

# BFT Consensus

From Academic Paper to Mainnet

Alberto Sonnino

# Alberto Sonnino

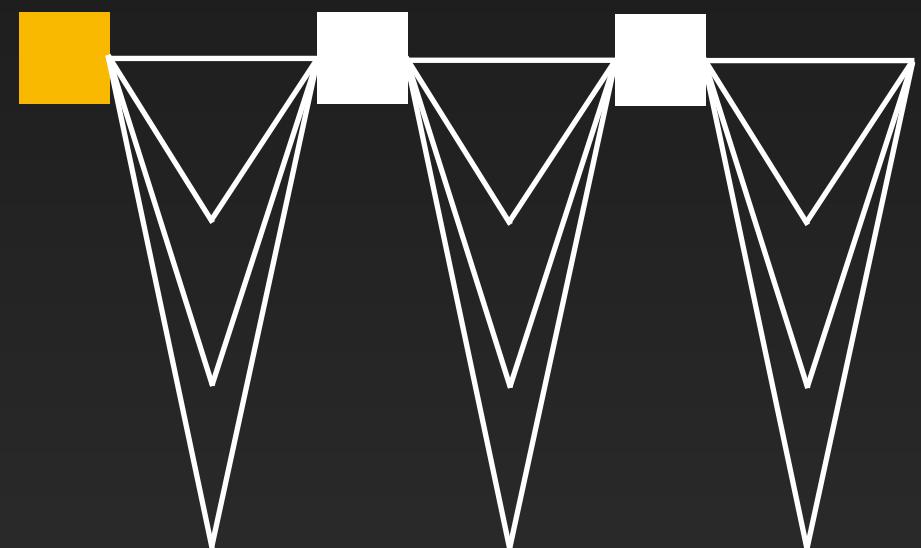
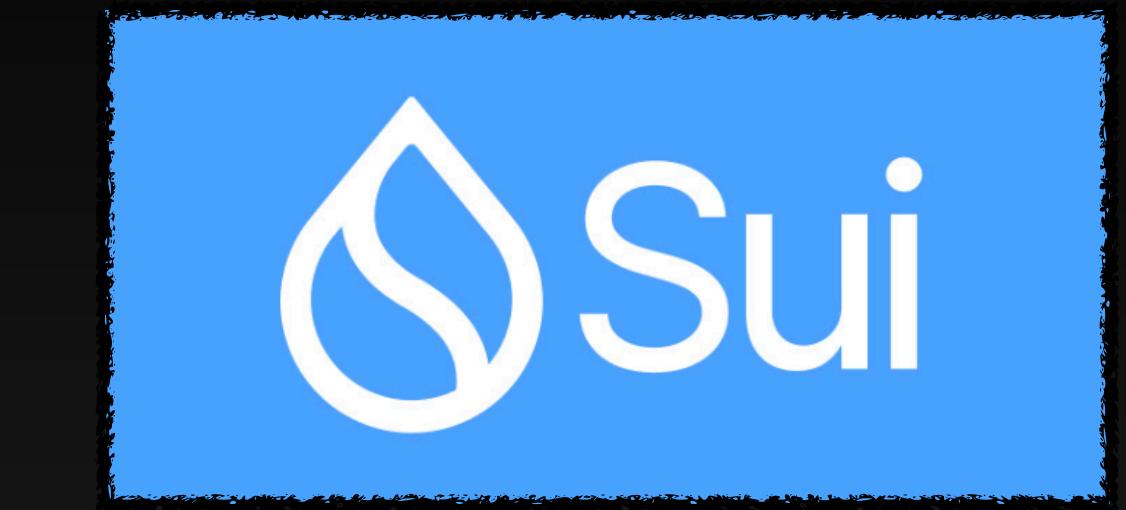
## Research Scientist

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

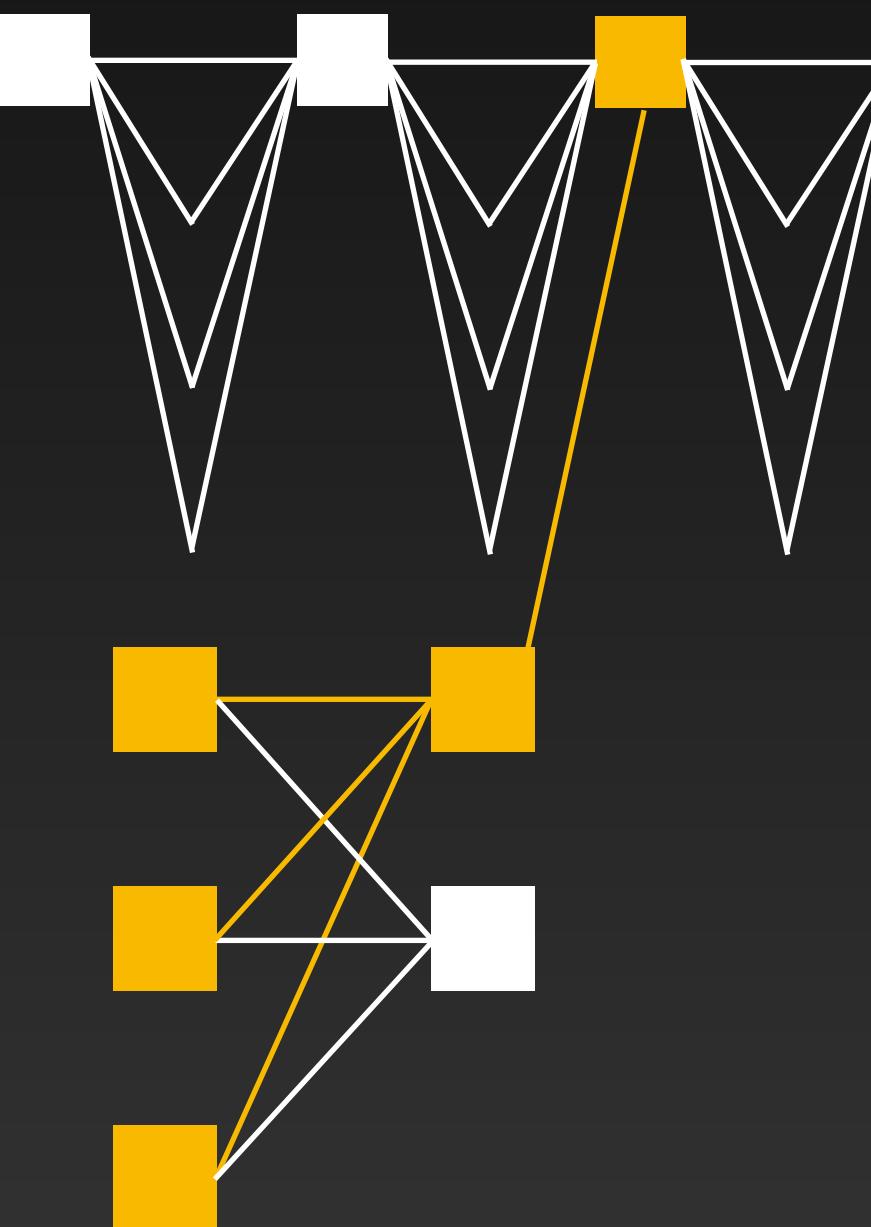
2019



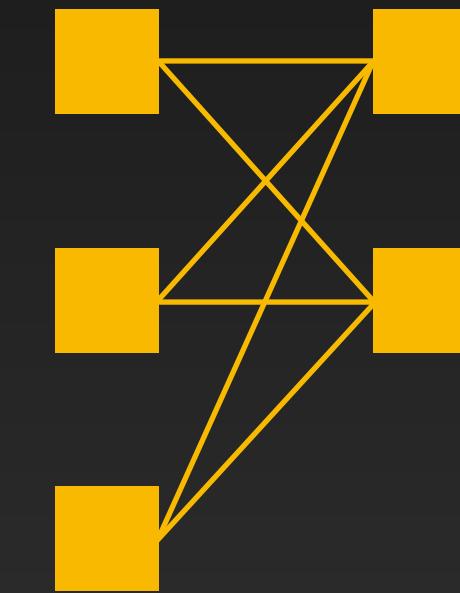
2024



**HotStuff**



**HotStuff + Mempool**

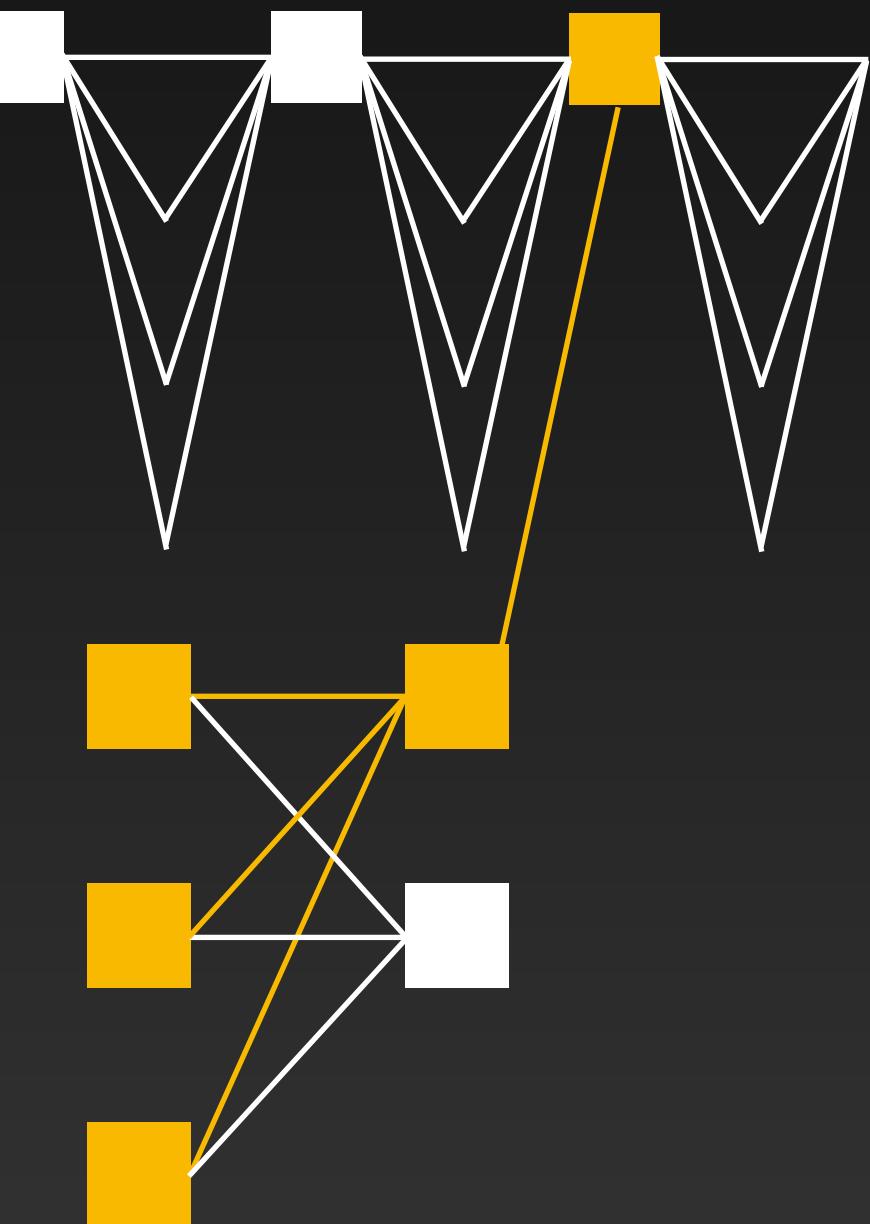


**Bullshark,  
Mysticeti**

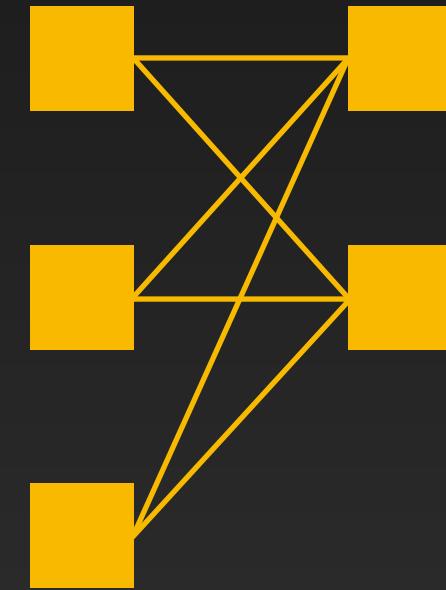
2019



2024



- Lessons learned
- Open research challenges



# Research Gifts



(please keep it short)

# Byzantine Fault Tolerance



# Byzantine Fault Tolerance



$> 2/3$



# Partial Synchrony

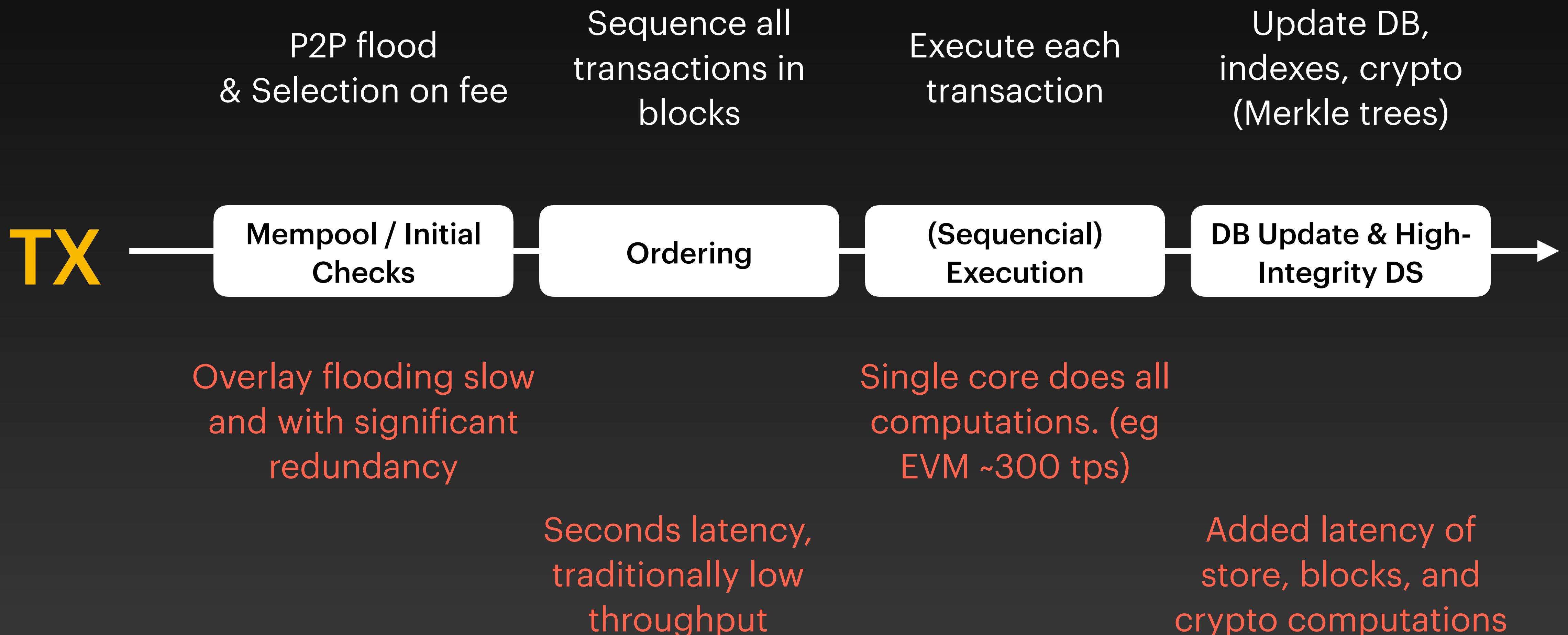


# Research Questions

1. Network model?

# Lessons Learned

# Typical Blockchain



# Typical Blockchain

**Sequence all  
transactions in  
blocks**

Ordering

**Seconds latency,  
traditionally low  
throughput**

# 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.

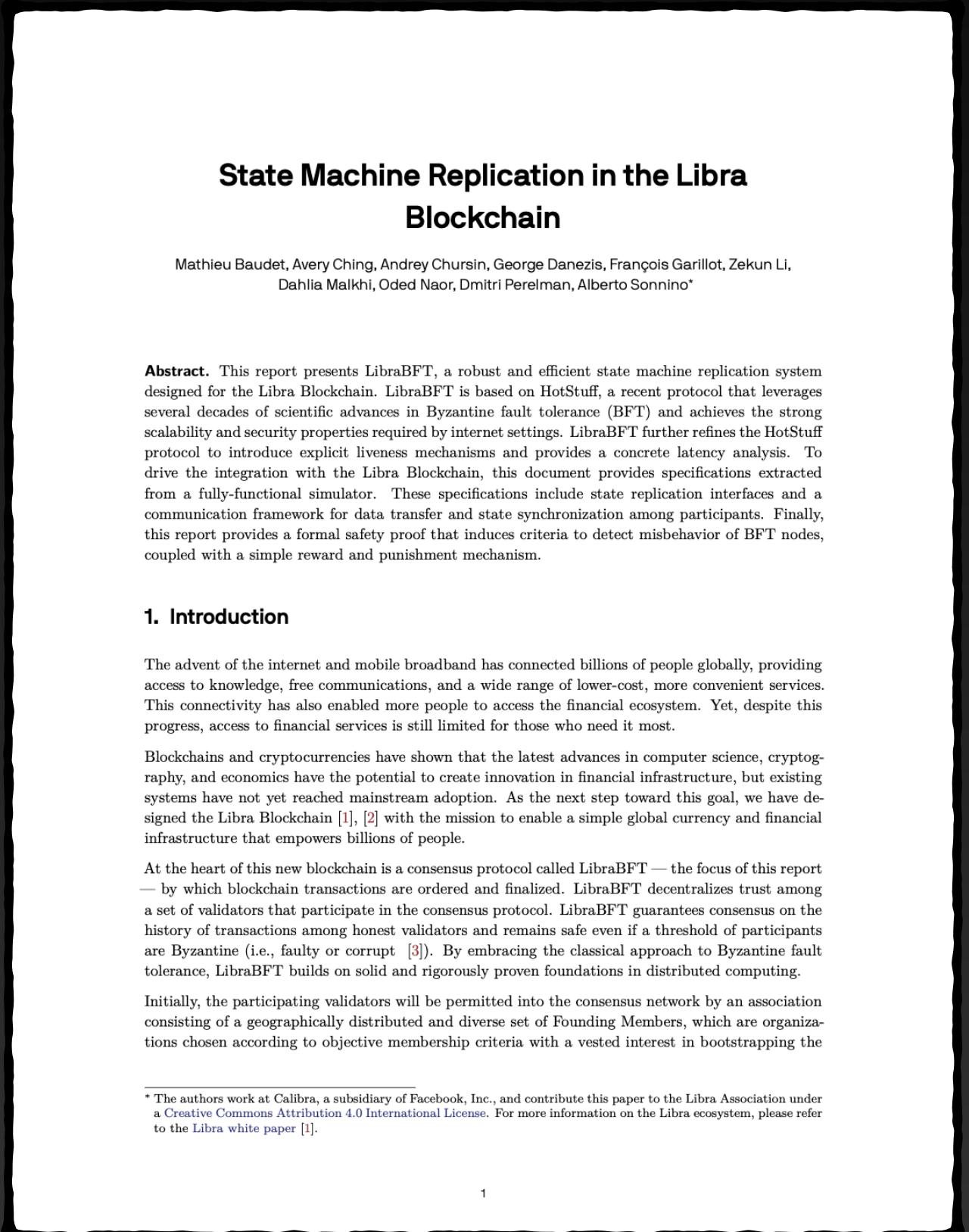
The algorithm for a new leader to collect information and propose it to replicas—called a *view-change*—is the epicenter of replication. Unfortunately, view-change based on the two-phase paradigm is far from simple [38], is bug-prone [4], and incurs a significant communication penalty for even moderate 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

# The first 6 months...



## SMR in the Libra Blockchain

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

### 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].

# Research Questions

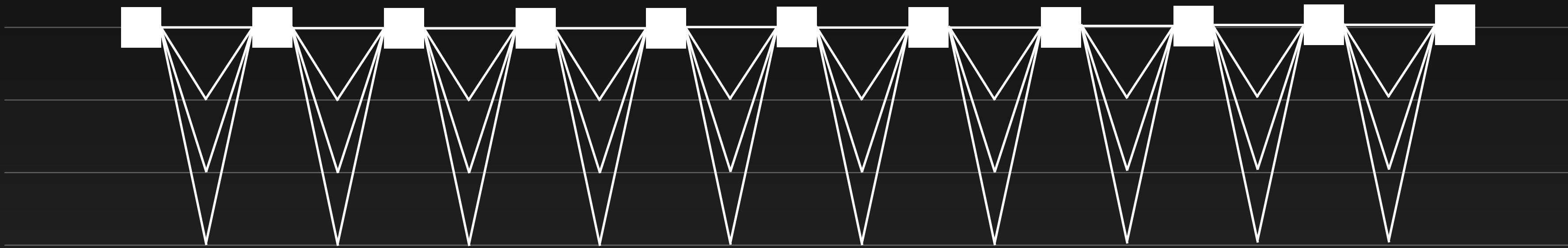
1. Network model?

# Lessons Learned

1. Modularisation is a design strategy

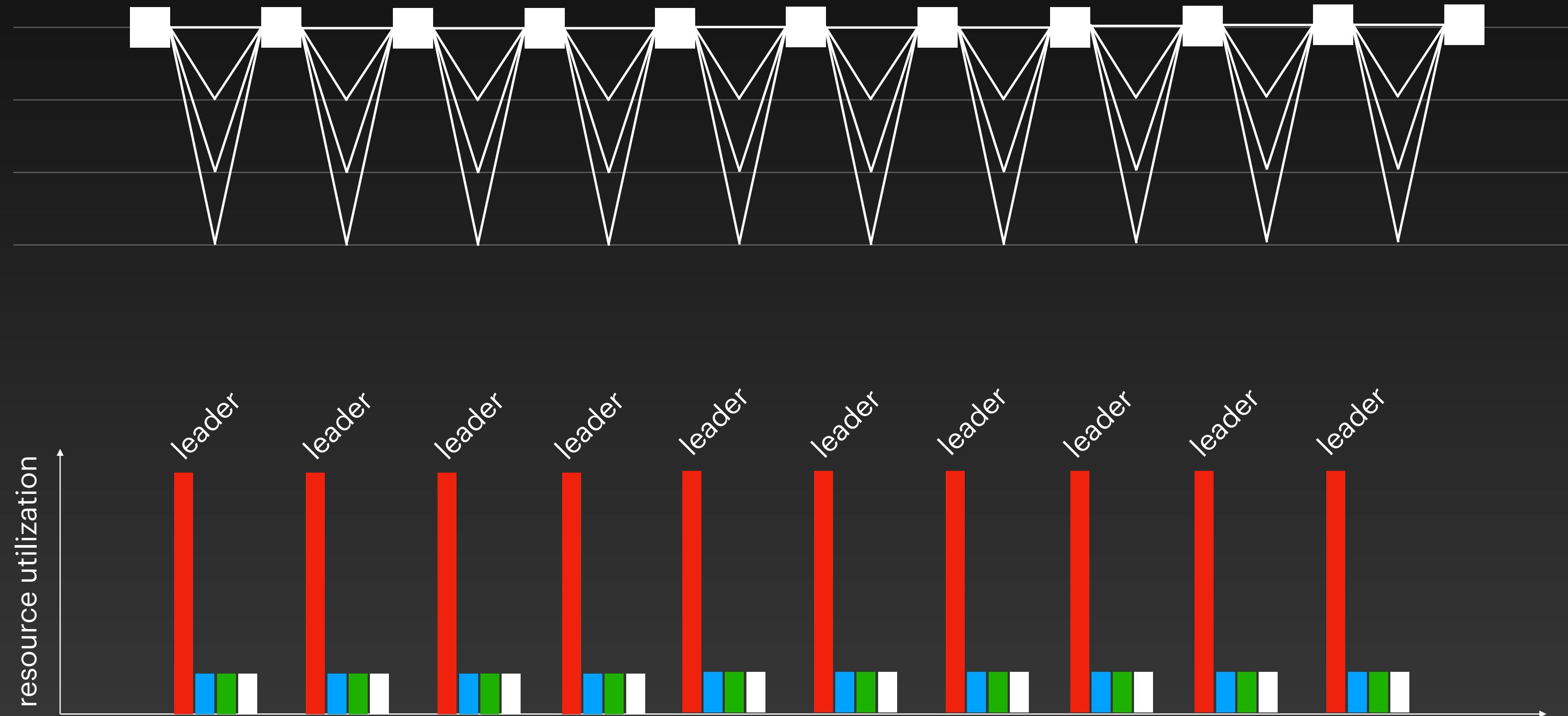
# HotStuff

## Typical leader-based protocols



# HotStuff

## Uneven resource utilisation



# Research Questions

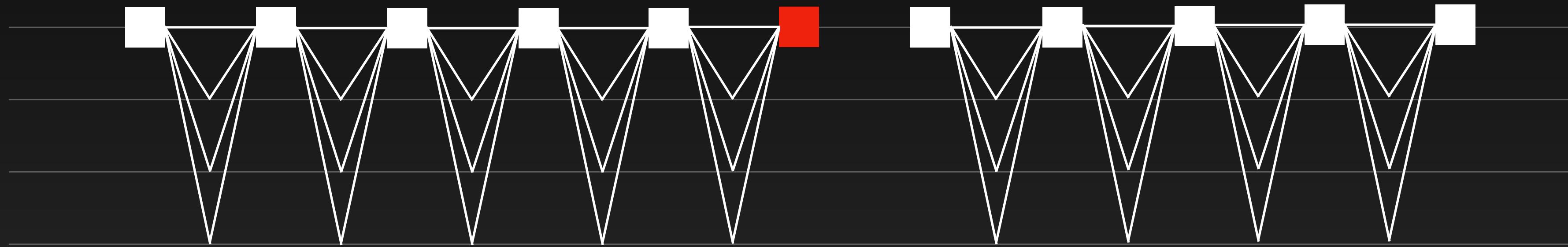
1. Network model?

# Lessons Learned

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

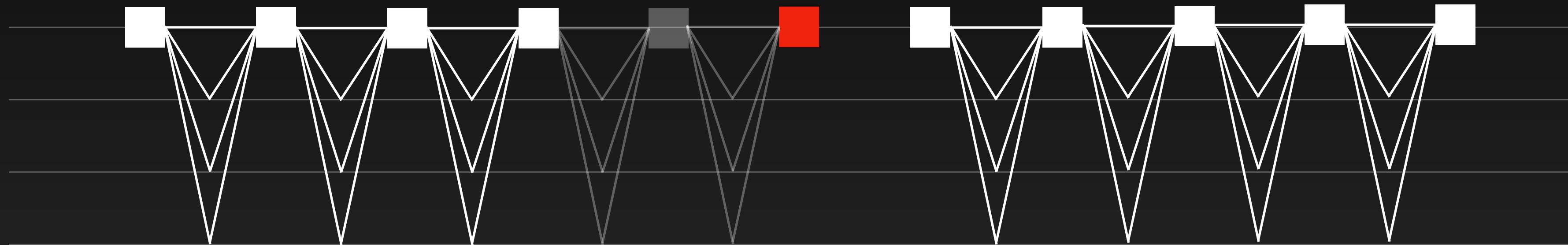
# HotStuff

## Fragility to faults and asynchrony

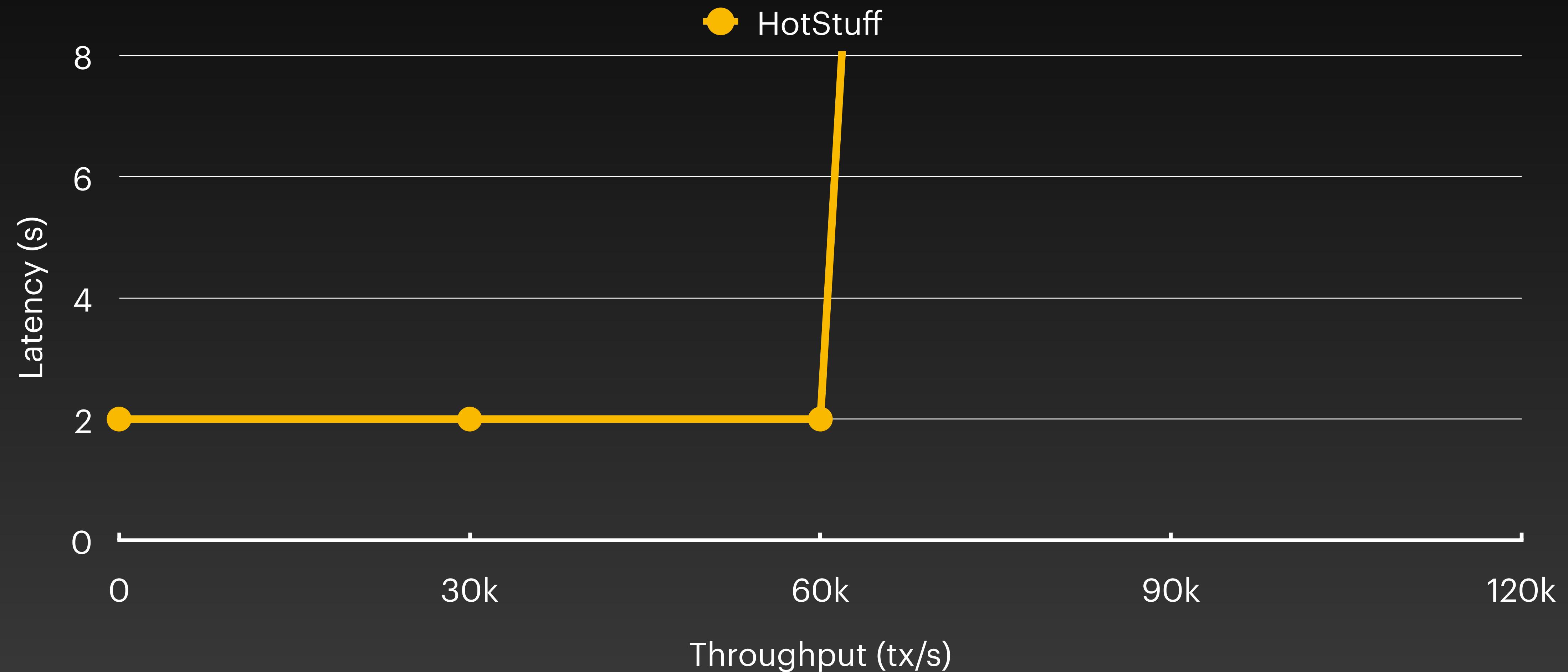


# HotStuff

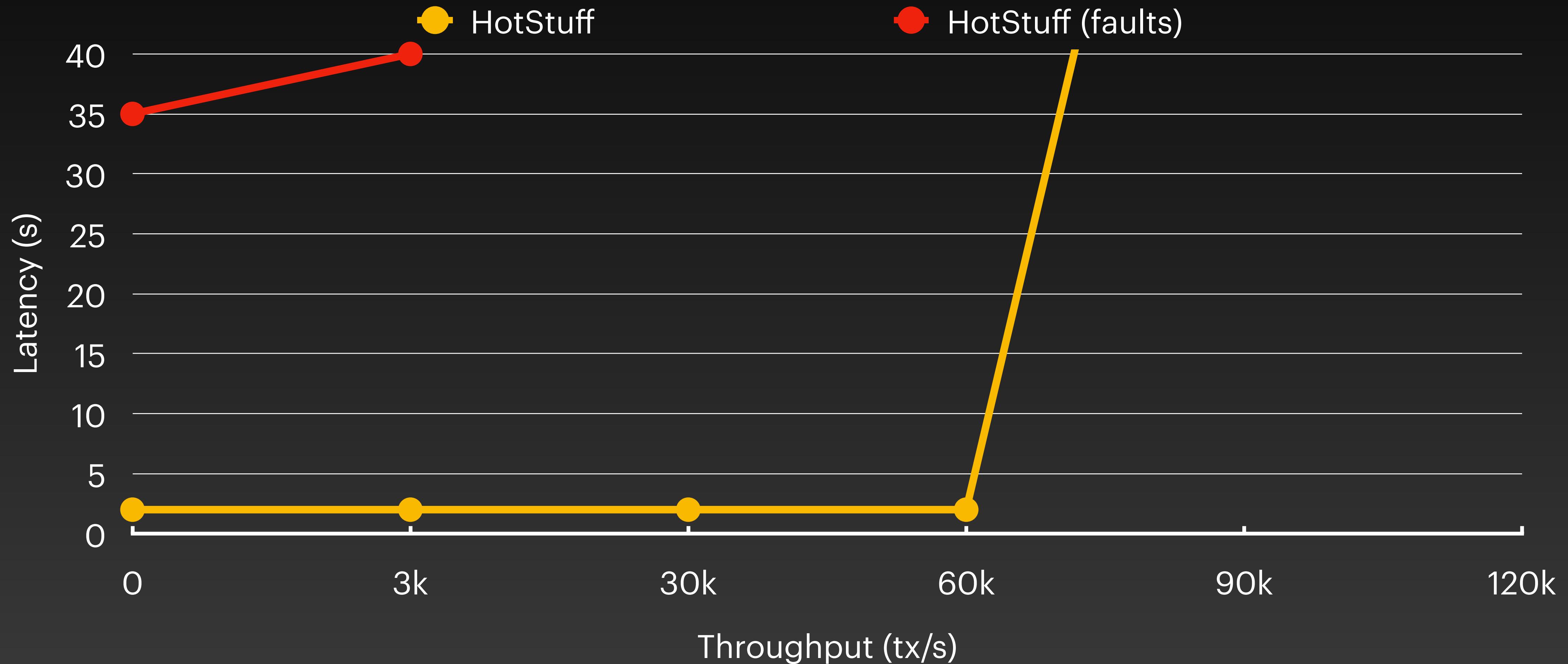
## Fragility to faults and asynchrony



# Performance



# Performance



# Research Questions

1. Network model?

# Lessons Learned

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

# 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
- 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).

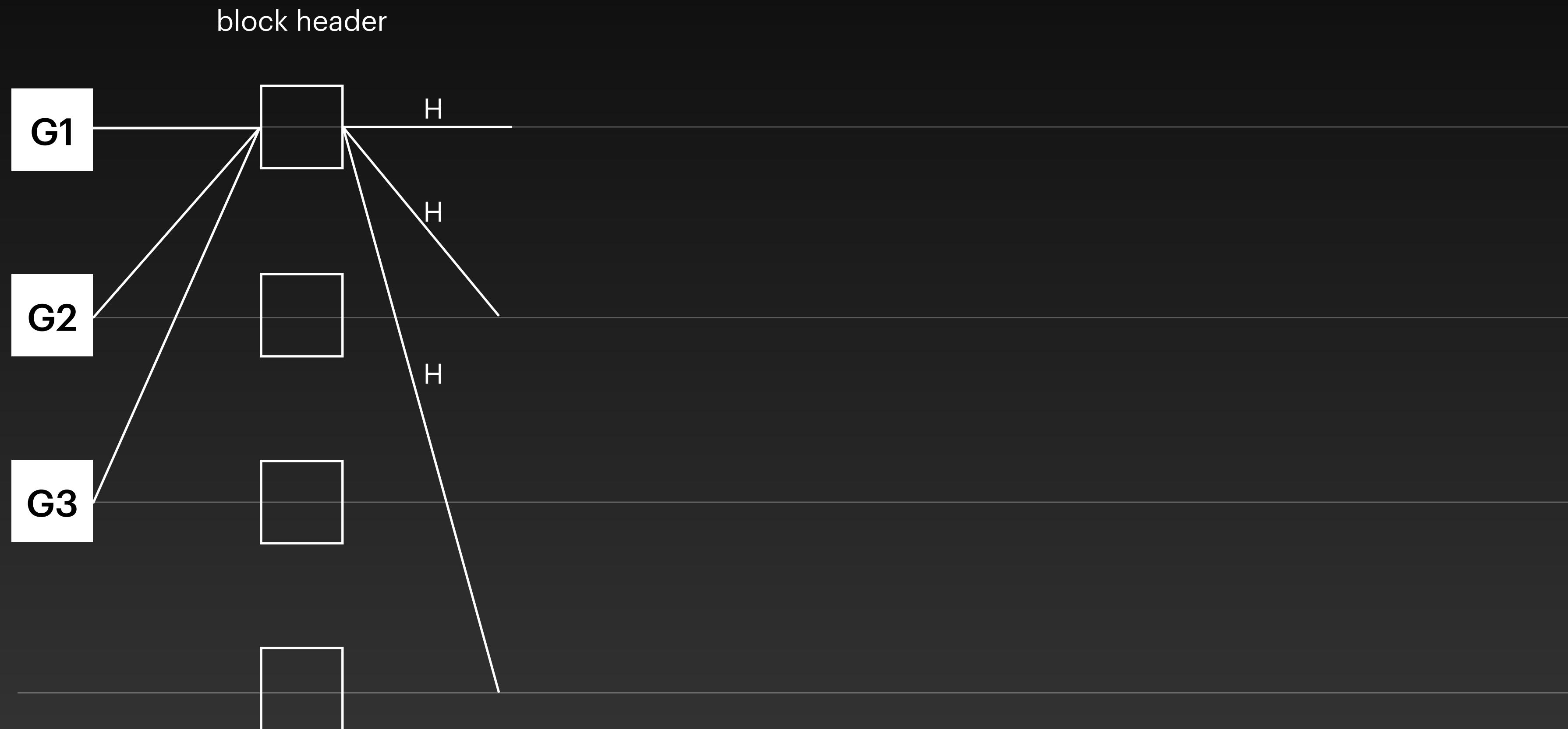
Additionally, consensus protocols have grouped a lot of functions into a monolithic protocol. In a typical distributed

## Narwhal

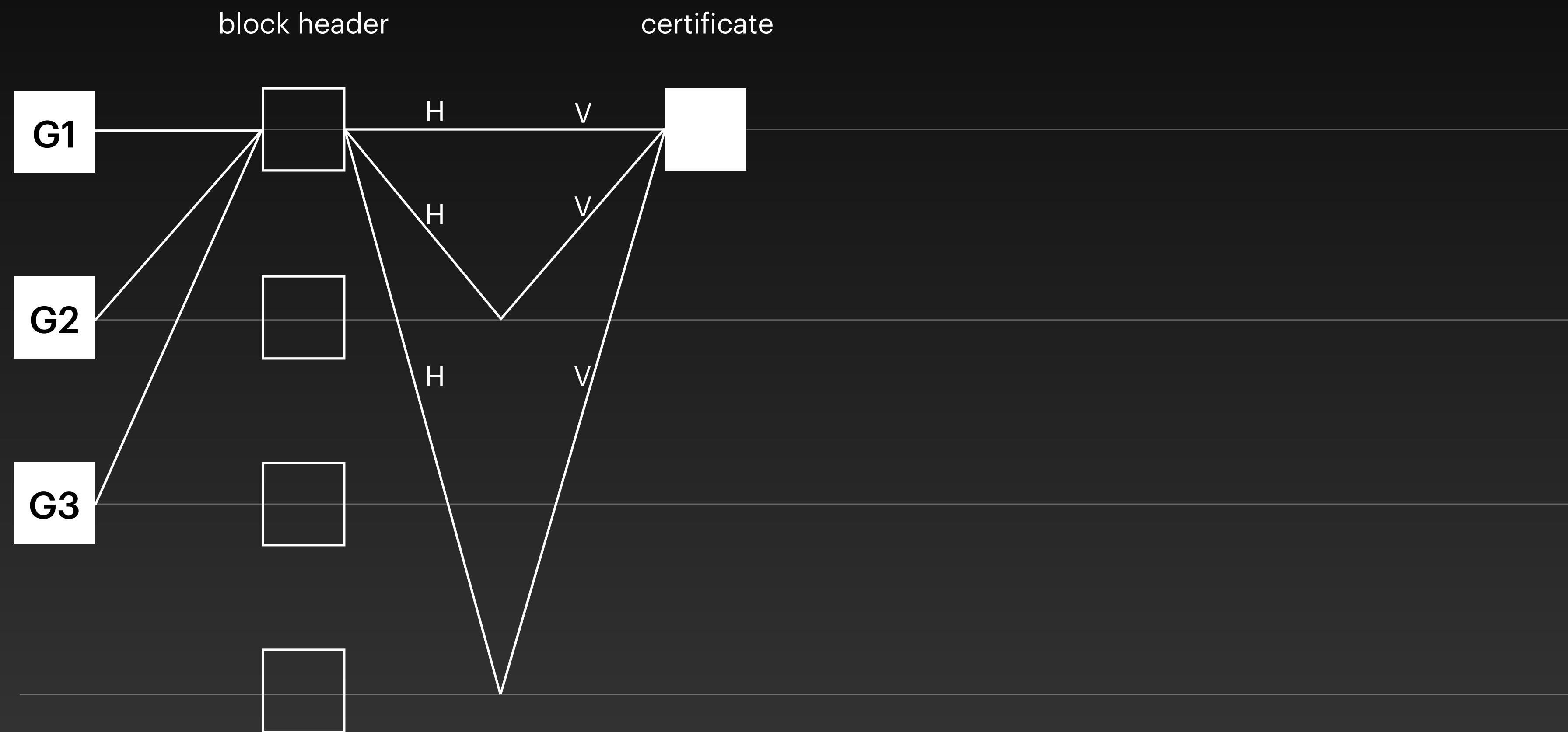
- Quadratic but even resource utilisation
- Separation between consensus and data dissemination

# Narwhal

## The primary machine

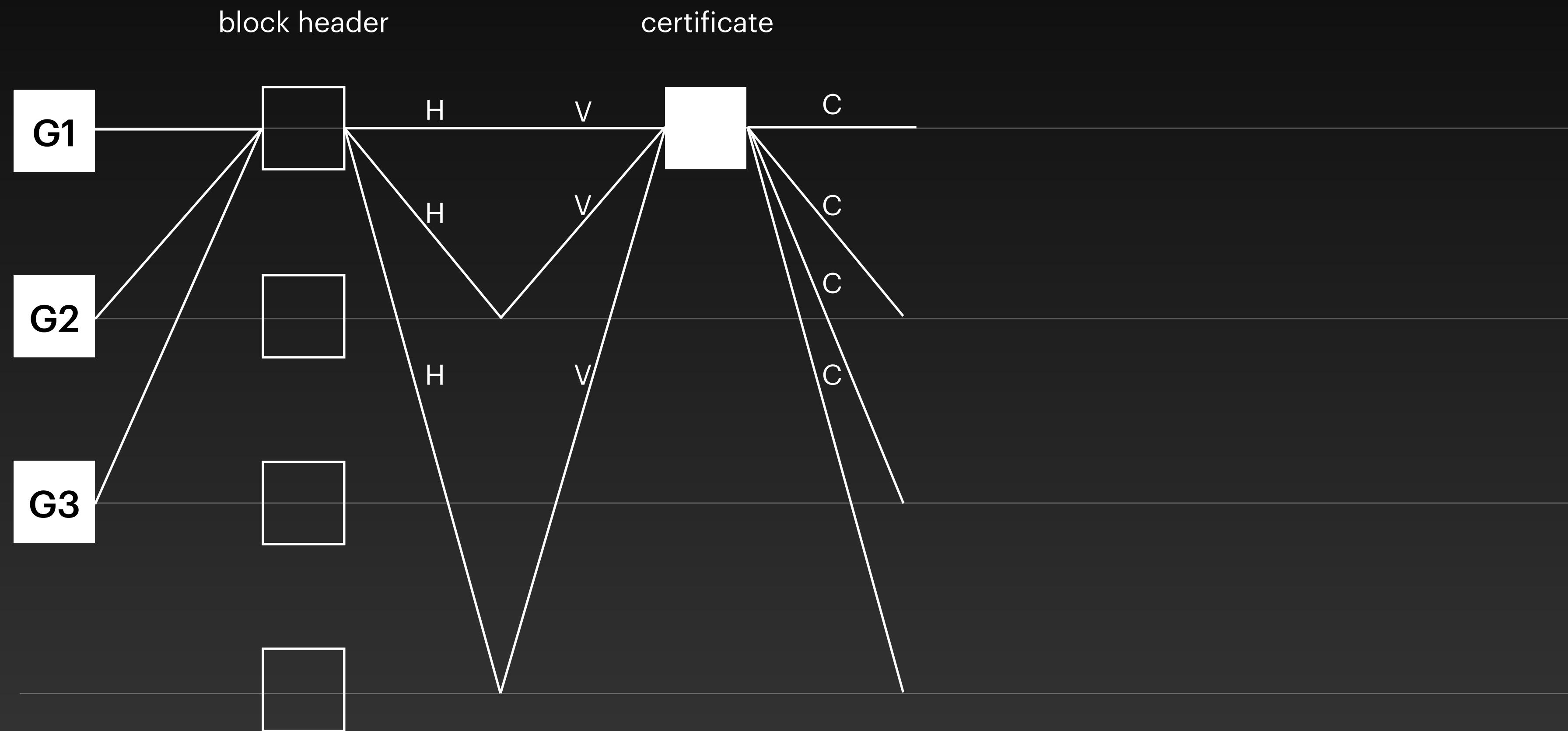


# Narwhal



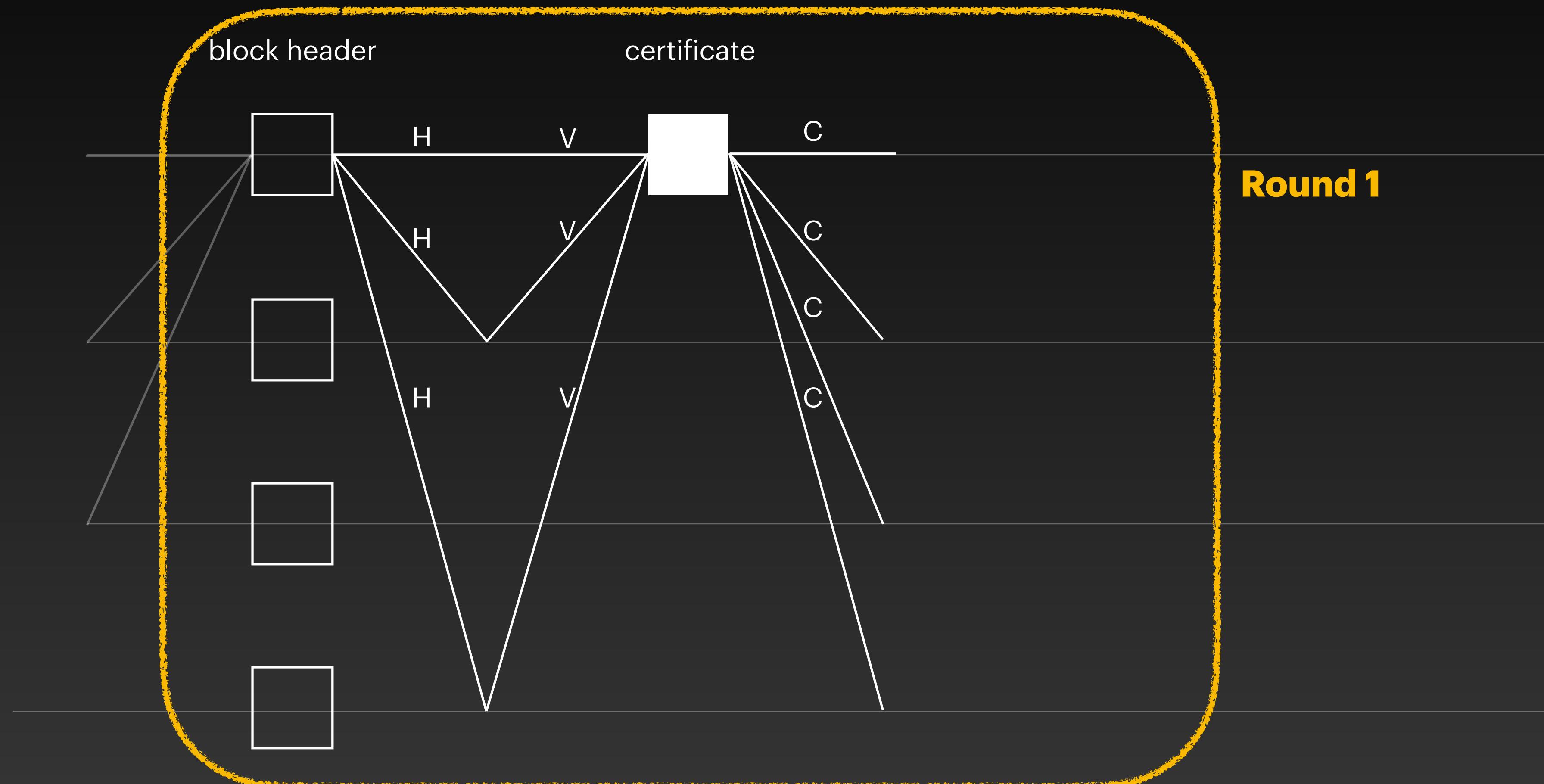
# Narwhal

## The primary machine



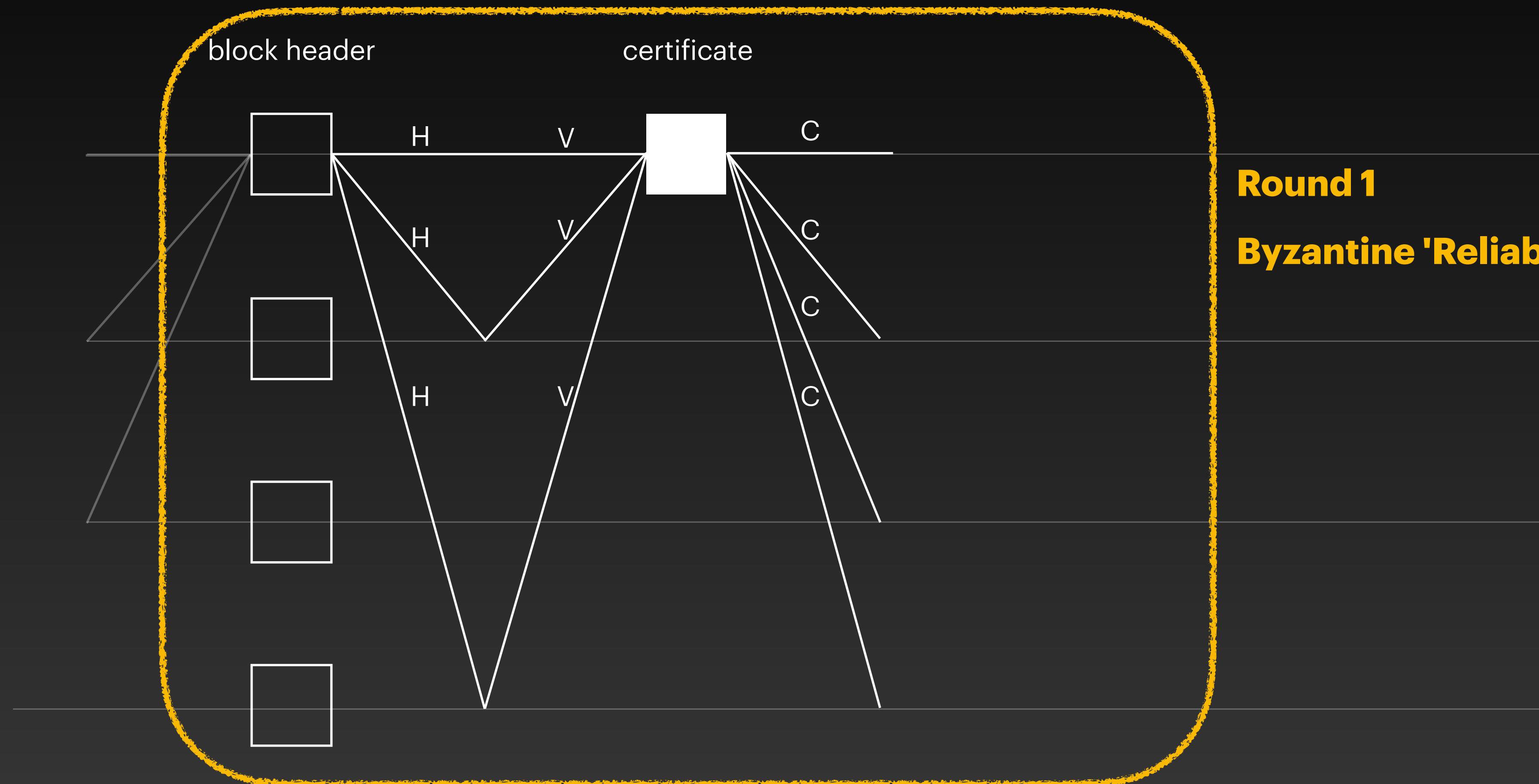
# Narwhal

## The primary machine



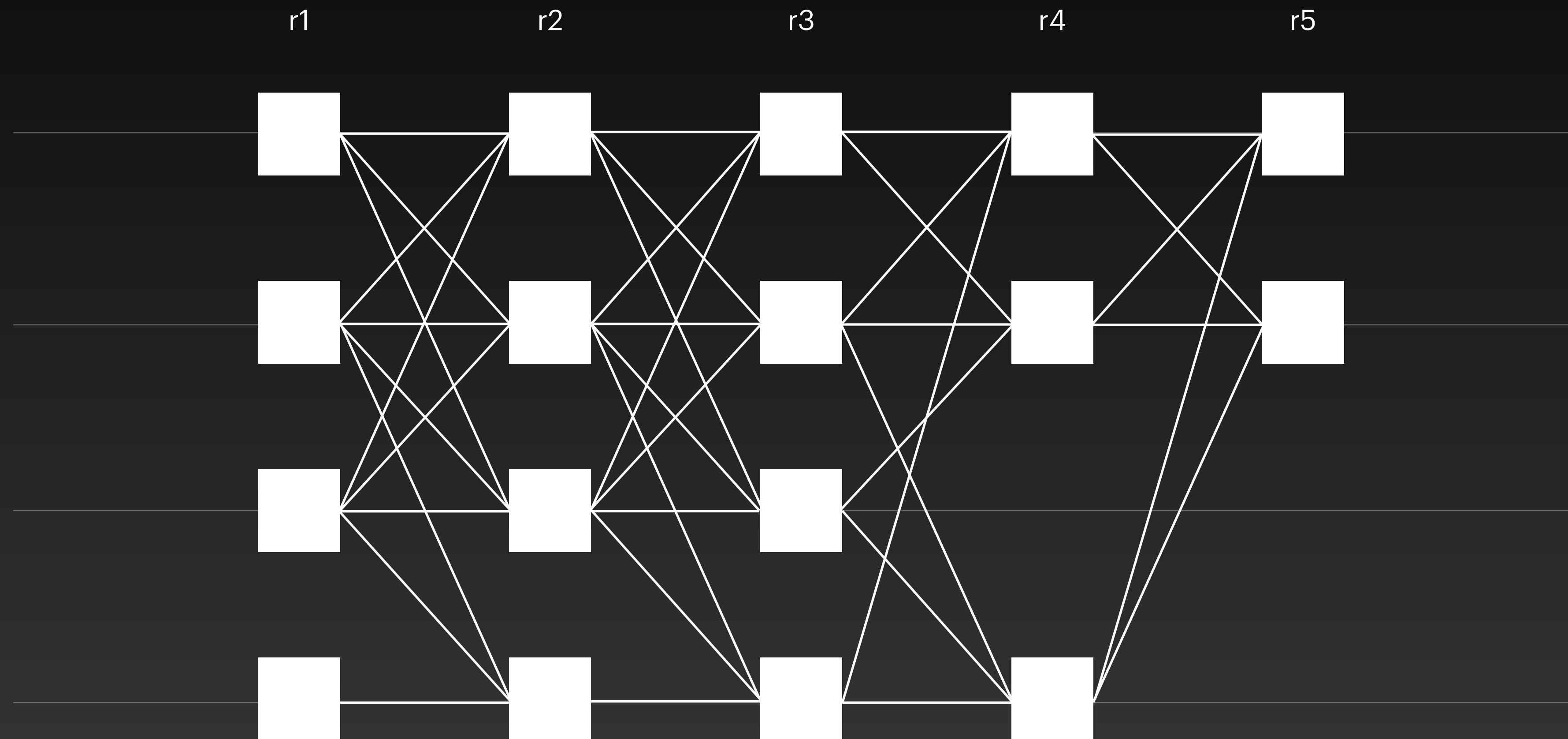
# Narwhal

## The primary machine



# Narwhal

## The primary machine



# Research Questions

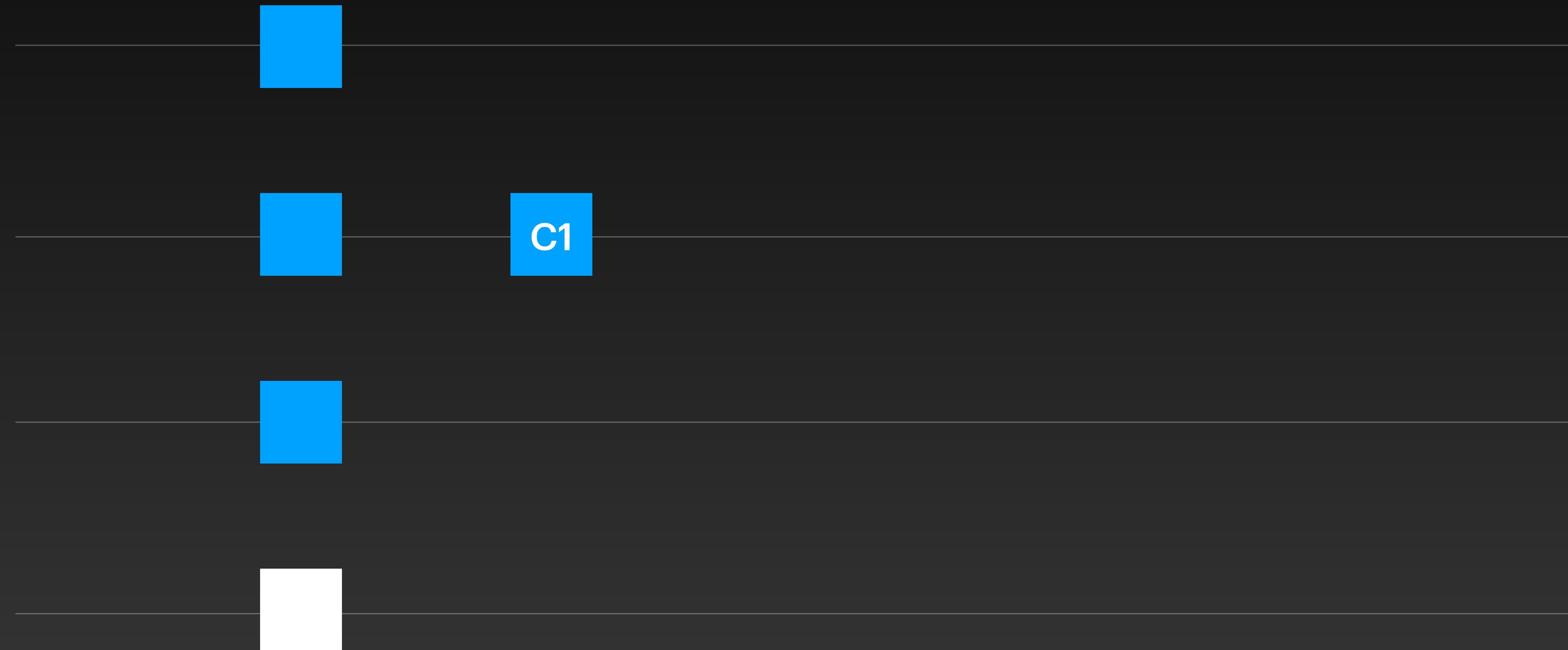
1. Network model?

# Lessons Learned

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

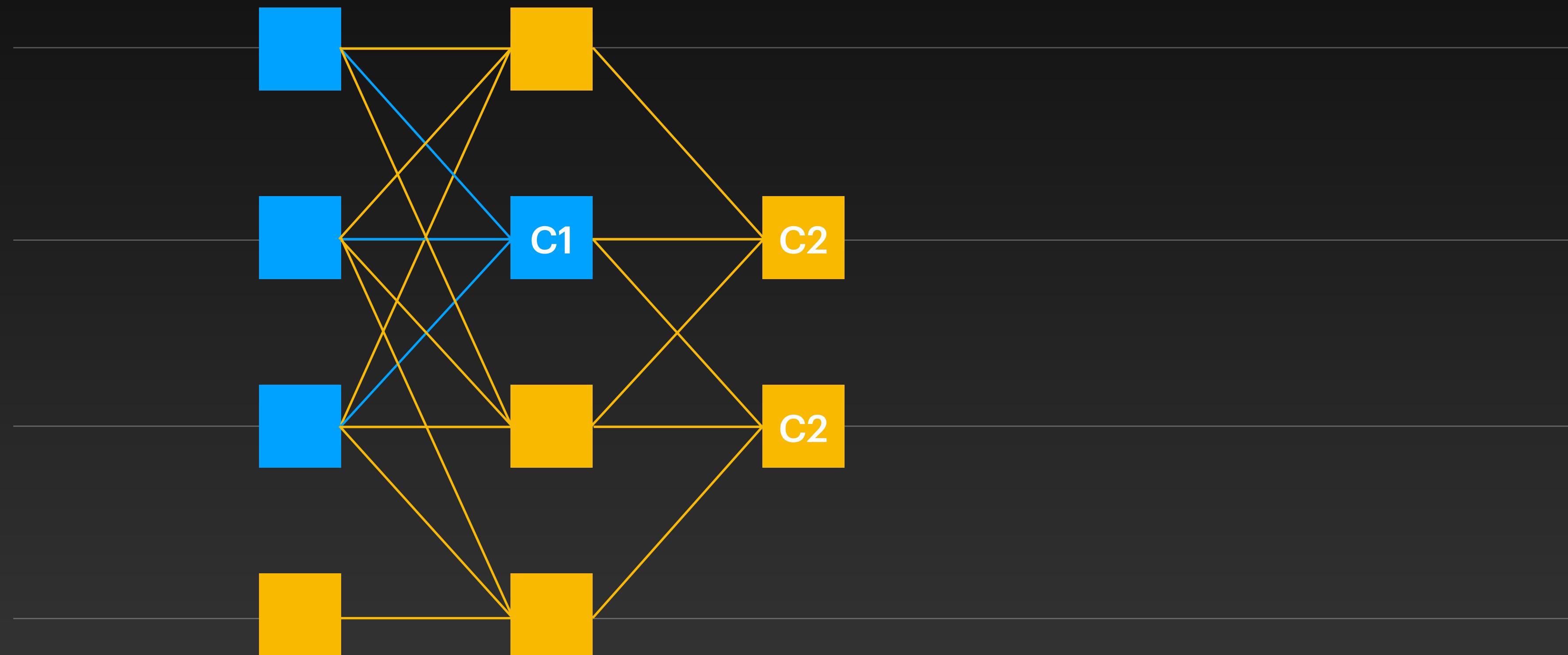
# HotStuff on Narwhal

## Enhanced commit rule



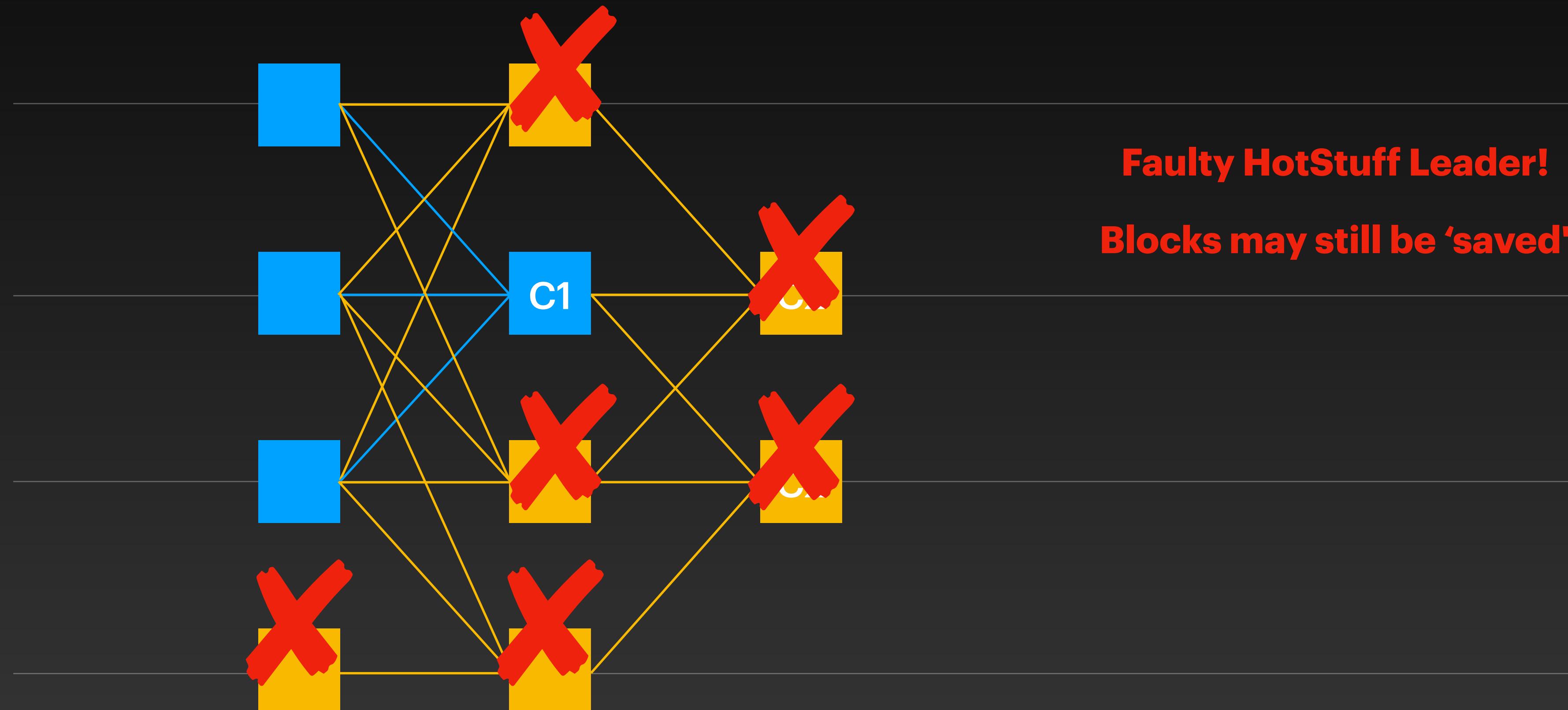
# HotStuff on Narwhal

## Enhanced commit rule



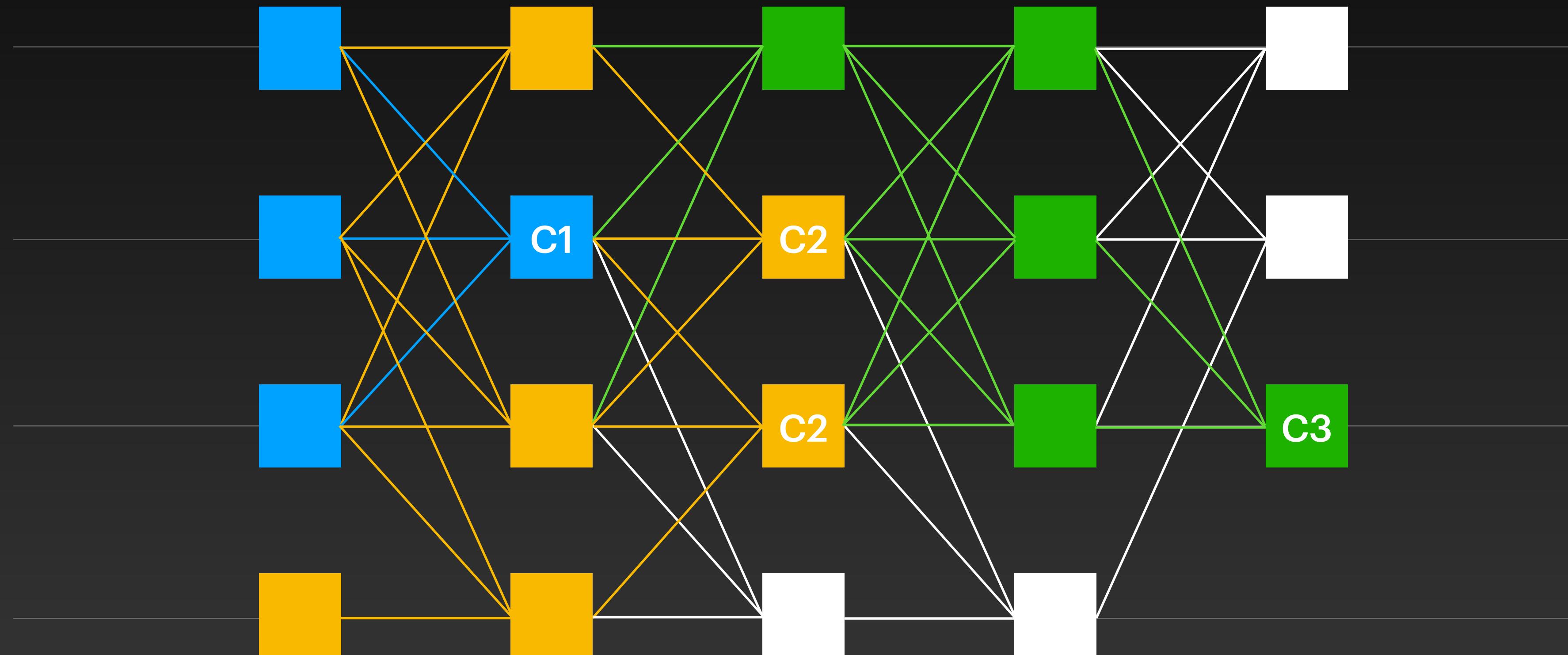
# HotStuff on Narwhal

## Enhanced commit rule



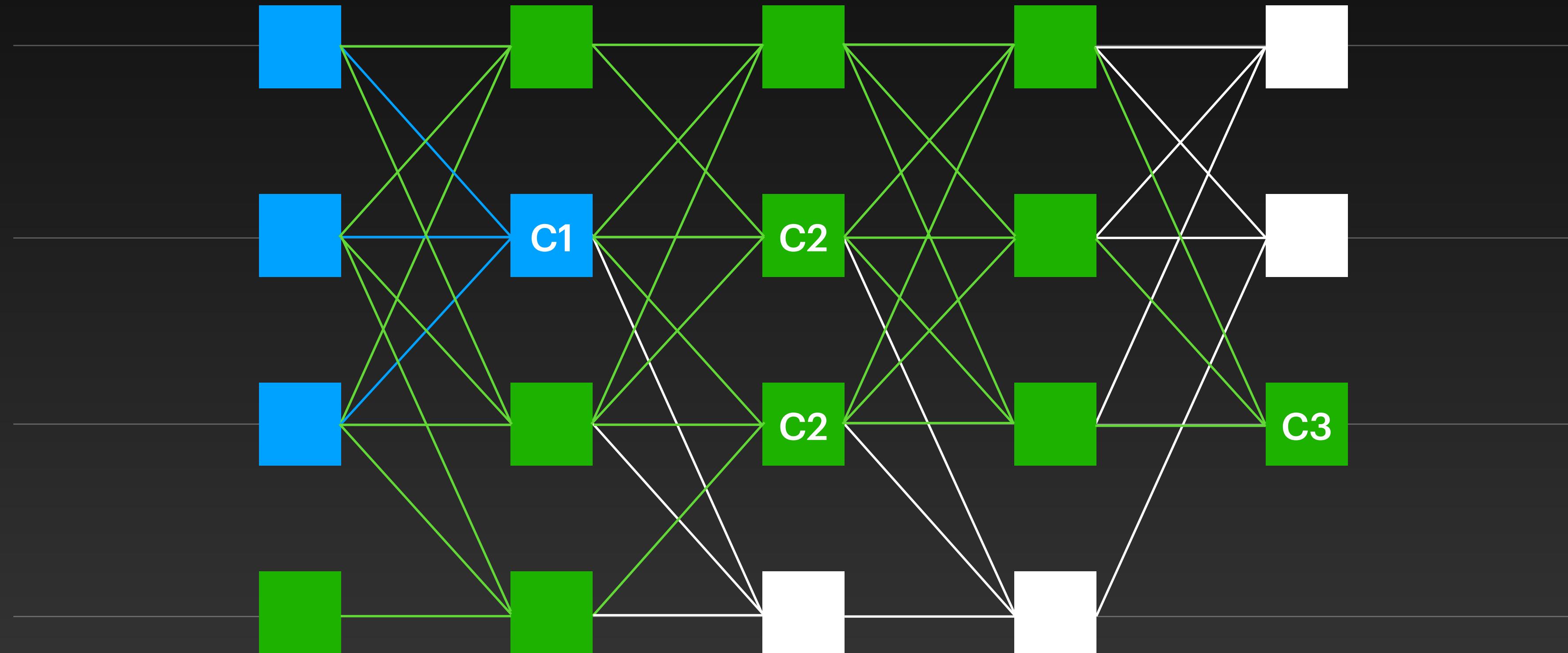
# HotStuff on Narwhal

## Enhanced commit rule

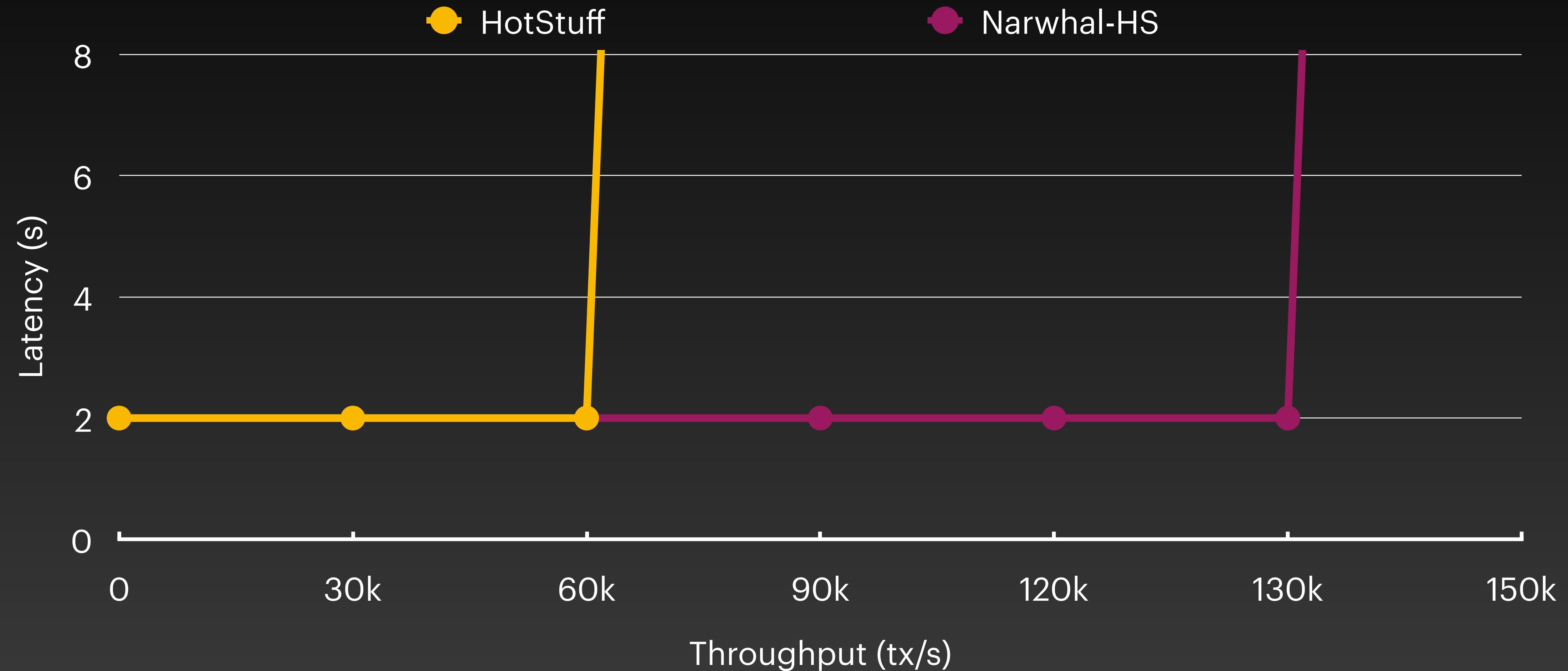


# HotStuff on Narwhal

## Enhanced commit rule



# Performance



# 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 relationship
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 first layer, we propose a novel protocol that does not rely on asymmetric cryptographic assumption. Therefore, when using a deterministic threshold-based coin implementation, their properties hold in a stronger sense. Directed Acyclic Graph (DAG) of the communication among them. In the second layer, we prove that DagRider preserves their DAGs and totally order all processes with no extra communication.

**ACM Reference Format:**  
Idit Keidar, Eleftherios Kokoris-Kogias, Oded Naor, and Alexander Spiegelman. 2021. DagRider: A Deterministic Byzantine Fault-Tolerant Consensus Protocol in Asynchronous Systems. In *Proceedings 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 trust in scalable geo-replicated Byzantine fault-tolerant consensus systems has motivated an enormous amount of study on the Byzantine Stateless Ledger (BSL) System [17, 20, 21] and the Byzantine Fault Tolerant (BFT) consensus [18, 22, 23]. In particular, the BSL system is simply a reliable broadcast. The agreement property of the reliable broadcast ensures that all correct processes eventually see the same value. Specifically, to implement BSL's [17] throughput of  $O(n^2)$  (twice as early as the best known protocols [23]) suggested consensus protocols. For higher throughput and lower latency committee-based consensus are required, and they are mainly based on the proof-of-stake mechanism. DagRider also guarantees that all proposals by correct processes are eventually included in the DAG. As a result, in contrast to the VARA and Dumbo protocols that return a single half of the consensus, DagRider returns a consensus value out of all proposals. DagRider does not waste any of the messages and all proposed values by correct processes are eventually included (i.e., there is no proposal rejection).

Byzantine 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 main reasons are usually cited to support this claim: (1) a single value is sent to all nodes, which costs 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 VARA and Dumbo consensus to independently agree on every slot. To this end, it applies in respect to our time complexity definition, we allow VARA and Dumbo based SMR to be up to a slot sequentially. Note, however, that the time complexity of the consensus is not necessarily in a sequential order (no gaps). Therefore, based on the proof in [6], the time complexity of VARA and Dumbo based SMR is  $O(\log(n))$ . This is in contrast to the time complexity of DagRider, which is  $O(1)$ .

Since our protocol uses a reliable broadcast abstraction as a basic building block, different instantiations yield different complexities. For example, if we use the classic BSL broadcast [11] to propose a

\*Oded Naor is grateful to the Technion Hsouki Fujisawa Cyber Security Research Center for providing a research grant. Part of Oded's work was done while at NUS Singapore.

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 to the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit 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](http://permissions@acm.org).  
© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-8548-2/21/07. \$15.00  
<https://doi.org/10.1145/3492211.3519594>

# 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 first layer, Narwhal is a mempool, Narwhal tolerates a high degree of unreliable dissemination and storage of causal histories of transactions. Narwhal tolerates an asynchronous network and maintains high performance despite multiple workers at each node. In the second layer, Tusk is the safety properties of our BFT protocol are post-quantum secure.

**ACM Reference Format:**  
George Danezis, Eleftherios 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 5–6, 2022, Rennes, France. ACM, New York, NY, USA, 17 pages. <https://doi.org/10.1145/3492211.3519594>

**1 Introduction**  
Composing Narwhal with a partially synchronous consensus protocol (e.g., PoS) yields significantly better performance than using a single consensus protocol. In fact, distributed systems have been studied for over 40 years. However, with the rise in popularity of blockchains there has been a renewed interest in understanding high-performance consensus protocols. Specifically, to implement Bitcoind's [13] throughput of  $O(n^2)$  (twice as early as the best known protocols [23]) suggested consensus protocols. For higher throughput and lower latency committee-based consensus are required, and they are mainly based on the proof-of-stake mechanism. DagRider also guarantees that all proposals by correct processes are eventually included in the DAG. As a result, in contrast to the VARA and Dumbo protocols that return a single half of the consensus, DagRider returns a consensus value out of all proposals. DagRider does not waste any of the messages and all proposed values by correct processes are eventually included (i.e., there is no proposal rejection).

As a summary of results, on a WAN, Narwhal-HotStuff achieves over 100K transaction per second (TPS), which is faster than VARD [1] and Speed-Dumbo [22] by a factor of 1,800 times 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 to the full citation on the first page. Copyright for components of this work owned by others than ACM must be honored. Abstracting with credit 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](http://permissions@acm.org).  
© 2021 Copyright held by the owner/author(s). Publication rights licensed to ACM. ISBN 978-1-4503-8548-2/21/07. \$15.00  
<https://doi.org/10.1145/3492211.3519594>

# 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 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 ignore the view of other parties in the DAG and update 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 standard BFT protocol.

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 round. Bullshark is a structured DAG, where each round is a threshold-based DAG and encodes a shared randomness in each round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor DAG-Rider, Bullshark is much more efficient. Specifically, it provides fairness and asynchronous liveness, and it is guaranteed to safely order them. The idea is simple. To propose a message, parties send them in a way that respects a causal order among them. They also keep track of the causal order of messages as well as references to previously received messages, which together form a *directed acyclic graph* (DAG). Interestingly, the structure of the DAG allows us to ignore the view of other parties in the DAG and update 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 standard BFT protocol.

**ACM Reference Format:**  
George Danezis, Eleftherios Kokoris-Kogias, Alberto Sonnino, and Alexander Spiegelman. 2022. Bullshark: DAG BFT Protocols Made Practical. In *Proceedings of ACM Conference, Los Angeles, CA, USA, November 2022 (CompoCNS '22)*, November 2022, Los Angeles, CA, USA, 17 pages. <https://doi.org/10.1145/3492211.3519594>

**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, the seminal impossibility proof [36] shows that no consensus mechanism can tolerate a single faulty node and still guarantee that the consensus value is reached in a sequential order (no gaps). Therefore, based on the proof in [6], the asymptotic message complexity is  $O(n^2)$  for consensus protocols designed for fixed mid-size committees (up to  $\sim 50$  nodes).

However, although DAG-based consensus protocols are deployable in practice, they all optimize for the worst case interpretation logic of DAG-Rider to totally order the DAG spans over less than 30 lines of pseudocode.

DAG-Rider [19] is a consensus protocol for the atomic broadcast (BAB), which achieves optimal amortized communication complexity ( $O(n)$  per transaction), post quantum safety, and some notion of fairness (asymptotic fairness). The main idea behind DAG-Rider is to delay an honest value until it is eventually delivered (ordered). To achieve optimal amortized communication DAG-Rider combines batching techniques with an efficient asynchronous verifiable consensus dispersed protocol [25] for the bottom of the DAG's building block. The protocol is post quantum safe because it does not rely on primitives that a quantum computer can break for the safety properties. That is, DAG-Rider is the first DAG-based consensus protocol that it cannot violate safety guarantees.

However, although DAG-based consensus protocols are deployable in practice, they all optimize for the worst case interpretation logic of DAG-Rider to totally order the DAG spans over less than 30 lines of pseudocode.

For example, the seminal impossibility proof [36] shows that no consensus mechanism can tolerate a single faulty node and still guarantee that the consensus value is reached in a sequential order (no gaps). Therefore, based on the proof in [6], the asymptotic message complexity is  $O(n^2)$  for consensus protocols designed for fixed mid-size committees (up to  $\sim 50$  nodes).

However, although DAG-based consensus protocols are deployable in practice, they all optimize for the worst case interpretation logic of DAG-Rider to totally order the DAG spans over less than 30 lines of pseudocode.

For example, the seminal impossibility proof [36] shows that no consensus mechanism can tolerate a single faulty node and still guarantee that the consensus value is reached in a sequential order (no gaps). Therefore, based on the proof in [6], the asymptotic message complexity is  $O(n^2)$  for consensus protocols designed for fixed mid-size committees (up to  $\sim 50$  nodes).

# Dumbo-NG

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

Yingqi Guo\*  
ISCA & UCAS

Yuan Liu\*  
ISCA

Zherui Lu\*  
USID

Qiang Tang\*  
UCAS

Jing Xu\*  
ISCA

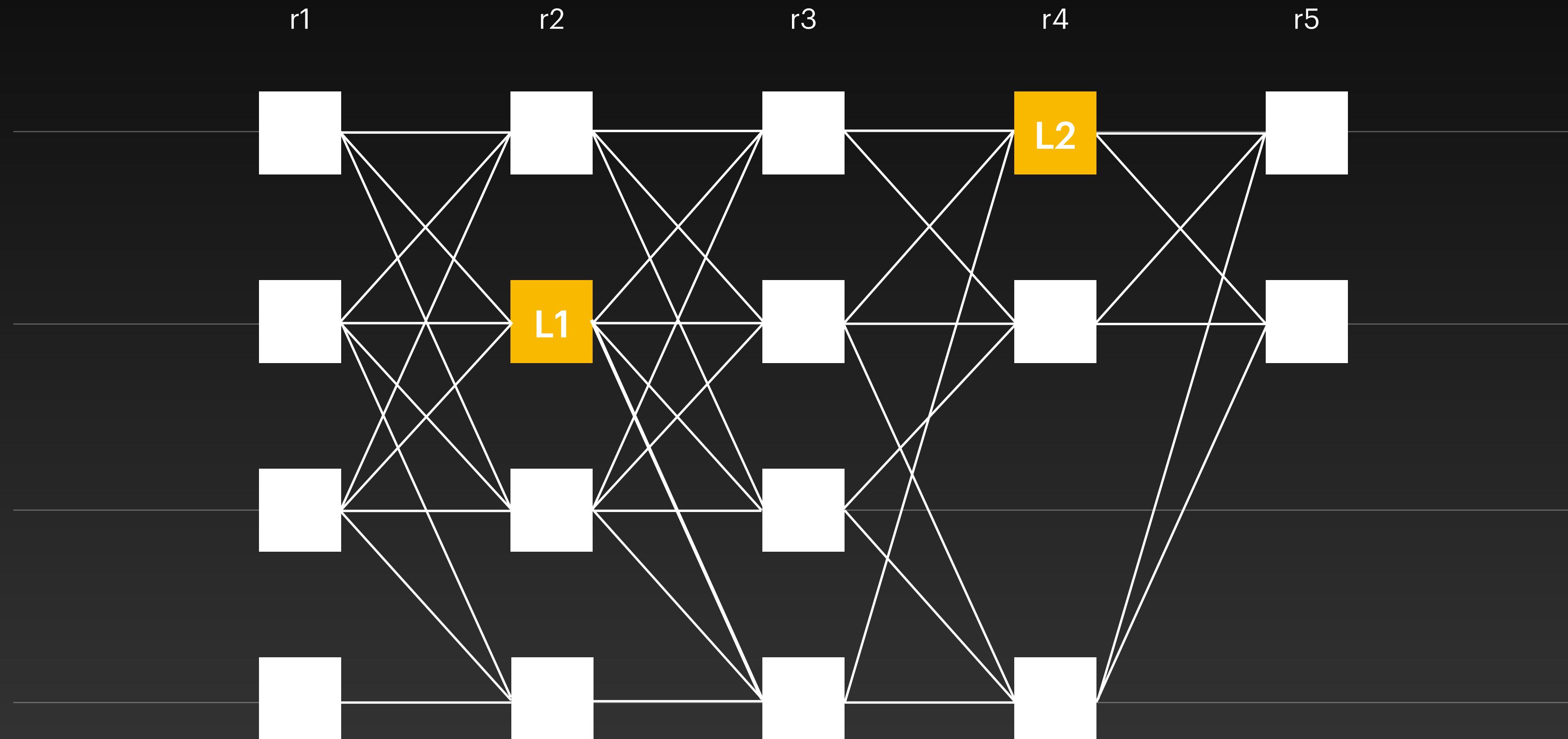
Zhenfeng Zhang\*  
ISCA

**Abstract**  
We present Dumbo-NG, the first asynchronous Byzantine-fault tolerant (BFT) consensus, the state-of-the-art designs still suffer from suboptimal performance. Particularly, to obtain maximum throughput, it is required to pay off the latency overhead. Specifically, to implement BFT consensus with throughput oblivious latency, it is required to pay off the latency overhead. This is, the latency overhead is proportional to the number of faults. Thus, it is difficult to achieve maximum throughput with minimum latency.

We present Dumbo-NG, a novel asynchronous BFT consensus (throughput oblivious latency) that is designed for the Internet. The technical core is a non-trivial direct delivery from asynchronous broadcast to multi-valued validated Byzantine agreement (MVBA) through a series of rounds. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that while maintaining the desired properties of its predecessor Dumbo, Dumbo-NG only needs to pay off the latency overhead. Thus, this conquering 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 round via a threshold signature scheme to achieve constant latency in expectation. 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 consensus blocks. Blocks are disseminated via reliable broadcast [11] to avoid equivocation and an honest party advances to the next round once it reliably sees  $\lceil \frac{n}{2} + 1 \rceil$  blocks from the previous round. Note that

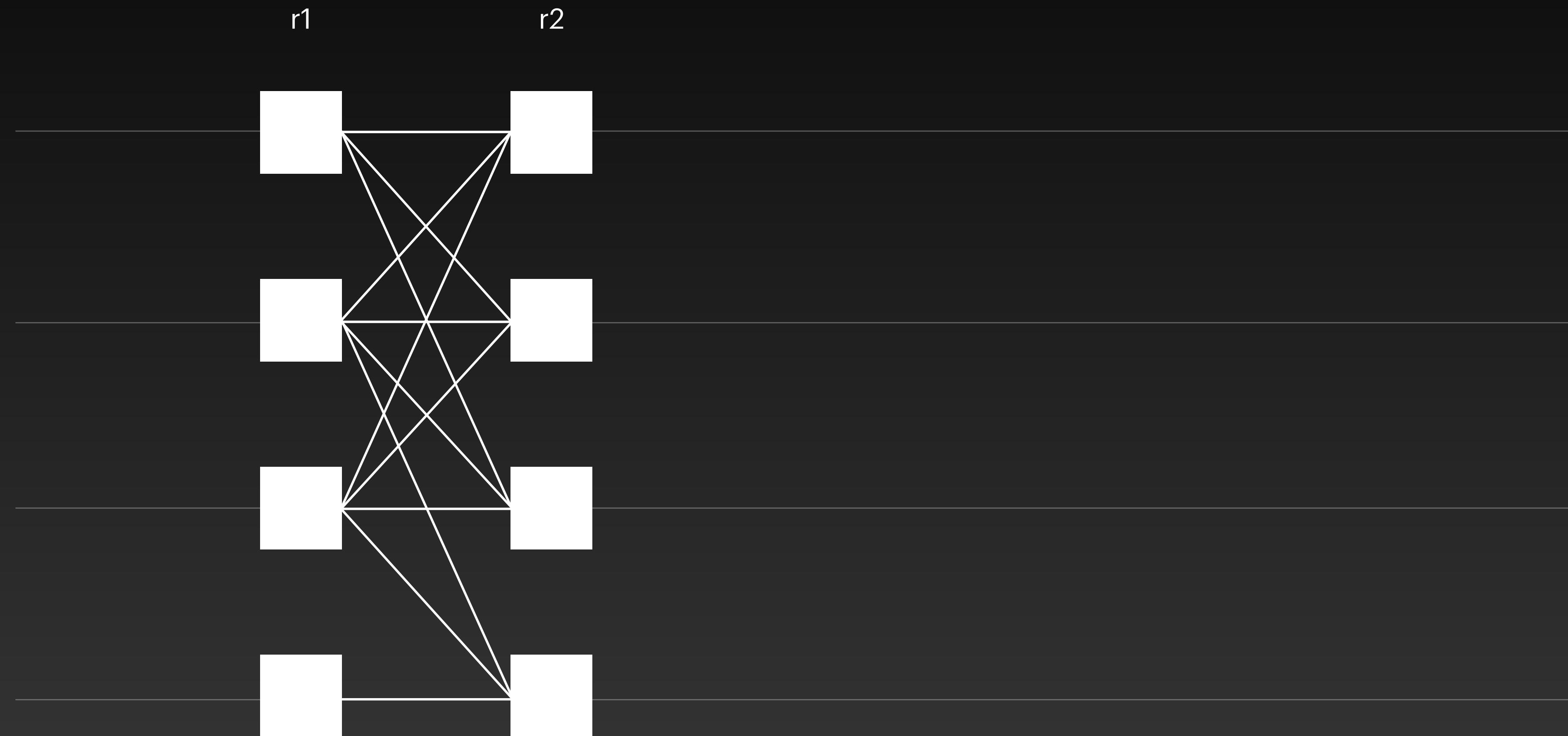
# Modified Narwhal

Even rounds: wait for the leader (or to timeout)



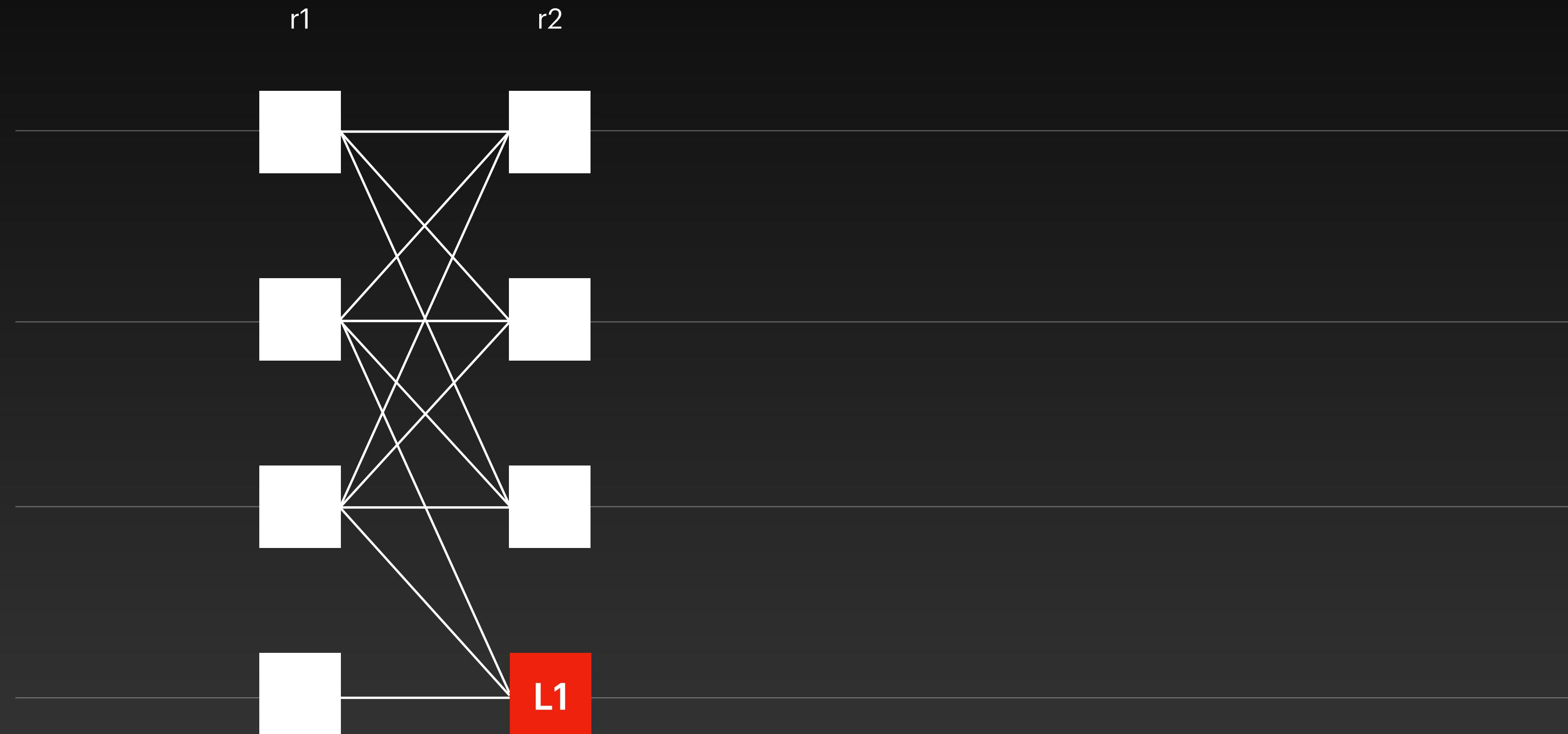
# Bullshark

## Just interpret the DAG



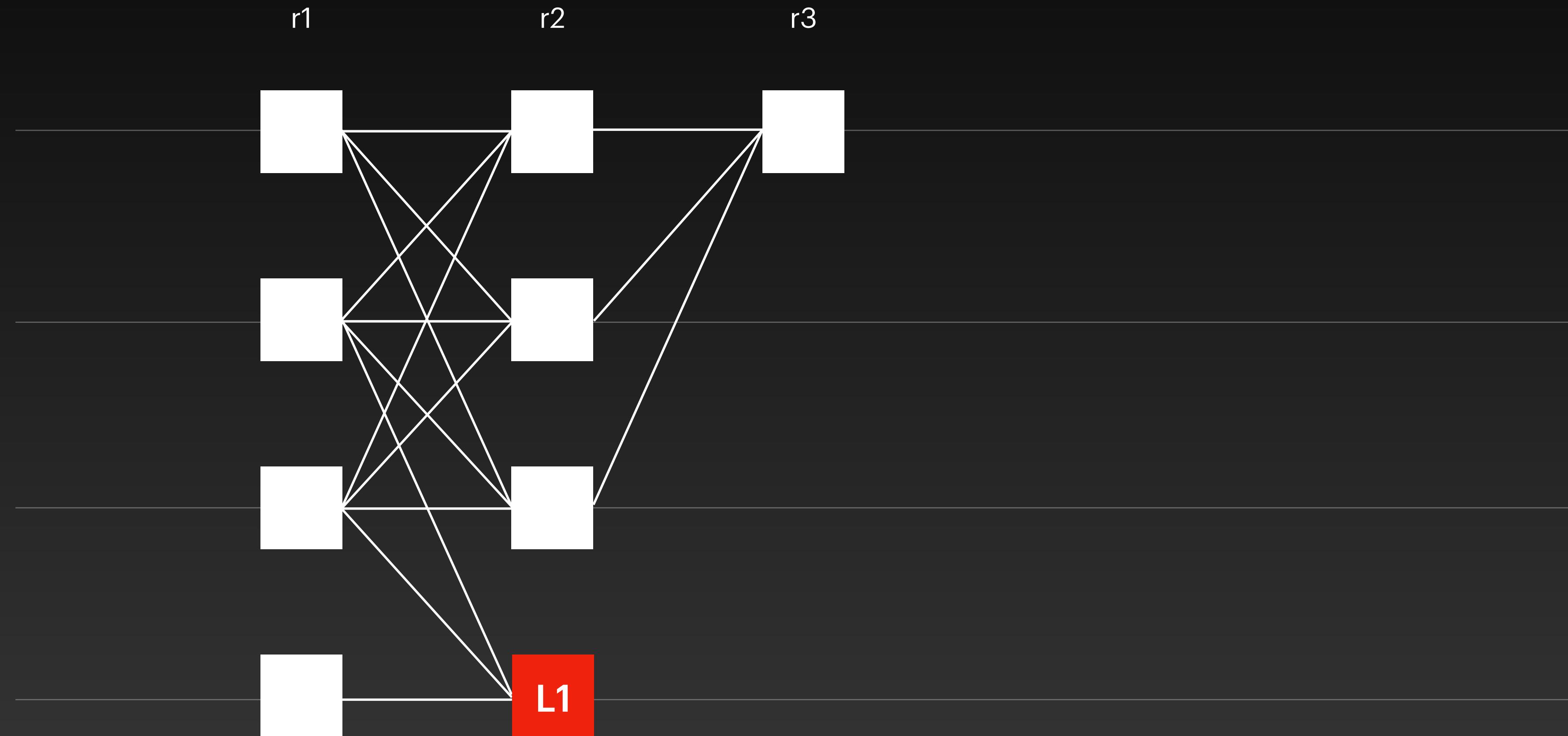
# Bullshark

## Deterministic leader every 2 rounds



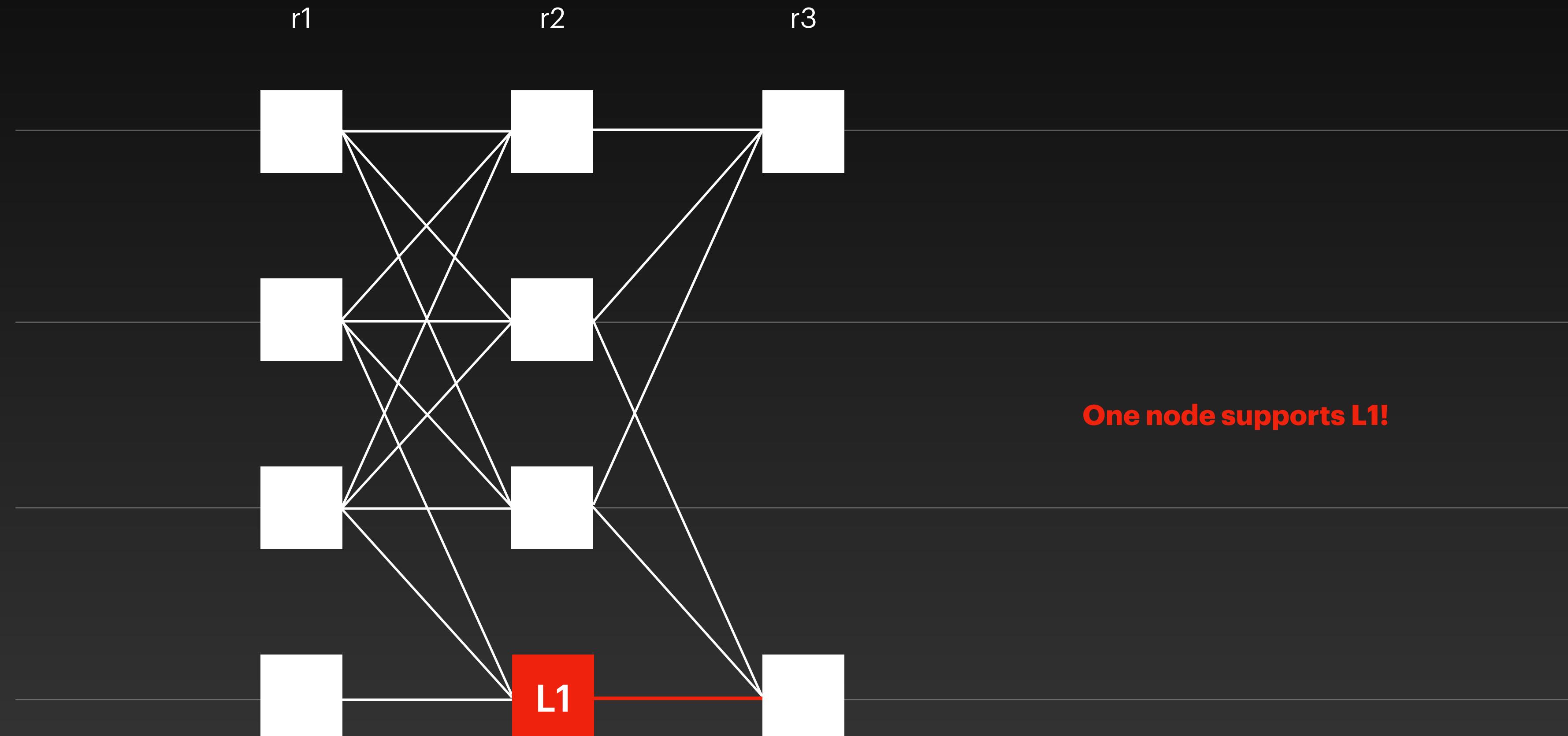
# Bullshark

## The leader needs $f+1$ links from round $r$



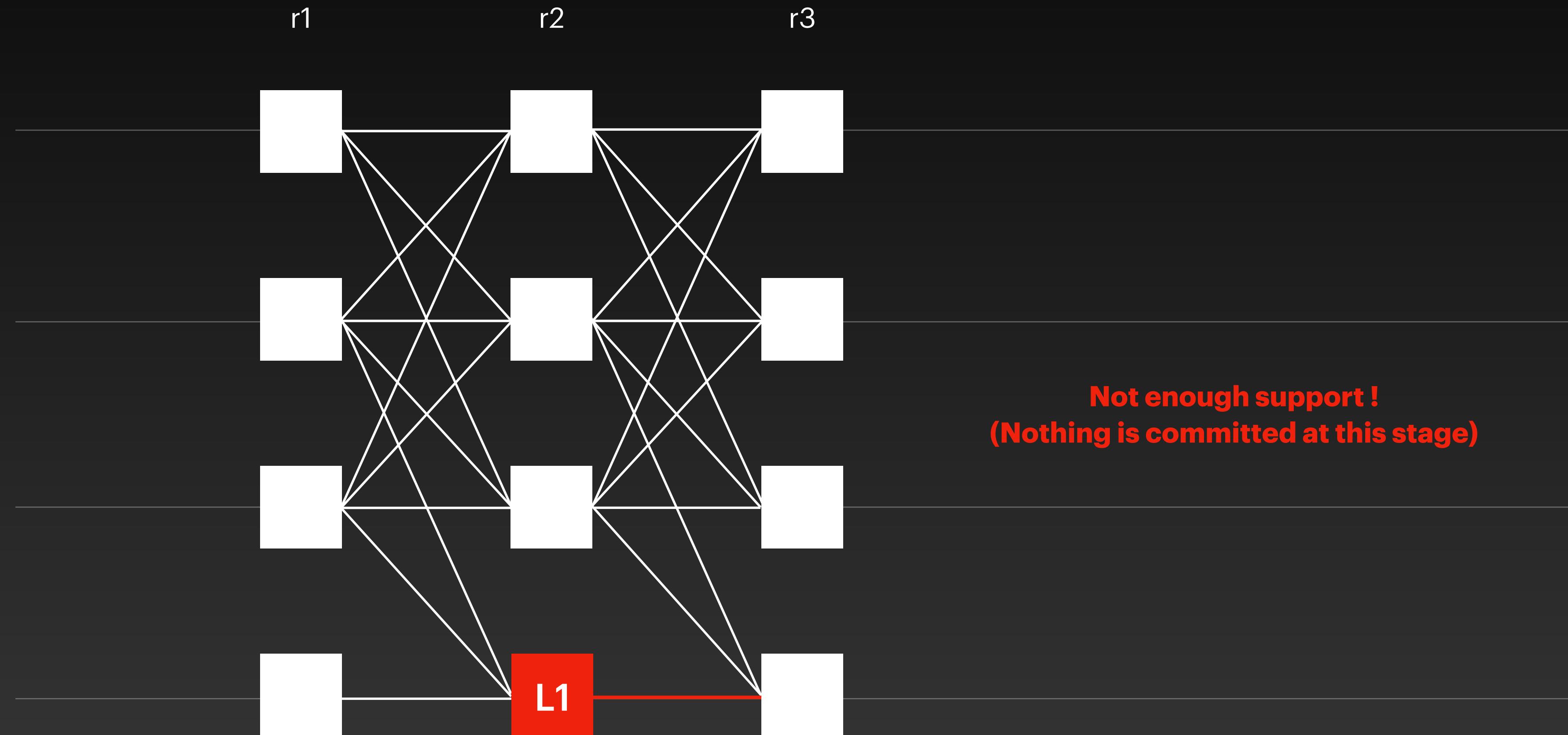
# Bullshark

## The leader needs $f+1$ links from round $r$



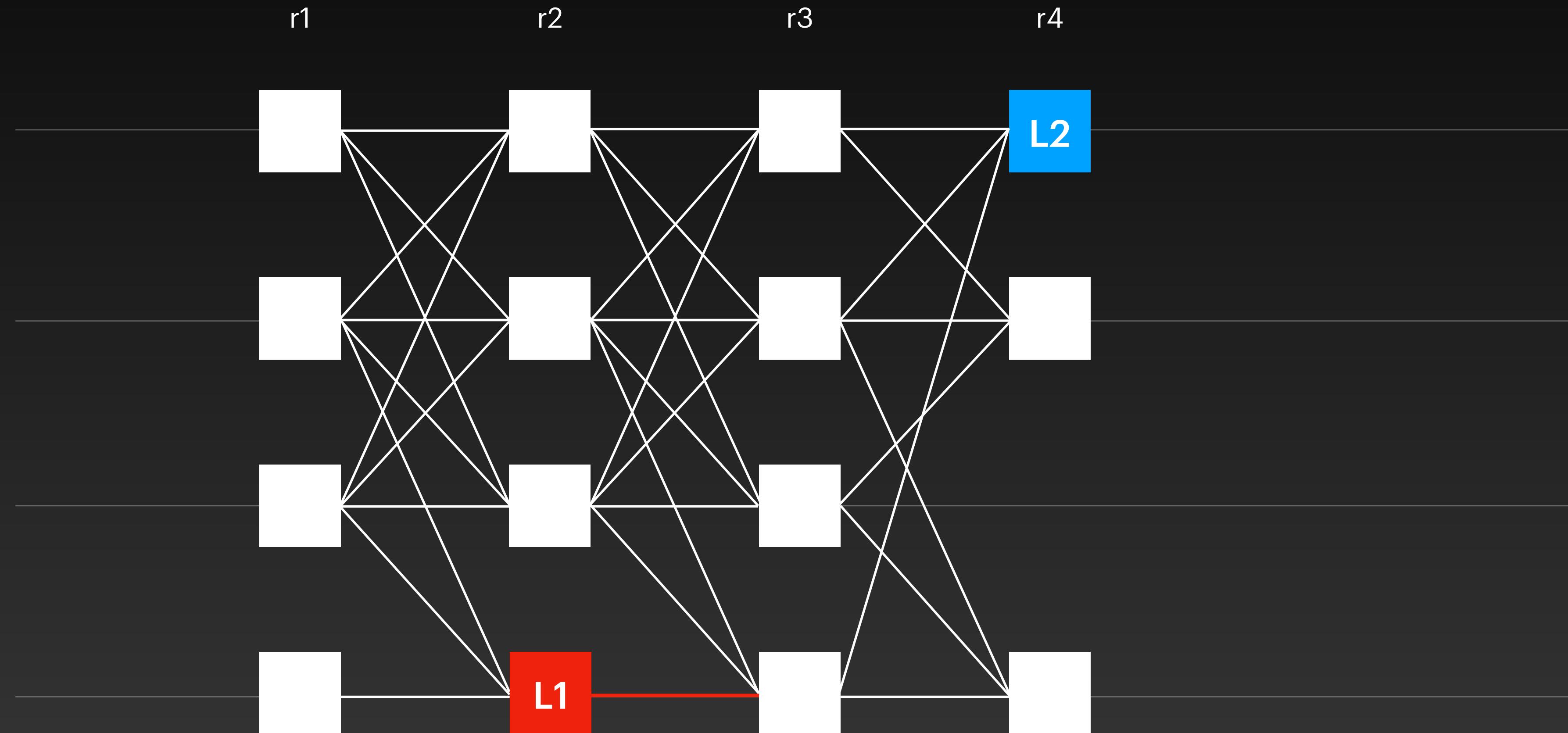
# Bullshark

## The leader needs $f+1$ links from round $r$



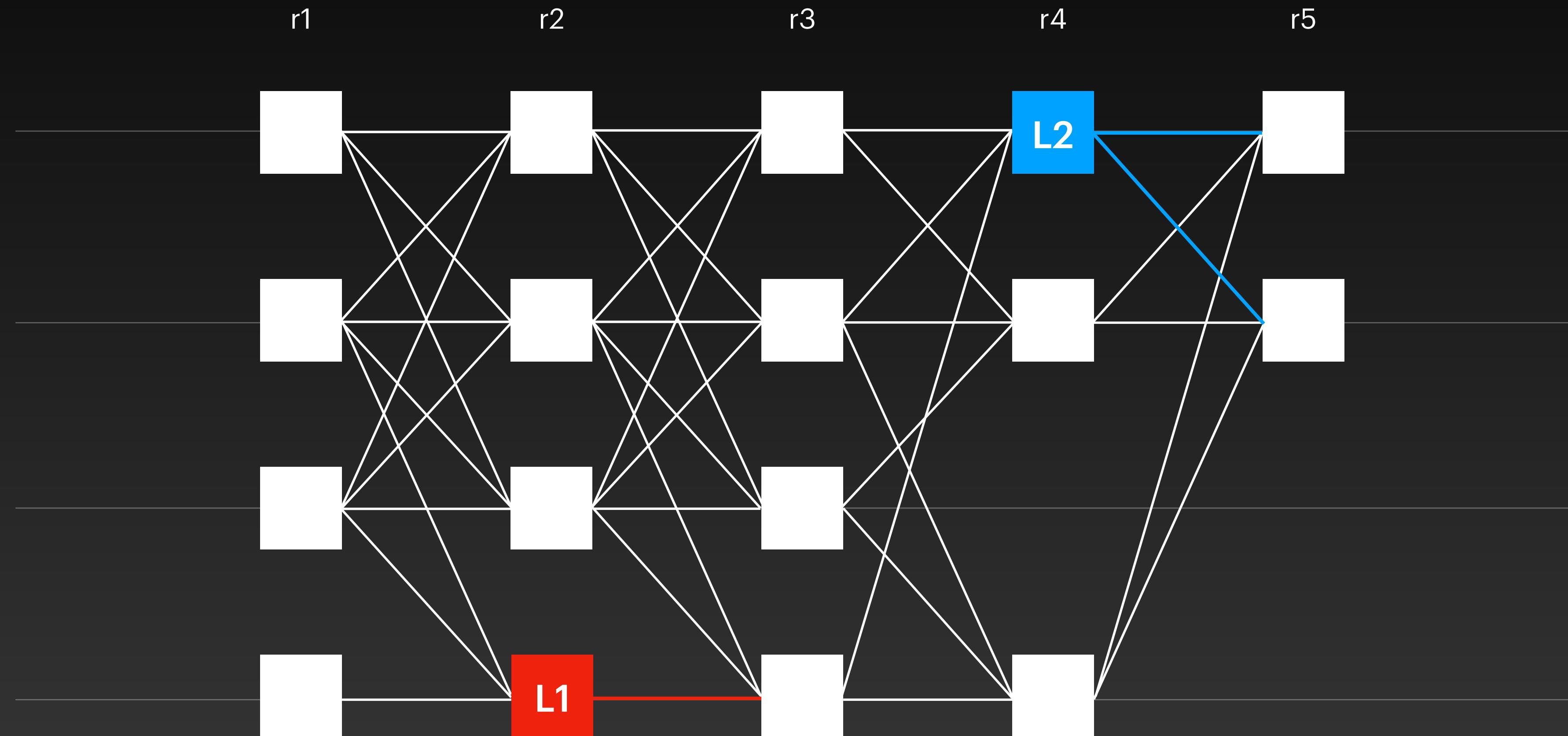
# Bullshark

## Elect the leader of r4



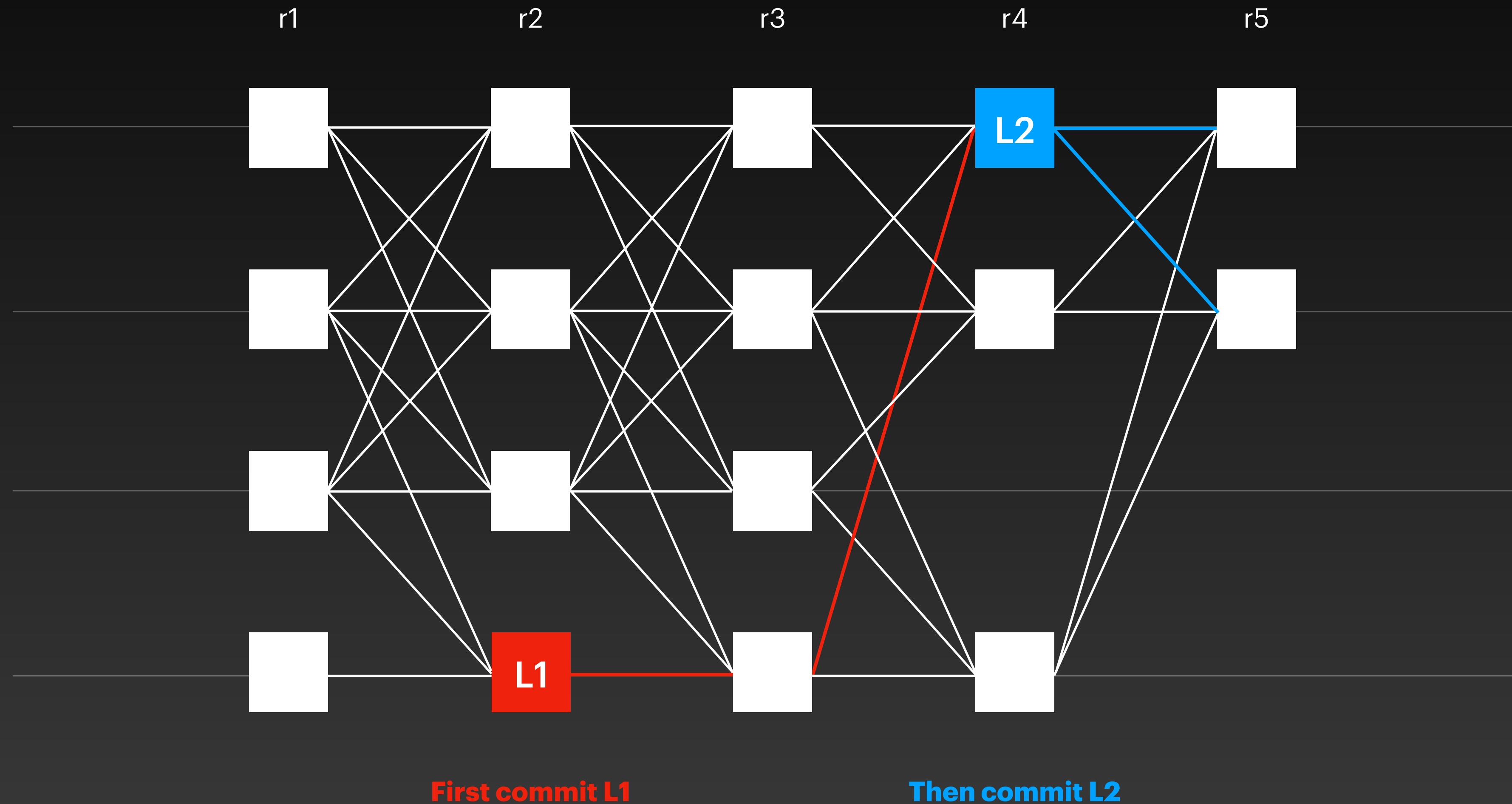
# Bullshark

## Leader L2 has enough support



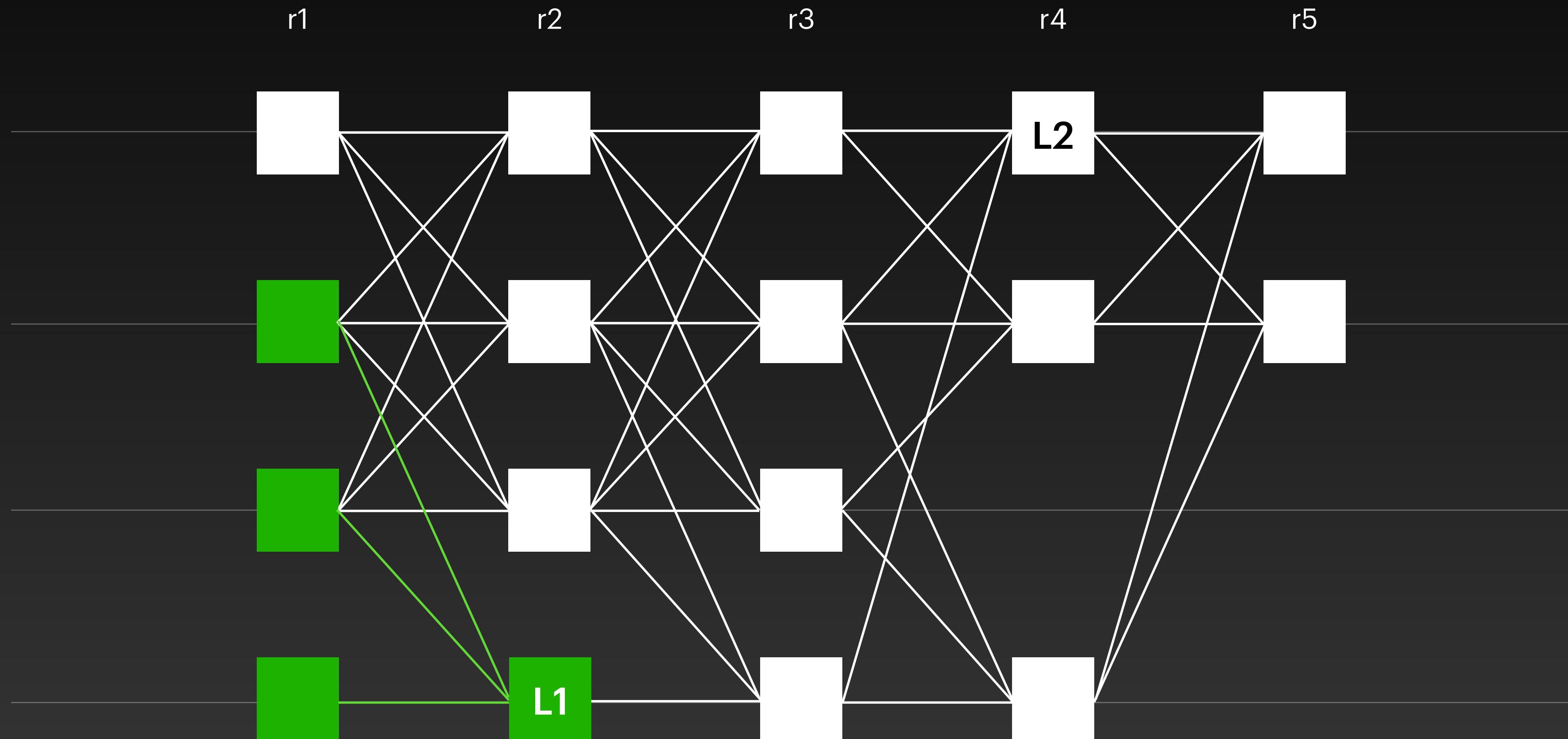
# Bullshark

## Leader L2 has links to leader L1



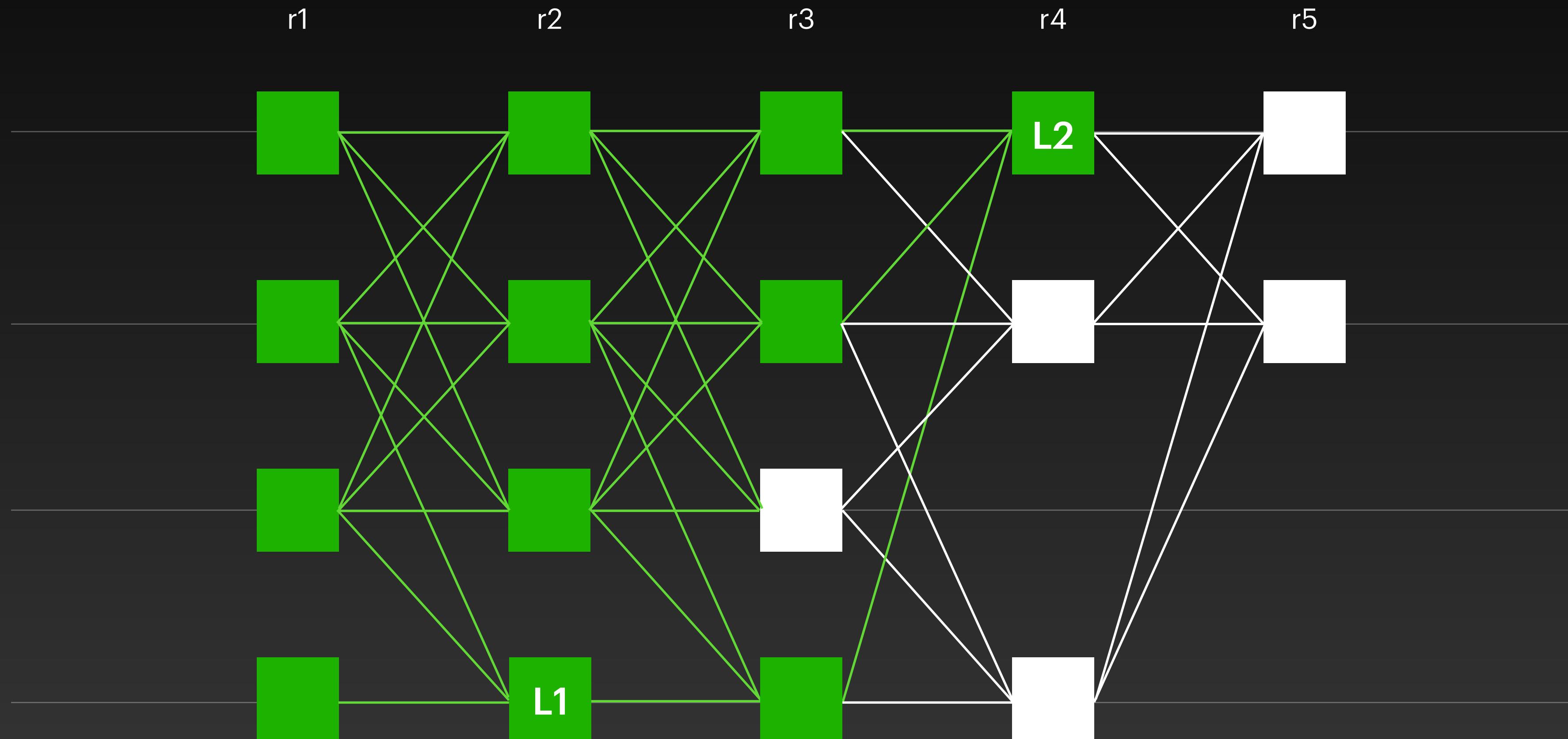
# Bullshark

## Commit all the sub-DAG of the leader

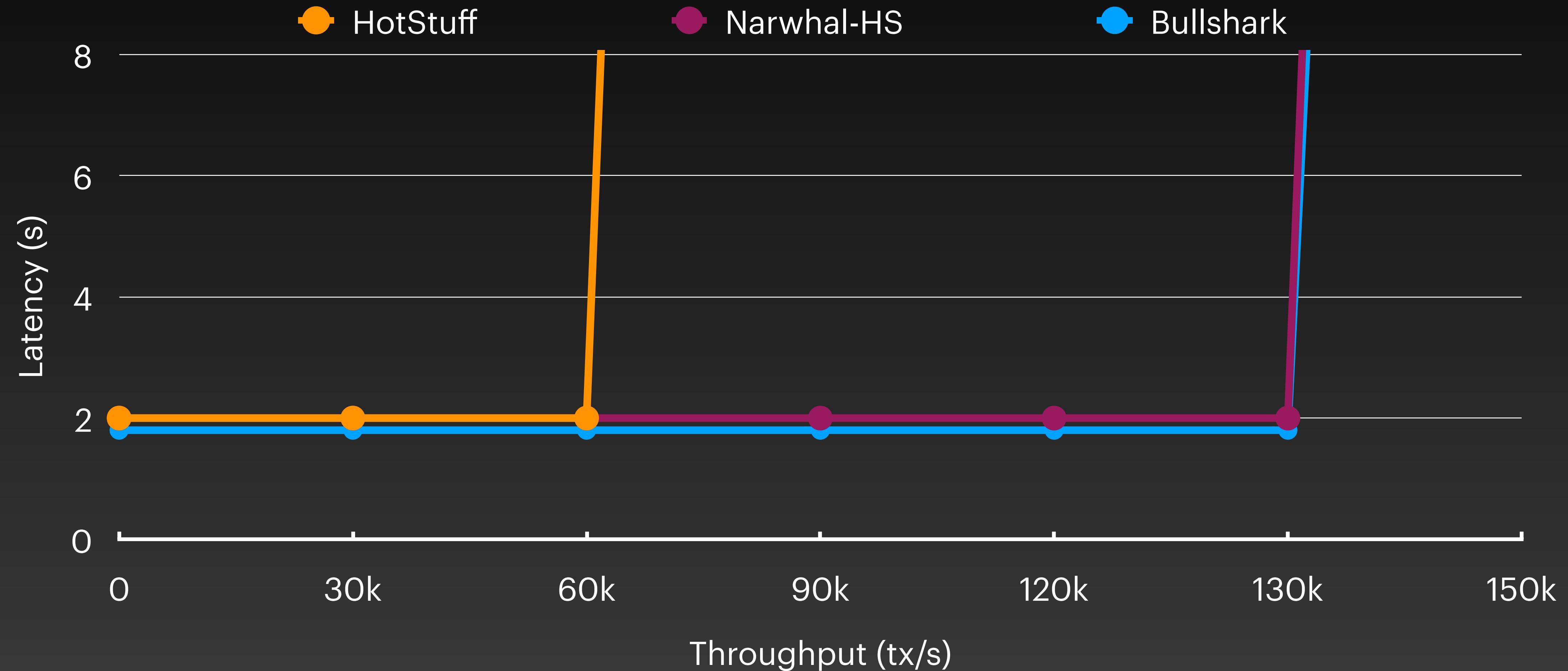


# Bullshark

## Commit all the sub-DAG of the leader



# Performance



# Research Questions

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

# 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

# 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 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

# 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 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 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.

# Sailfish

**Sailfish: Towards Improving the Latency of DAG-based BFT**

Nibesh Shrestha  
*n.shrestha@supraoracle.com*  
Supra Research  
Balaji Arun  
Aptos  
Zekun Li  
Aptos

**Abstract**

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. Instead, the current blockchain network targets. For example, the state-of-the-art Hotstuff [46] protocol in this line of work 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.

A recent breakthrough, however, stemmed from the realization that data dissemination is the primary bottleneck for leader-based protocols, and it can be mitigated by utilizing a novel approach. Specifically, the proposed protocol, called Sailfish, separates data dissemination from the core consensus logic and proposed an architecture where all validators simultaneously process data, while only one validator performs a smaller number of leader-related tasks. A notable advantage of this architecture is that not only it delivers impressive throughput on a single machine, but also naturally supports scaling out each blockchain validator for scaling more nodes. The paper proposes to evaluate this idea in a geo-replicated environment 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 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.

# 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

**Abstract**

Cordial Miners are a family of efficient Byzantine Atomic Broadcast protocols, with instances achieving the lower bounds of latency. Our main contribution is that Cordial Miners also achieve high resource efficiency and censorship resistance. MYSTICETI-C achieves its latency improvement by avoiding the need for a consensus DAG, and instead uses a novel commit rule such that every block can be committed without delay, resulting in a fast commit path in both normal and under crash failures. We further extend MYSTICETI-C to MYSTICETI-FPC, which incorporates a fast commit path that allows for a fast commit even if a validator fails. Similar to prior fast commit path protocols, MYSTICETI-FPC minimizes the number of signatures and messages by weaving the fast path into the leader vertex. We prove that the safety of Cordial Miners is robust to the leader vertex failing, and subsequently result in better performance. We prove the safety of Cordial Miners in a Byzantine setting. We evaluate both MYSTICETI and Cordial Miners with their own consensus and fast path protocols to demonstrate their low latency and resource efficiency. We also extend Cordial Miners to handle crash failures. MYSTICETI-C is the first Byzantine consensus protocol to achieve WAN latency of 4.5 s for consensus commit while simultaneously maintaining a throughput of 100,000 TPS. Finally, we report on integrating MYSTICETI-C as the consensus layer in the Sul blockchain [67], resulting in over 4x latency reduction.

**1 Introduction**

The problem of ordering transactions in a permissioned Byzantine distributed system, also known as *Byzantine Atomic Broadcast (BAB)*, has been investigated for four decades [39], and in the last decade, has attracted renewed attention due to the emergence of cryptocurrencies. At a high level, a BAB protocol enables a group of  $n$  parties to agree on a sequence of values, even if a bound of  $f$  to these parties are Byzantine (arbitrarily malicious). Overall, the design of a consensus protocol has been a lot of progress in improving the key efficiency metrics, namely, latency, communication complexity, and throughput. While several consensus protocols have been proposed 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), Ethereum (the state-of-the-art blockchain) [30, 31, 43, 44] and now engaged in a race to deliver a scalable BFT system with the utmost throughput and minimal latency.

# 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†

**Abstract**

We introduce MYSTICETI-C, the first DAG-based Byzantine consensus protocol to achieve the lower bounds of latency. Our main contribution is that MYSTICETI-C also achieves high resource efficiency and censorship resistance. MYSTICETI-C achieves its latency improvement by avoiding the need for a consensus DAG, and instead uses a novel commit rule such that every block can be committed without delay, resulting in a fast commit path in both normal and under crash failures. We further extend MYSTICETI-C to MYSTICETI-FPC, which incorporates a fast commit path that allows for a fast commit even if a validator fails. Similar to prior fast commit path protocols, MYSTICETI-FPC minimizes the number of signatures and messages by weaving the fast path into the leader vertex. We prove that the safety of Cordial Miners is robust to the leader vertex failing, and subsequently result in better performance. We prove the safety of Cordial Miners in a Byzantine setting. We evaluate both MYSTICETI and Cordial Miners with their own consensus and fast path protocols to demonstrate their low latency and resource efficiency. We also extend Cordial Miners to handle crash failures. MYSTICETI-C is the first Byzantine consensus protocol to achieve WAN latency of 4.5 s for consensus commit while simultaneously maintaining a throughput of 100,000 TPS. Finally, we report on integrating MYSTICETI-C as the consensus layer in the Sul blockchain [67], resulting in over 4x latency reduction.

**1 Introduction**

Several recent blockchains, such as Sul [67], [12], have adopted a consensus mechanism based on a distributed acyclic graph (DAG) of blocks [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. In addition, they support 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 a family of uncertified DAG-based consensus protocols. In each round, a potential DAG vertex is signed by every party that can create a potential DAG vertex containing transactions that do not yet appear in previous rounds. These protocols rely on committing a designated "leader vertex" and order other non-leader vertices according to the DAG structure. The DAG is distributively constructed from messages of miners running the consensus protocol. While building the DAG, the miners must ensure that the order in which leaders are designated and how fast the leader vertices emerge directly influences the commit latency.

The two state-of-the-art protocols in this context are DAG-Rider [21] and Bullshark [33]. DAG-Rider works in the asynchronous setting, in which the adversary controls the finite delay on message delivery between miners, and Bullshark works in the Eventual Synchrony (ES) model, in which eventually all messages between correct miners are delivered within a known time-bound.

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-C is the first DAG-based consensus protocol based on a threshold logical clock [29] of  $2f+1$  blocks. 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 ever round
- Uncertified DAG

# Discussion

**Certified DAG**

**Uncertified DAG**



**Shoal/shoal++**

- Low latency
- Easier synchroniser
- Leverage existing code

**Sailfish/BBCA**

- Lower latency
- Easier synchroniser
- Flexible

**CM/Mysticeti**

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

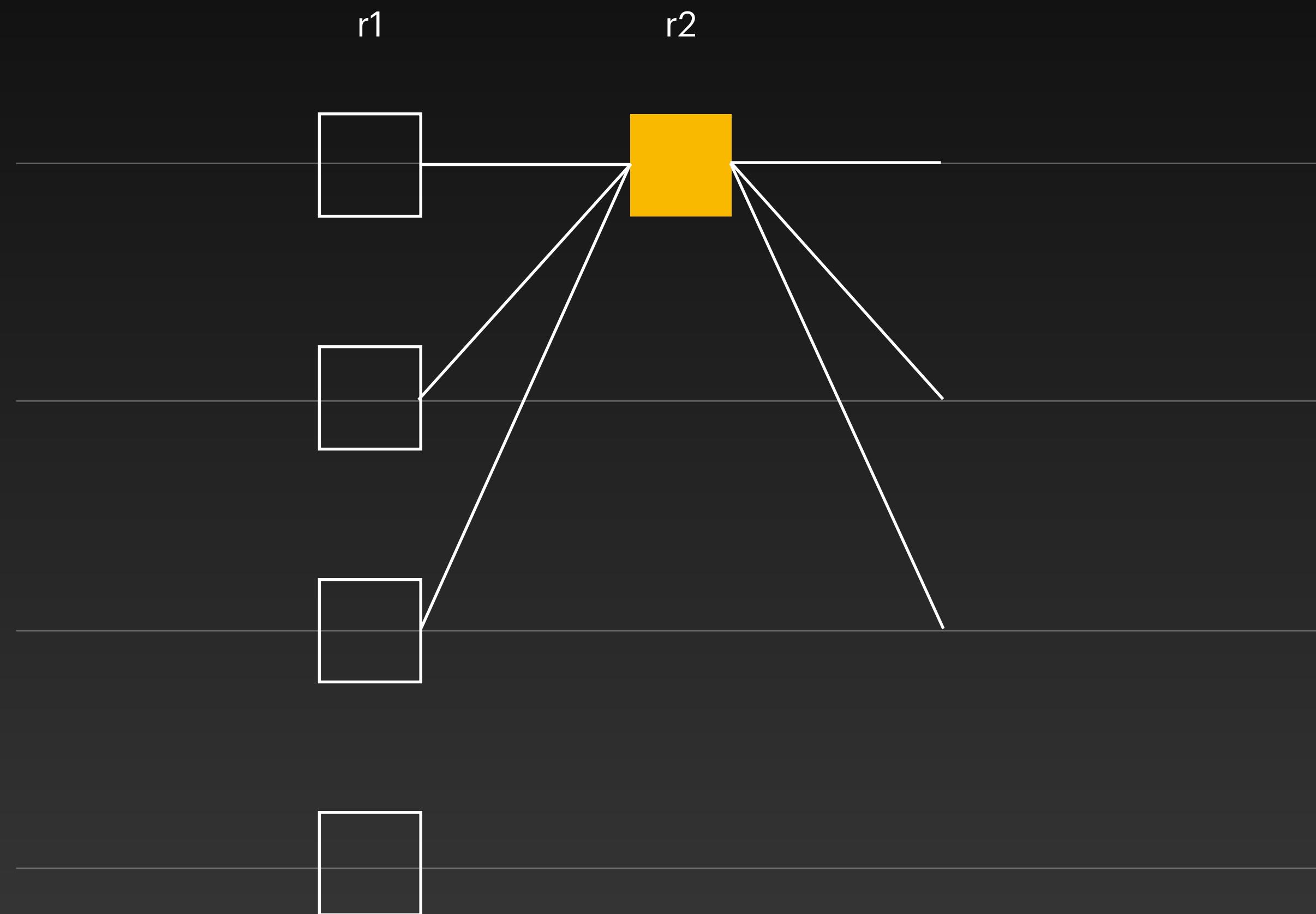
# Research Questions

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

# Lessons Learned

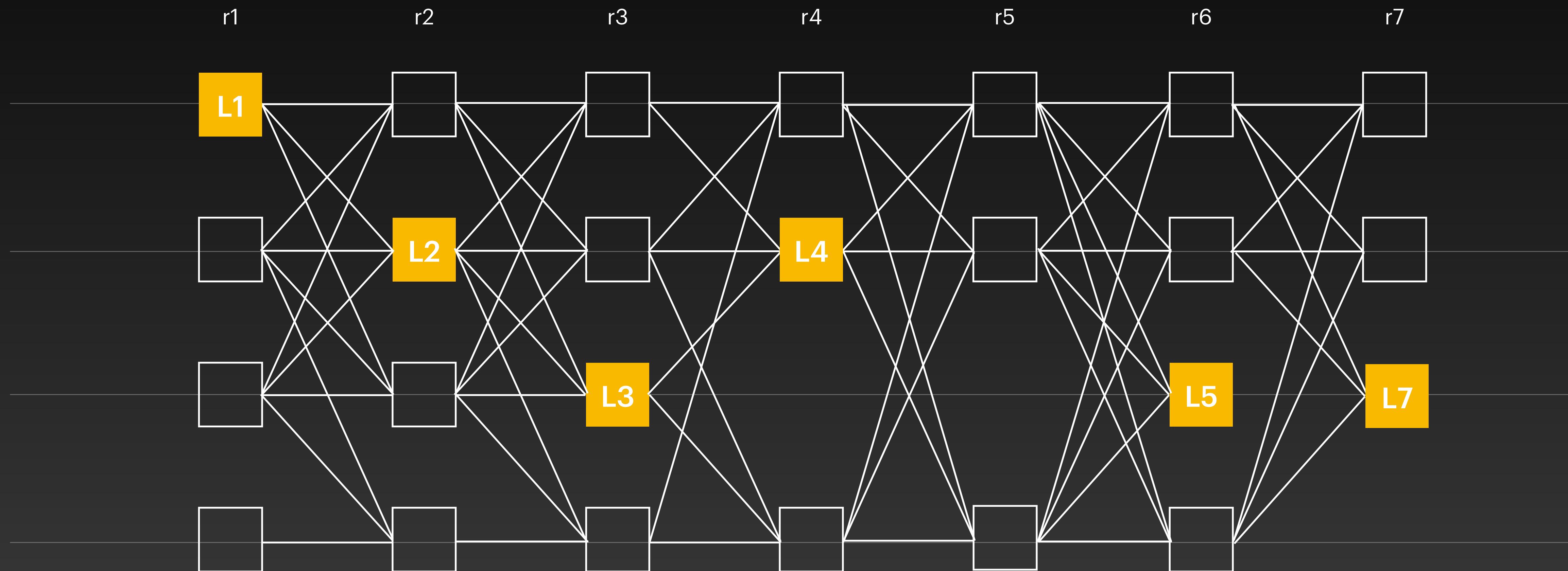
1. Modularisation is a design strategy
2. Tasks-threads relationship
3. Benchmark early
4. Codesign with mem. and storage
5. Core is hard, consensus is easy
6. Epoch change is not an add-on

# Uncertified DAG

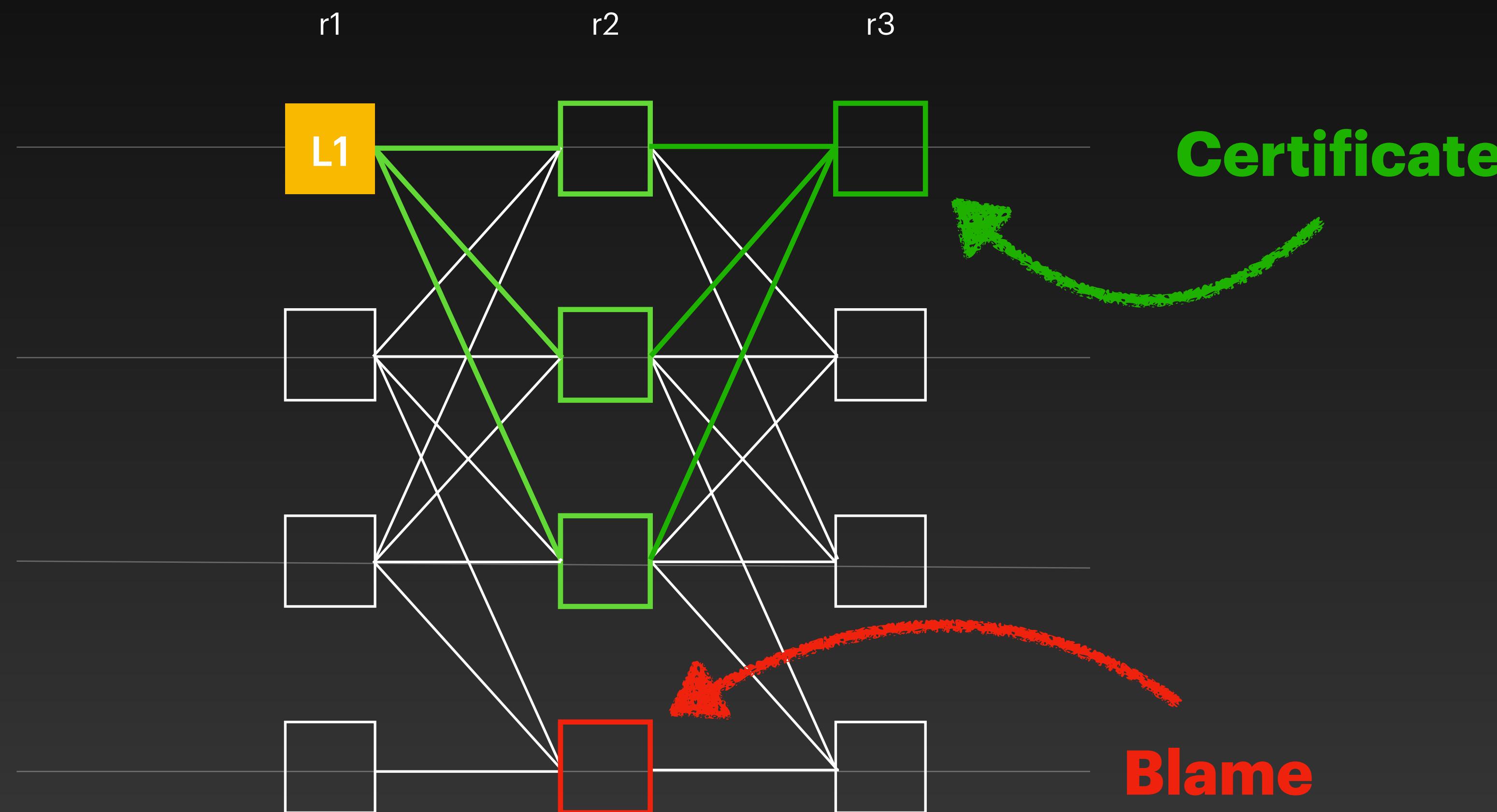


- Round number
- Author
- Payload (transactions)
- Signature

# Uncertified DAG



# Interpreting DAG Patterns



# Direct Decision Rule

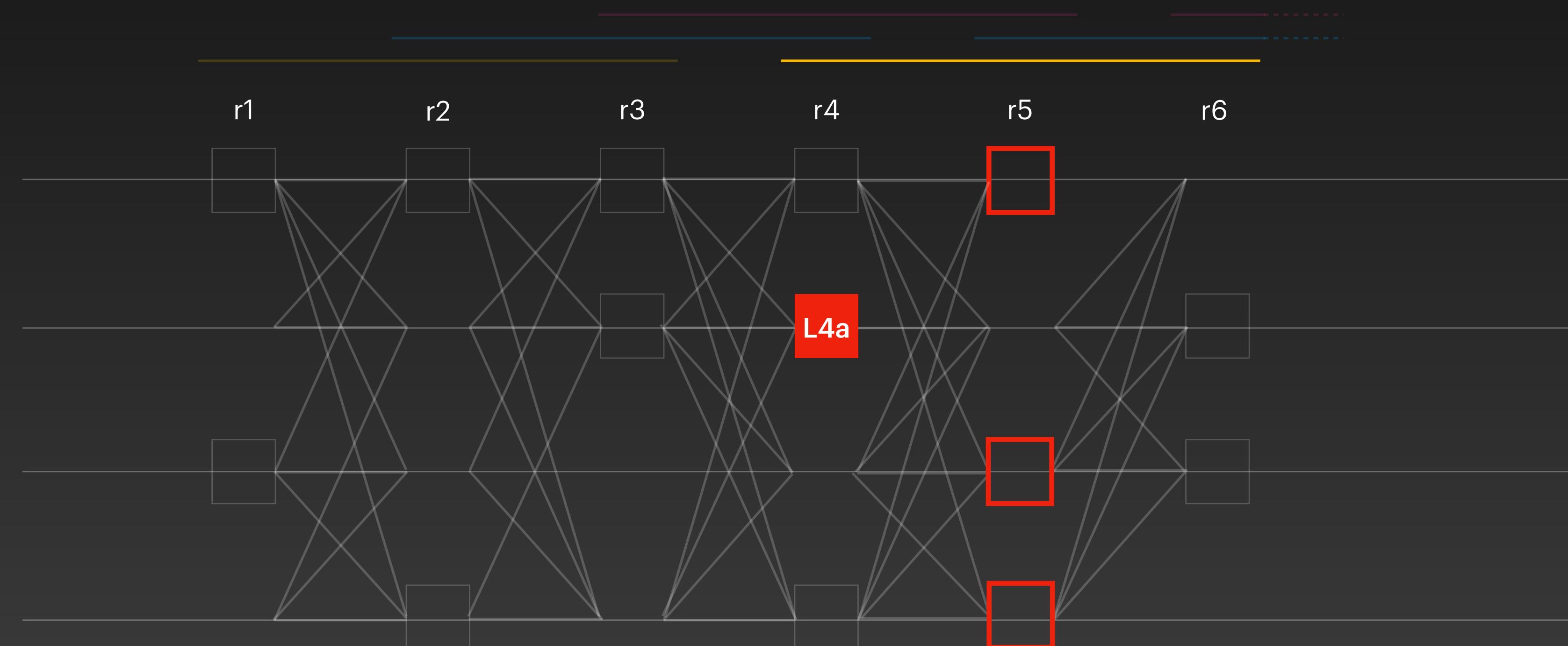
On each leader starting from highest round:

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

# Direct Decision Rule

On each leader starting from highest round:

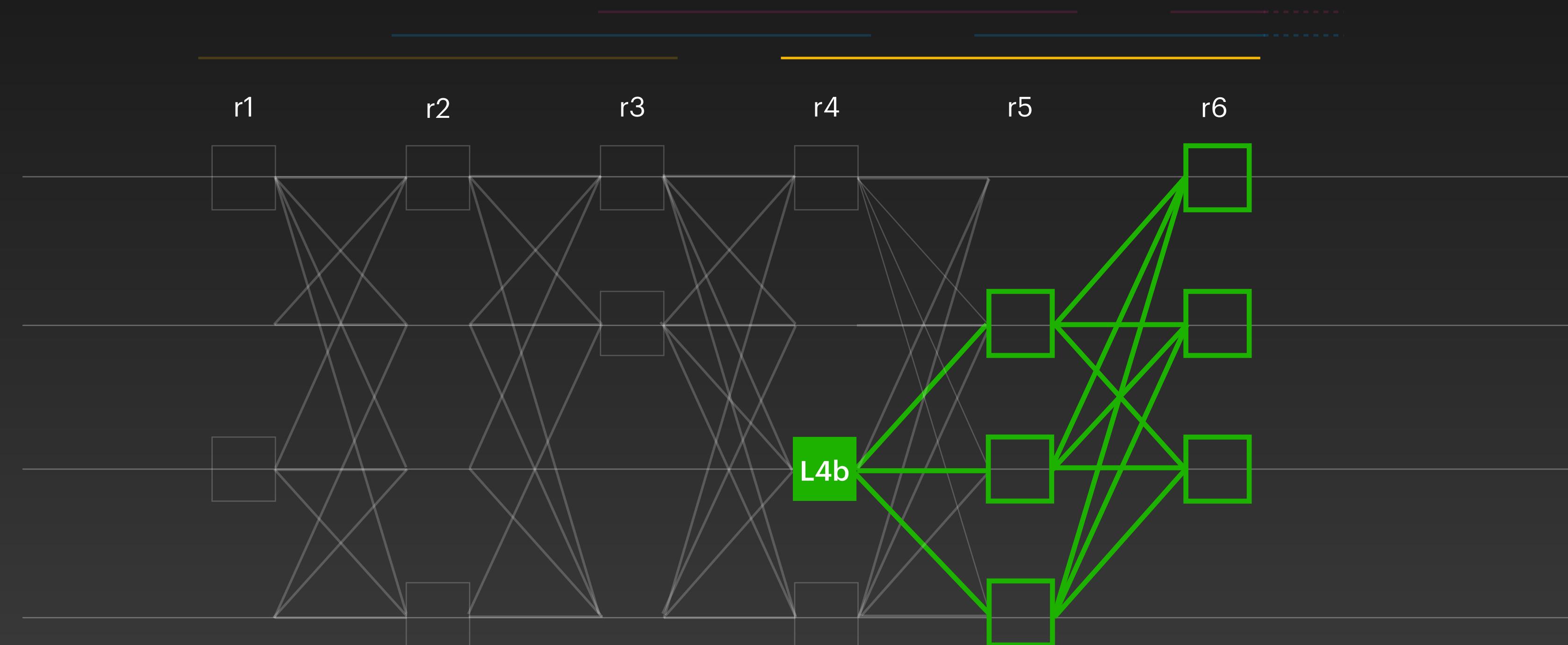
- **Skip** if  $2f+1$  blames
- **Commit** if  $2f+1$  certificates
- **Undecided** otherwise



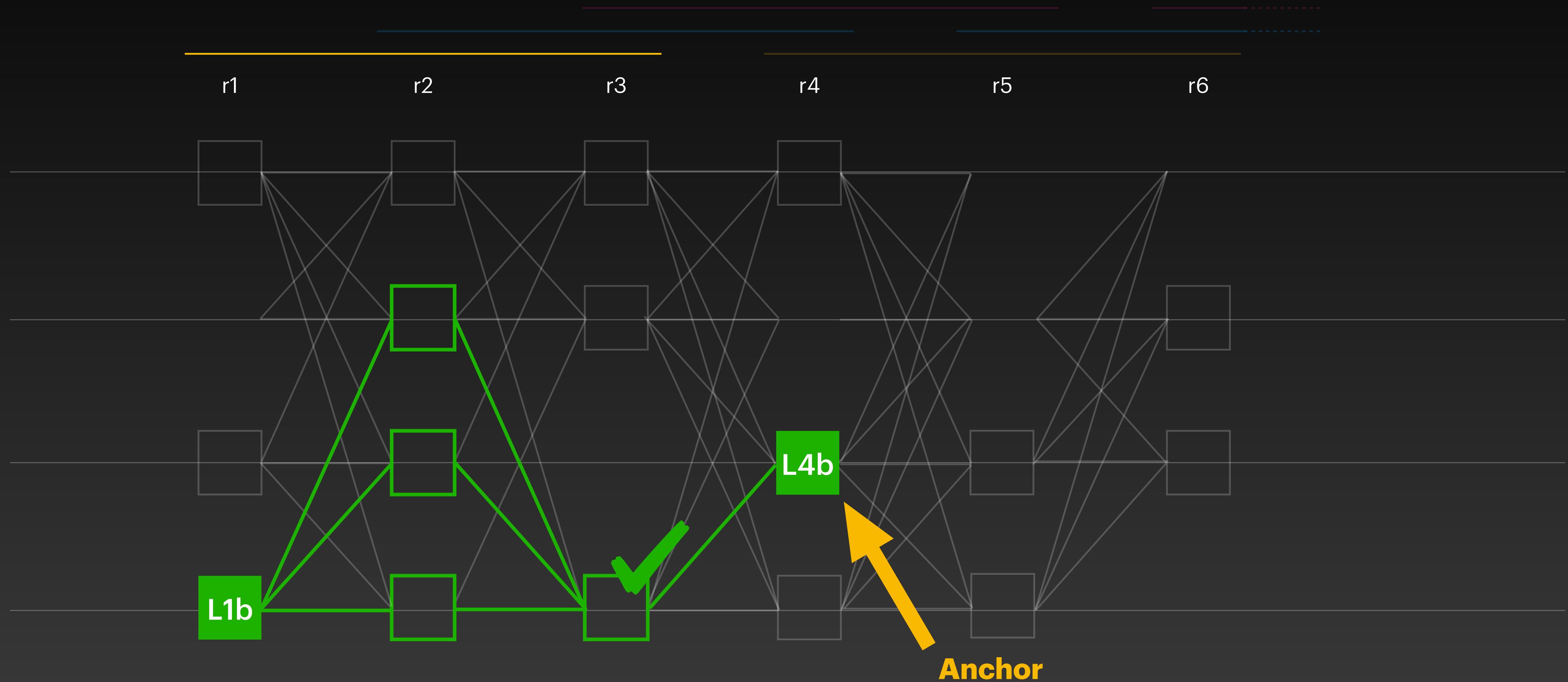
# Direct Decision Rule

On each leader starting from highest round:

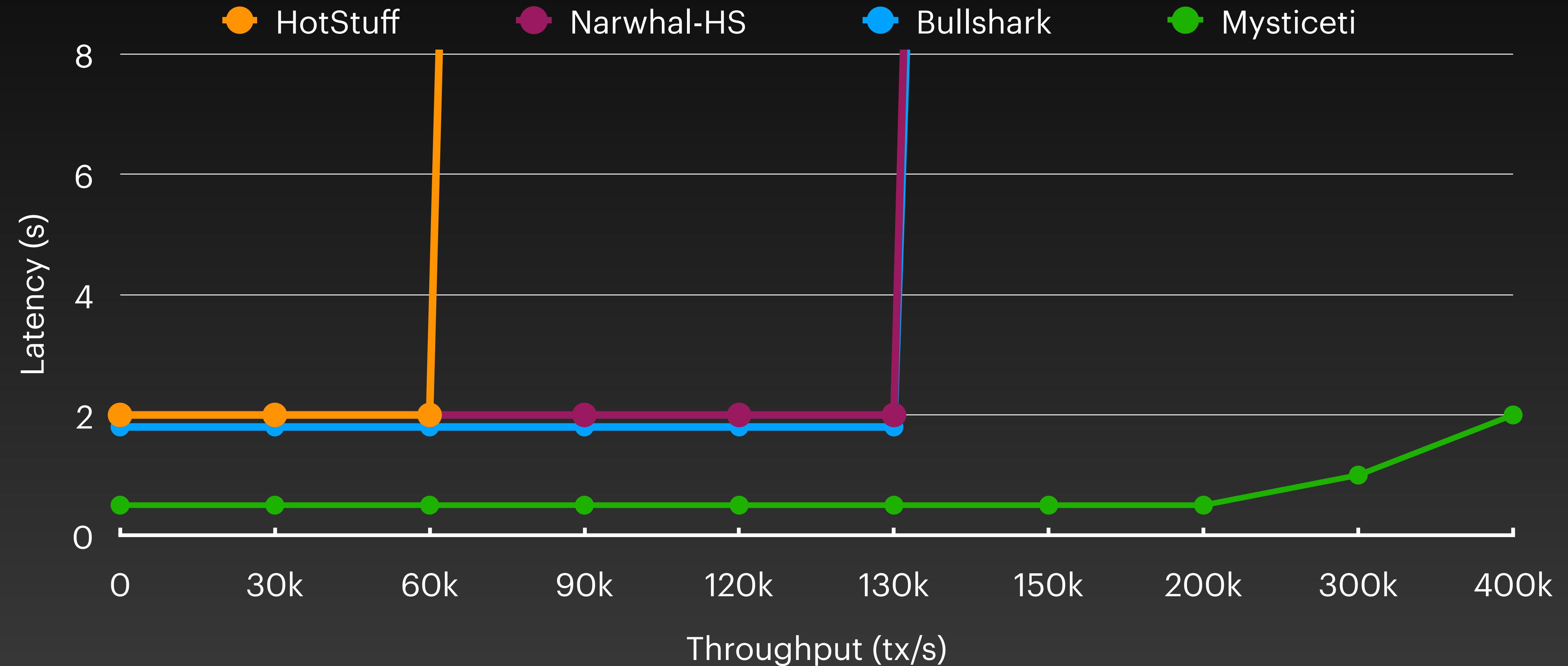
- **Skip** if  $2f+1$  blames
- **Commit** if  $2f+1$  certificates
- **Undecided** otherwise



# 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 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

# Testing Strategy



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

# 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 

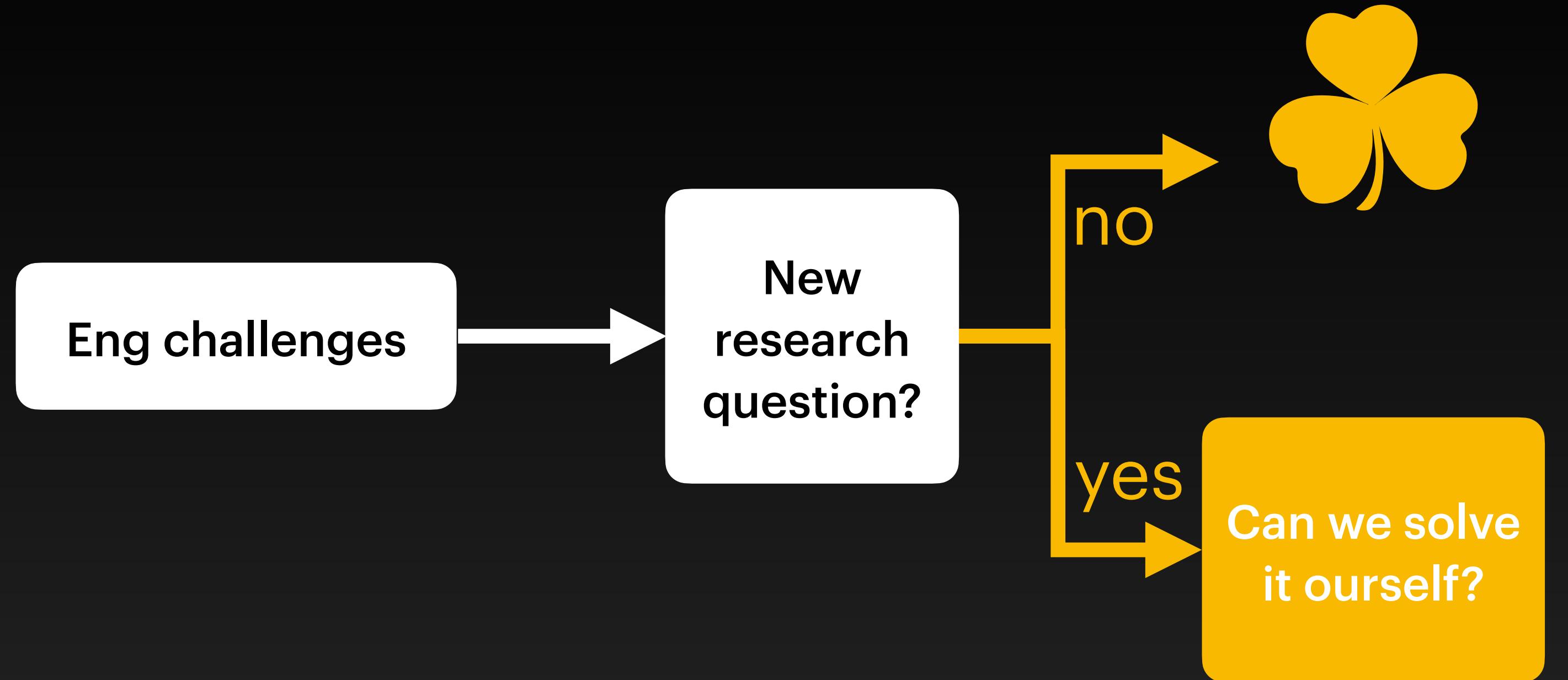
 **Dmitri Perelman** Oct 18th at 7:36 AM  
You're in the invite list!

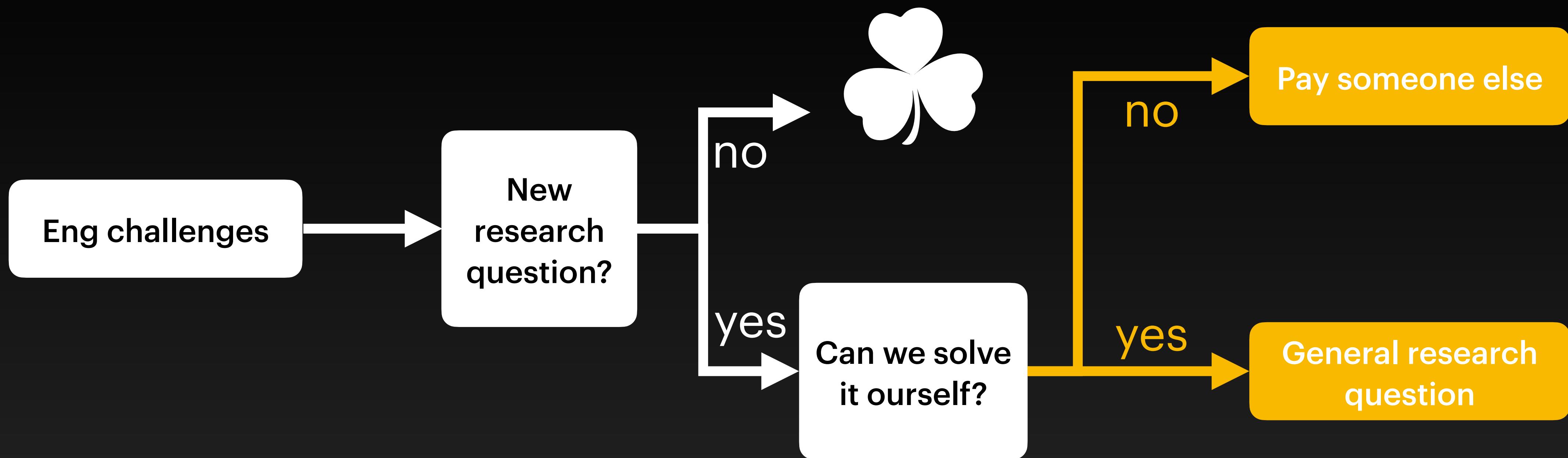
 1 

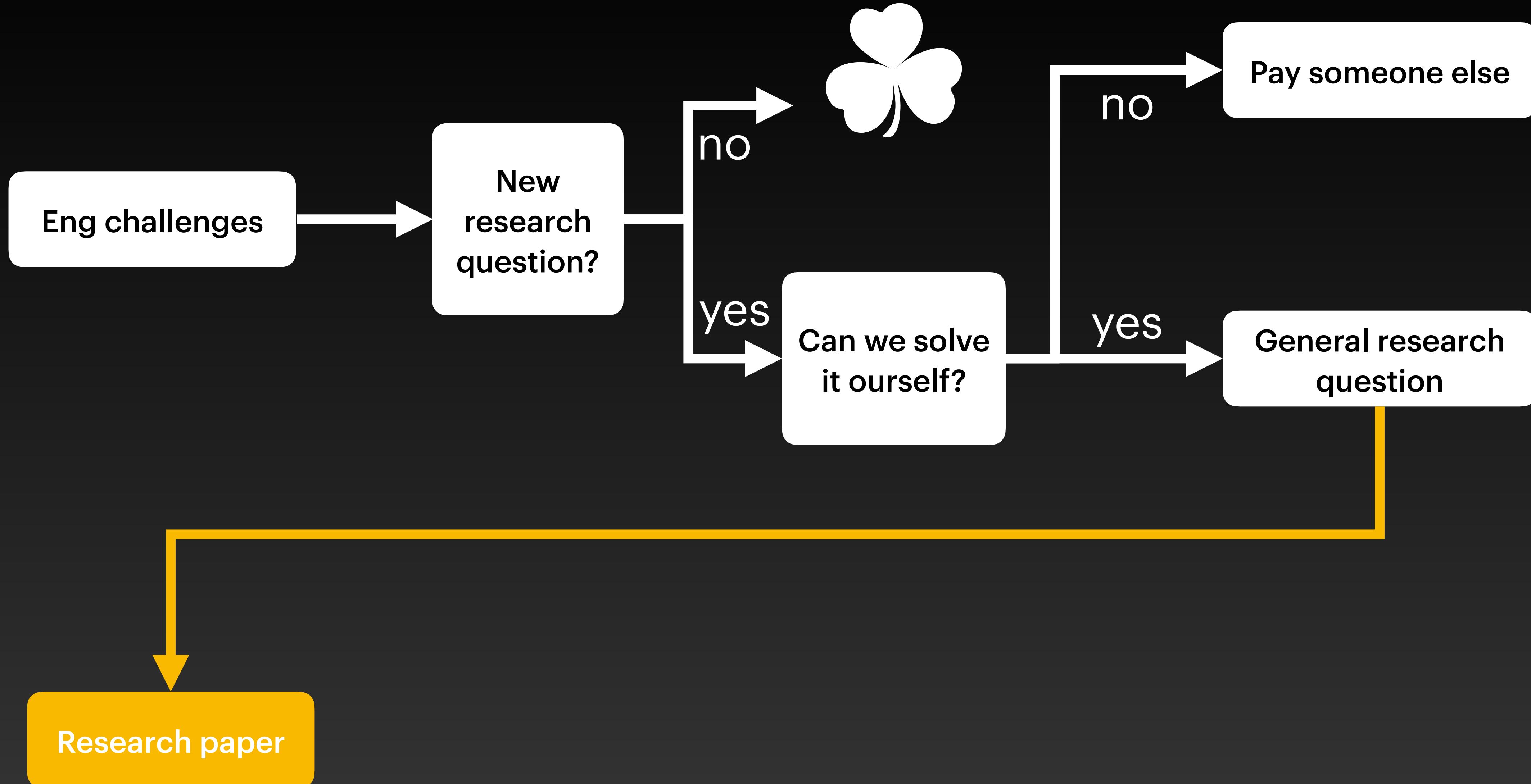
Eng challenges

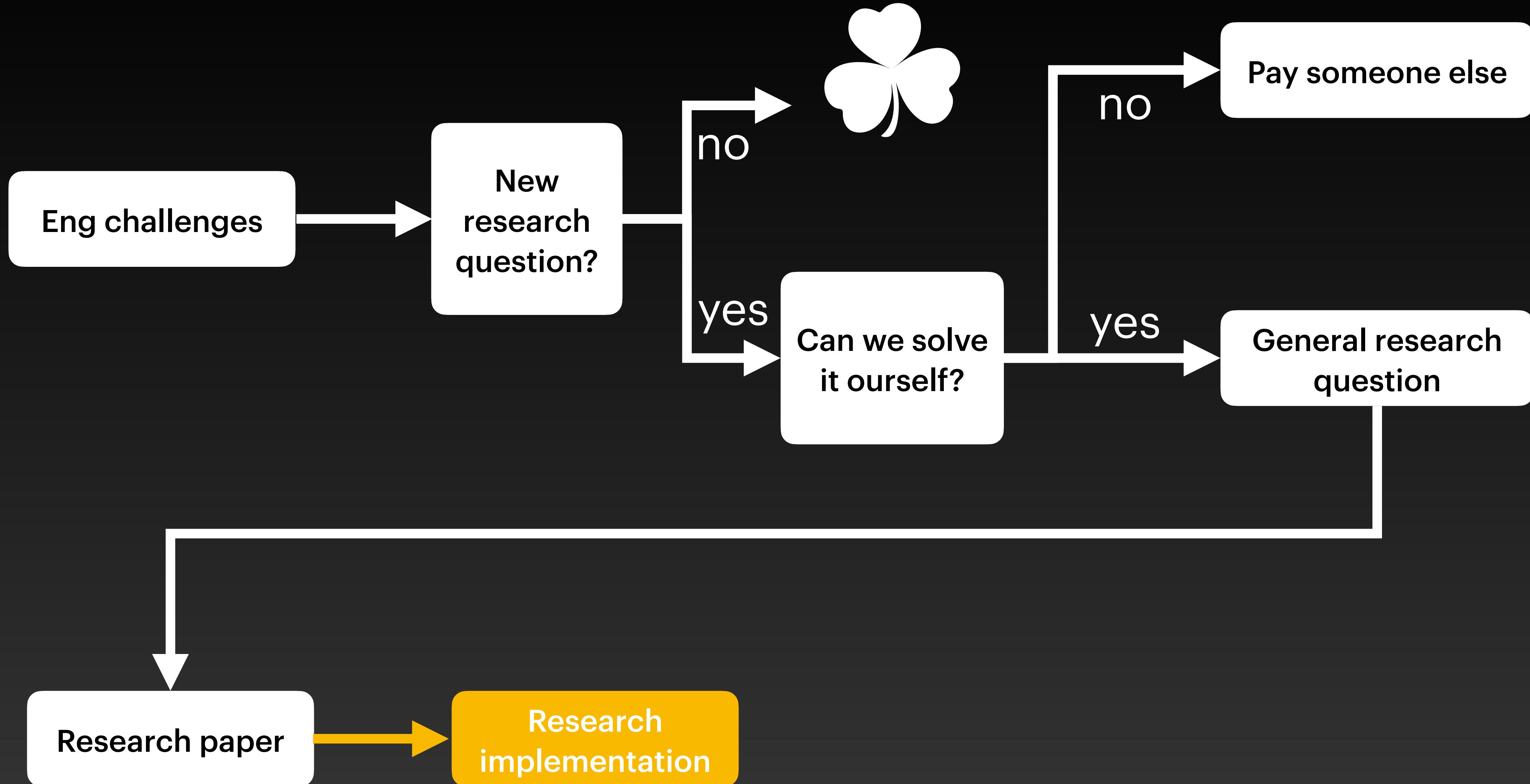


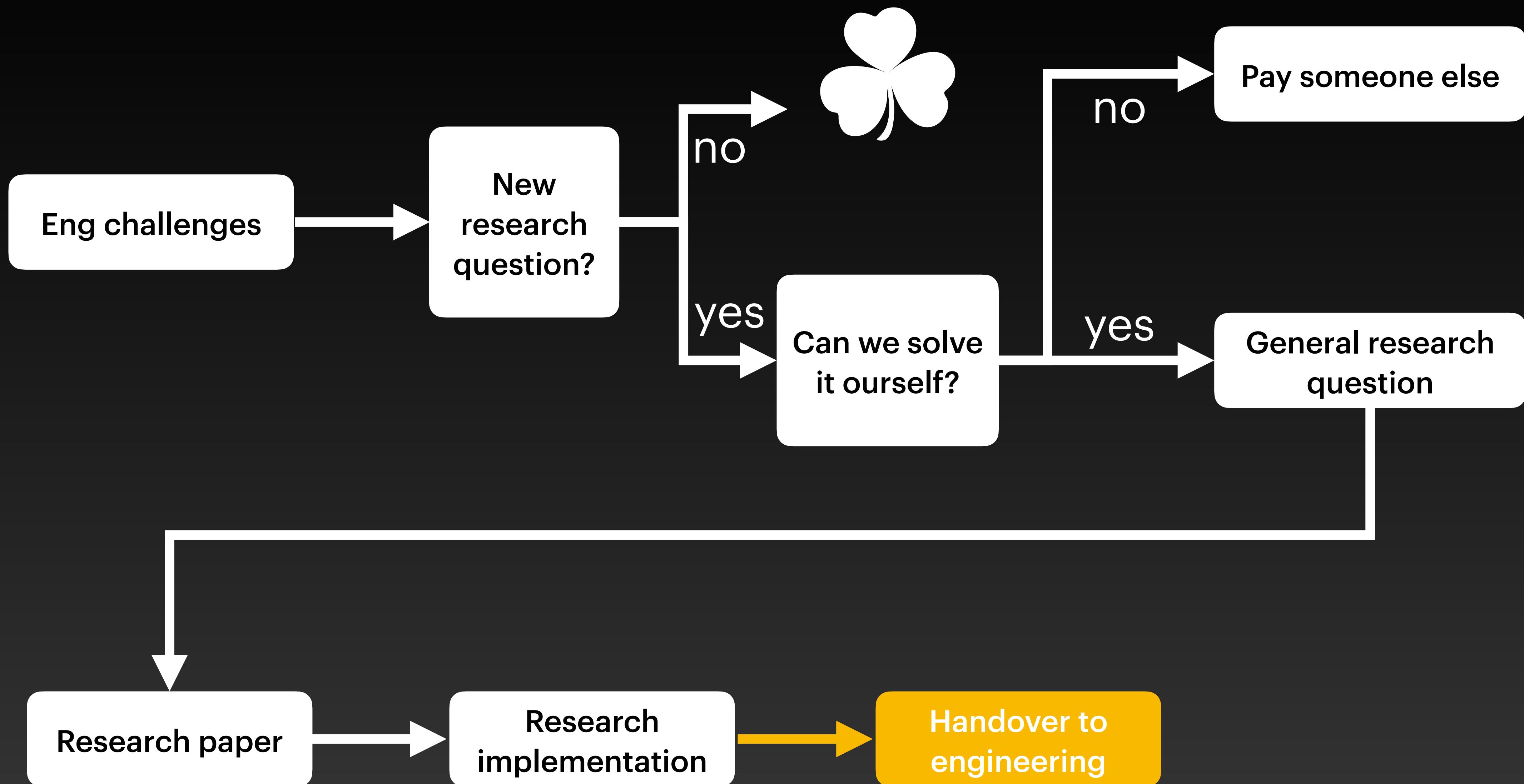
New  
research  
question?

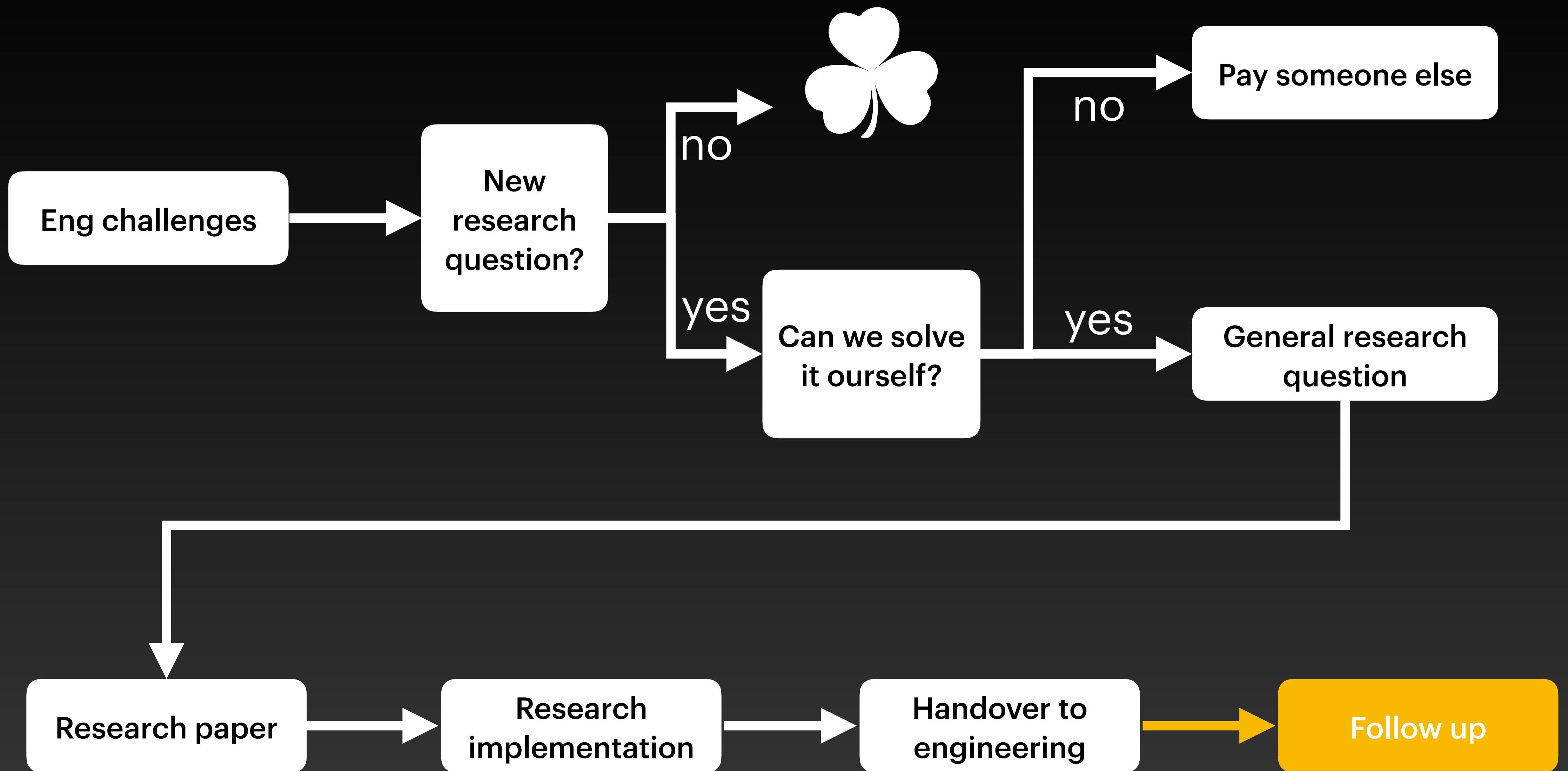












[Chainspace: A Sharded Smart Contracts Platform](#)

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

alberto@mystenlabs.com

**EXTRA**

# Implementation

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

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

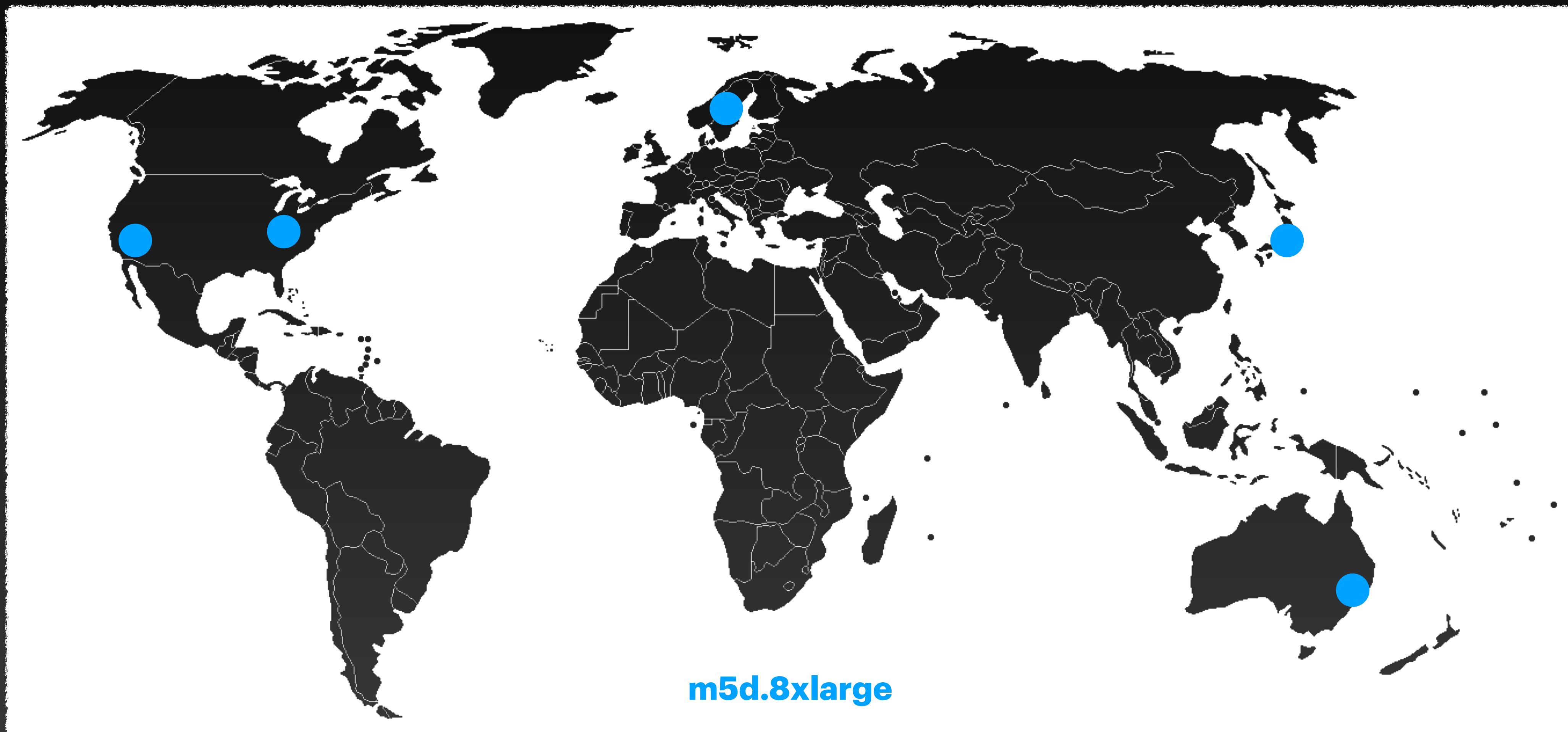
# Implementation

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

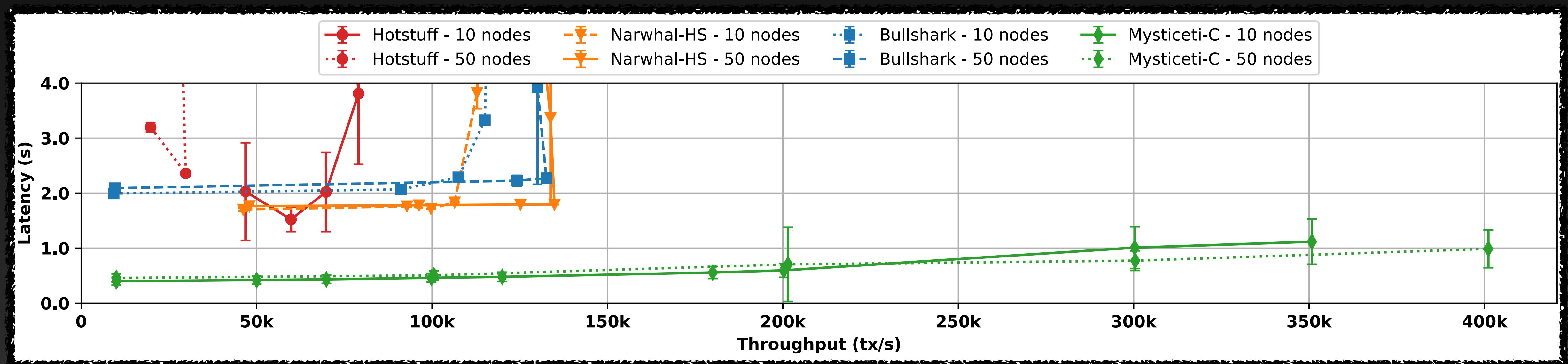
<https://github.com/mystenlabs/mysticeti>

# Evaluation

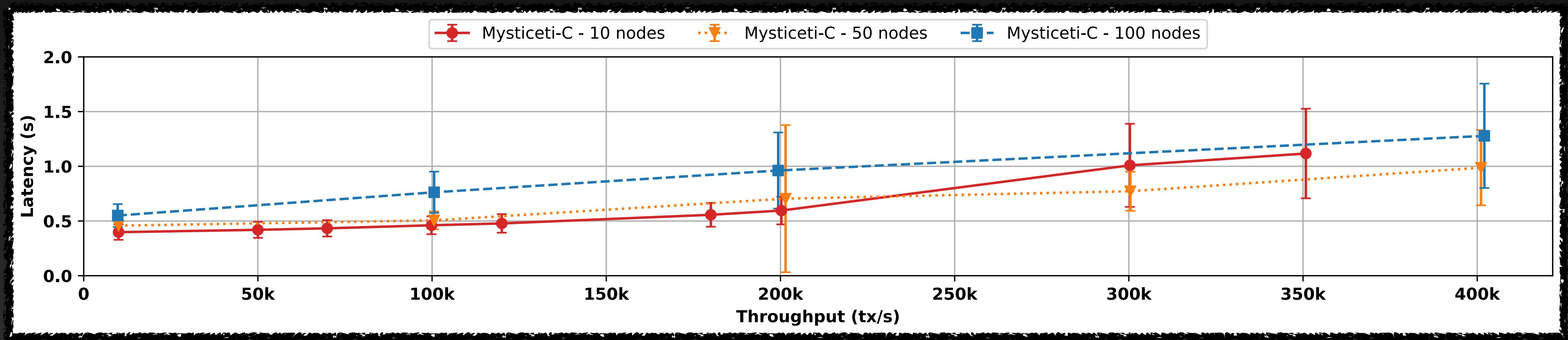
## Experimental setup on AWS



# Prototype Benchmarks



# Prototype Benchmarks



# 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