# The Evolution of Sui
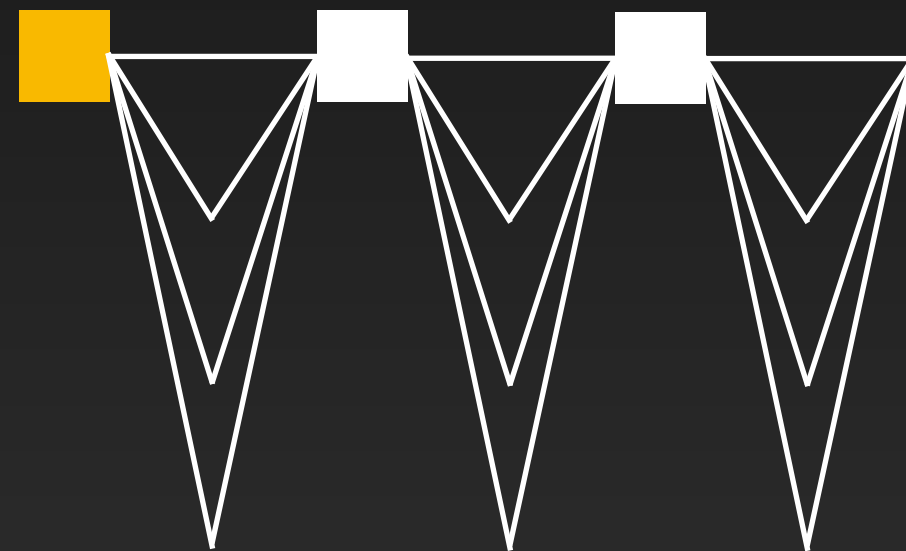
## From Academic Paper to Mainnet
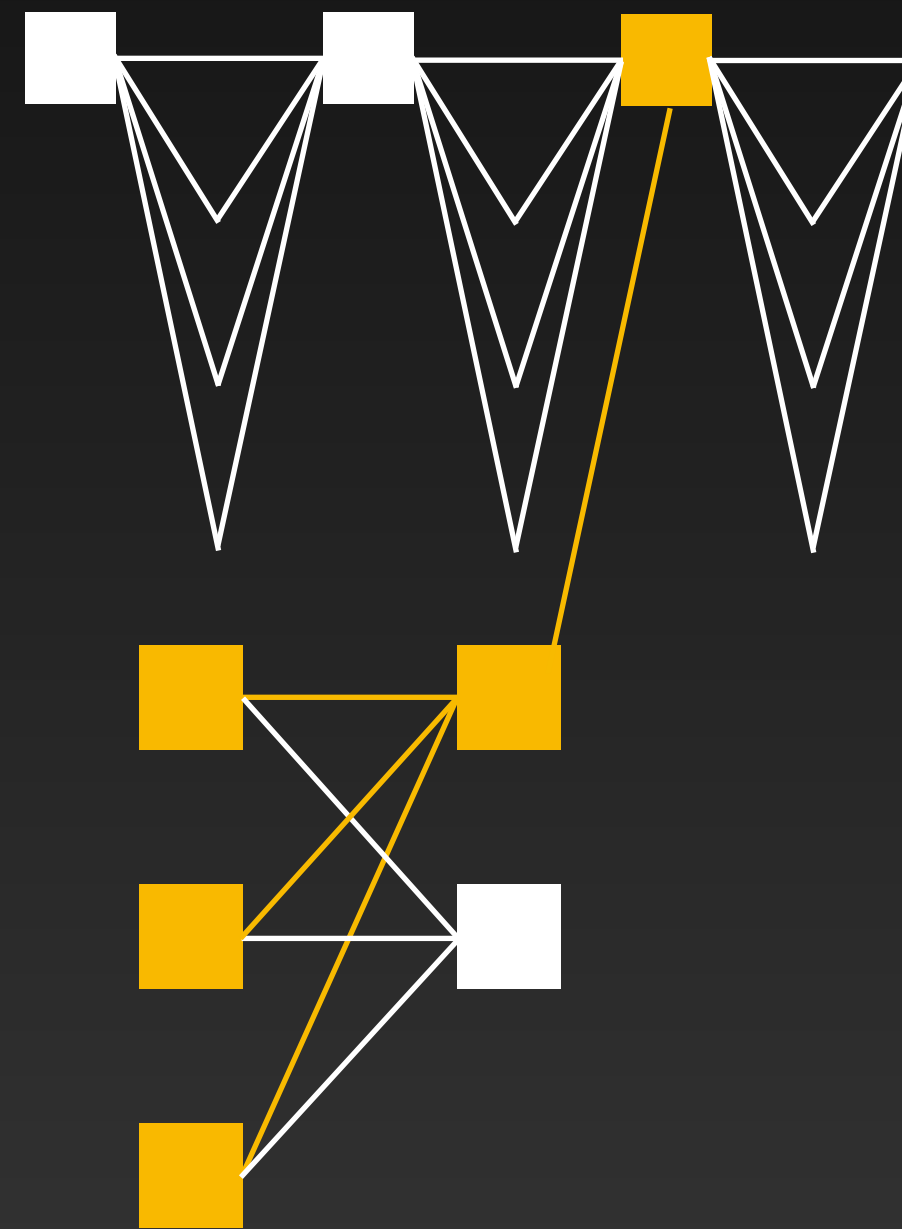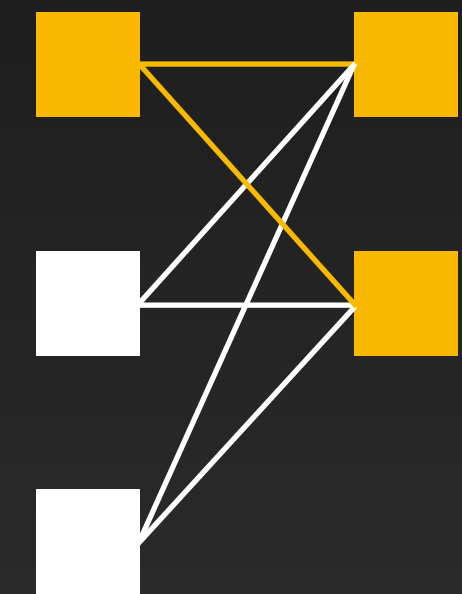
**2019**

**2024**



**HotStuff**

**HotStuff + Mempool**

**Bullshark, Mysticeti**

# Libra, 2019

## HotStuff

## HashGraph

?

# Libra, 2019

## HotStuff

✓ Linear

✓ Clearly isolated components

## HashGraph

✗ Hard 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

leader leader leader leader leader leader leader leader leader leader

resource utilization

# Research Questions

# Lessons Learned

1. Network model?

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

# Research Questions

# Lessons Learned

1. Network model?

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

# Libra, 2019



HotStuff: BFT Consensus in the Lens of Blockchain

Maofan Yin[1,2], Dahlia Malkhi[2], Michael K. Reiter[2,3], Guy Golan Gueta[2], and Ittai Abraham[2]

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

## HotStuff (naive mempool)

- Linear

- Clearly isolated components

- Uneven resource utilisation

- Fragile to faults and asynchrony

- Unspecified components (pacemaker)

# Libra, 2021



## Narwhal

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

# Narwhal

# Narwhal

# Narwhal



block       certificate

B    V    C

B    V    C

B    V    C

**Round 1**

# Narwhal



block

certificate

B    V    C

B    V    C

C

B    V

C

**Round 1**

**Byzantine 'Reliable' Broadcast**

Narwhal

# Research Questions

1. Network model?

# Lessons Learned

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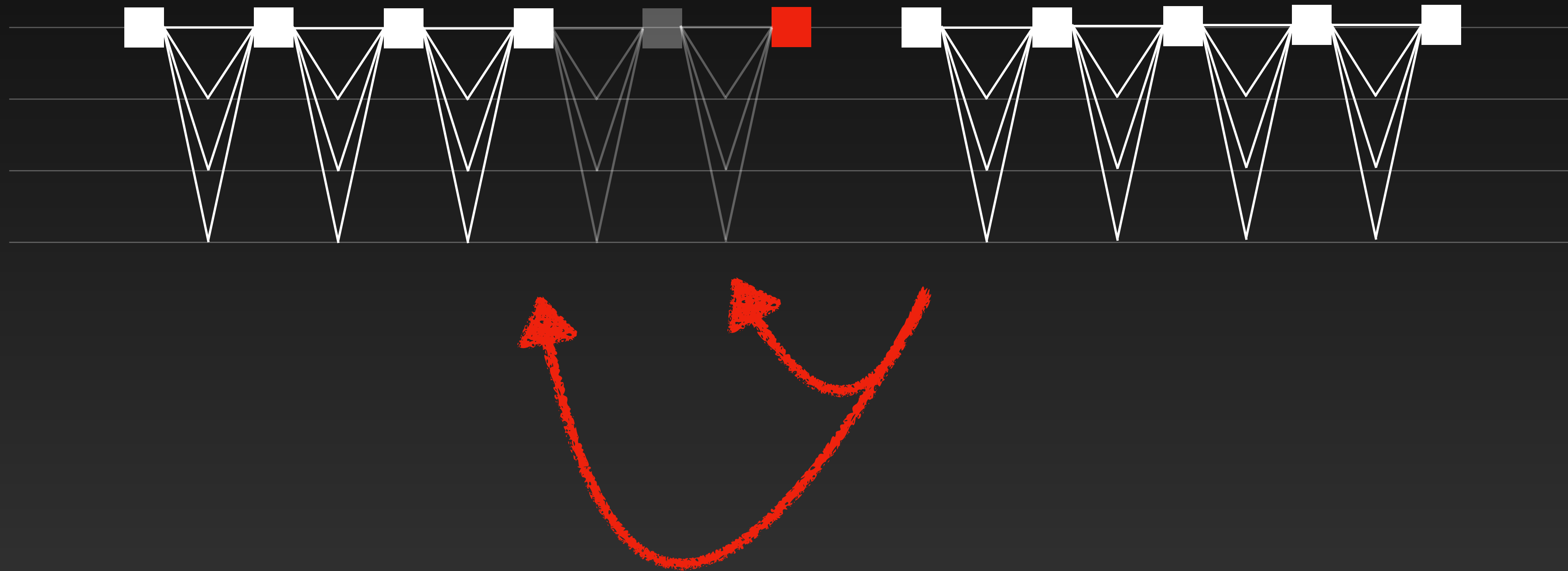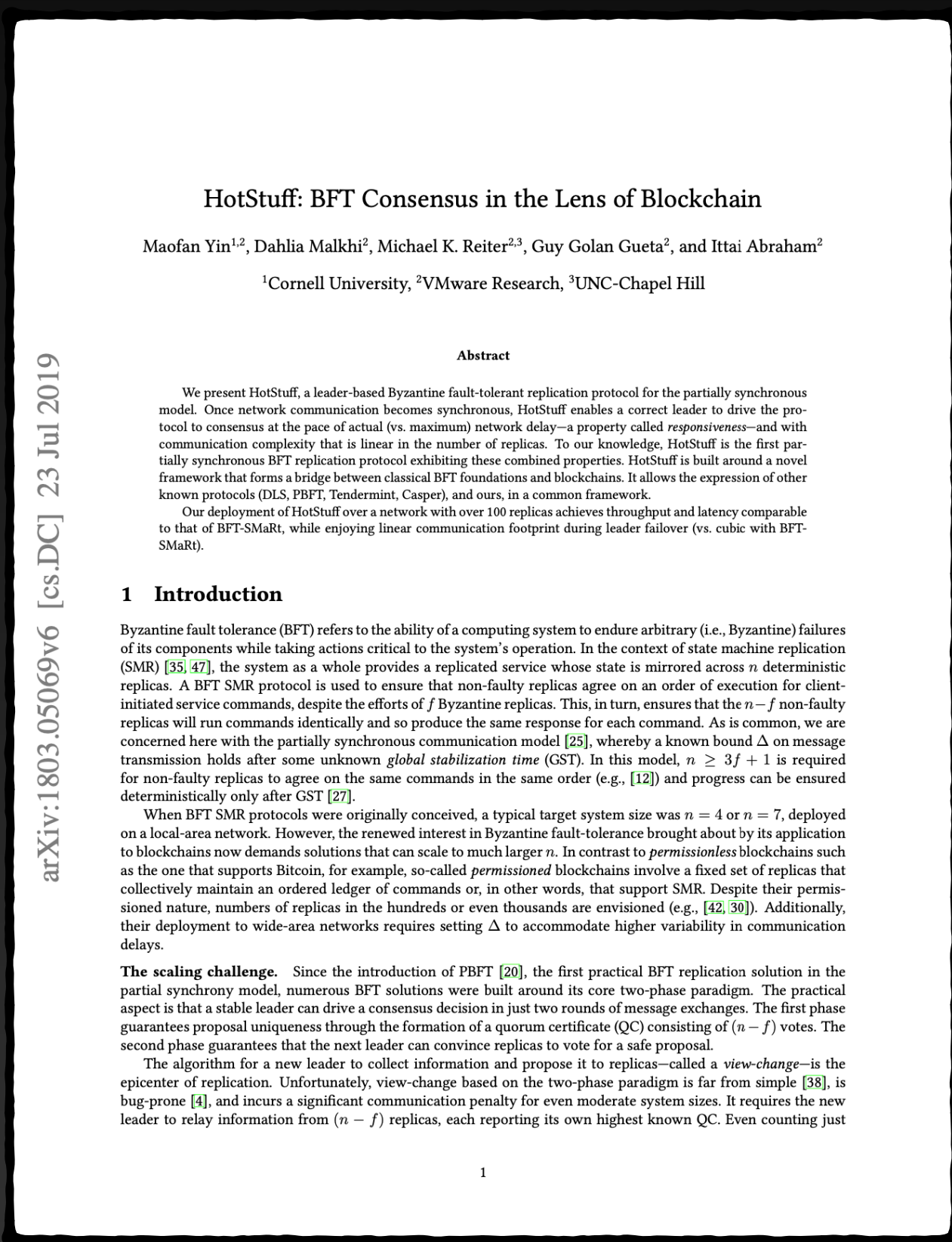4. Codesign with mem. and storage

# HotStuff on Narwhal
## Enhanced commit rule

# HotStuff on Narwhal
## Enhanced commit rule

# HotStuff on Narwhal
## Enhanced commit rule



**Faulty HotStuff Leader!**

**Blocks may still be 'saved'**

C1

# HotStuff on Narwhal
## Enhanced commit rule

# Performance

# Performance

# Libra, 2021



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

## 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
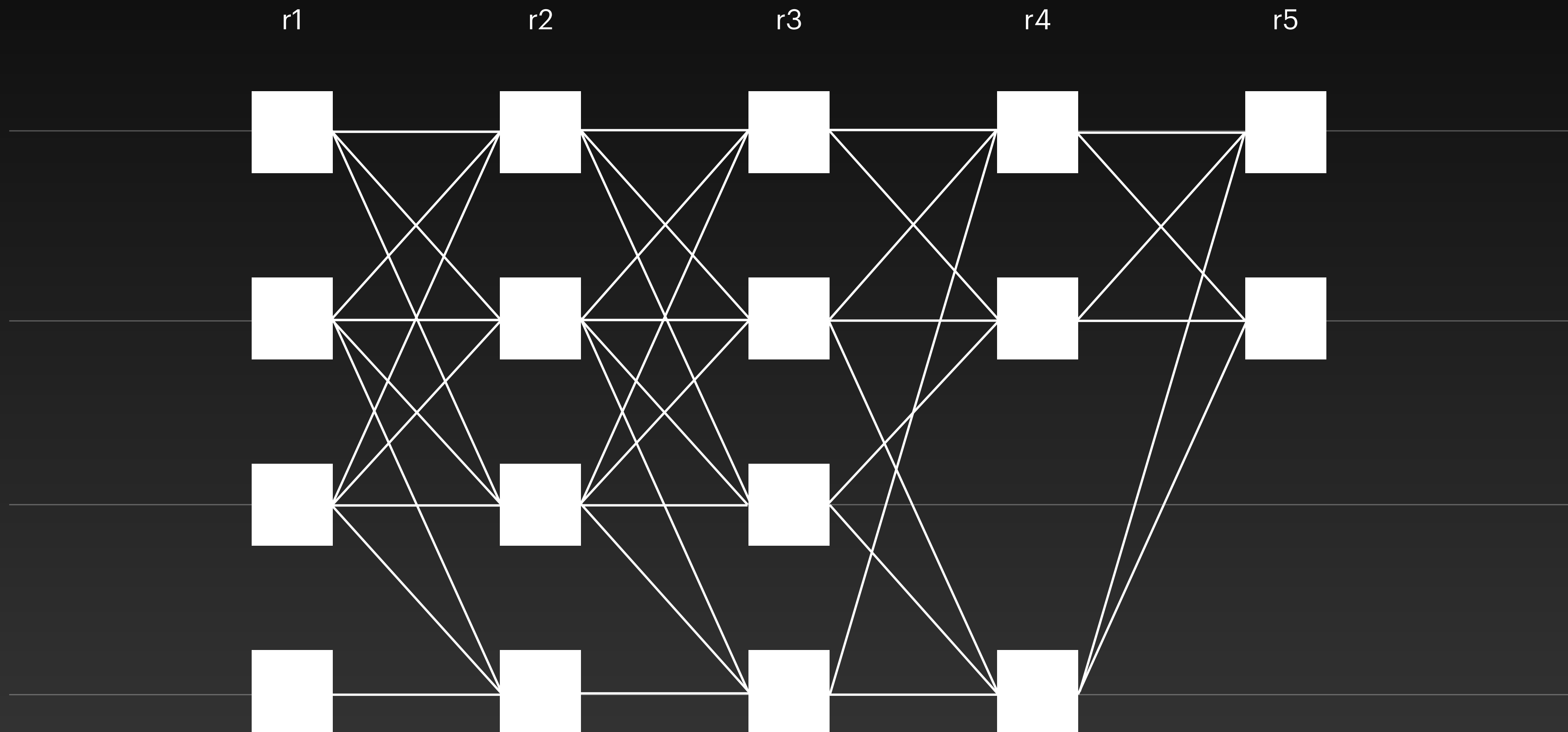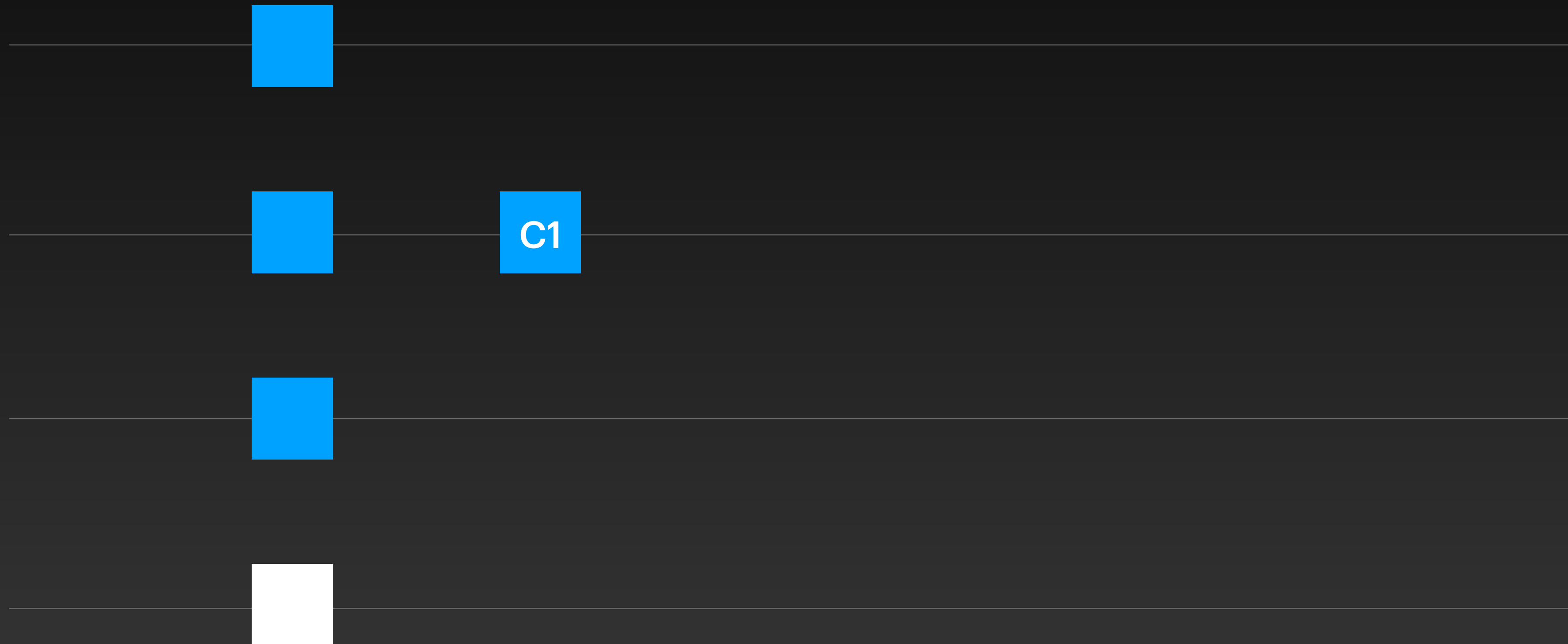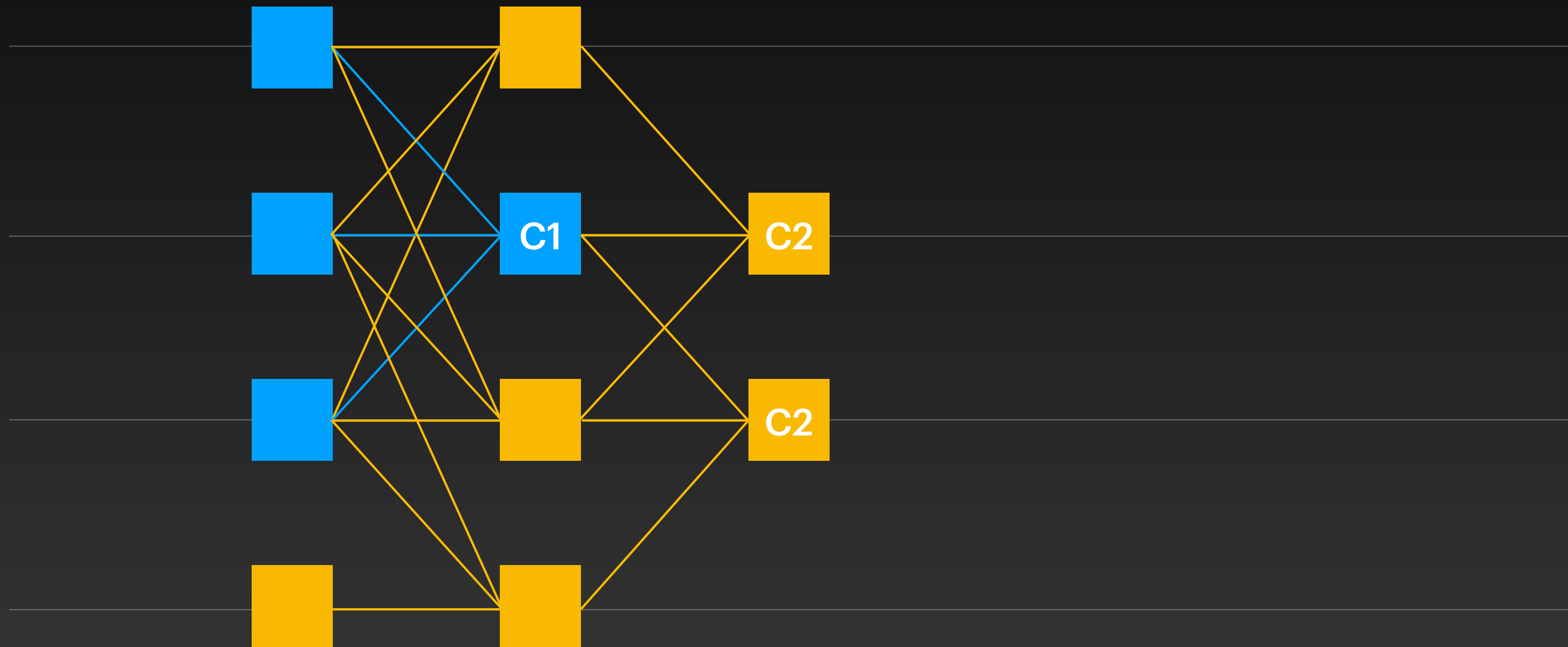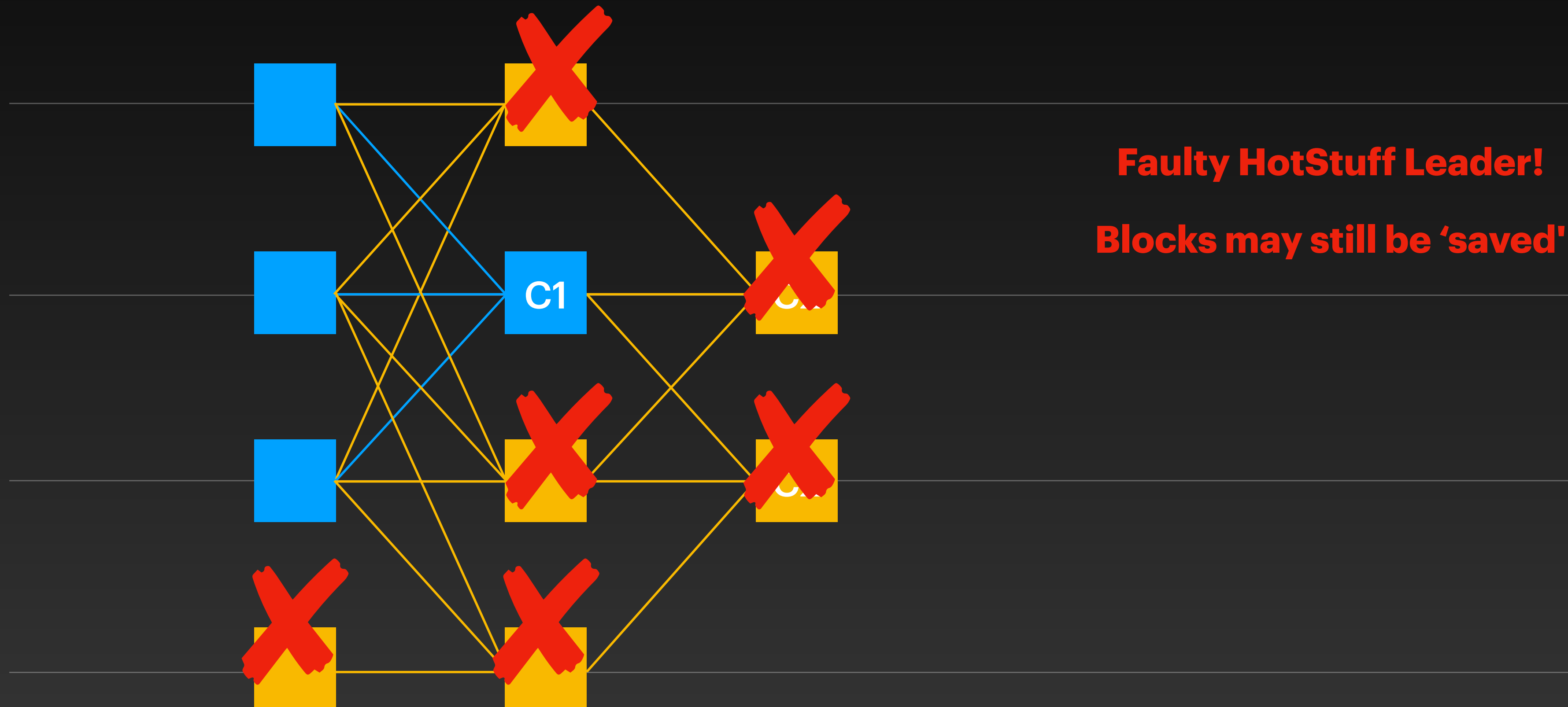4. Codesign with mem. and storage

# DagRider

# Tusk

# Bullshark

# Dumbo-NG

## Data Dissemination

- Hard to make efficient

- 99% of the code

## Consensus

- Error prone

- Isolated, easy to maintain

# Bullshark

## Bullshark: DAG BFT Protocols Made Practical

Alexander Spiegelman
sasha@aptoslabs.com
Aptos

Neil Giridharan
giridhn@berkeley.edu
UC Berkeley

Alberto Sonnino
alberto@mysternlabs.com
Mysten Labs

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

# Performance

# Research Questions

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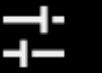
# Lessons Learned

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...

📌 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

# Sailfish

# CM

# Mysticeti

## 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

# Mysticeti

## Shoal: Improving DAG-BFT Latency And Robustness

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

## Cordial Miners: Fast and Efficient Consensus for Every Eventuality

## MYSTICETI: Reaching the Latency Limits with Uncertified DAGs

# Uncertified DAG



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

# Interpreting DAG Patterns

# Performance

# 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

# The Sui Mainnet

## p50 consensus latency ⓘ

2024-05-17 09:27:15

━ p50      **397 ms**

# The Roadmap

**2019**

naive consensus

→

**2020-2021**

mempool ❤️ consensus

→

**2022-2023**

Fold into DAG

→

**2024**

Remove overhead

# 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!

TY 1

# 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
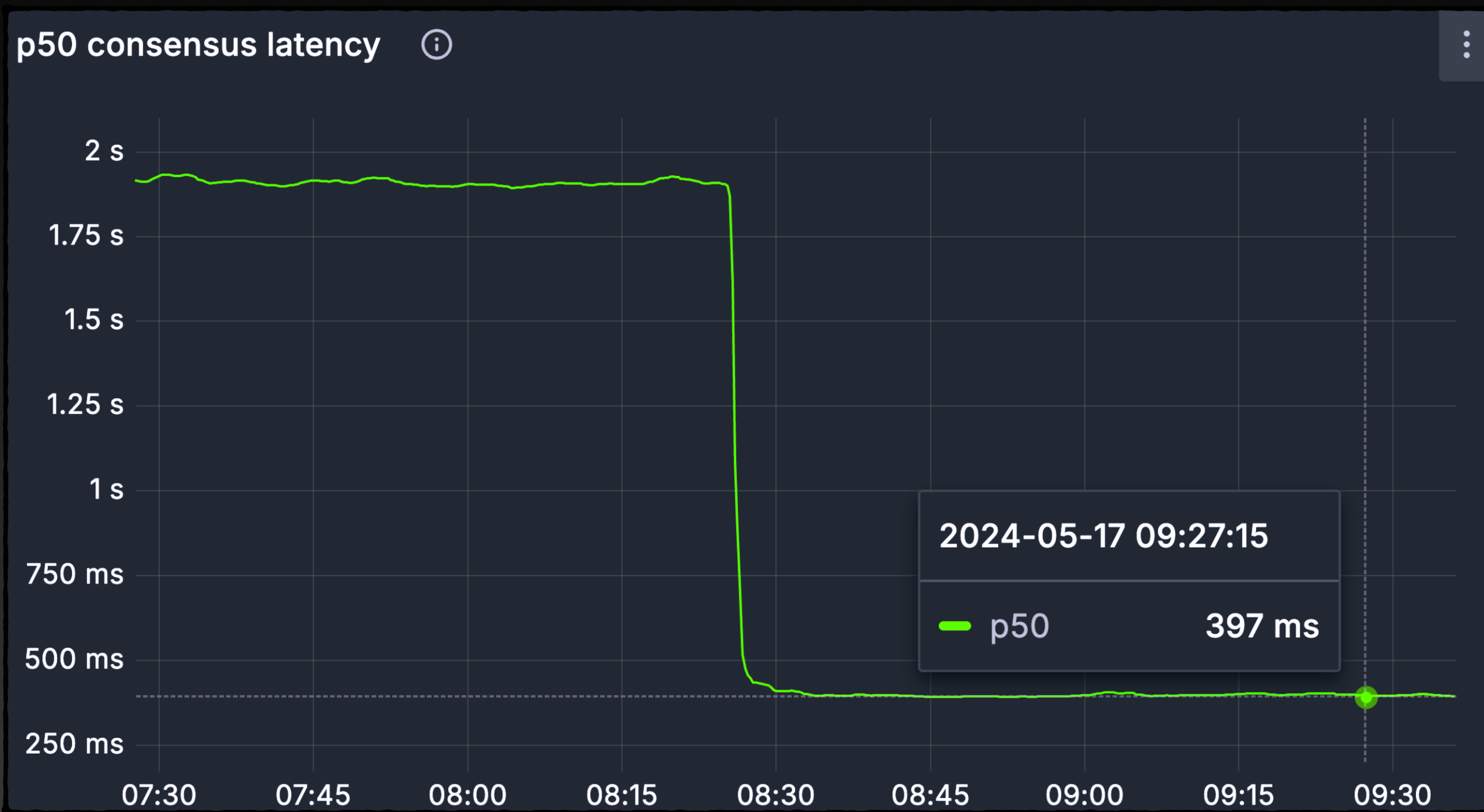5. Core is hard, consensus is easy
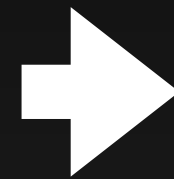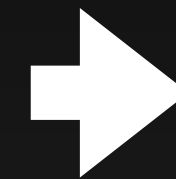6. Epoch change is not an add-on
7. Writing papers to explore designs

# 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



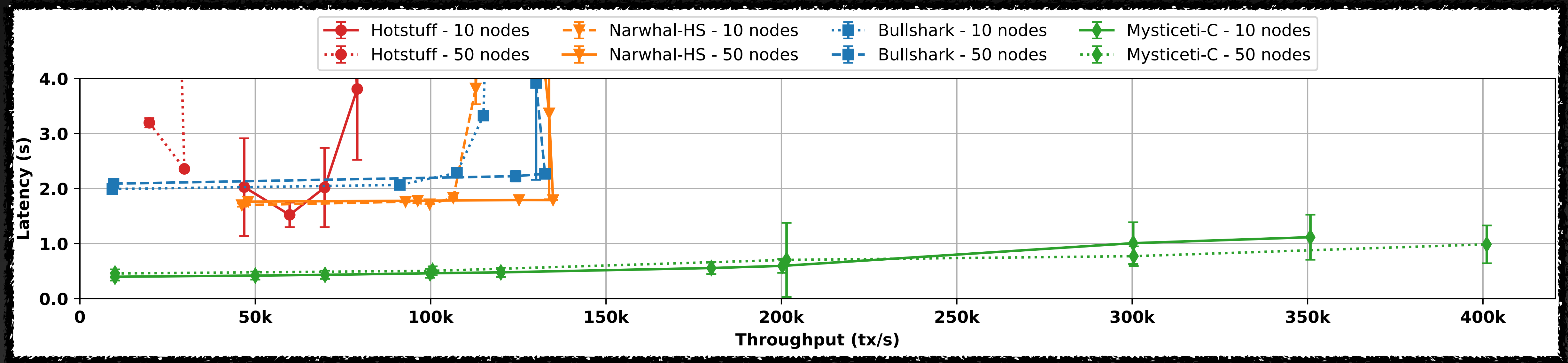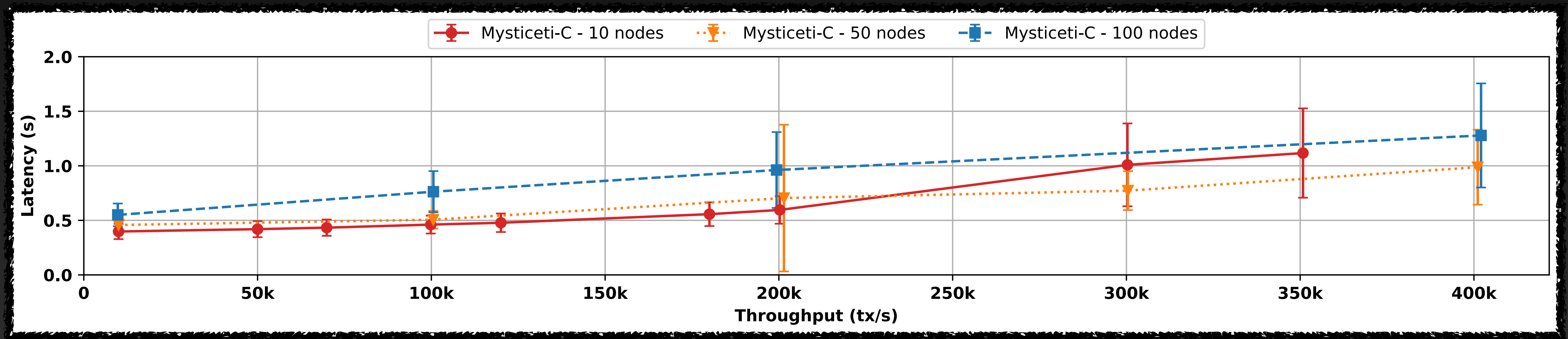m5d.8xlarge

# Prototype Benchmarks

# Prototype Benchmarks

# Mysticeti
## Key Limitations

- Block Synchroniser

- Parallelise block creation and synchronisation

- Rigid DAG structure?