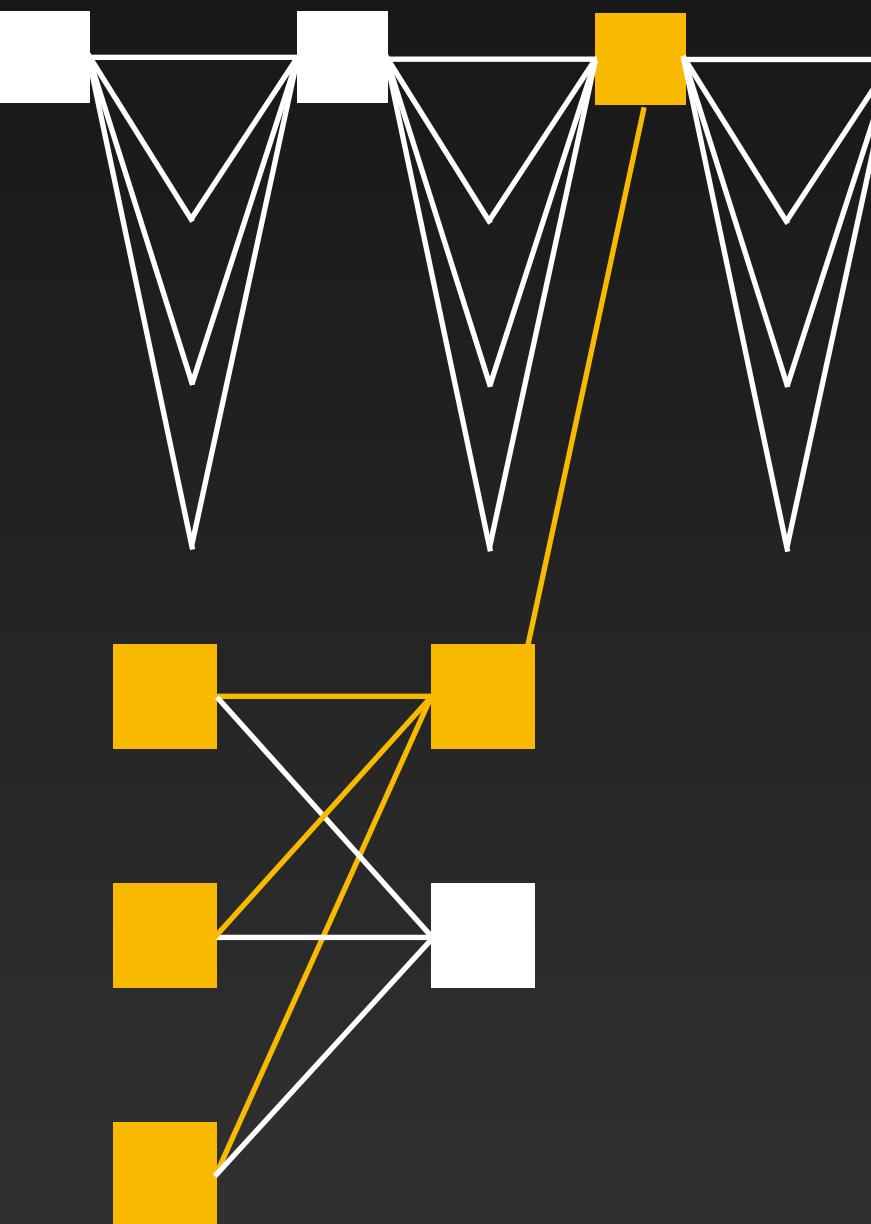
2019



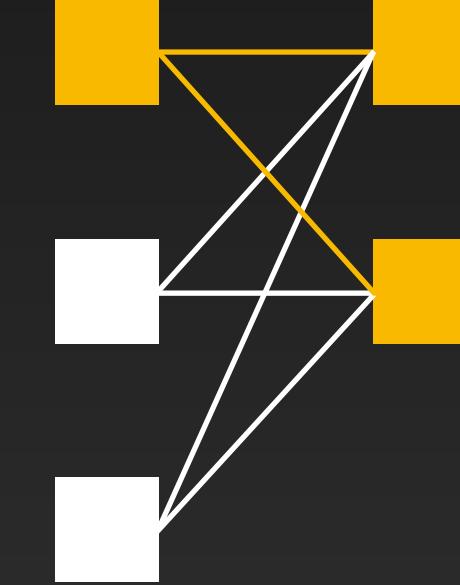
2024



**HotStuff**



**HotStuff + Mempool**

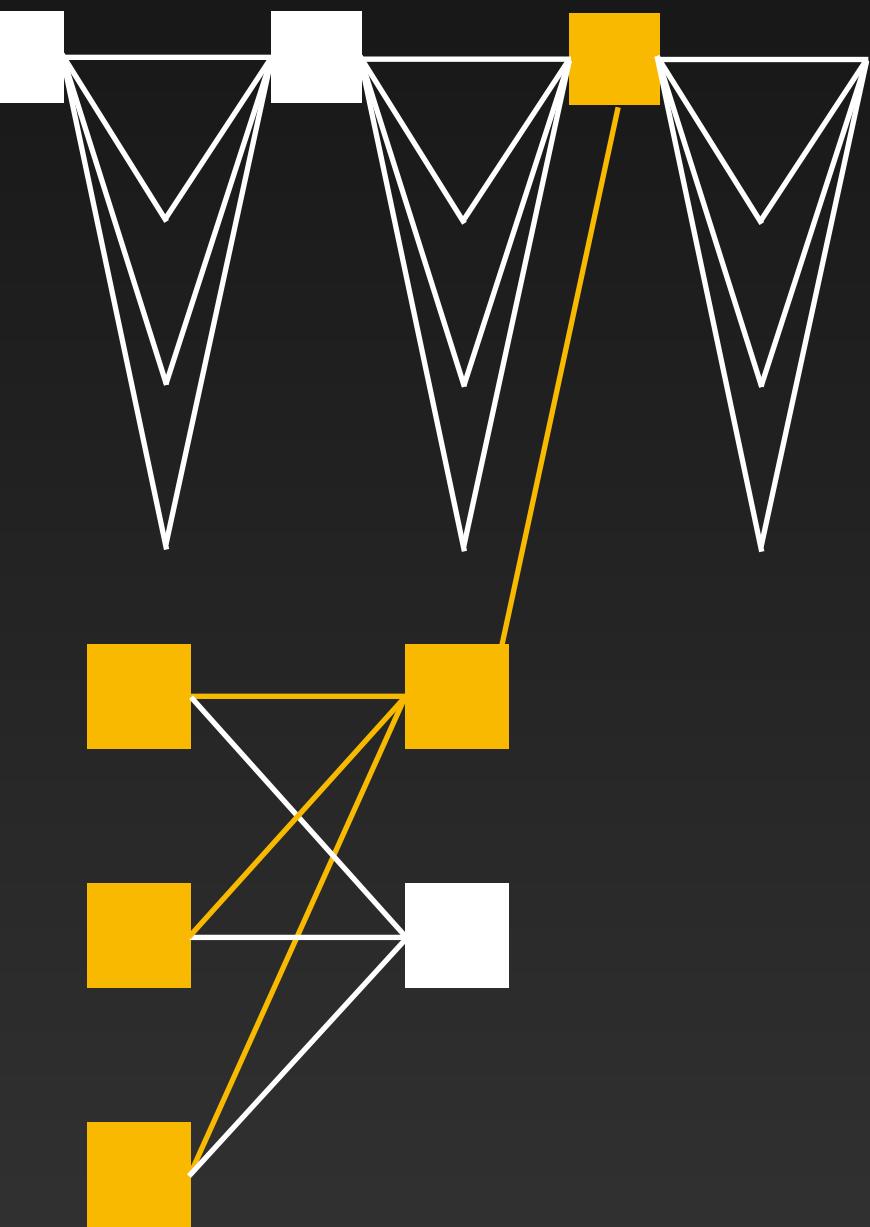


**Bullshark,  
Mysticeti**

2019

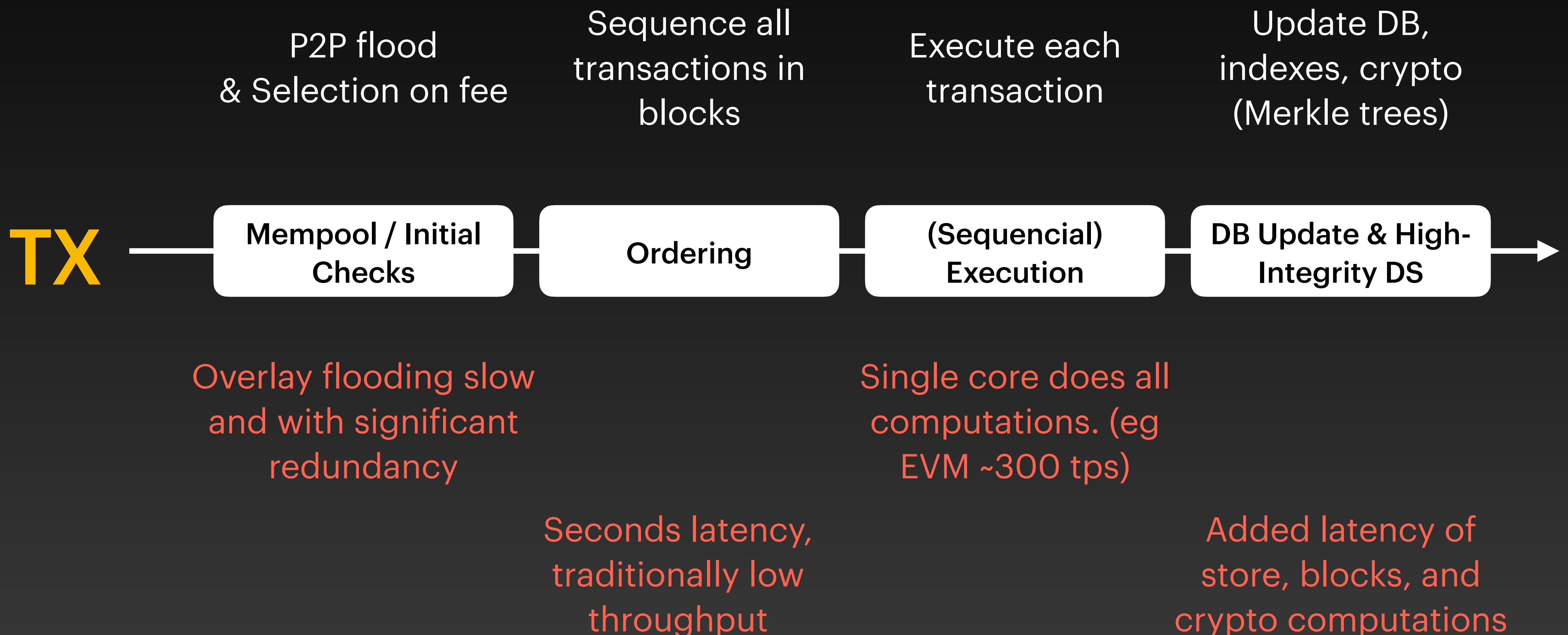


2024



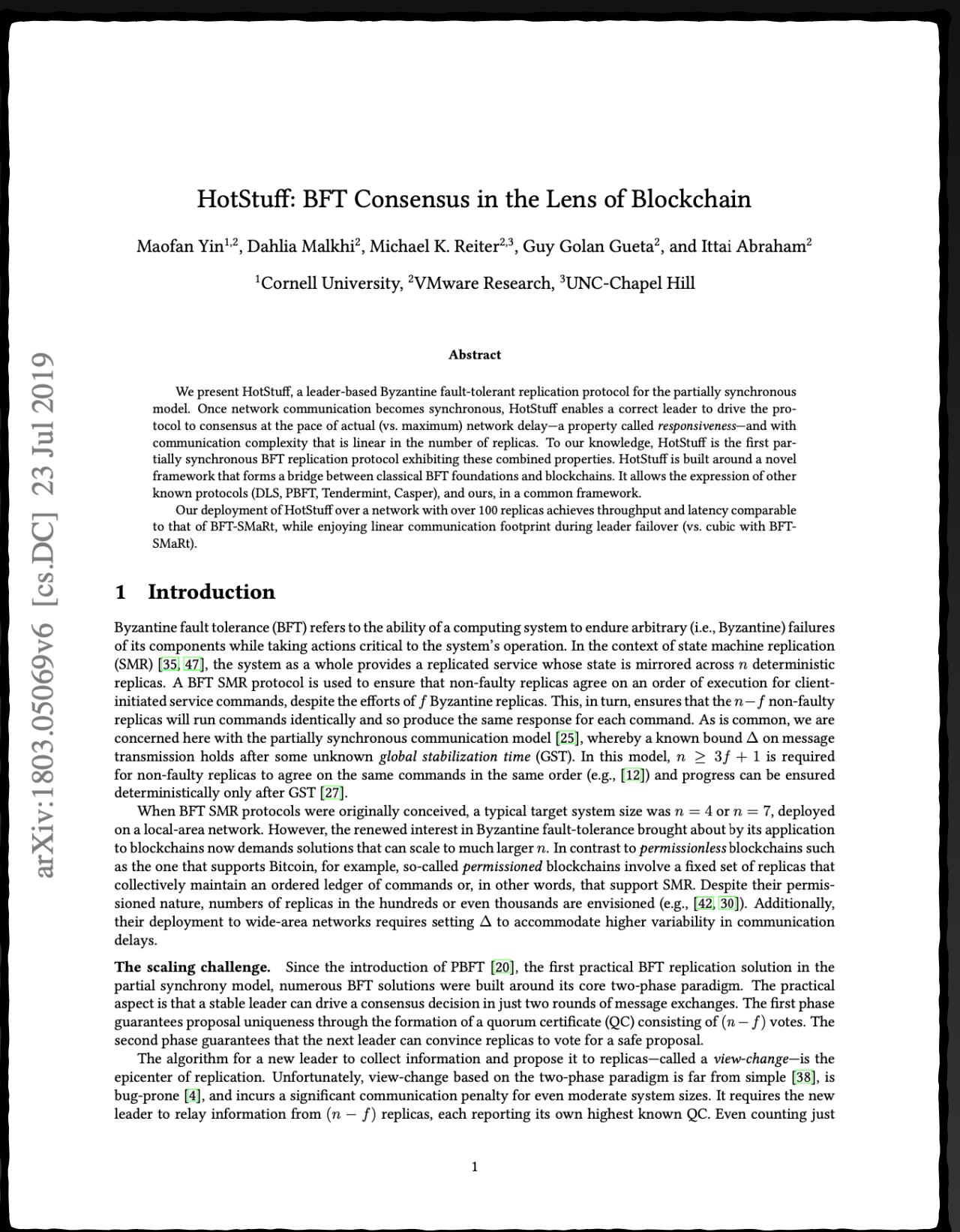
- Lessons learned
- Open research challenges

# Typical Blockchain



# Libra, 2019

## HotStuff



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

### HotStuff: BFT Consensus in the Lens of Blockchain

Maofan Yin<sup>1,2</sup>, Dahlia Malkhi<sup>2</sup>, Michael K. Reiter<sup>2,3</sup>, Guy Golan Gueta<sup>2</sup>, and Ittai Abraham<sup>2</sup>

<sup>1</sup>Cornell University, <sup>2</sup>VMware Research, <sup>3</sup>UNC-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  $\Delta$  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  $\Delta$  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

# Libra, 2019

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

HotStuff: BFT Consensus in the Lens of Blockchain

Maofan Yin<sup>1,2</sup>, Dahlia Malkhi<sup>2</sup>, Michael K. Reiter<sup>2,3</sup>, Guy Golan Gueta<sup>2</sup>, and Ittai Abraham<sup>2</sup>

<sup>1</sup>Cornell University, <sup>2</sup>VMware Research, <sup>3</sup>UNC-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  $\Delta$  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  $\Delta$  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

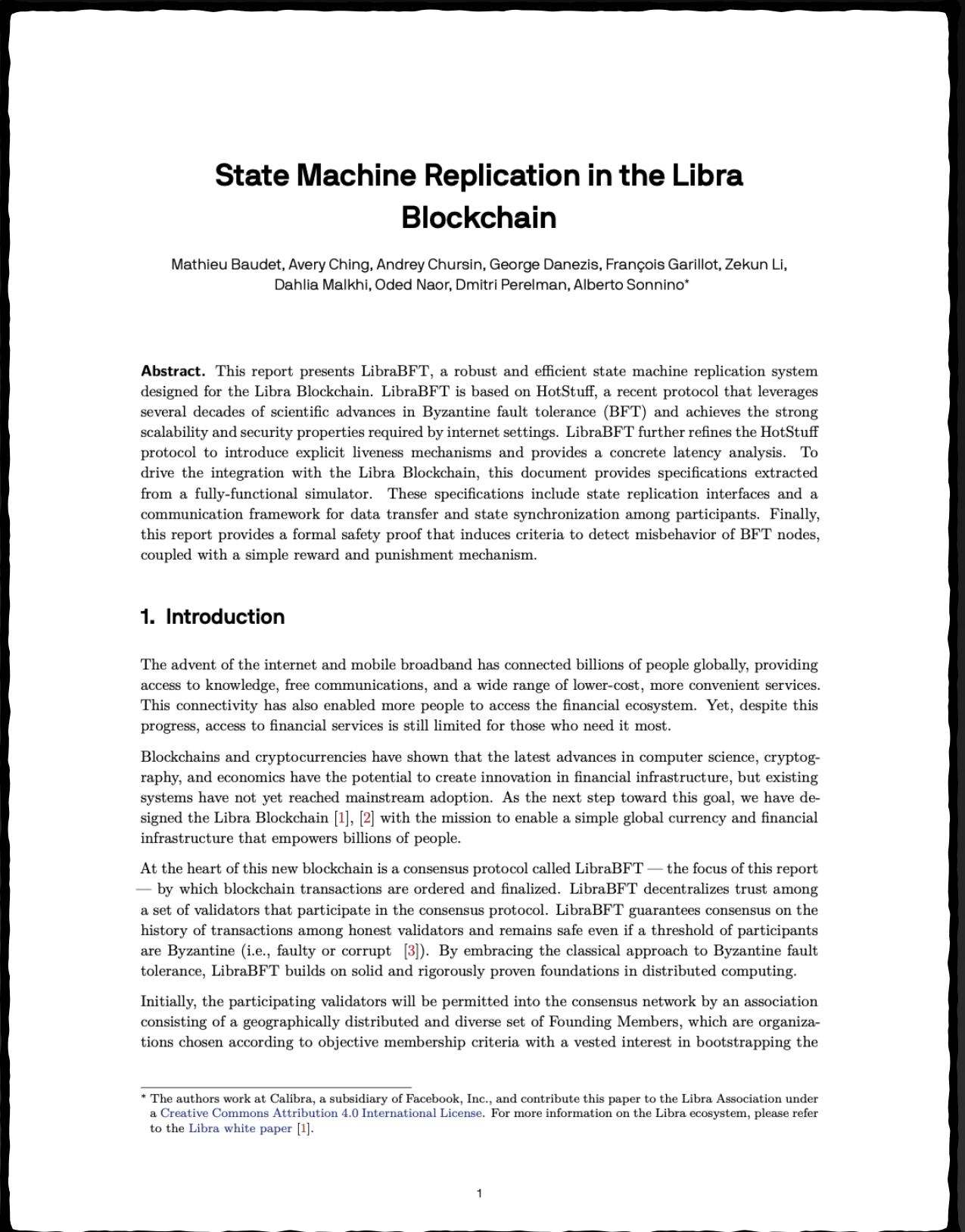
✓ Clearly isolated components

## HashGraph

✗ Impossible to garbage collect

✗ Unclear block synchroniser

# The first 6 months...



## SMR in the Libra Blockchain

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

# Research Questions

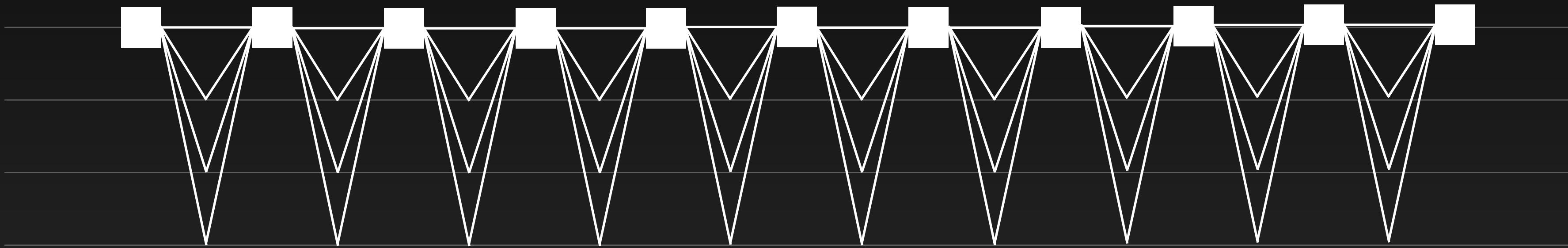
1. Network model?

# Lessons Learned

1. Modularisation is a design strategy

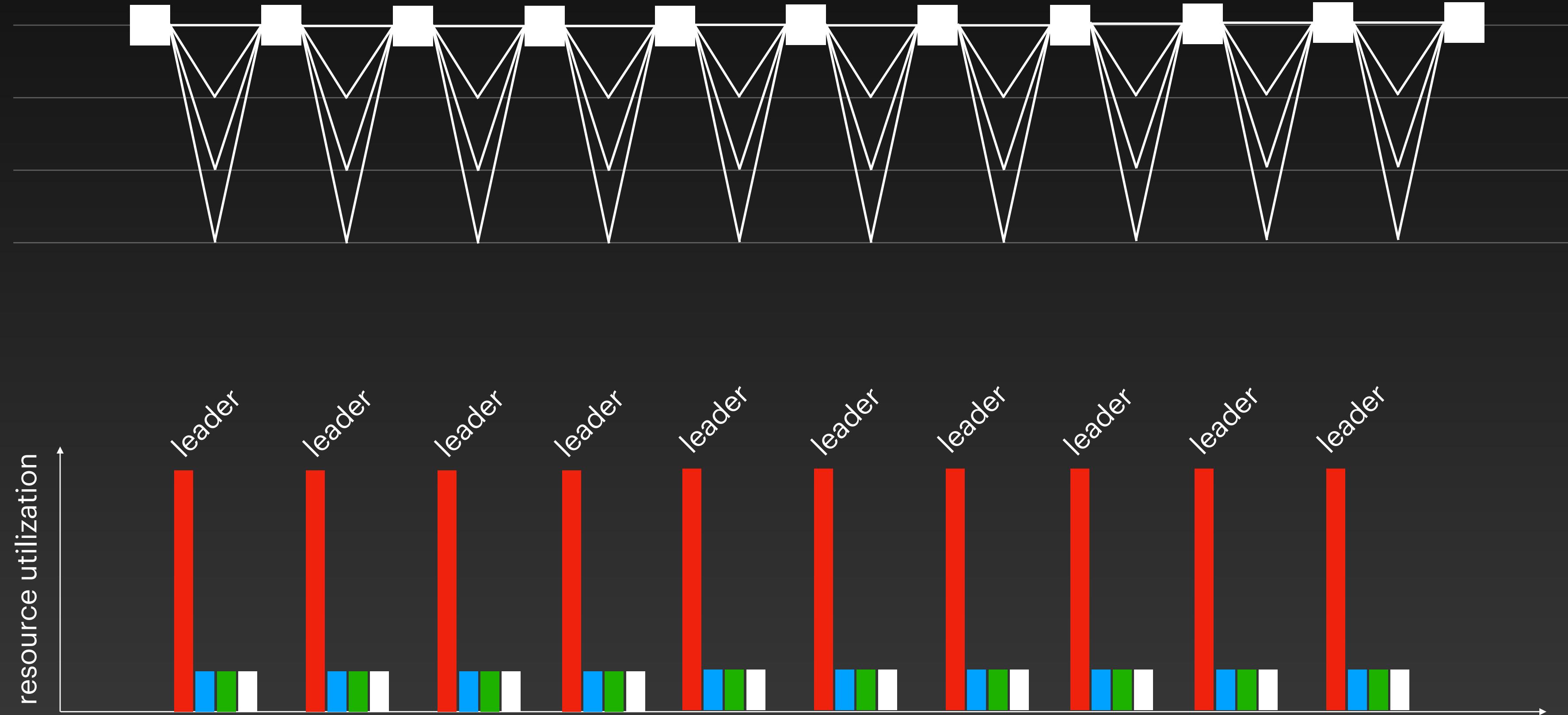
# HotStuff

## Typical leader-based protocols



# Naive Implementation

## Uneven resource utilisation



# Research Questions

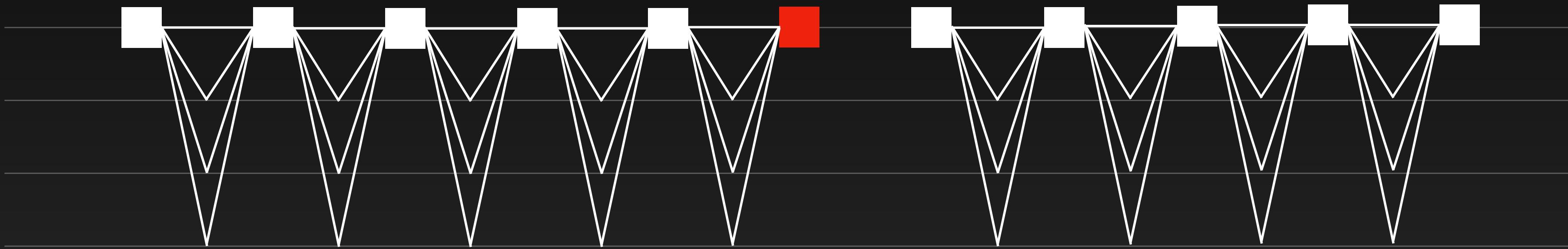
1. Network model?

# Lessons Learned

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

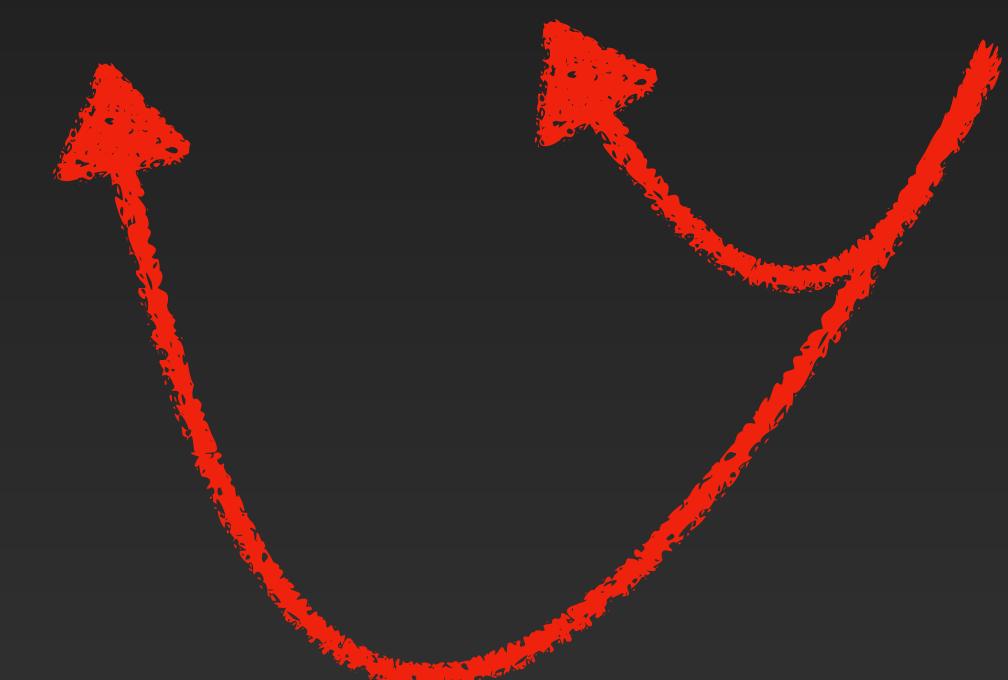
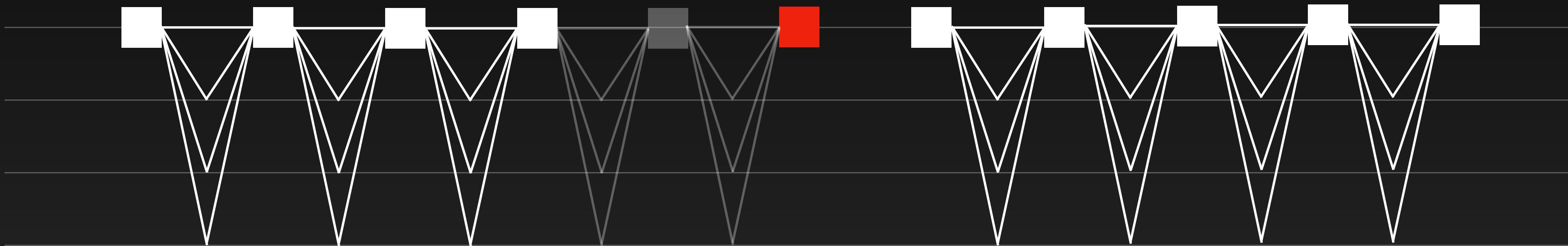
# Leader-Driven Consensus

## Fragility to faults and asynchrony

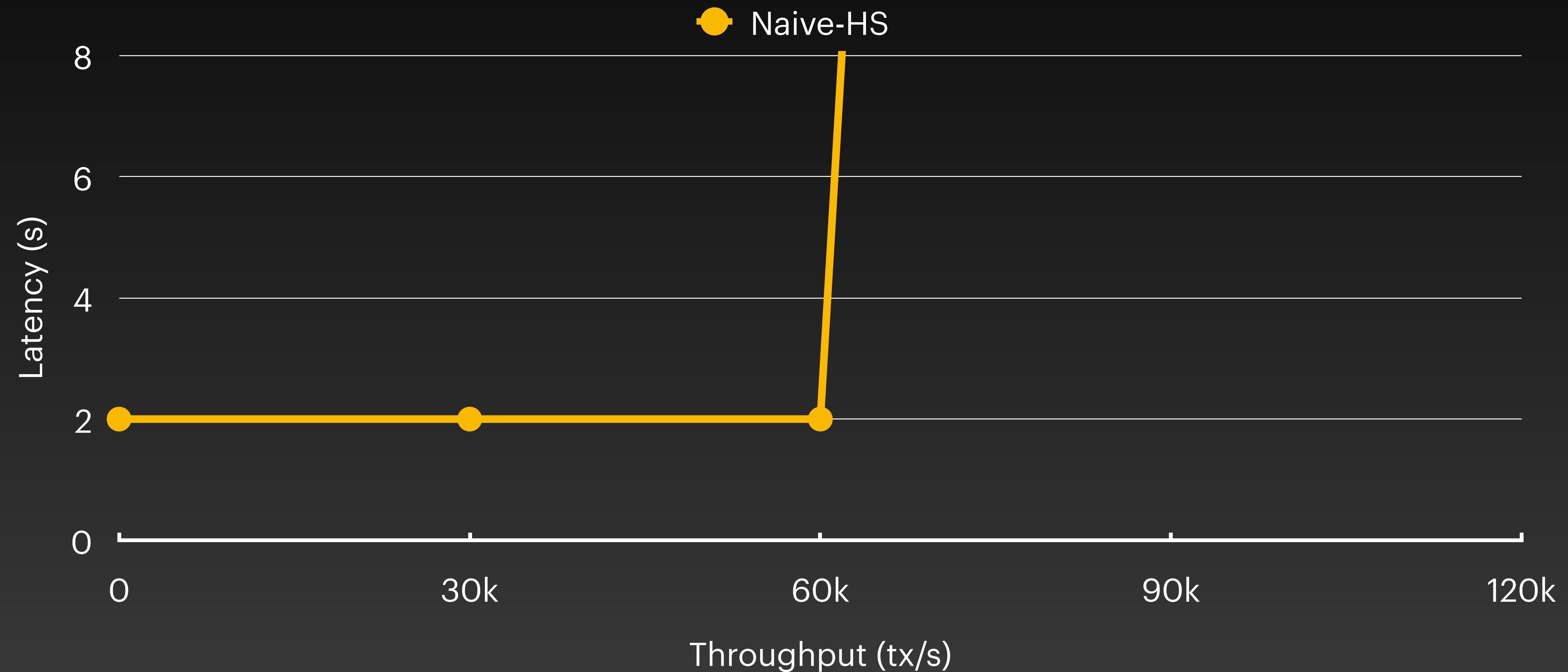


# Leader-Driven Consensus

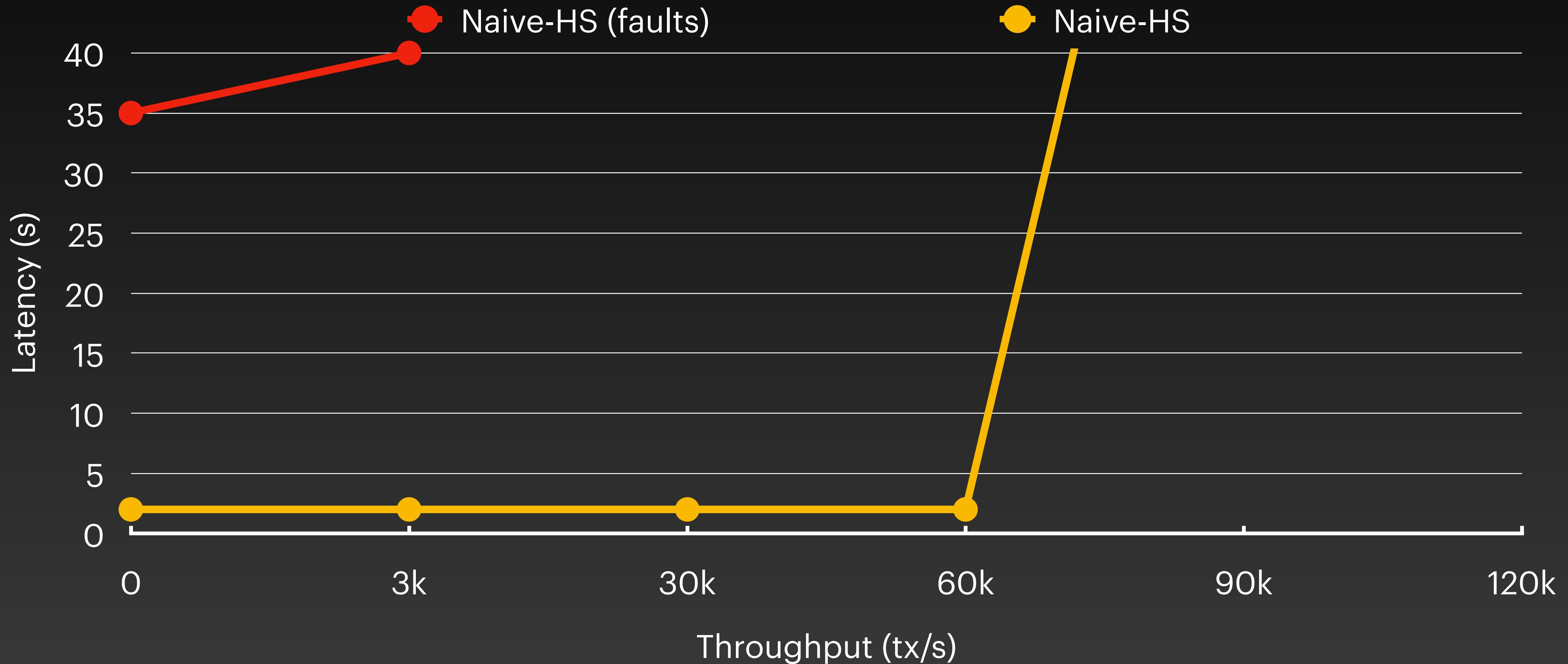
## Fragility to faults and asynchrony



# 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<sup>1,2</sup>, Dahlia Malkhi<sup>2</sup>, Michael K. Reiter<sup>2,3</sup>, Guy Golan Gueta<sup>2</sup>, and Ittai Abraham<sup>2</sup>

<sup>1</sup>Cornell University, <sup>2</sup>VMware Research, <sup>3</sup>UNC-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  $\Delta$  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  $\Delta$  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](mailto: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:

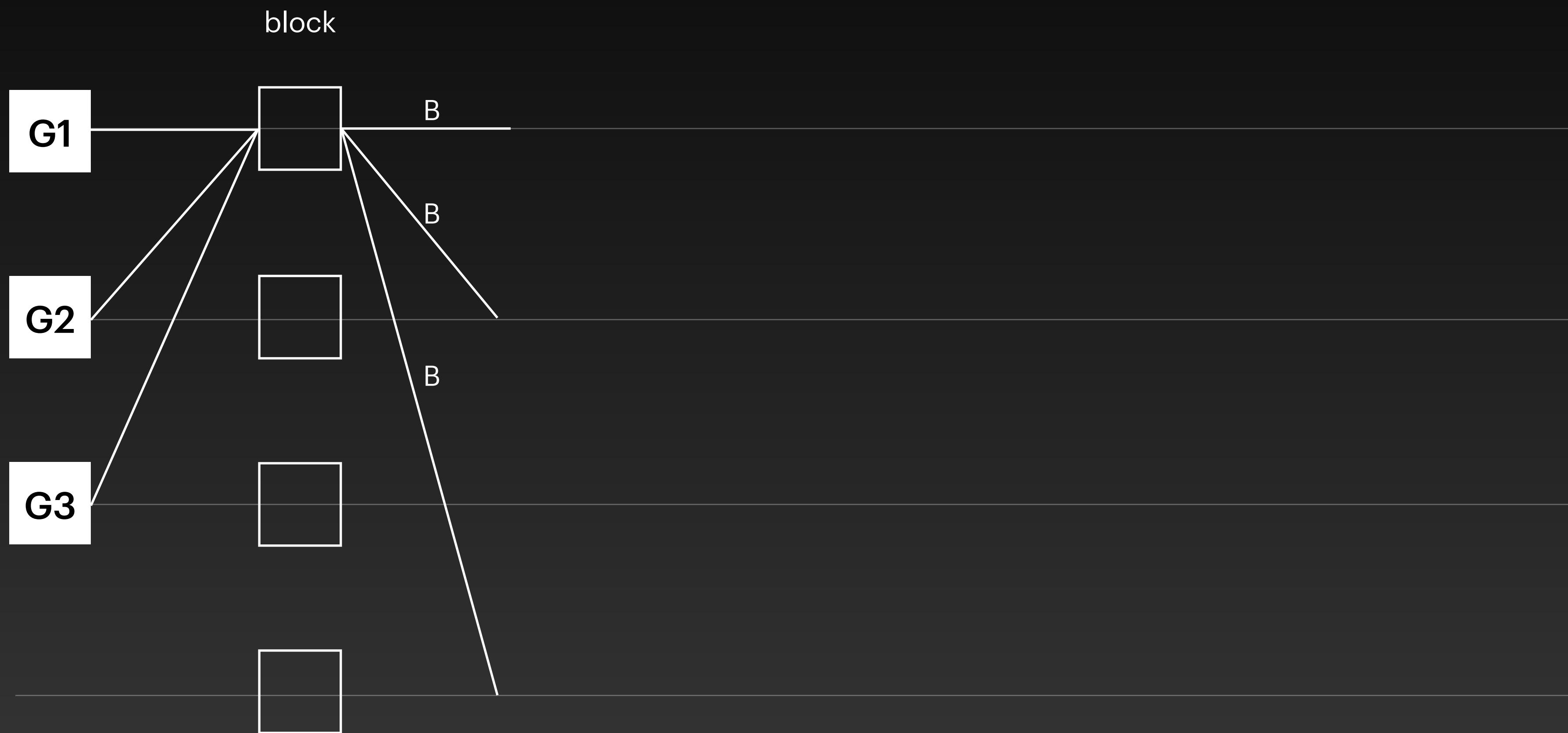
- Any (partially-synchronous) protocol that minimizes overall message number, but relies on a leader to produce proposals and coordinate consensus, fails to capture the high load this imposes on the leader who inevitably becomes a bottleneck.
- Message complexity counts the number of *metadata* messages (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

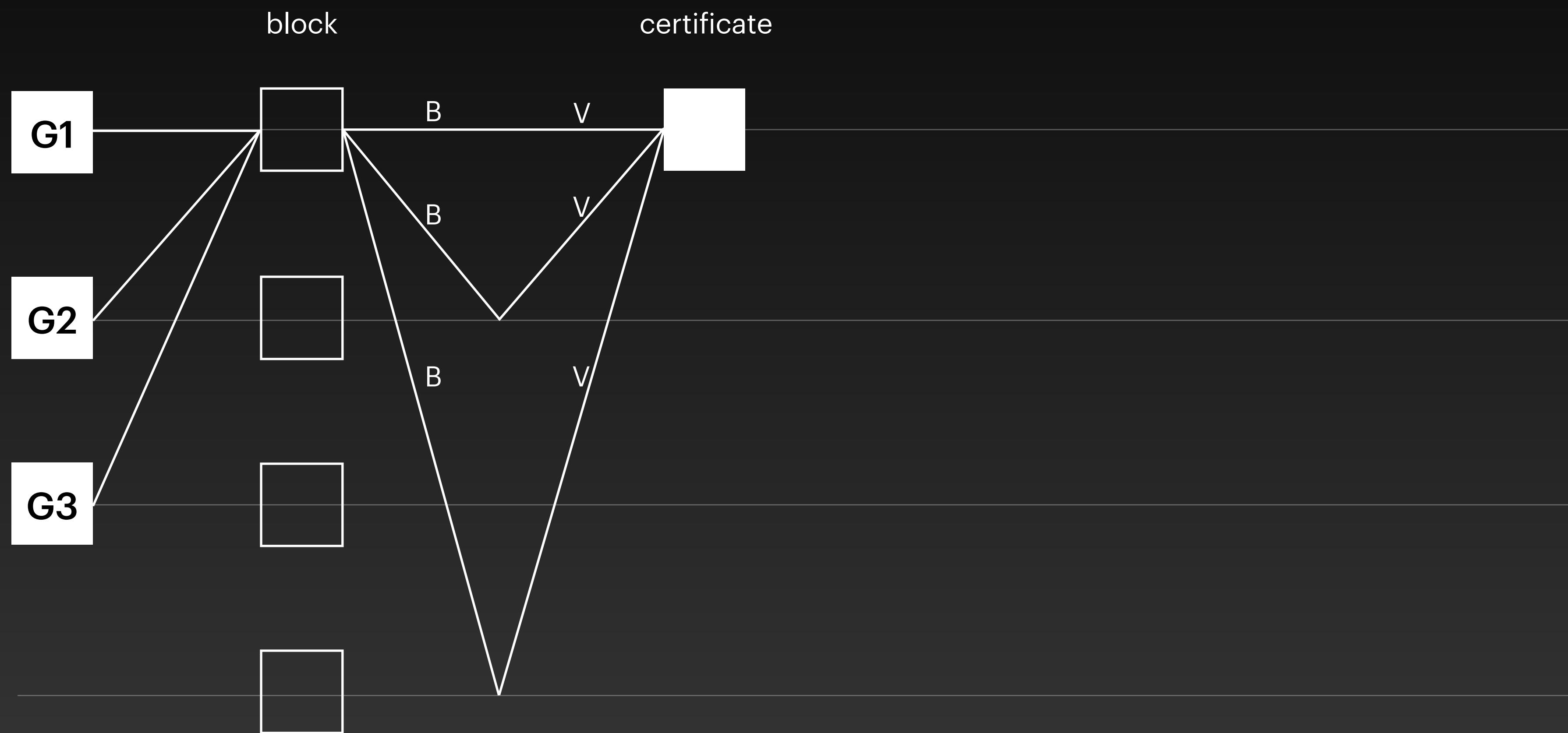
## Narwhal

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

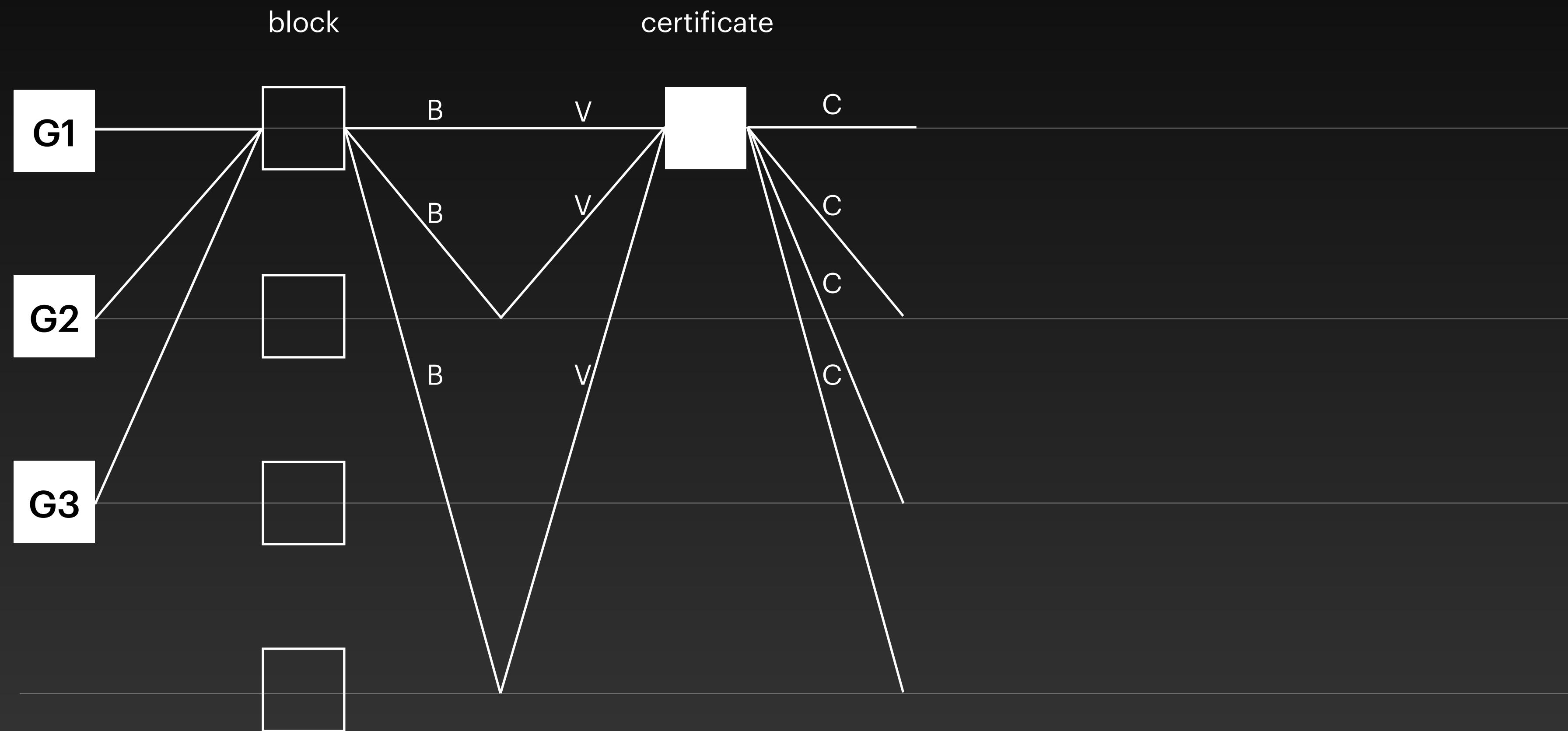
# Narwhal



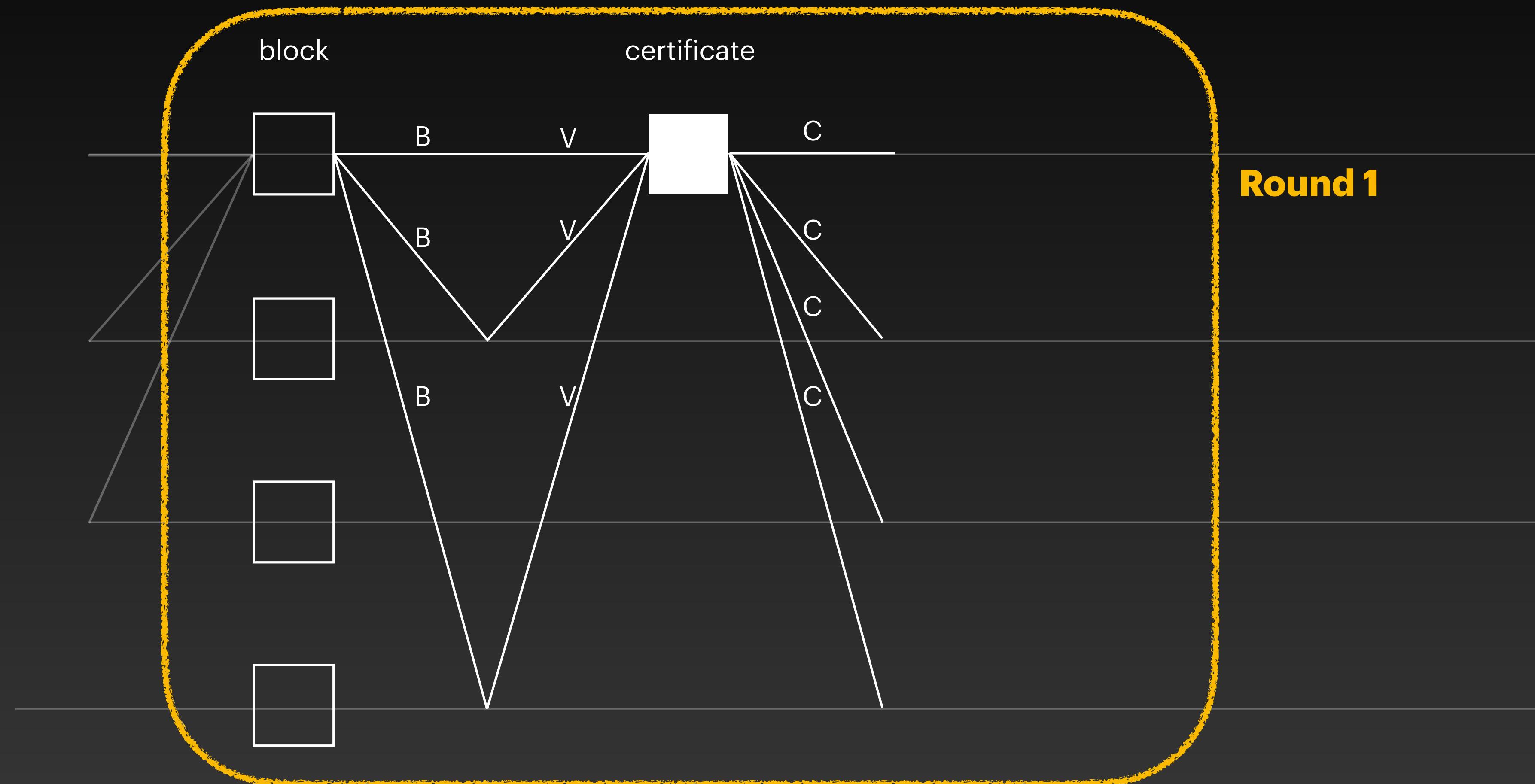
# Narwhal



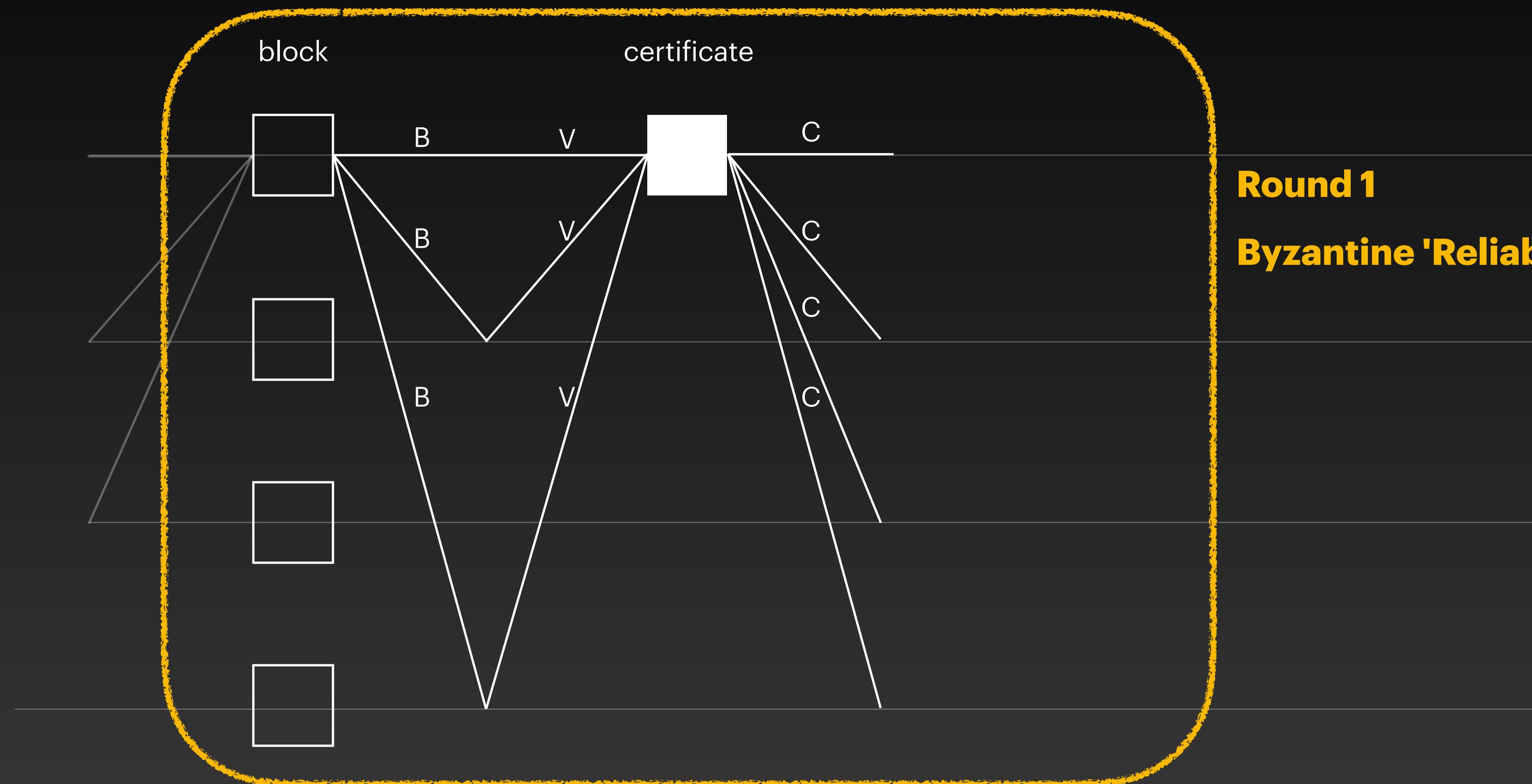
# Narwhal



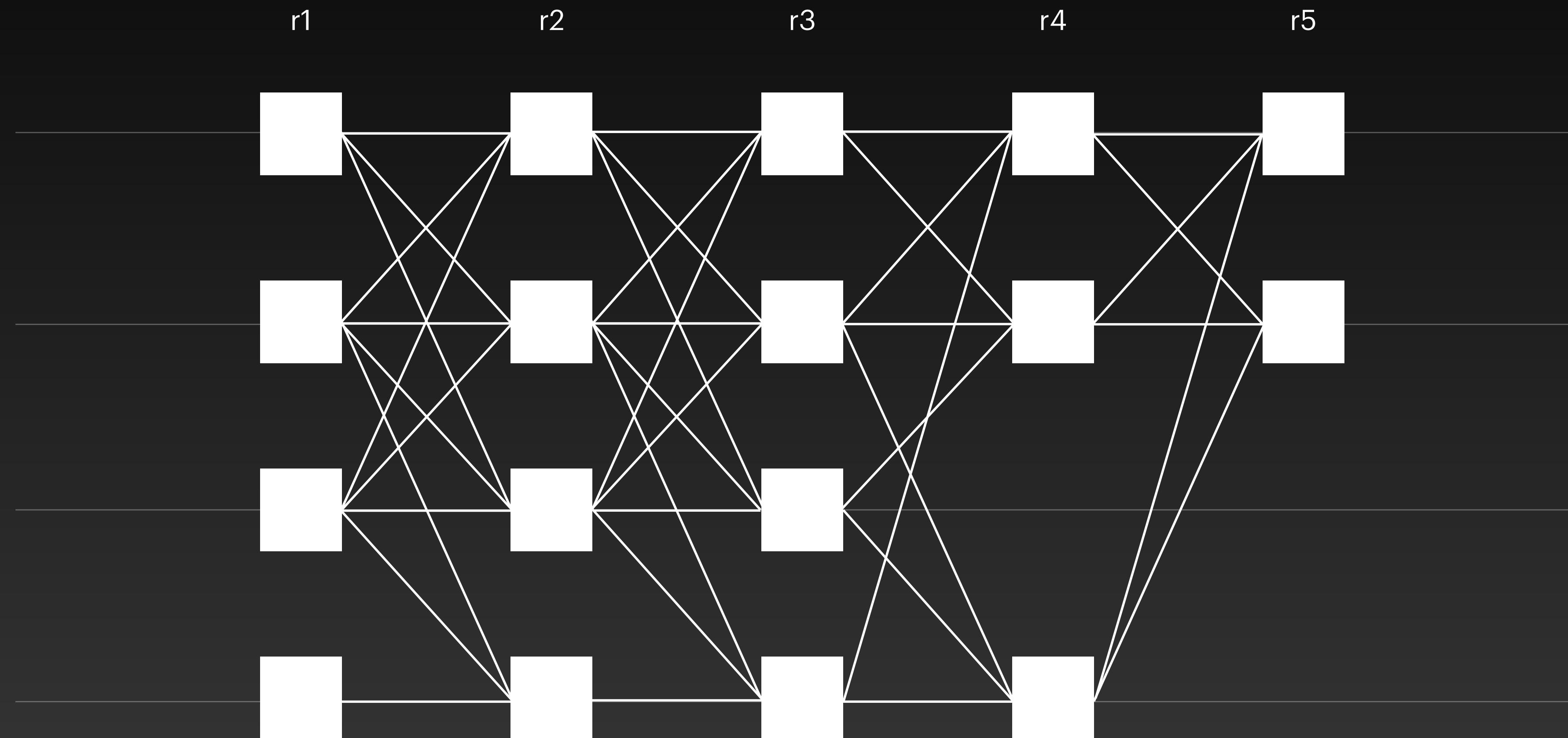
# Narwhal



# Narwhal



# Narwhal



# Research Questions

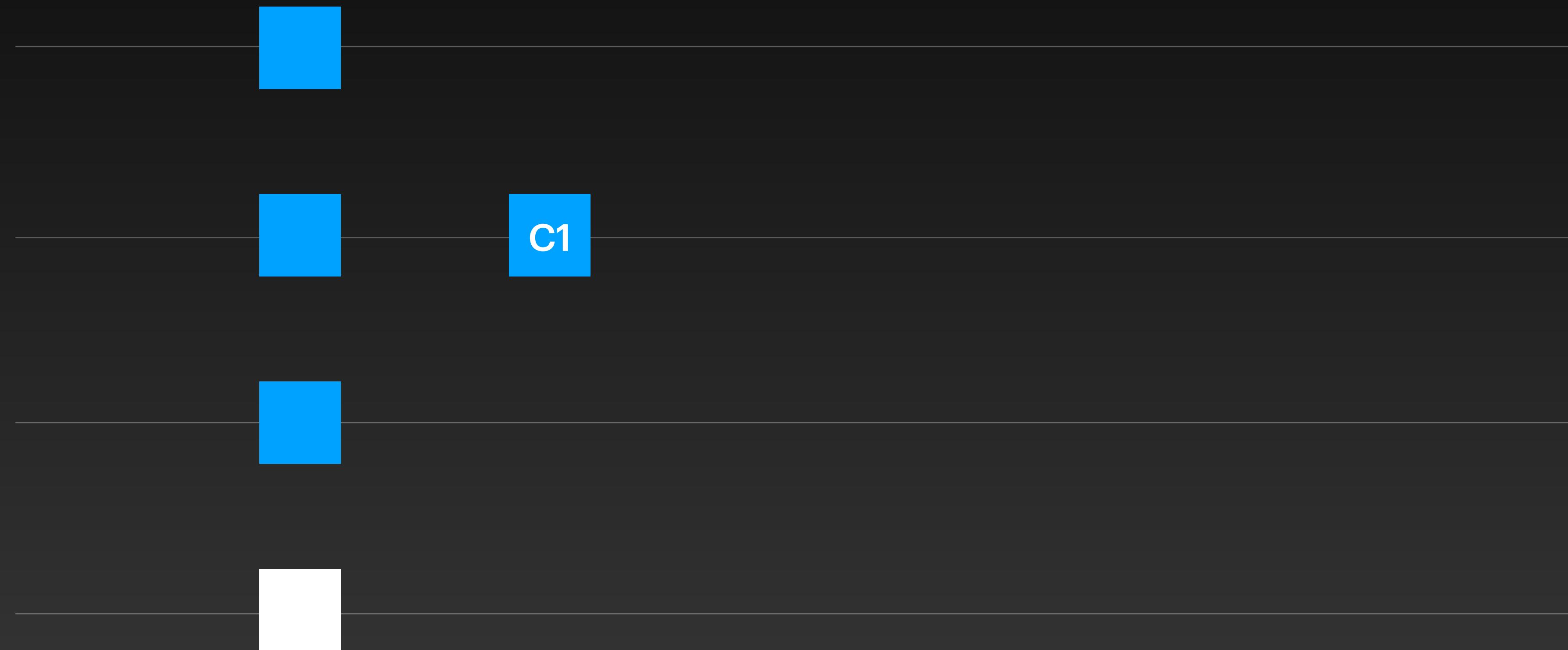
1. Network model?

# Lessons Learned

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

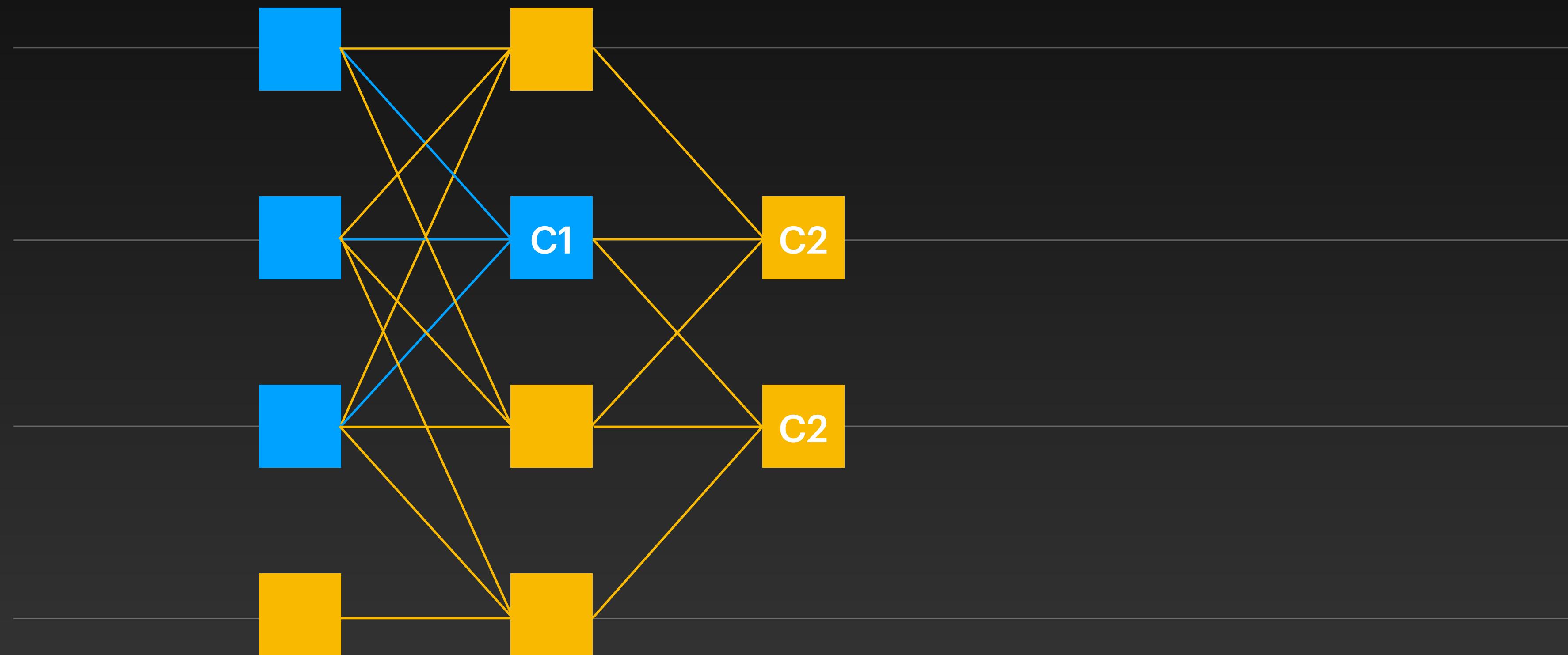
# HotStuff on Narwhal

## Enhanced commit rule



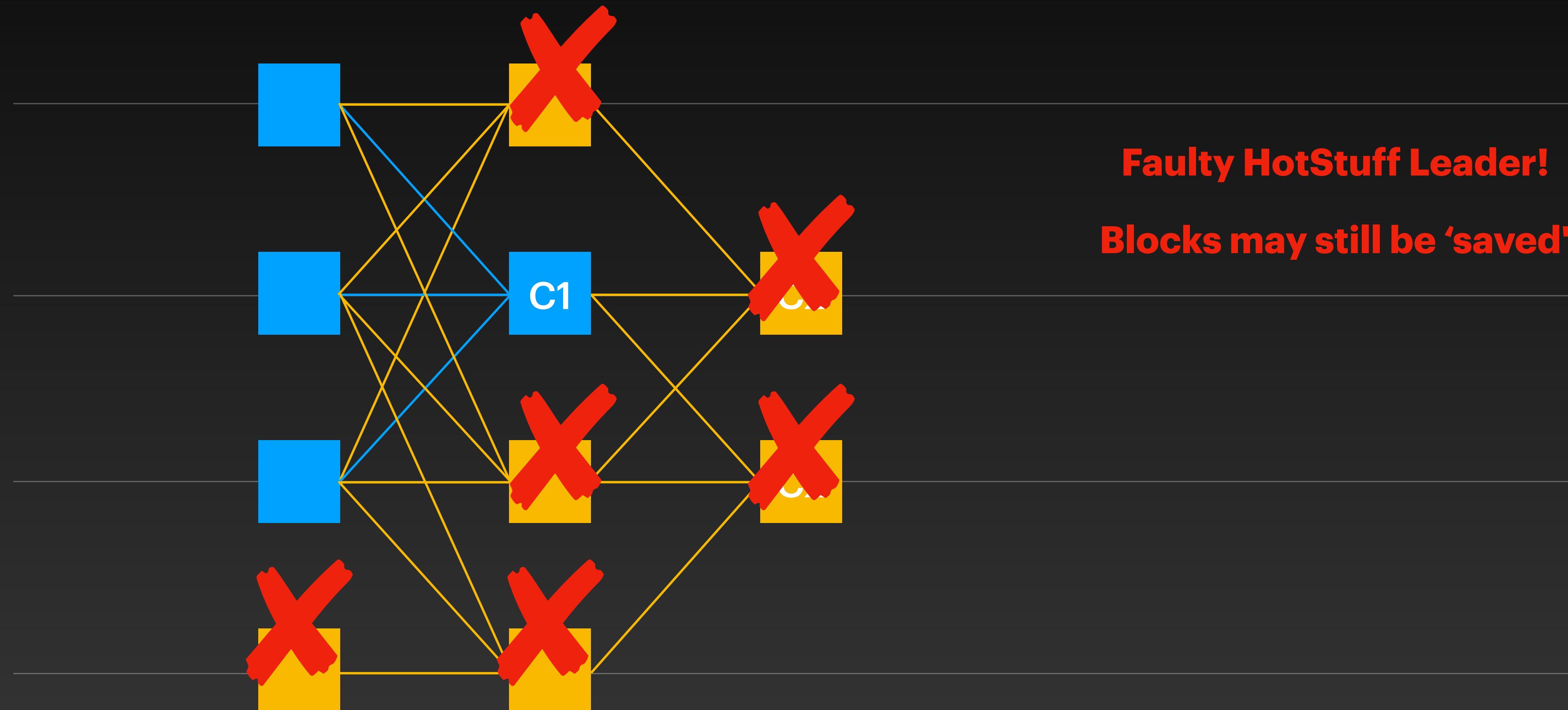
# HotStuff on Narwhal

## Enhanced commit rule



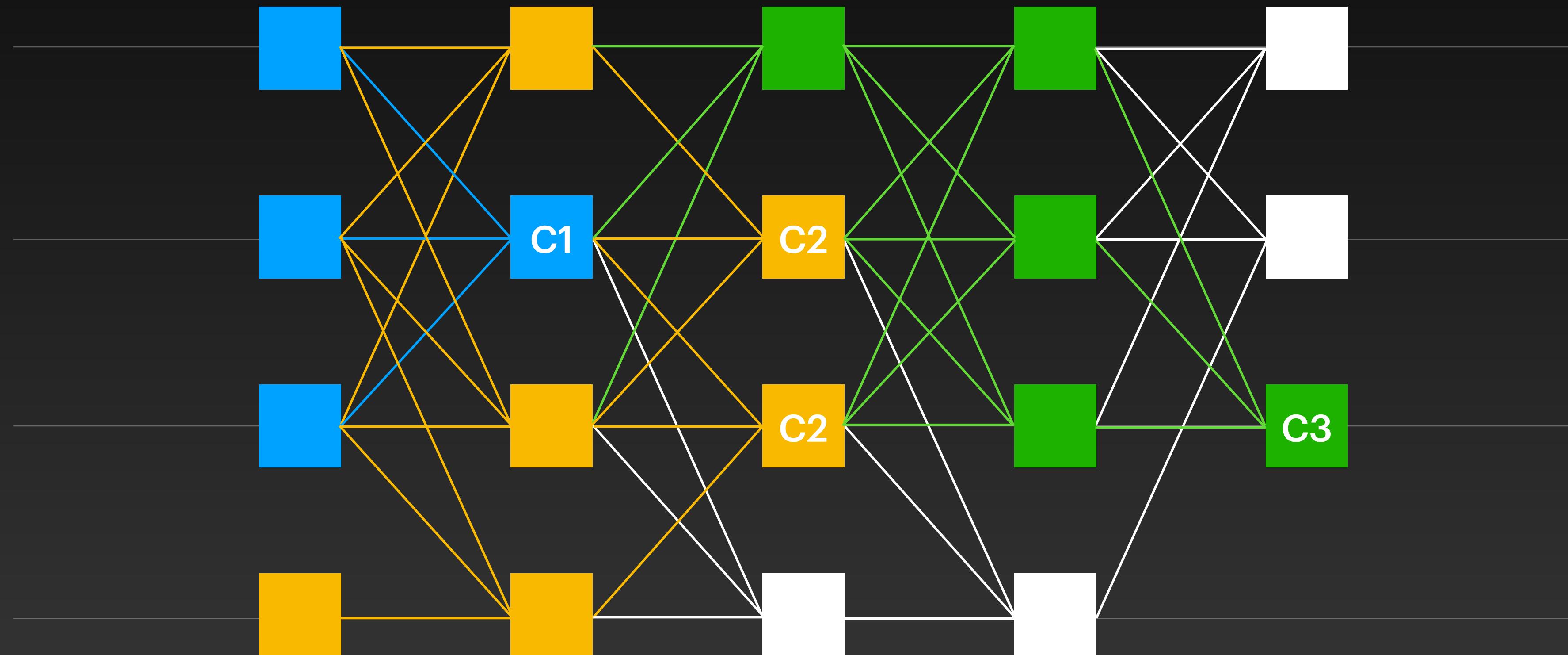
# HotStuff on Narwhal

## Enhanced commit rule



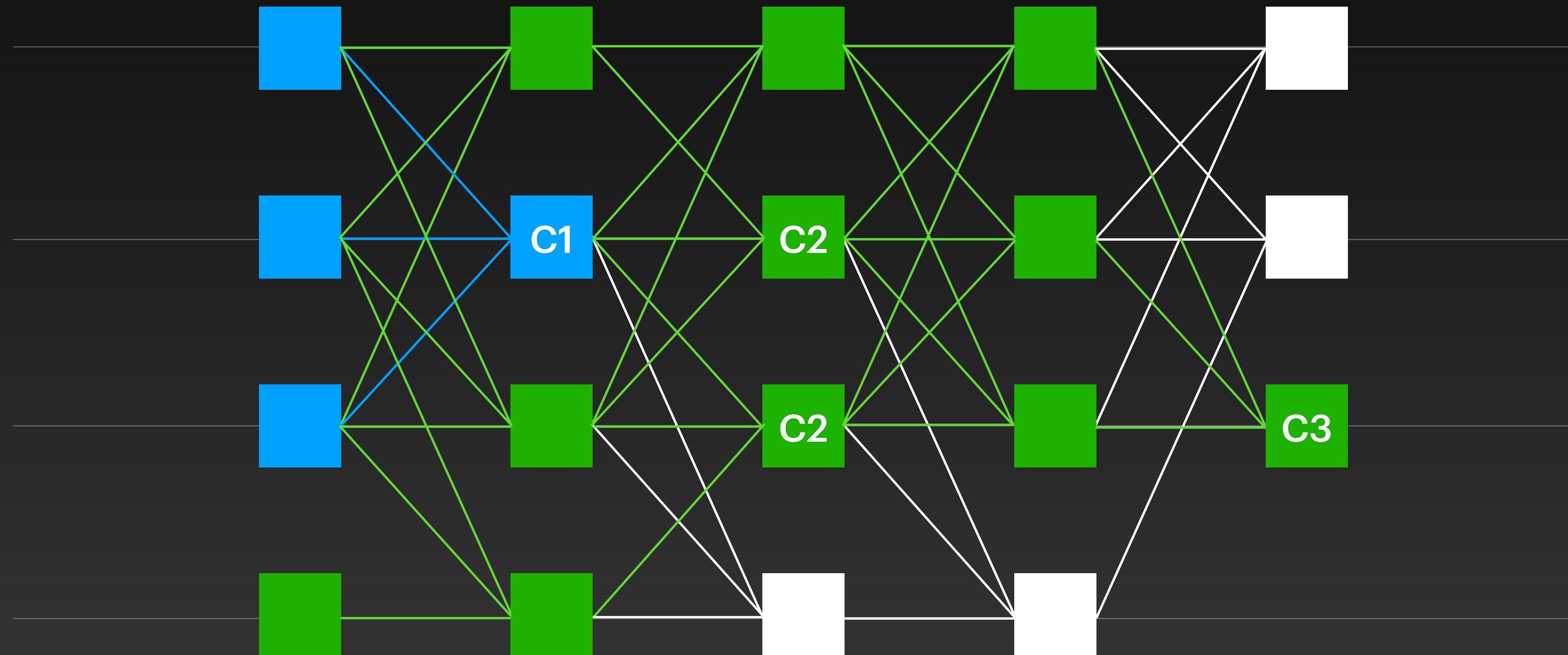
# HotStuff on Narwhal

## Enhanced commit rule

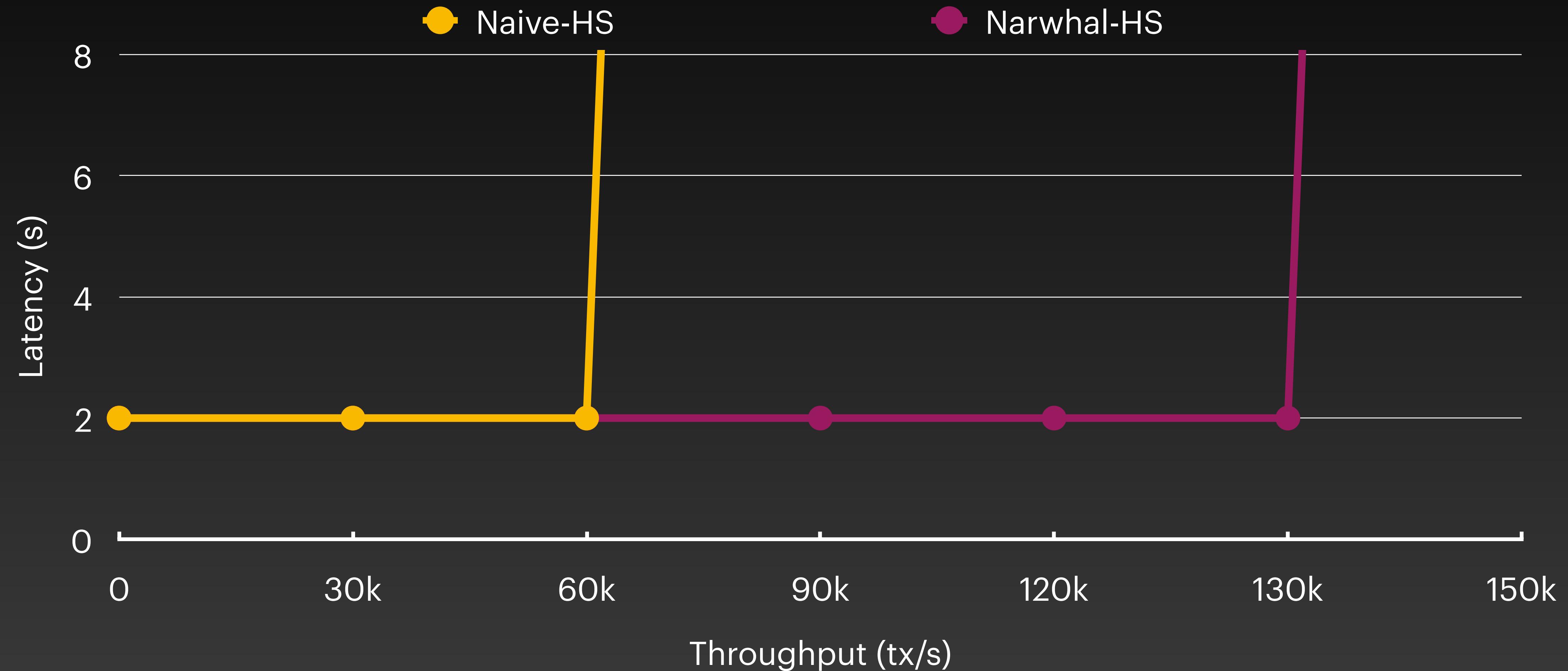


# HotStuff on Narwhal

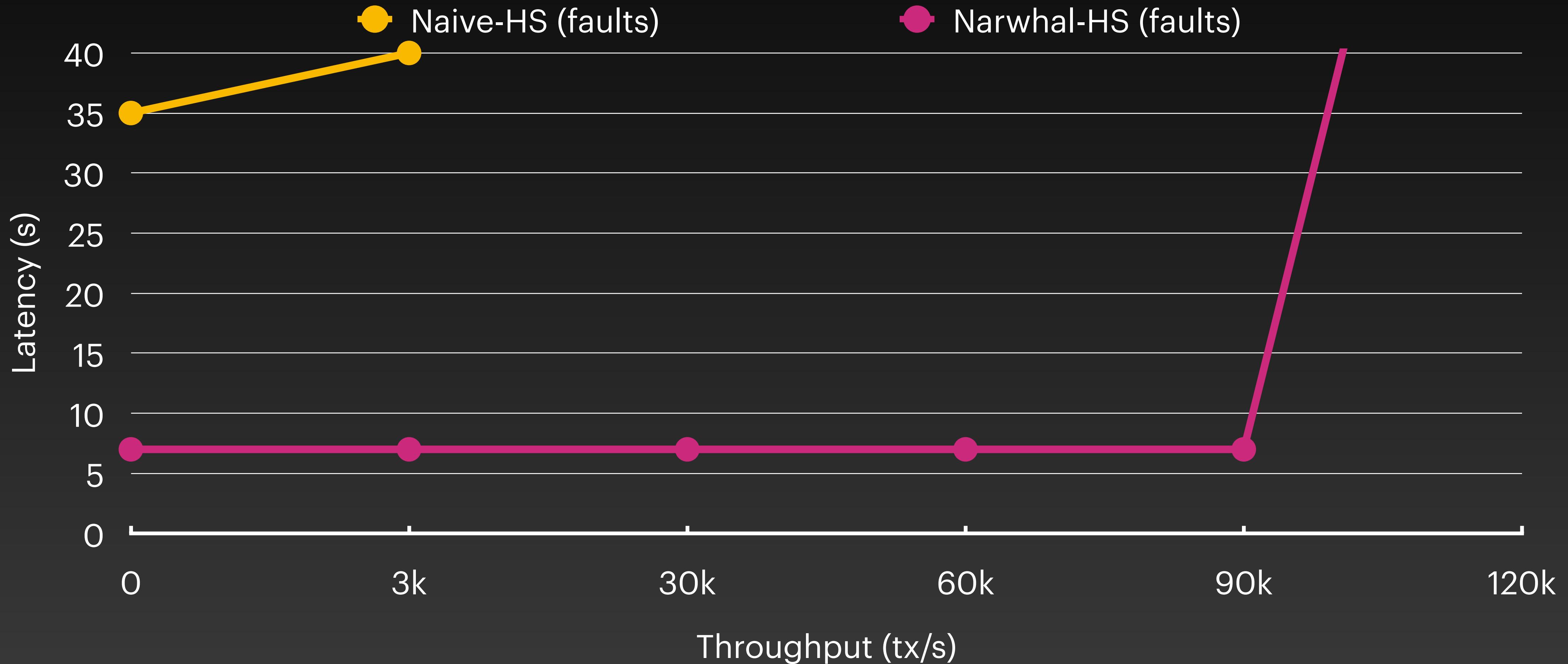
## Enhanced commit rule



# Performance

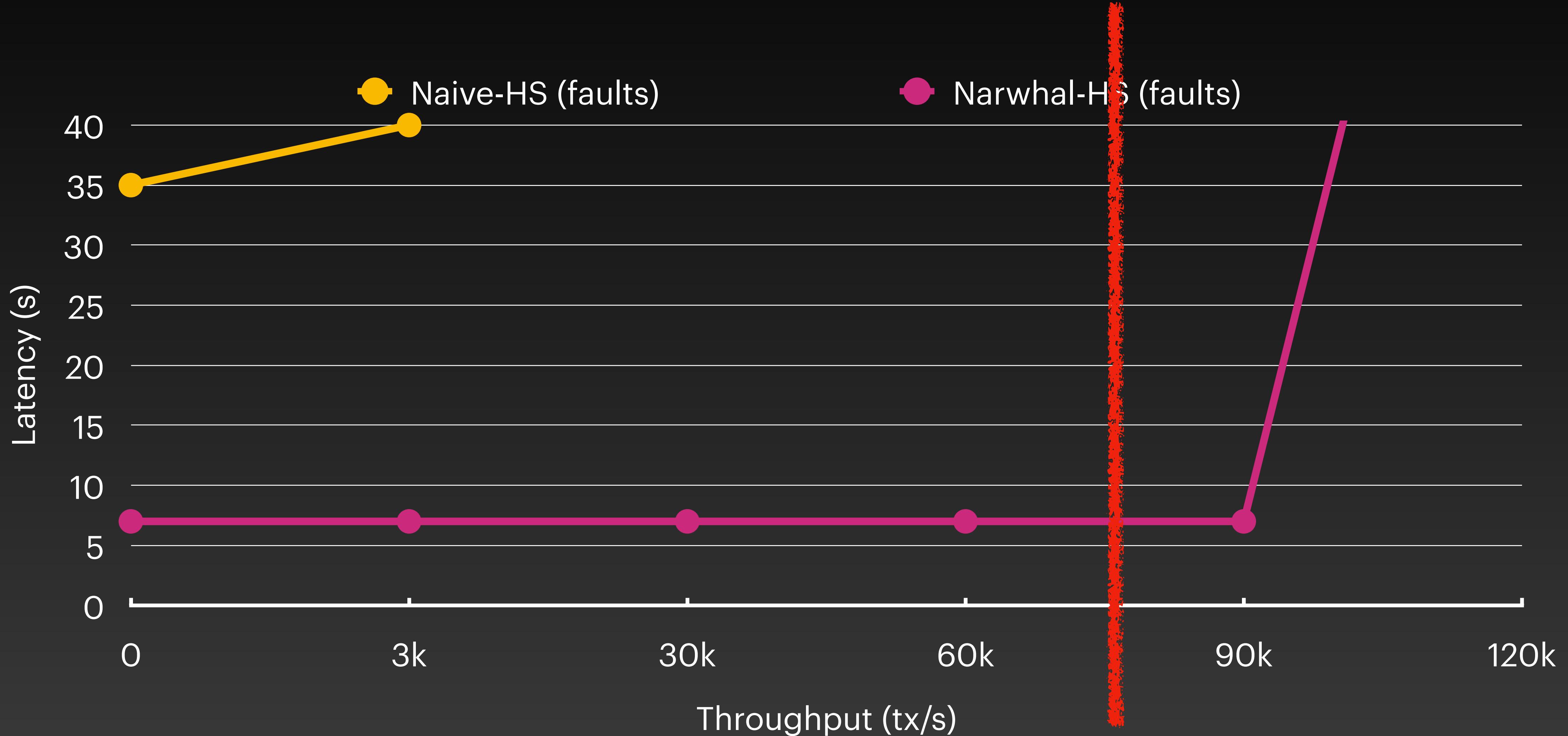


# Performance



# Performance

visa+mastercard



# 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](mailto: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

# DagRider

**All You Need is DAG**

Idit Keidar  
Technion

Oded Naor\*  
Technion

**ABSTRACT**  
We present *DagRider*, the first asynchronous Byzantine Atomic Broadcast protocol that is post-quantum secure, optimal communication complexity, and optimal time complexity. DagRider is correct processes eventually get delivered. We construct two layers. In the lower layer, *Narwhal* tolerates a single faulty node and does 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. All You Need is DAG. In *Proceedings of the 2021 ACM SIGPLAN 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>

# 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 Narwhal as a mempool-based protocol, Narwhal tolerates at most two faulty nodes in the network and stores a copy of each transaction in its mempool. Tusk is a stateless consensus protocol that casts their processes and has a timestamped Directed Acyclic Graph (DAG) of the communication among them. In the second layer, Tusk preserves the safety properties of their DAGs and totally order all processes 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/10.1145/3451959.3459994>

# Bullshark

**Bullshark: DAG BFT Protocols Made Practical**

Alexander Spiegelman  
sasha@septangle.com  
Aptos

George Danezis  
IST Austria

Alexander Spiegelman  
Aptos

**Abstract**  
We present BullShark, the first directed acyclic graph (DAG) based asynchronous Byzantine atomic broadcast 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 safely order messages sent by other parties in the bottom half of 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 reading off the top of the DAG.

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 totally order the DAG in expectation. BullShark follows a similar approach, but uses a structured round-based DAG and encodes a shared randomness in each round via a threshold signature scheme to achieve constant latency in expectation. BullShark is the first DAG-based BFT protocol built on previous ideas. Every round 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 round's messages in a previous round. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party adversarial to the next round once it reliably receives  $n - f$  blocks from the previous round. 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 transitive, so it is not safe to totally order all the DAG's vertices. So 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's structure, BullShark can implement a more efficient edge interpretation logic of Dag-Rider 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 security, and some notion of fairness (e.g., no double spend) have been around for a long time. The main issue with these protocols is that they impose a high load on the leader who inevitably has to handle all the communication.

• **Message complexity** counts the number of *metadata* messages (e.g., voter, signatures, hashes) which take minimal bandwidth compared to the size of the transaction blocks (transaction data). Since blocks are orders of magnitude larger (10MB) than a typical consensus message (100B), the asymptotic message complexity is proportional to the size of the consensus protocol, but it cannot be optimized for fixed mid-size committees (up to ~50 nodes).

• **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 has to handle all the communication.

• **Message complexity** counts the number of *metadata* messages (e.g., voter, signatures, hashes) which take minimal bandwidth compared to the size of the transaction blocks (transaction data). Since blocks are orders of magnitude larger (10MB) than a typical consensus message (100B), the asymptotic message complexity is proportional to the size of the consensus protocol, but it cannot be optimized for fixed mid-size committees (up to ~50 nodes).

• **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 has to handle all the communication.

**ACM Reference Format:**  
Alexander Spiegelman, Neil Giridharan, Alberto Sonnino, and George Danezis. 2022. Bullshark: DAG BFT Protocols Made Practical. In *Proceedings of ACM Conference, Los Angeles, CA, USA, November 2022 (Comsense '22)*, November 2022, Los Angeles, CA, USA, 17 pages. <https://doi.org/10.1145/3492211.3492220>

# Dumbo-NG

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

Yingqi Guo\*  
ISCA & UCAS

Yuan Liu\*  
ISCA

Zherlang Lu\*  
USID

Neil Giridharan  
giridharan@berkeley.edu  
UC Berkeley

Jing Xu\*  
ISCA

Zhenfeng Zhang\*  
ISCA

Alberto Sonnino  
alberto@mystenlabs.com  
Mysten Labs

Lefteris Kokoris-Kogias  
ekokoris@ist.ac.at  
IST Austria

**Abstract**  
We propose Dumbo-NG, a novel asynchronous BFT consensus algorithm that is throughput oblivious. Specifically, Dumbo-NG achieves the same throughput of order  $\sqrt{n}$  as the well-known Dumbo [3] (with  $n$  the number of nodes). It is throughput oblivious, meaning that it only needs to receive one message per node per round. Thus, Dumbo-NG is able to tolerate a large number of faulty nodes while maintaining the same throughput as Dumbo. Moreover, Dumbo-NG is asynchronous, meaning that it does not require any synchronization between nodes. This makes Dumbo-NG suitable for use in distributed systems where nodes may have different clock speeds or may be located in different time zones. Finally, Dumbo-NG is fast, meaning that it can reach consensus quickly even in the presence of many faulty nodes. This makes Dumbo-NG suitable for use in real-world applications such as blockchain and distributed ledger technologies.

We present Dumbo-NG, a novel asynchronous BFT consensus algorithm that is throughput oblivious. Specifically, Dumbo-NG achieves the same throughput of order  $\sqrt{n}$  as the well-known Dumbo [3] (with  $n$  the number of nodes). It is throughput oblivious, meaning that it only needs to receive one message per node per round. Thus, Dumbo-NG is able to tolerate a large number of faulty nodes while maintaining the same throughput as Dumbo. Moreover, Dumbo-NG is asynchronous, meaning that it does not require any synchronization between nodes. This makes Dumbo-NG suitable for use in distributed systems where nodes may have different clock speeds or may be located in different time zones. Finally, Dumbo-NG is fast, meaning that it can reach consensus quickly even in the presence of many faulty nodes. This makes Dumbo-NG suitable for use in real-world applications such as blockchain and distributed ledger technologies.

**ACM Reference Format:**  
Yingqi Guo, Yuan Liu, Zherlang Lu, Neil Giridharan, Alberto Sonnino, Jing Xu, and Zhenfeng Zhang. 2022. Dumbo-NG: Fast Asynchronous BFT Consensus with Throughput-Oblivious Latency. In *Proceedings of the 2022 ACM SIGPLAN Conference on Principles of Distributed Computing (PODC '22)*, July 26–30, 2022, Virtual Event, Italy. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3492211.3492220>

# Data Dissemination

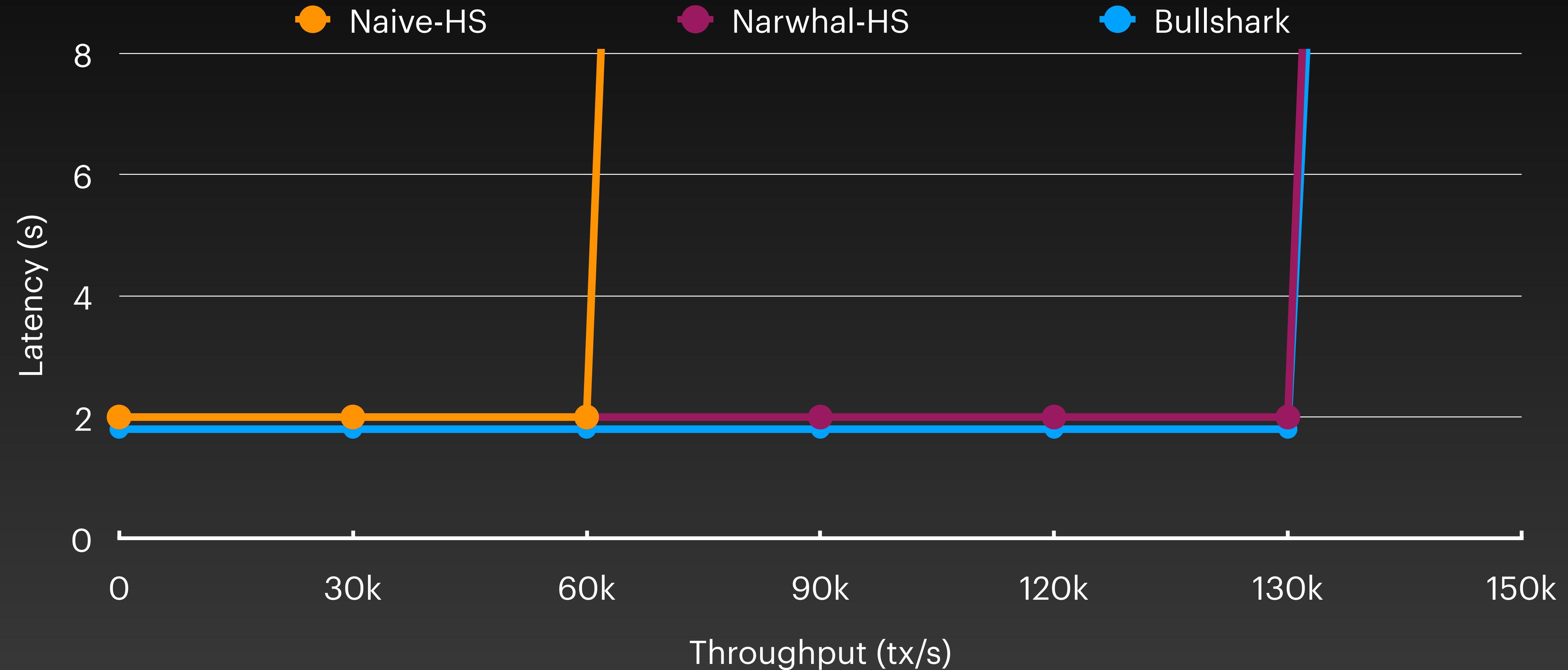
- Hard to make efficient
- 99% of the code

# Consensus

- Error prone
- Isolated, easy to maintain



# 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**

...

# 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 the predominant faults
- 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](mailto: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 DAG, each validator is assigned a unique vertex and each vertex has up to  $n$  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 details of consensus, the DAG can be constructed without contending with complex mechanisms like view-change or view-synchronization.

Due to periods of network asynchrony, each validator may observe a slightly different portion of the DAG at any given point in time. To handle this, each validator maintains a set of pending transactions, which are ordered by their arrival time. When a validator receives a transaction, it checks if it is valid and if so, it adds it to its set of pending transactions. If the transaction is invalid, it is discarded. This ensures that the DAG remains consistent across all validators.

# 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

Historically, the proposed values and ensure that the leader keeps them safe. Finally, miners commit the values. 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 memory are used to store the values and move them 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 full set of parties to provide fast consensus and fast path protocols to demonstrate their low latency and resource efficiency. Our paper proposes a novel DAG-based consensus protocol with 50 validators and reported a throughput of 160,000 TPS with one machine per validator, which further increased to 600,000 TPS with 10 machines per validator compared to existing DAG-based protocols, with similar throughput.

**Related Work**

Related Version: Cordial Miners: Fast and Efficient Consensus for Every Eventuality

Full Version: <https://arxiv.org/abs/2205.09174v6>

**Abstract**

Cordial Miners are a family of efficient Byzantine Atomic Broadcast protocols, with instances of latency 3 rounds. They achieve high resource efficiency and censorship resistance. MYSTICETI-C achieves high latency improvement by avoiding consensus on every block. Instead, Cordial Miners use a novel commit rule such that every block can be committed without delay, resulting in a much faster consensus than the state-of-the-art. MYSTICETI-C and under crash failures. We further extend MYSTICETI-C to MYSTICETI-FPC, which incorporates a fast commit path that allows for fast consensus even in the presence of a primary fast commit path protocols. MYSTICETI-FPC minimizes the number of signatures and messages by weaving the fast path protocols with the slow path protocols. This subsequently results in better performance. We prove the safety and liveness of our protocols in a BFT setting. We evaluate both MYSTICETI-C and MYSTICETI-FPC with both consensus and fast path protocols to demonstrate their low latency and resource efficiency. Our work shows that Cordial Miners can achieve a wide distribution of workload. Additionally, because each party is responsible for disseminating its own transactions, the burden of maintaining a single round and committing all blocks is shifted to the leader. Consequently, these protocols have demonstrated improved throughput compared to their DAG-based counterparts under the same framework [19], [26]. 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 BFT protocols have a problem: each round, every party can create a potential DAG vertex containing transactions from edges pointing to vertices from previous rounds. These protocols rely on committing a designated "leader vertex" and order other non-leader vertices around it. The question is: what is the order in 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 require a large amount of overhead or message delays to commit a proposal (instead of the 6 in Bullshark) and facilitates the pipeline of proposals to commit one block every round [38]. They, however, require a high number of messages to be exchanged between validators, a lot 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, which focuses on low-latency and low-CPU overhead. MYSTICETI-C is a consensus protocol that is a consensus protocol based on a threshold logical clock [29] of  $2\delta + 1$  blocks. Finally, (3) since all certified blocks need to

# 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 3 rounds. MYSTICETI-C achieves high resource efficiency and censorship resistance. MYSTICETI-C achieves high latency improvement by avoiding consensus on every block. Instead, Cordial Miners use a novel commit rule such that every block can be committed without delay, resulting in a much faster consensus than the state-of-the-art. MYSTICETI-C and under crash failures. We further extend MYSTICETI-C to MYSTICETI-FPC, which incorporates a fast commit path that allows for fast consensus even in the presence of a primary fast commit path protocols. MYSTICETI-FPC minimizes the number of signatures and messages by weaving the fast path protocols with the slow path protocols. This subsequently results in better performance. We prove the safety and liveness of our protocols in a BFT setting. We evaluate both MYSTICETI-C and MYSTICETI-FPC with both consensus and fast path protocols to demonstrate their low latency and resource efficiency. Our work shows that Cordial Miners can achieve a wide distribution of workload. Additionally, because each party is responsible for disseminating its own transactions, the burden of maintaining a single round and committing all blocks is shifted to the leader. Consequently, these protocols have demonstrated improved throughput compared to their DAG-based counterparts under the same framework [19], [26]. 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.

Recently, a novel approach known as DAG-based BFT has emerged [5], [18], [28], [33], [34], [40], [46]. These protocols enable a full set of parties to provide fast consensus and fast path protocols to demonstrate their low latency and resource efficiency. Our paper proposes a novel DAG-based consensus protocol with 50 validators and reported a throughput of 160,000 TPS with one machine per validator, which further increased to 600,000 TPS with 10 machines per validator compared to existing DAG-based protocols, with similar throughput.

**Related Work**

Related Version: Cordial Miners: Fast and Efficient Consensus for Every Eventuality

Full Version: <https://arxiv.org/abs/2205.09174v2>

**Abstract**

MYSTICETI-C is the first DAG-based consensus protocol to achieve the lower bounds of latency 3 rounds. We evaluate both MYSTICETI-C and MYSTICETI-FPC to demonstrate their low latency and resource efficiency. Our work shows that Cordial Miners can achieve a wide distribution of workload. Additionally, because each party is responsible for disseminating its own transactions, the burden of maintaining a single round and committing all blocks is shifted to the leader. Consequently, these protocols have demonstrated improved throughput compared to their DAG-based counterparts under the same framework [19], [26]. 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.

Several recent blockchains, such as Sui [67], [12], have adopted DAG-based consensus [30], [51], [56], [34], [30], [17], [12], [58], [44]. By design, these consensus protocols scale well in terms of throughput, with a performance of 100K TPS or more. However, they are not able to achieve a low latency. This is mainly due to the fact that they are partially synchronous, existing SMR protocols can commit with a latency overhead of  $3\delta$  (where  $\delta$  represents the number of validators) [23]. 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 a family of uncertified DAG-based consensus protocols where each vertex is delivered through consistent commitment [14], have high latency for three main reasons: (1) the consensus process requires three rounds to commit a proposal block between validators; (2) the consensus process requires a large amount of CPU on each validator, which grows with the number of validators [42]; (3) the consensus process requires a large amount of memory to store the data of the blocks that have been committed [16]. This burden is particularly heavy for a crash-recovered validator that typically needs to verify thousands of signatures before it can commit a proposal. At the same time, consensus certification seems to have the benefit that in adversarial cases one can advance the DAG without needing to synchronize the full block chain. Our experiments show that this benefit is negated when needing to execute the committed transactions. As a result, the certification benefit of consensus is not used if not used for the common case of powering a State Machine Replication system (e.g., a blockchain).

This comes in contrast to Cordial Miners [26], which require only 3 messages to commit a proposal (instead of the 6 in Bullshark) and facilitates the pipeline of proposals to commit one block every round [38]. They, however, require a high number of messages to be exchanged between validators, a lot 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, which focuses on low-latency and low-CPU overhead. MYSTICETI-C is a consensus protocol that is a consensus protocol based on a threshold logical clock [29] of  $2\delta + 1$  blocks. Finally, (3) since all certified blocks need to

# Techniques

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

# Discussion

**Certified DAG**

**Uncertified DAG**



**Shoal/shoal++**

- Low latency
- Easier synchroniser
- Leverage existing code

**Sailfish/BBCA**

- Lower latency
- Easy synchroniser
- Flexible

**CM/Mysticeti**

- Lowest latency
- Graceful crash faults
- Simpler, less CPU

# Research Questions

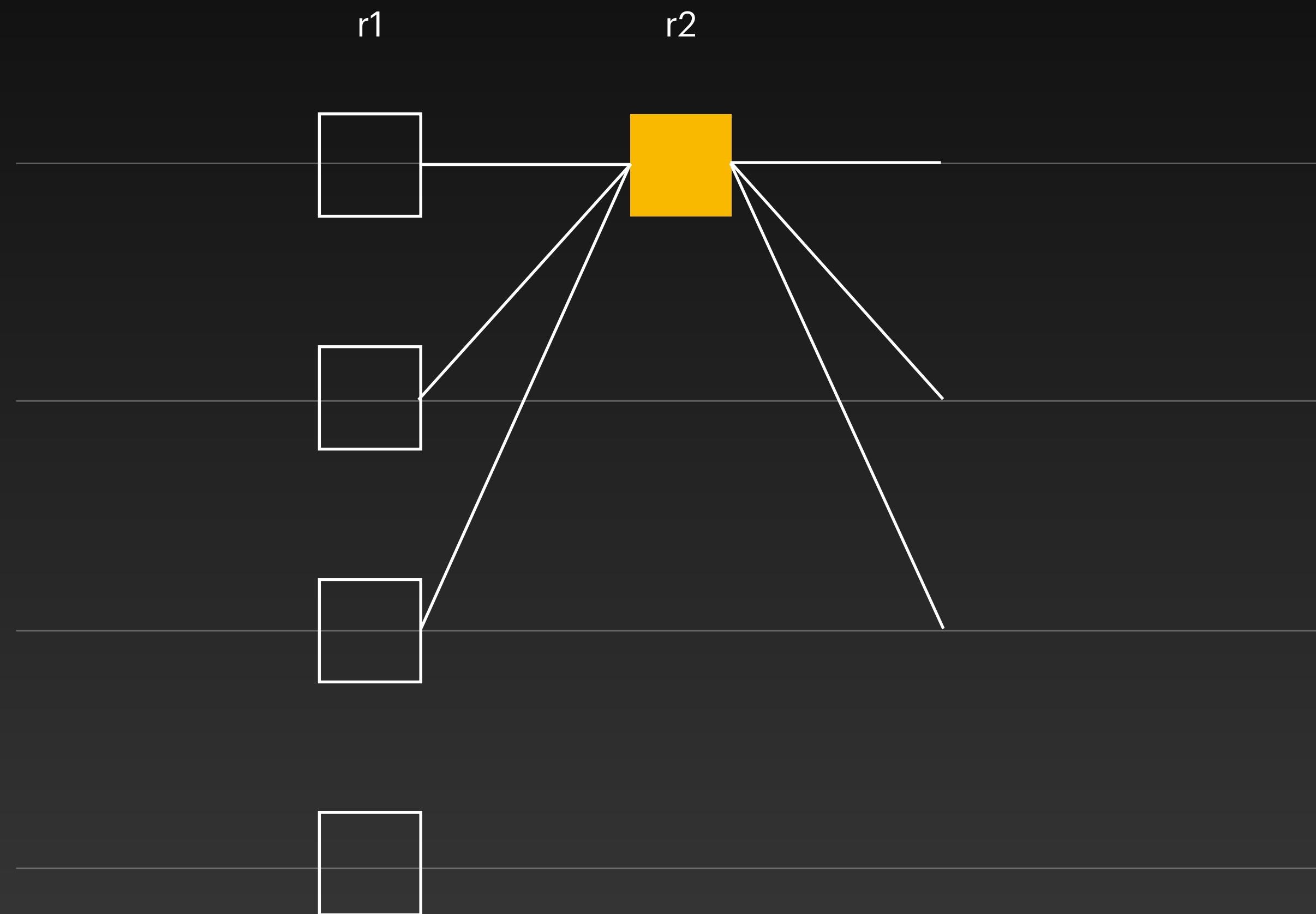
1. Network model?
2. BFT testing?
3. Consensus-exec interface?
4. Storage architecture?
5. Block synchroniser?

# 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

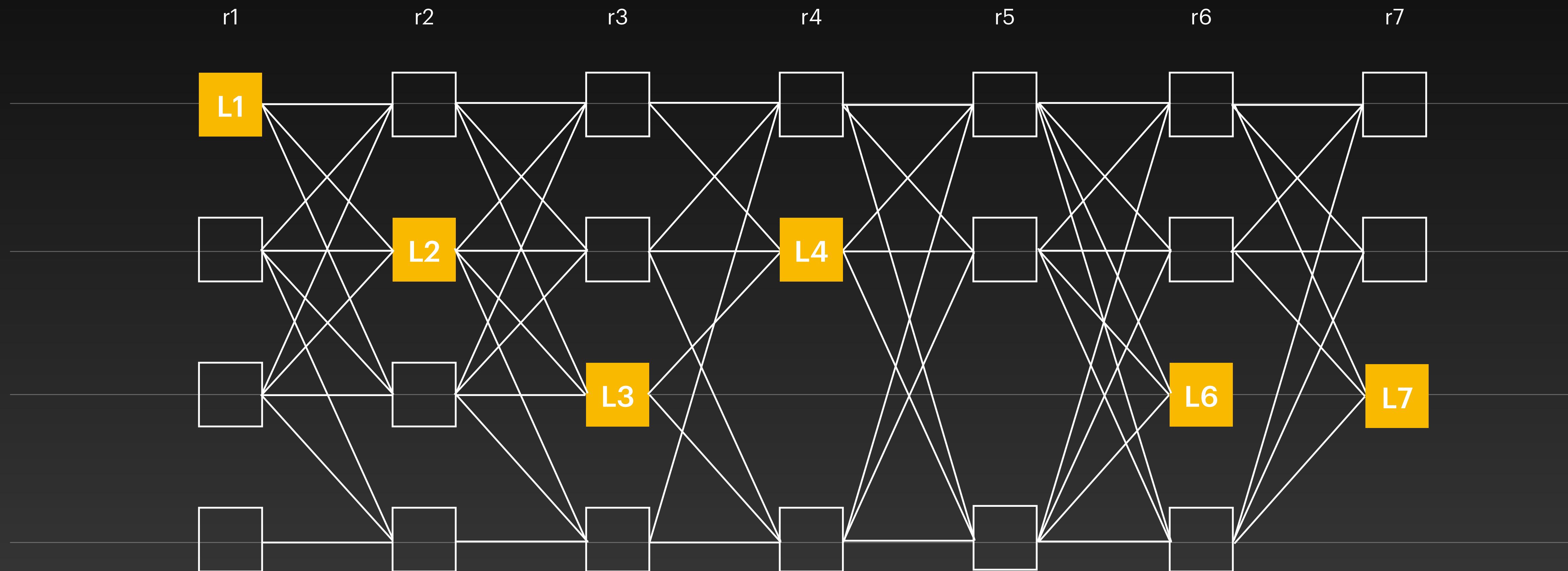


# Uncertified DAG

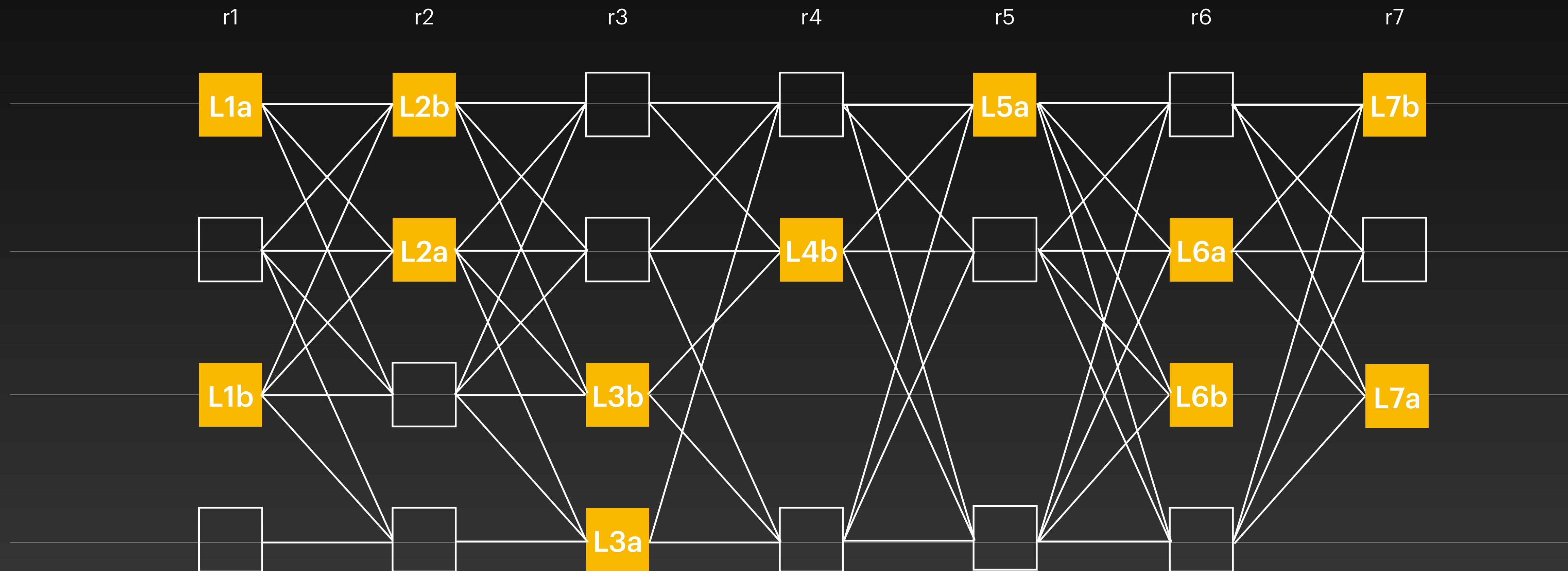


- Round number
- Author
- Payload (transactions)
- Signature

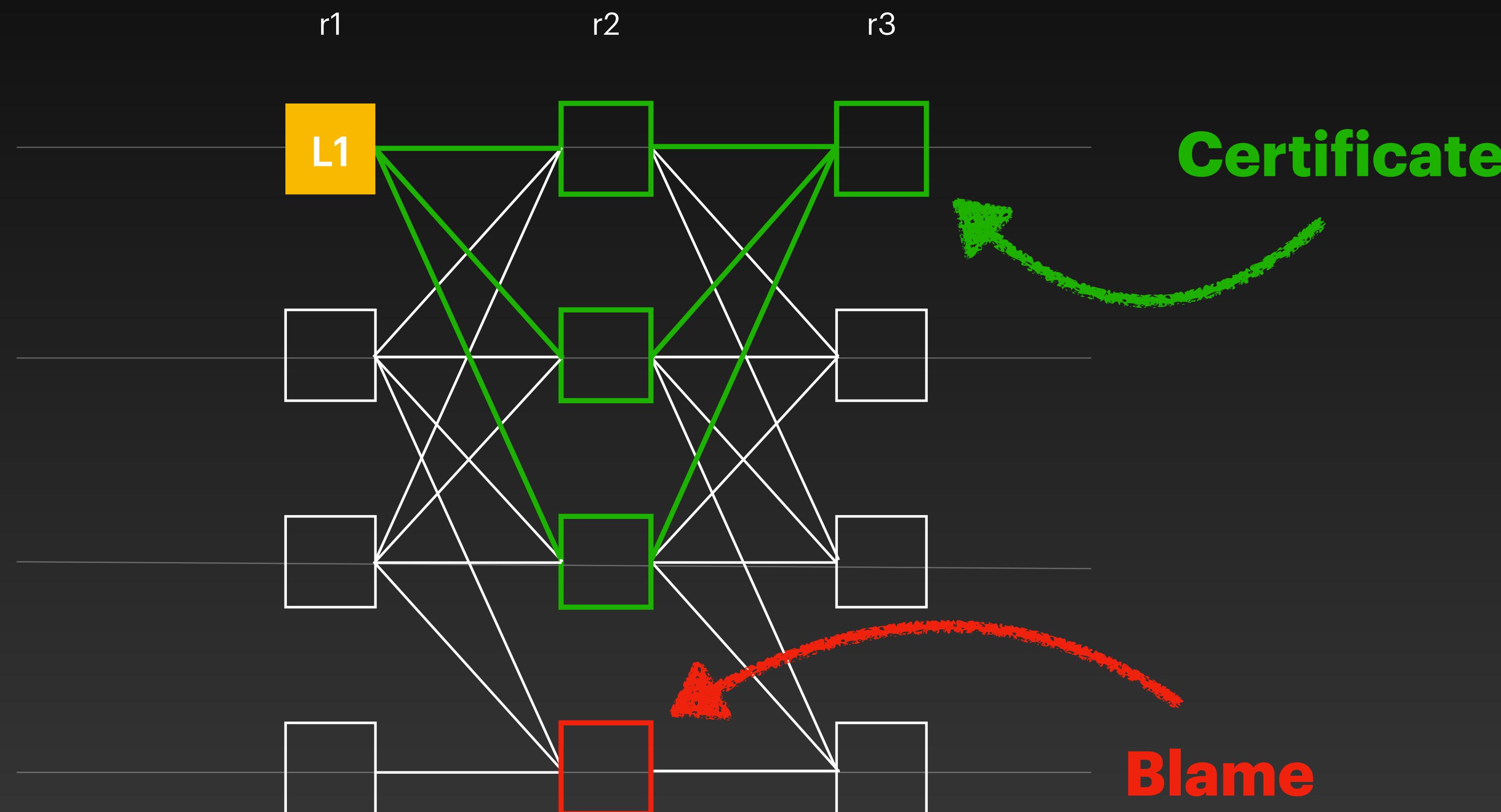
# Uncertified DAG



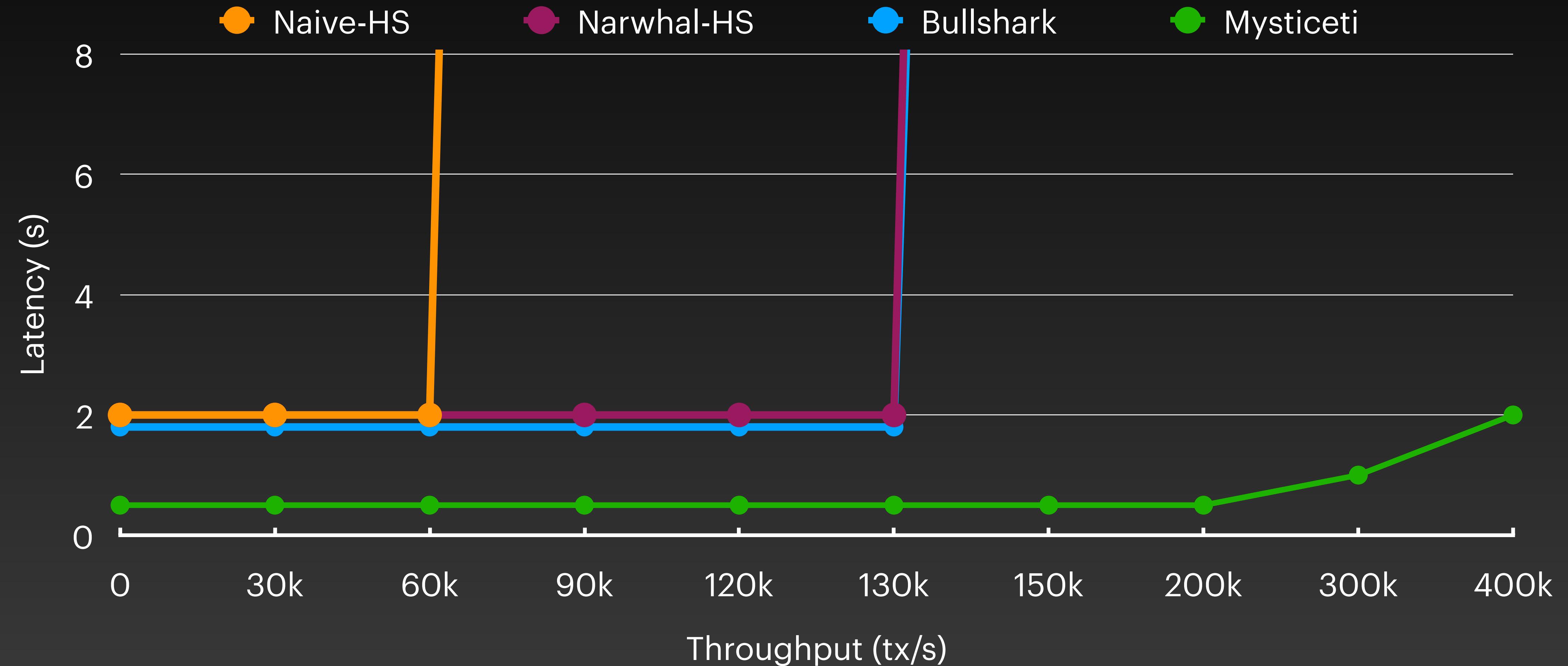
# Uncertified DAG



# Interpreting DAG Patterns



# 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

# 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

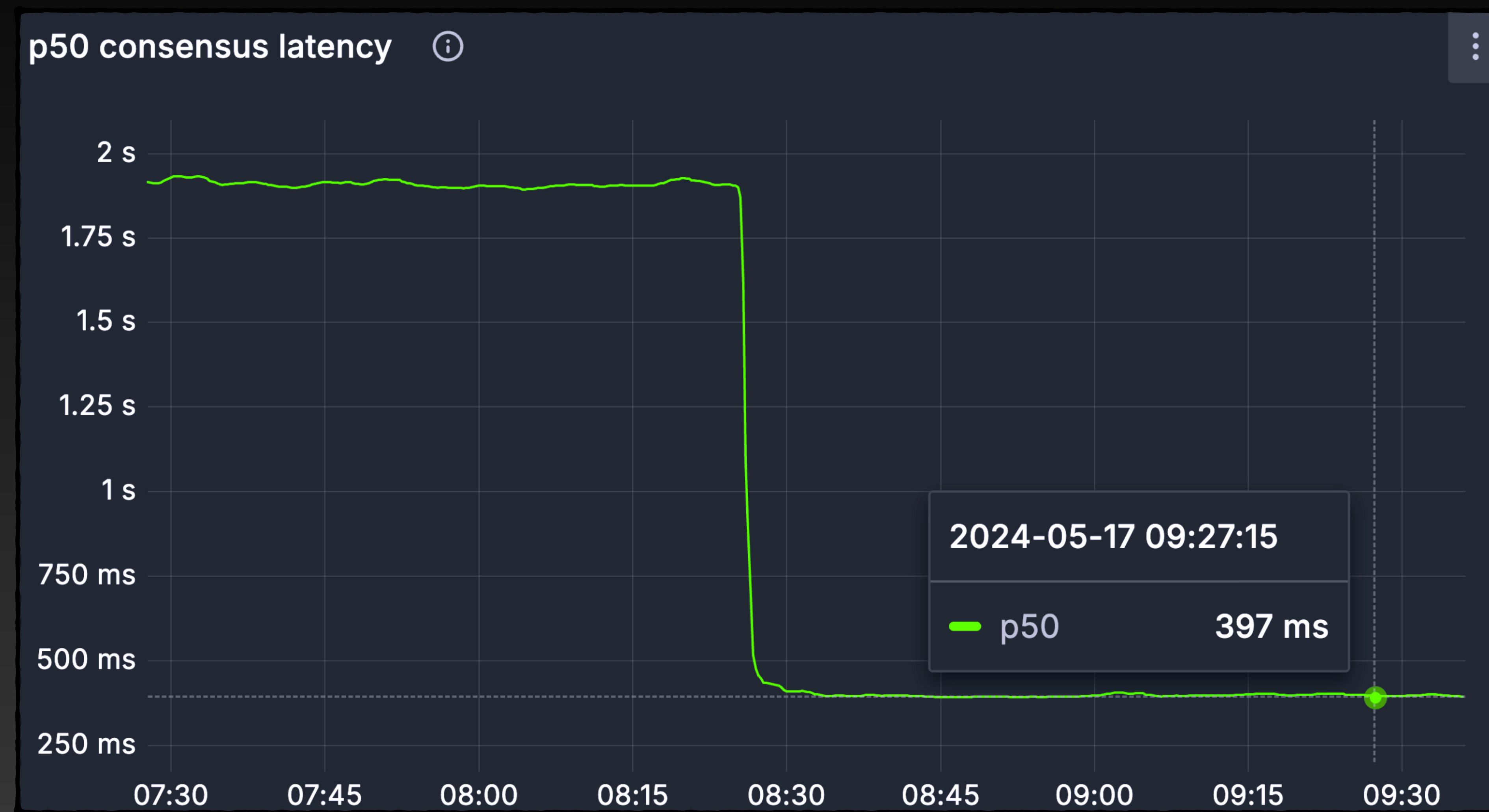
1. Modularisation is a design strategy
2. Tasks-threads allocation
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy
6. Epoch change is not an add-on

# Testing Strategy

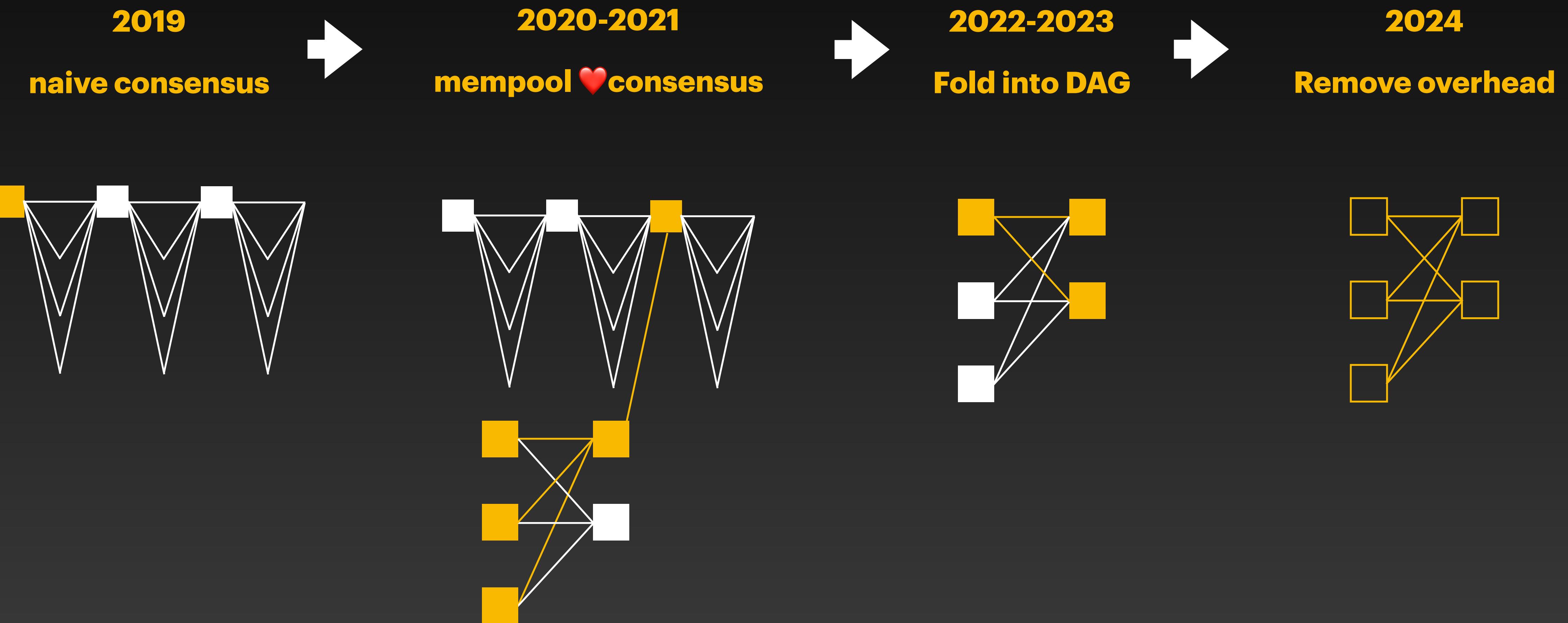


- Compare performance & robustness
- Test mainnet change bullshark -> mysticeti
- Prepare for the worst mysticeti -> bullshark

# The Sui Mainnet



# The Roadmap



# **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 

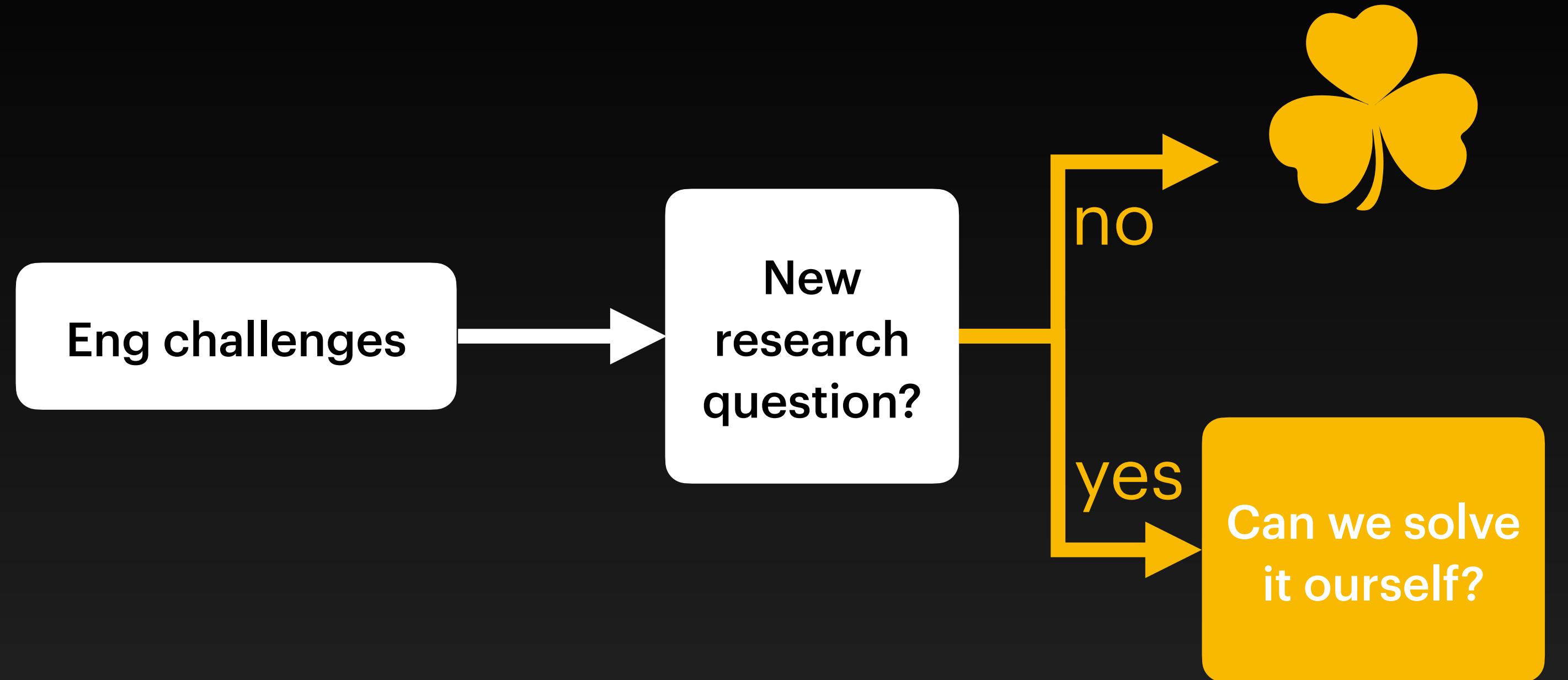
 **Dmitri Perelman** Oct 18th at 7:36 AM  
You're in the invite list!

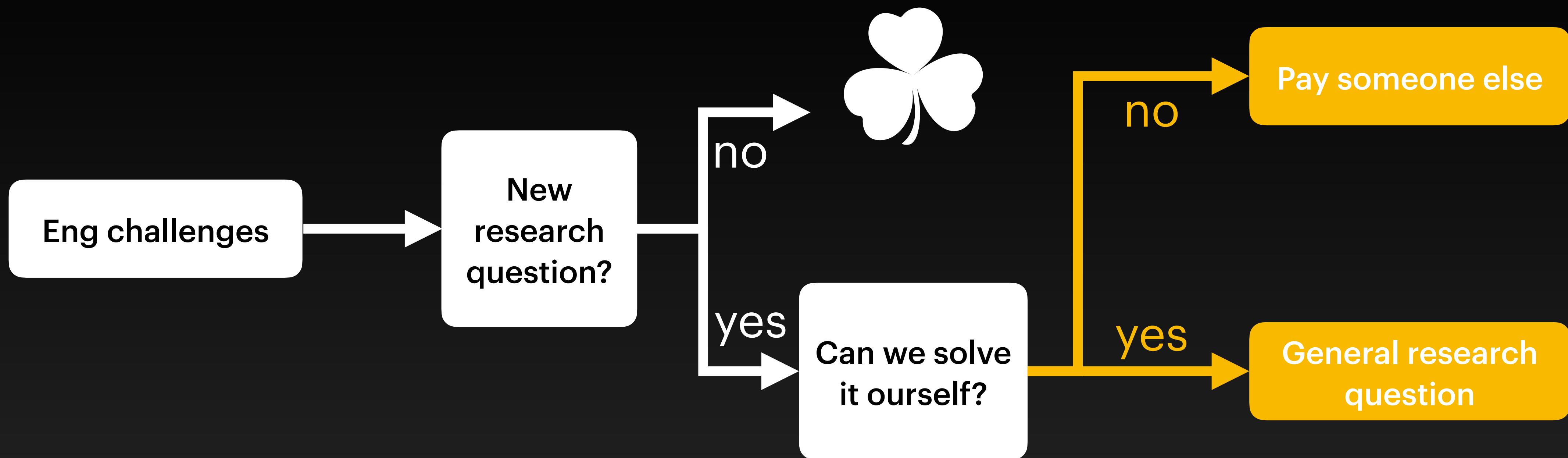
 1 

Eng challenges



New  
research  
question?

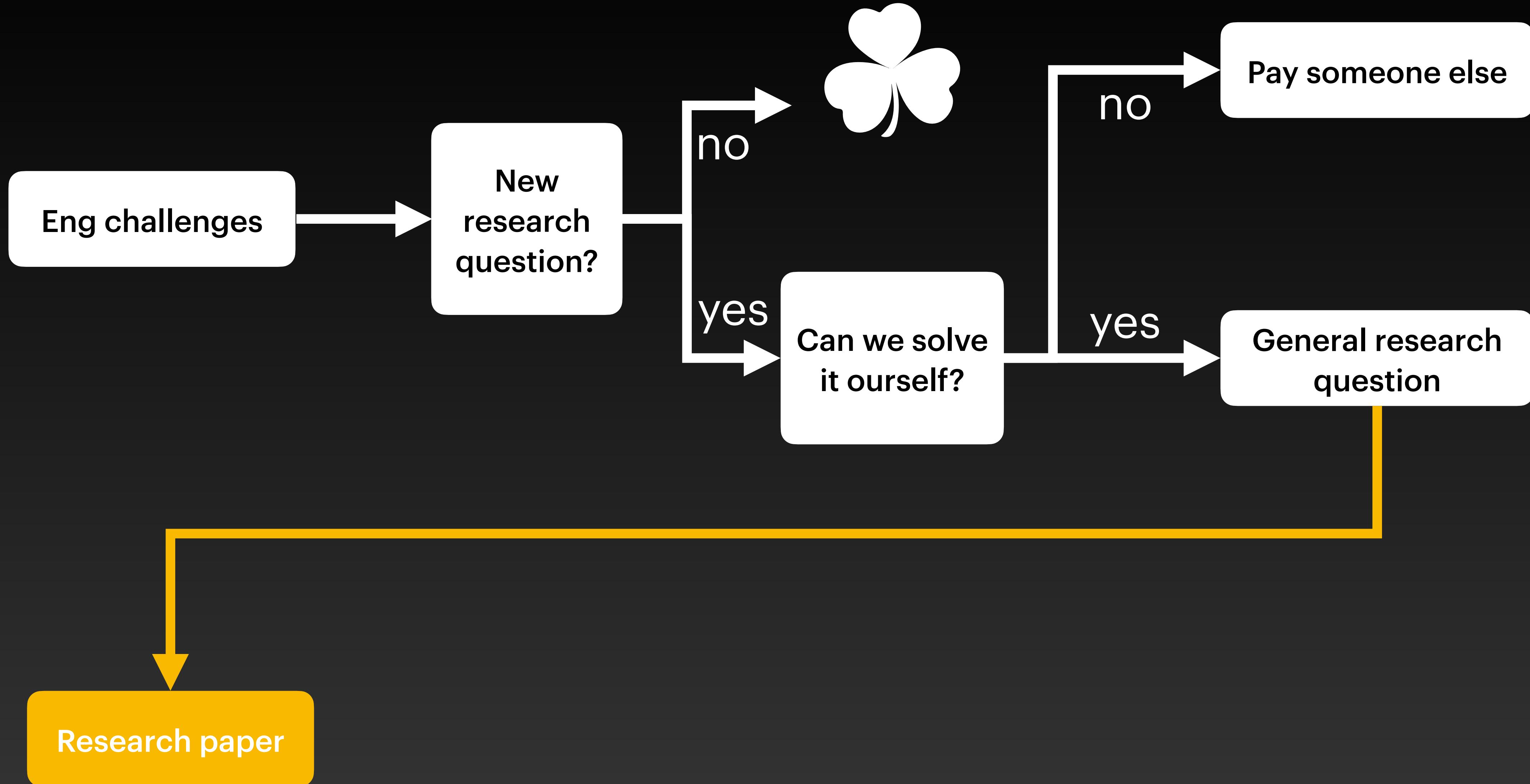


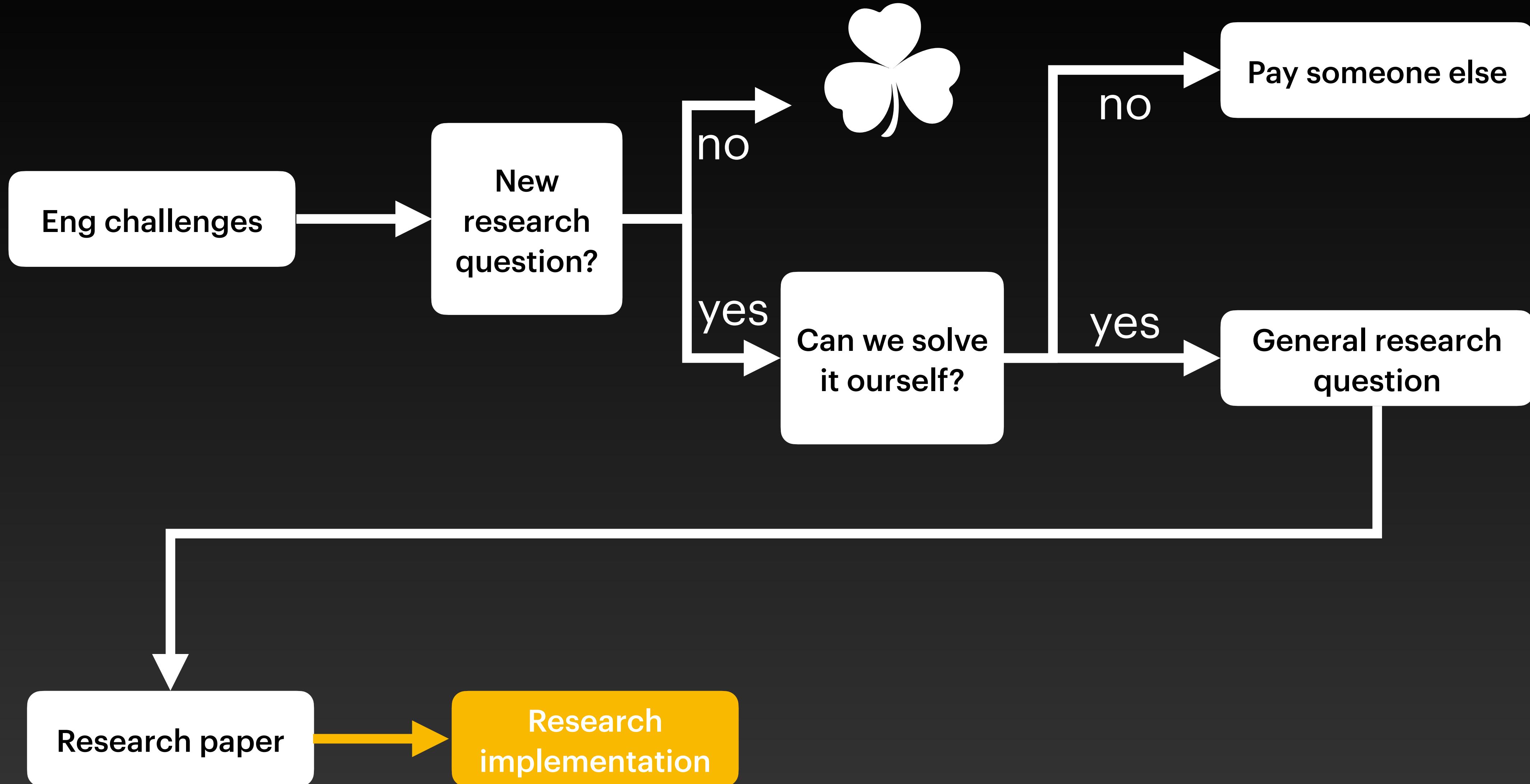


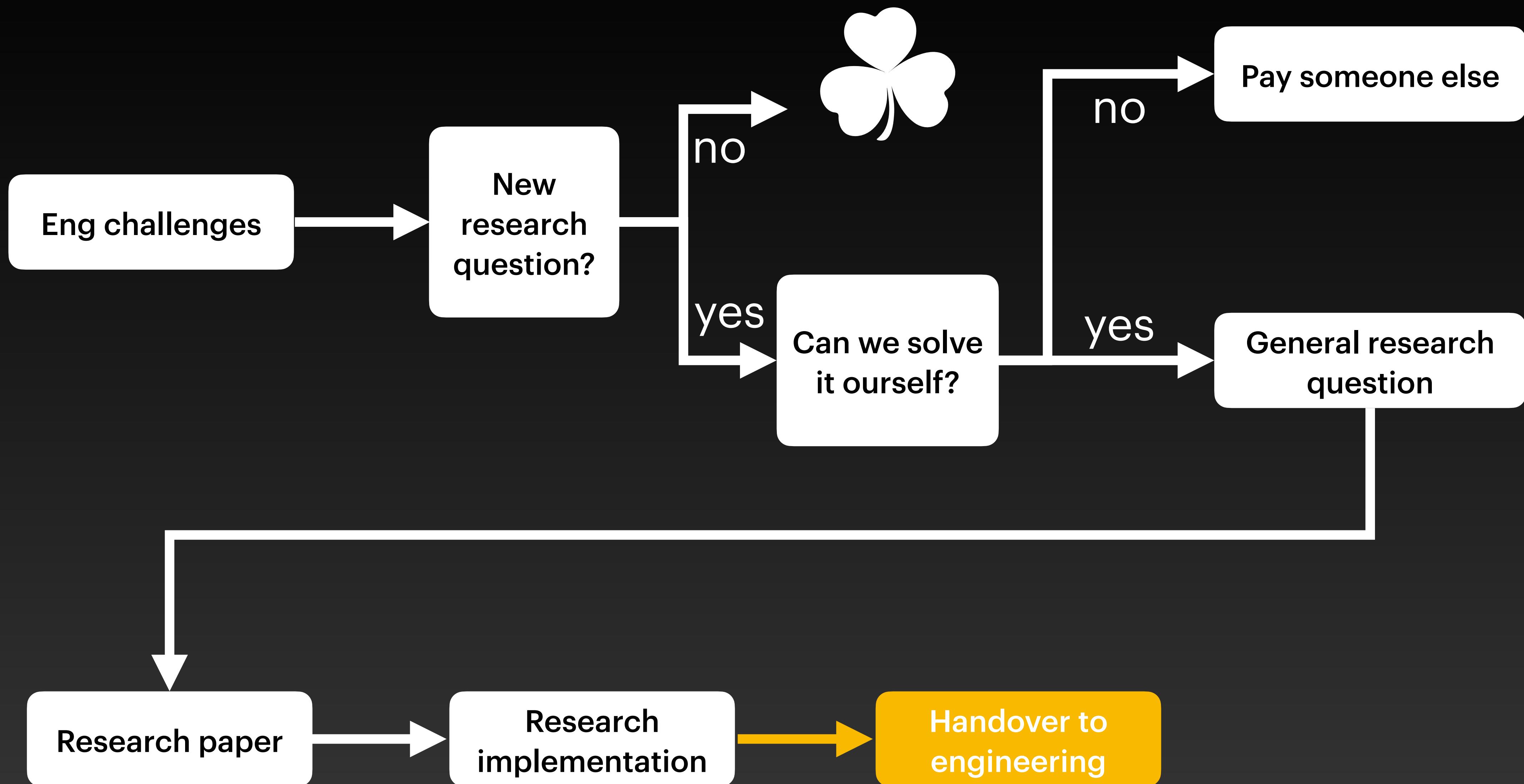
# Research Gifts

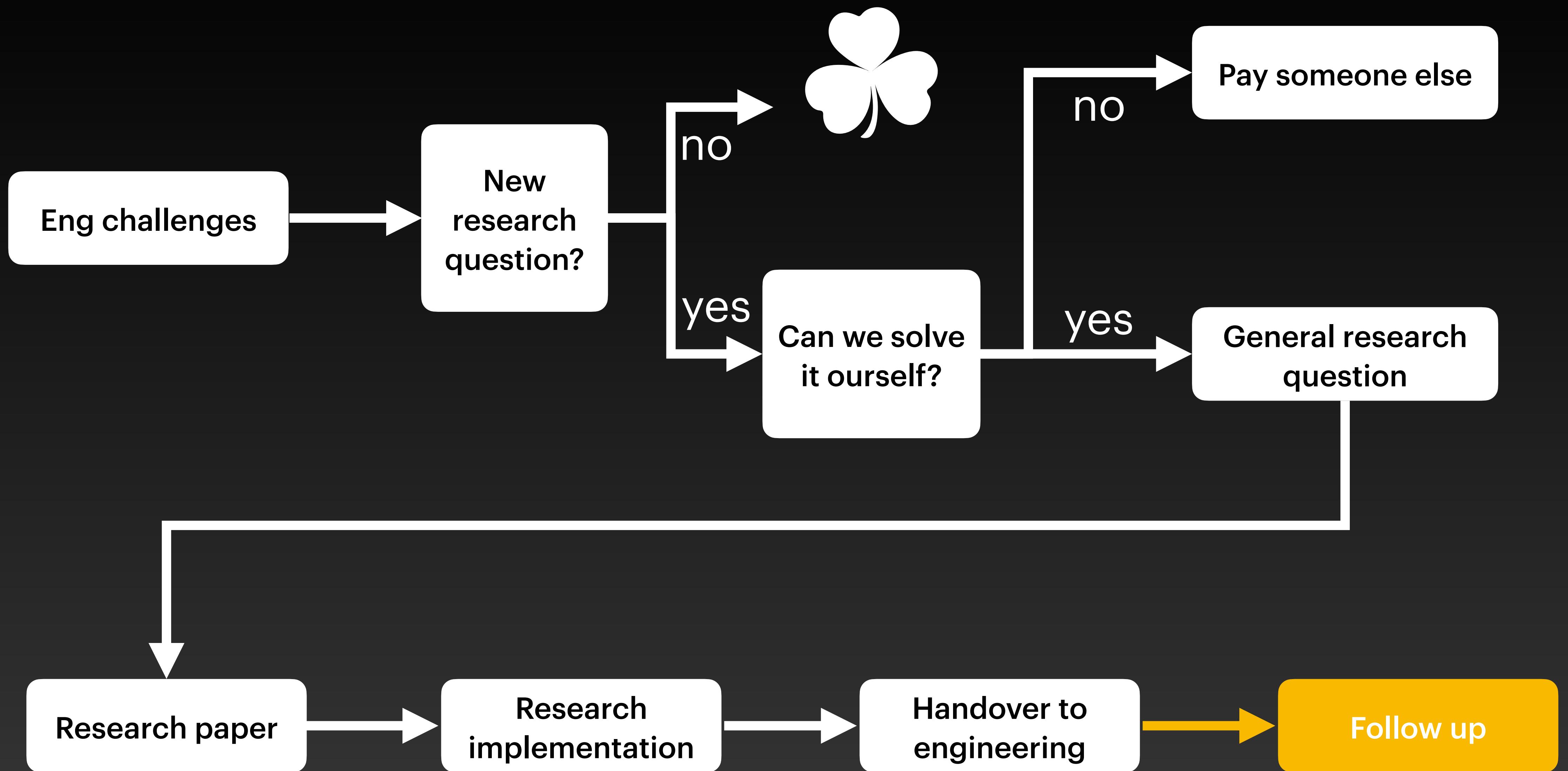


(please keep it short)









# 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>

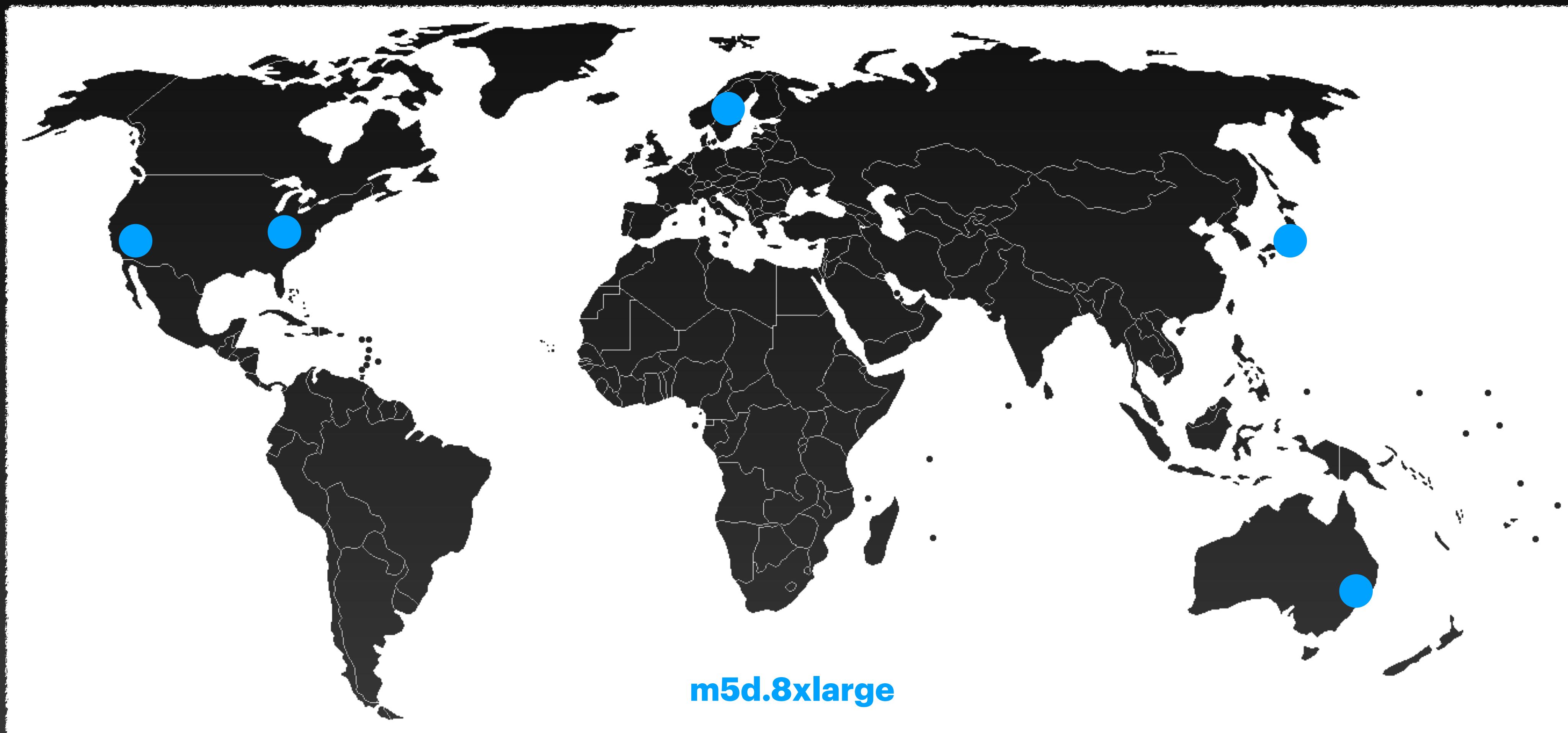
# Implementation

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

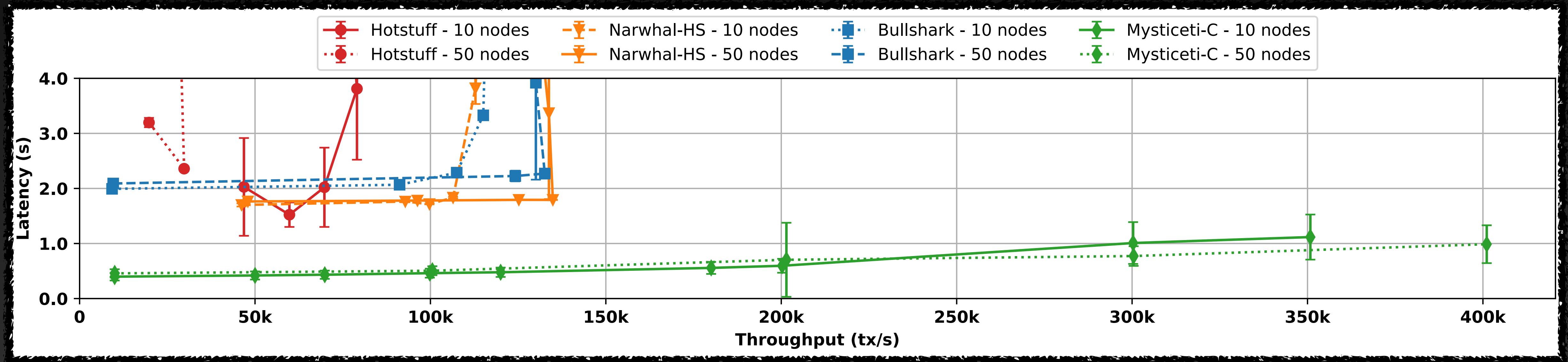
<https://github.com/mystenlabs/mysticeti>

# Evaluation

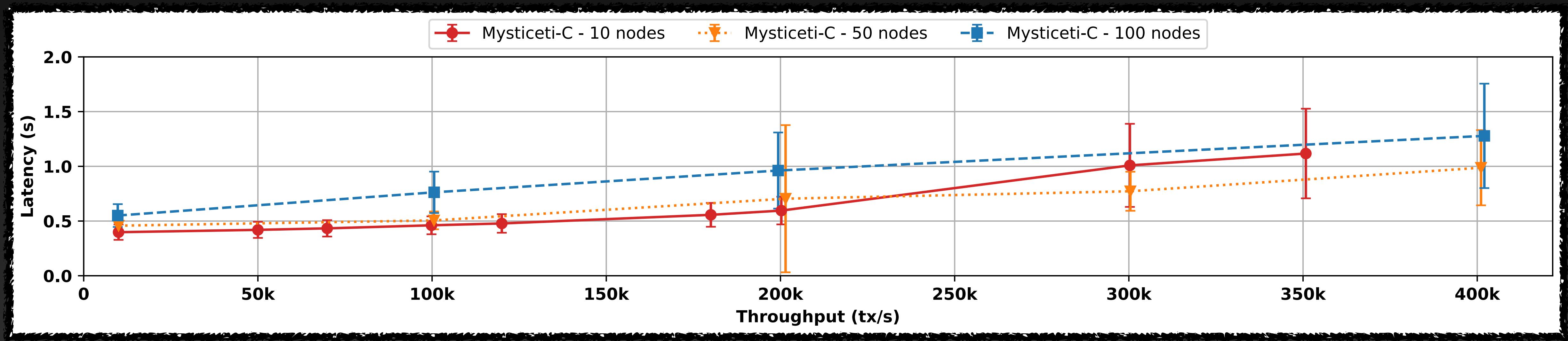
## Experimental setup on AWS



# Prototype Benchmarks



# Prototype Benchmarks



# Mysticeti

## Key Limitations

- Block Synchroniser
- Parallelise block creation and synchronisation
- Rigid DAG structure?