

## Ladon: High-Performance Multi-BFT Consensus via Dynamic Global Ordering

Hanzheng Lyu<sup>1</sup>, Shaokang Xie<sup>2</sup>, Jianyu Niu<sup>2</sup>, Chen Feng<sup>1</sup>,  
Yinqian Zhang<sup>2</sup>, and Ivan Beschastnikh<sup>1</sup>

<sup>1</sup>University of British Columbia, Canada

<sup>2</sup>Southern University of Science and Technology, China



THE UNIVERSITY  
OF BRITISH COLUMBIA



**SUSTech**

Southern University  
of Science and  
Technology



# BFT Consensus



THE UNIVERSITY  
OF BRITISH COLUMBIA

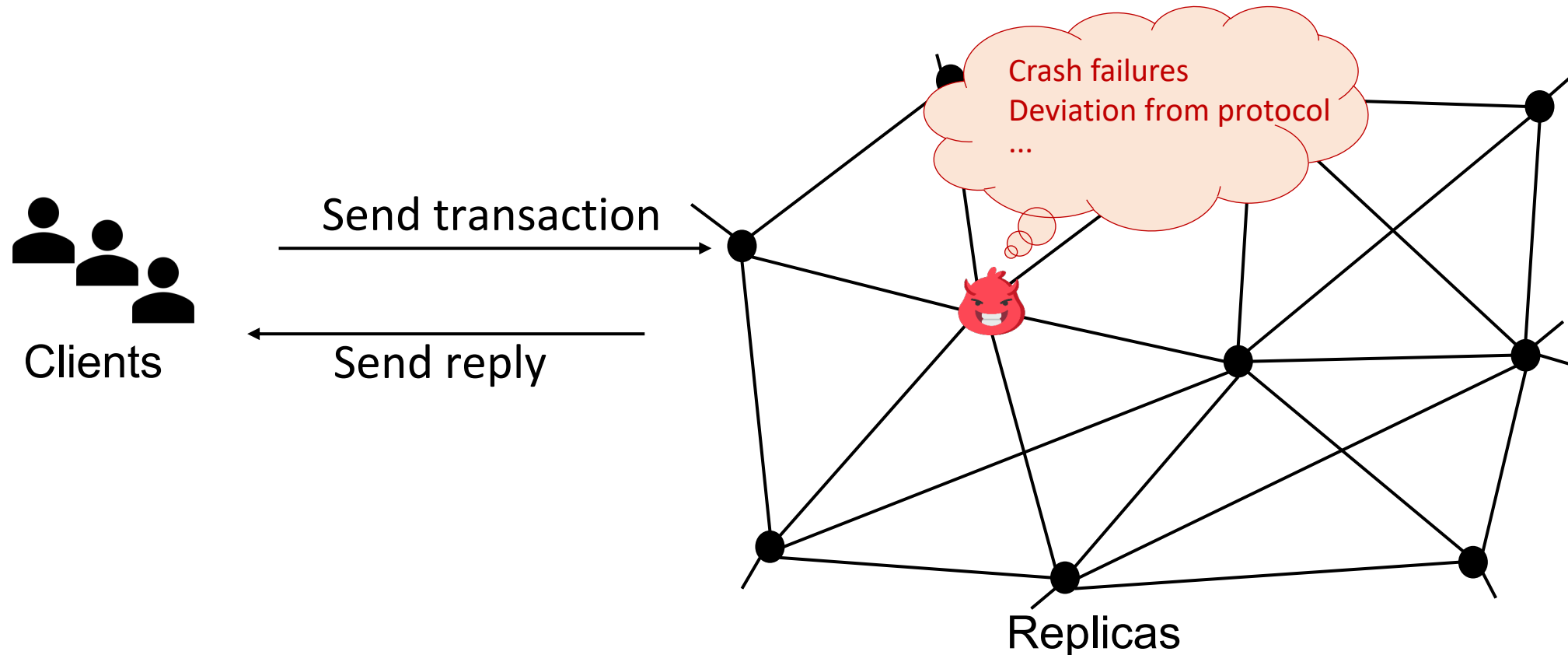


SUSTech

Southern University  
of Science and  
Technology

## Byzantine Fault Tolerant (BFT) Consensus

- ❑ allows a decentralized system to reach agreement among replicas
- ❑ despite Byzantine failures





# Leader-based BFT<sup>[1]</sup>

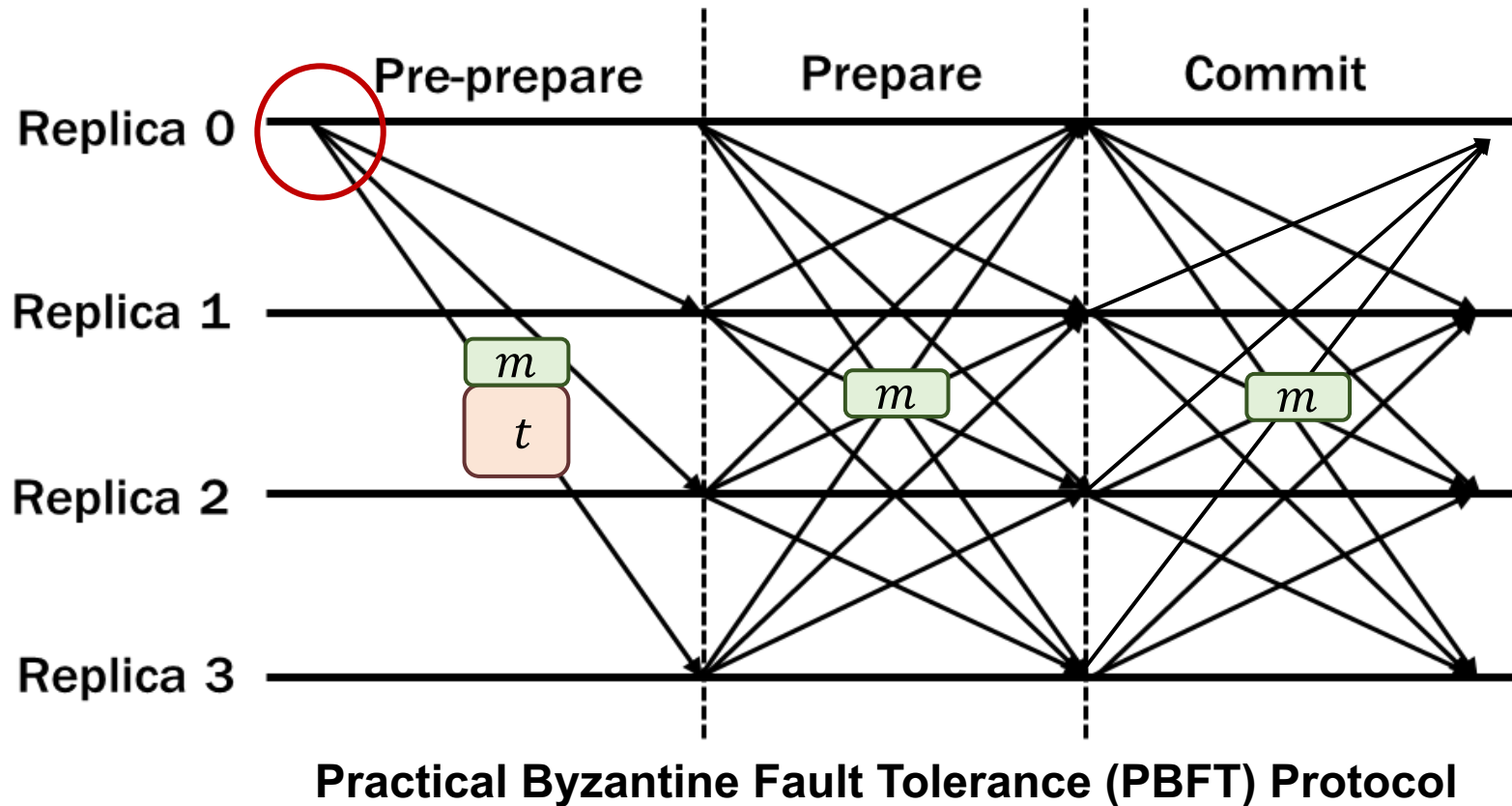


THE UNIVERSITY  
OF BRITISH COLUMBIA



SUSTech  
Southern University  
of Science and  
Technology

The leader limits the overall performance of the BFT protocol.



$$size(t) \gg size(m)$$

$T_{\max}$

Maximum throughput

$B$

Leader's  
bandwidth

$$\approx \frac{B}{(n-1)size(t)}$$

number of replicas

[1] Miguel Castro and Barbara Liskov. Practical Byzantine Fault Tolerance. In OSDI, 1999.



# Multi-BFT Consensus<sup>[1-3]</sup>

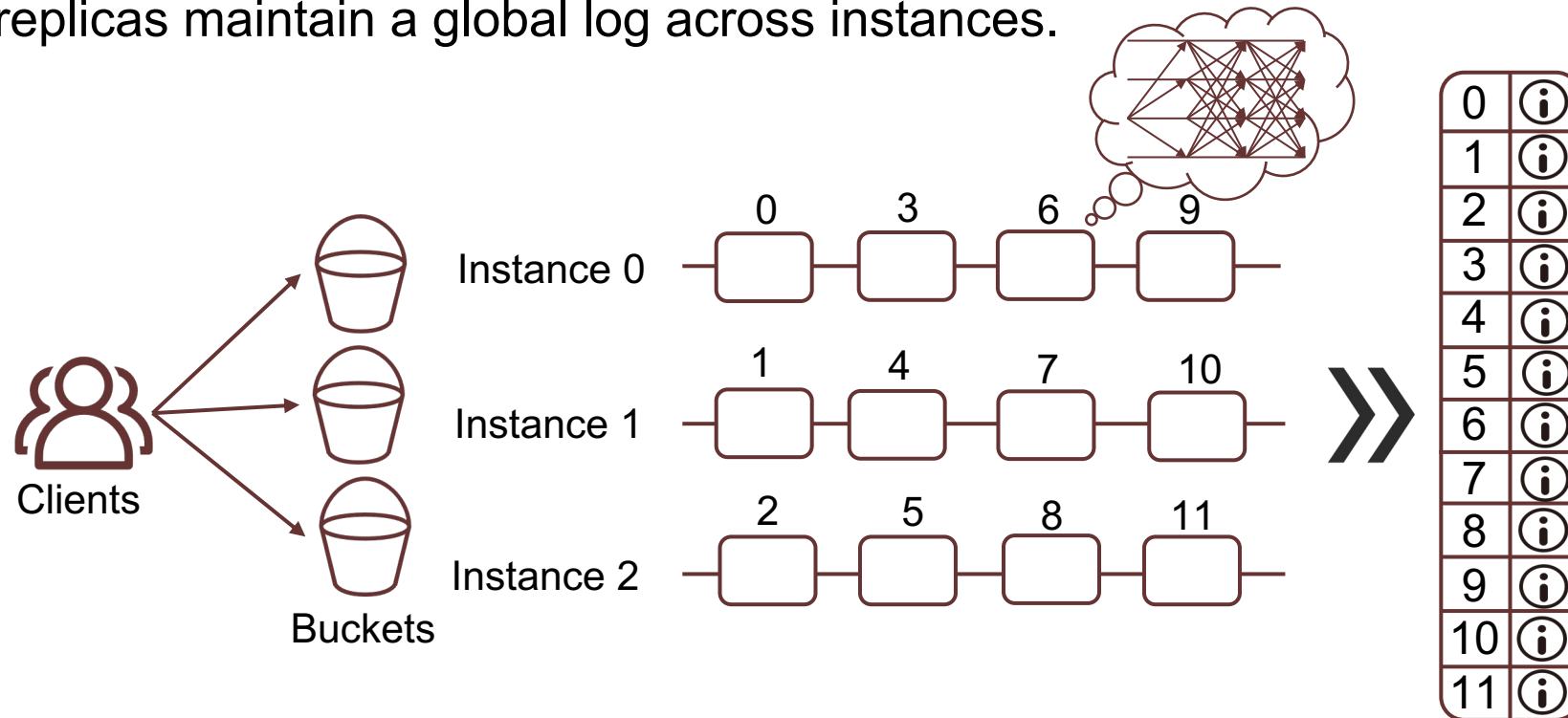


THE UNIVERSITY  
OF BRITISH COLUMBIA



SUSTech  
Southern University  
of Science and  
Technology

- ❑ Each replica act as the leader of one instance and backups of other instances.
- ❑ Each instance independently outputs a sequence of committed blocks.
- ❑ All replicas maintain a global log across instances.



- [1] Chrysoula Stathakopoulou, Matej Pavlovic, and Marko Vukolić. State Machine Replication Scalability Made Simple, in ACM EuroSys'22.
- [2] Chrysoula Stathakopoulou, Tudor David, and Marko Vukolic. Mir BFT: Scalable and Robust BFT for Decentralized Networks, in Jsyz'22.
- [3] Suyash Gupta, Jelle Hellings, and Mohammad Sadoghi. RCC: Resilient Concurrent Consensus for High-Throughput Secure Transaction Processing, in IEEE ICDE'21.



# Stragglers in Network



THE UNIVERSITY  
OF BRITISH COLUMBIA



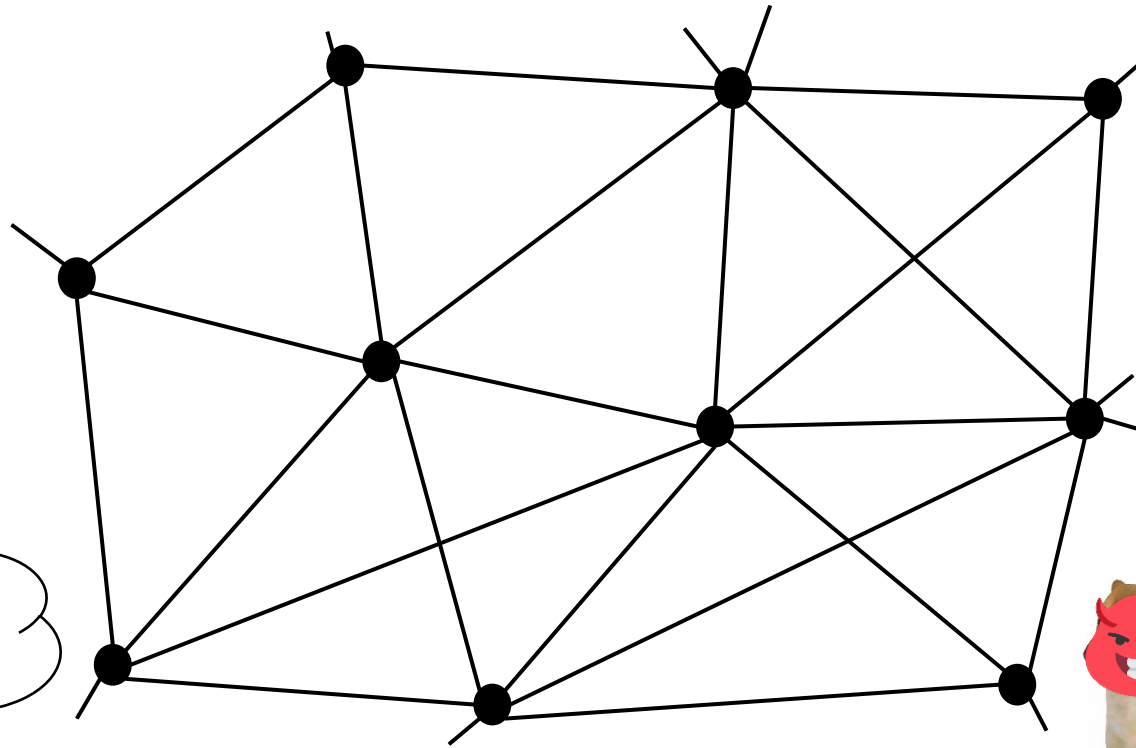
SUSTech

Southern University  
of Science and  
Technology

In distributed systems, disparities in CPU, bandwidth, and other resources, along with potential malicious behavior, can lead to some replicas being slower, causing reduced consensus efficiency.



I have abundant  
resources—  
I am fast.



I have limited  
resources—  
I am slow.



I am malicious—  
I will deliberately slow  
things down.



# The Impact of Stragglers



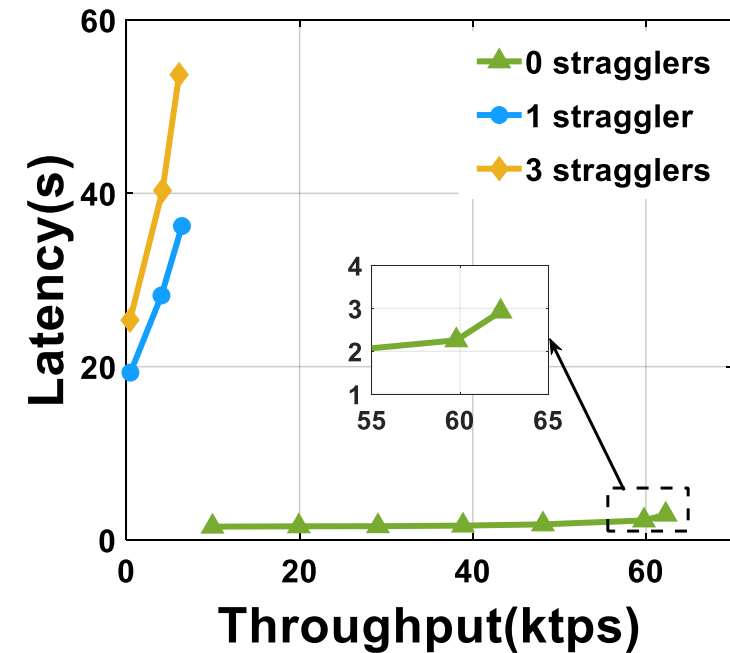
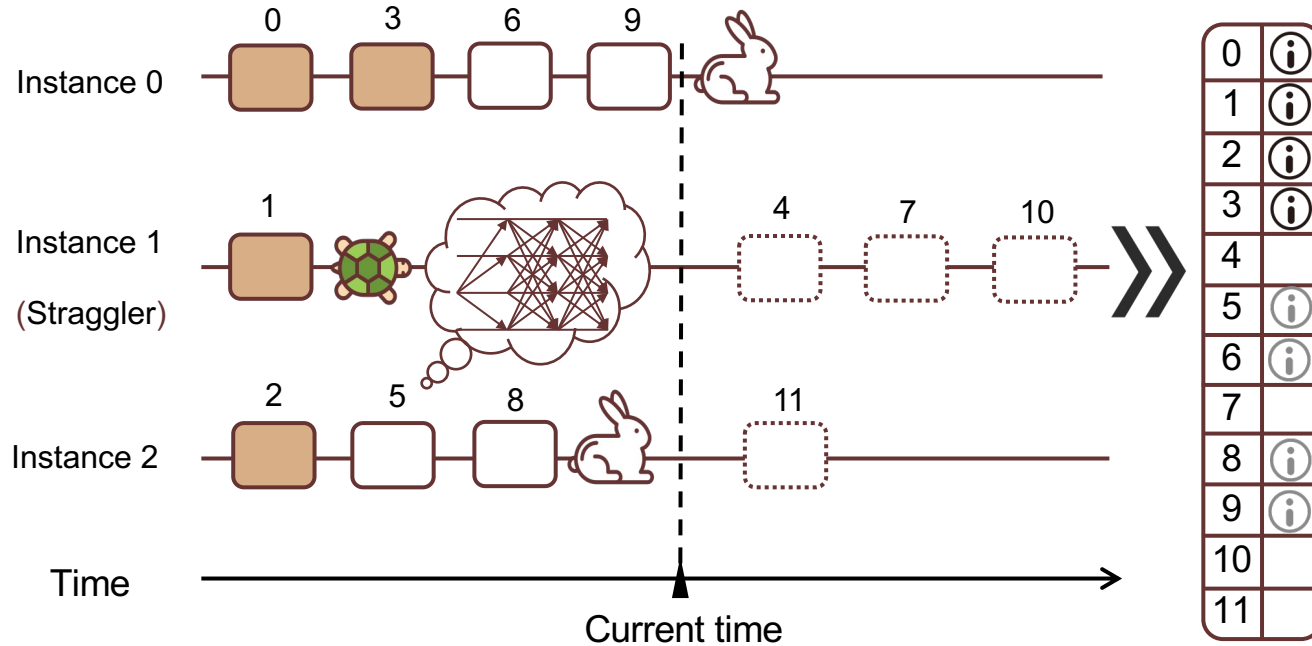
THE UNIVERSITY  
OF BRITISH COLUMBIA



SUSTech

Southern University  
of Science and  
Technology

Performance degradation of pre-determined global ordering with stragglers





# Key Observation



THE UNIVERSITY  
OF BRITISH COLUMBIA

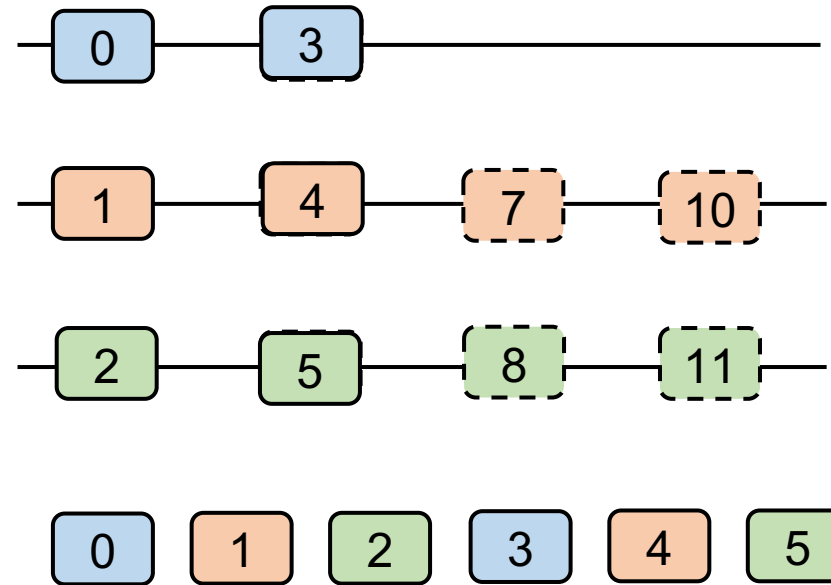


SUSTech

Southern University  
of Science and  
Technology

## Limitations of Pre-Determined Ordering

- A slow instance prevents other instances from confirming their blocks.





# Key Observation



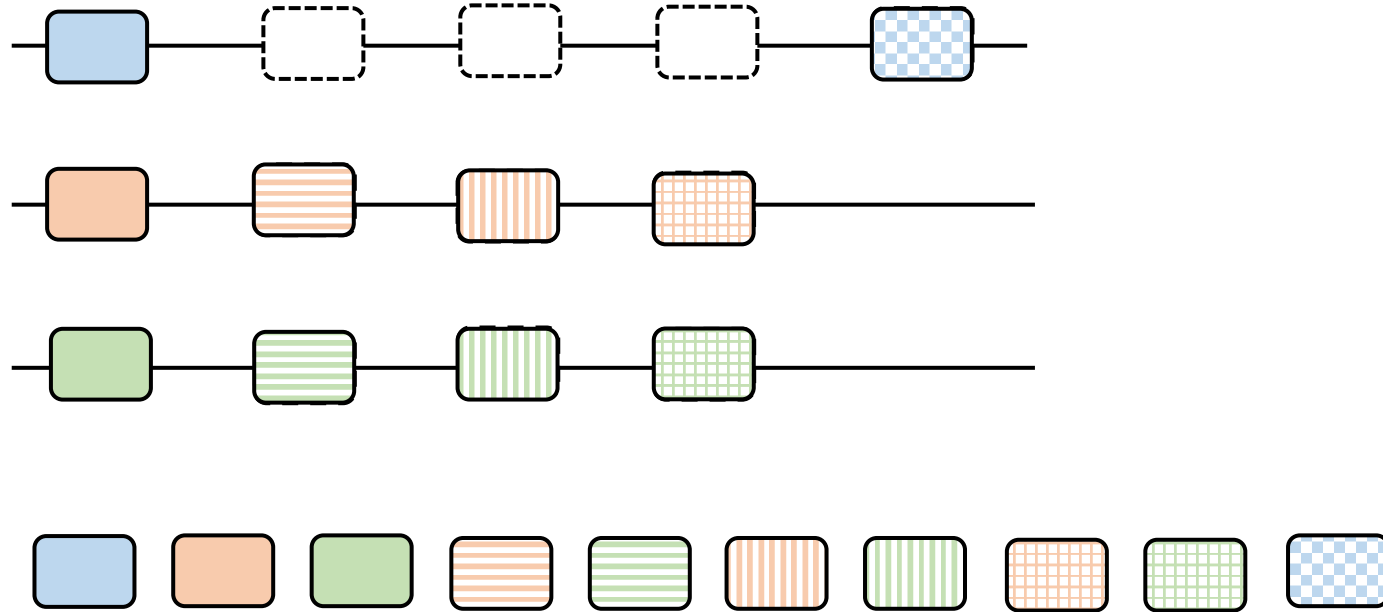
THE UNIVERSITY  
OF BRITISH COLUMBIA



SUSTech

Southern University  
of Science and  
Technology

If a new block could jump directly to the front, it would no longer obstruct global ordering.







# Dynamic Global Ordering



THE UNIVERSITY  
OF BRITISH COLUMBIA



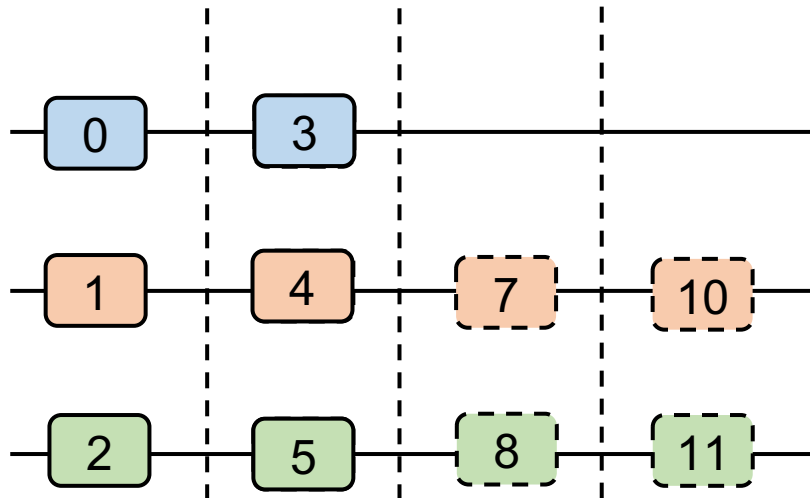
SUSTech

Southern University  
of Science and  
Technology

We use a dynamically generated parameter *rank* to mark the position of a new block.

- The *rank* of a new block = the highest *rank* + 1

## Pre-determined global ordering



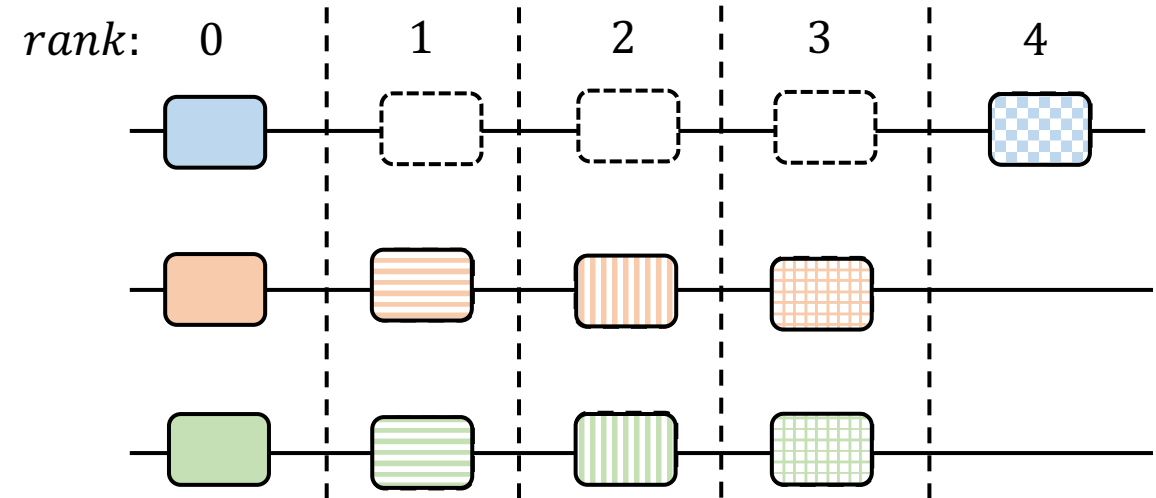
Pre-determined global ordering



Dynamic global ordering



## Dynamic global ordering



the highest *rank*  
=  $\max(0, 3, 3) = 3$



We use a dynamically generated parameter *rank* to mark the position of a new block.

- The *rank* of a new block = the highest *rank* + 1

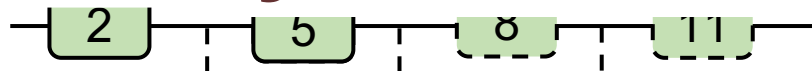
Pre-determined global ordering

Dynamic global ordering

*rank*: 0 1 2 3 4



How to generate a consistent and protocol-compliant *rank* in a distributed system with Byzantine replicas?



Pre-determined global ordering



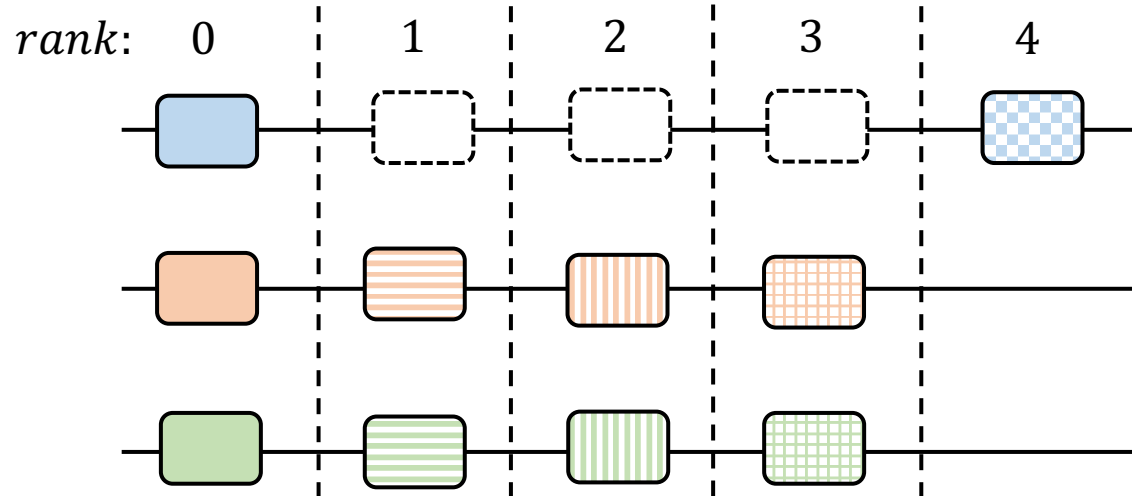
Dynamic global ordering



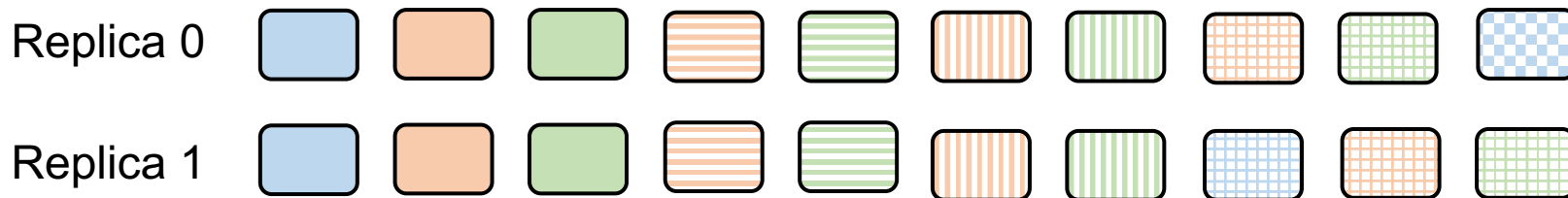
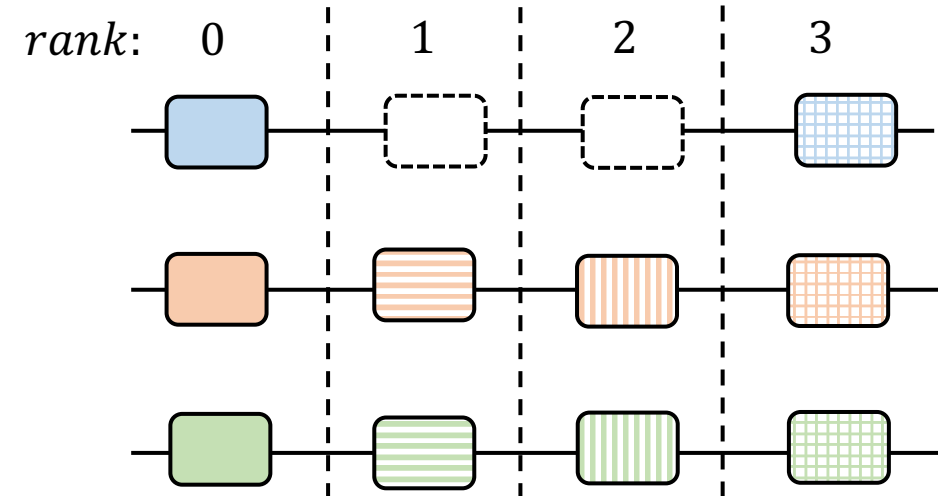


Challenge 1: Different replicas may have different views on system state

Replica 0

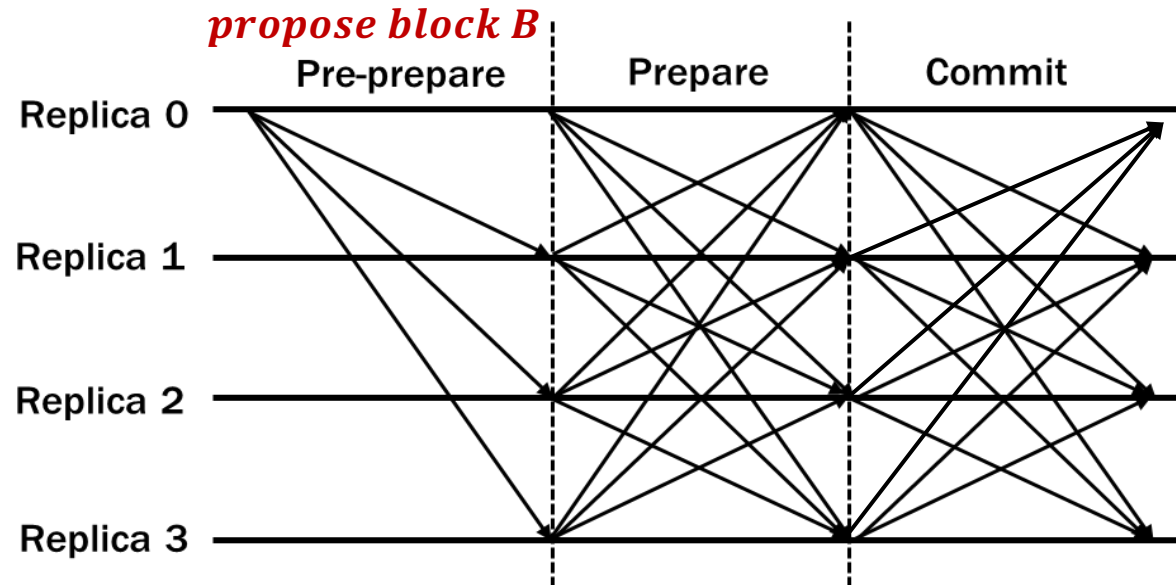


Replica 1





Solution: Make rank as a parameter of the block when the leader proposes



$B = (txs, index, round, rank)$

*txs*: a batch of clients' transactions

*index*: the index of the consensus instance

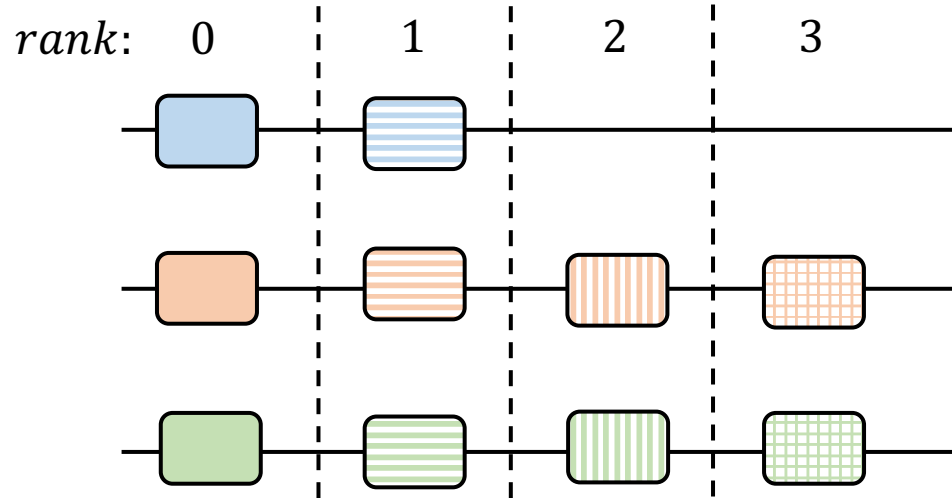
*round*: the consensus round

**Agreement:** All honest replicas have the same rank for a partially committed block.

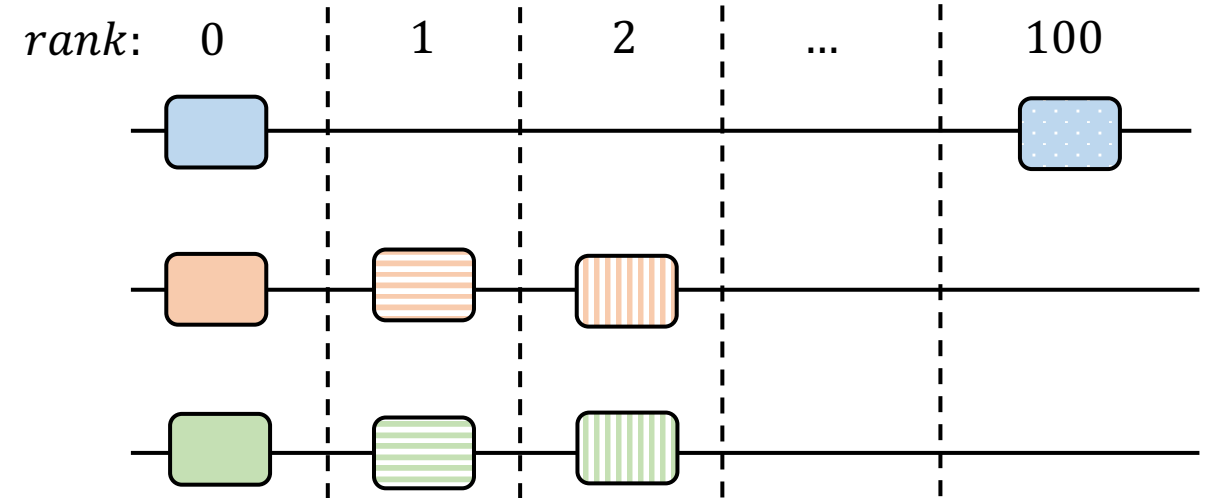


## Challenge 2: A malicious leader might generate a fake rank

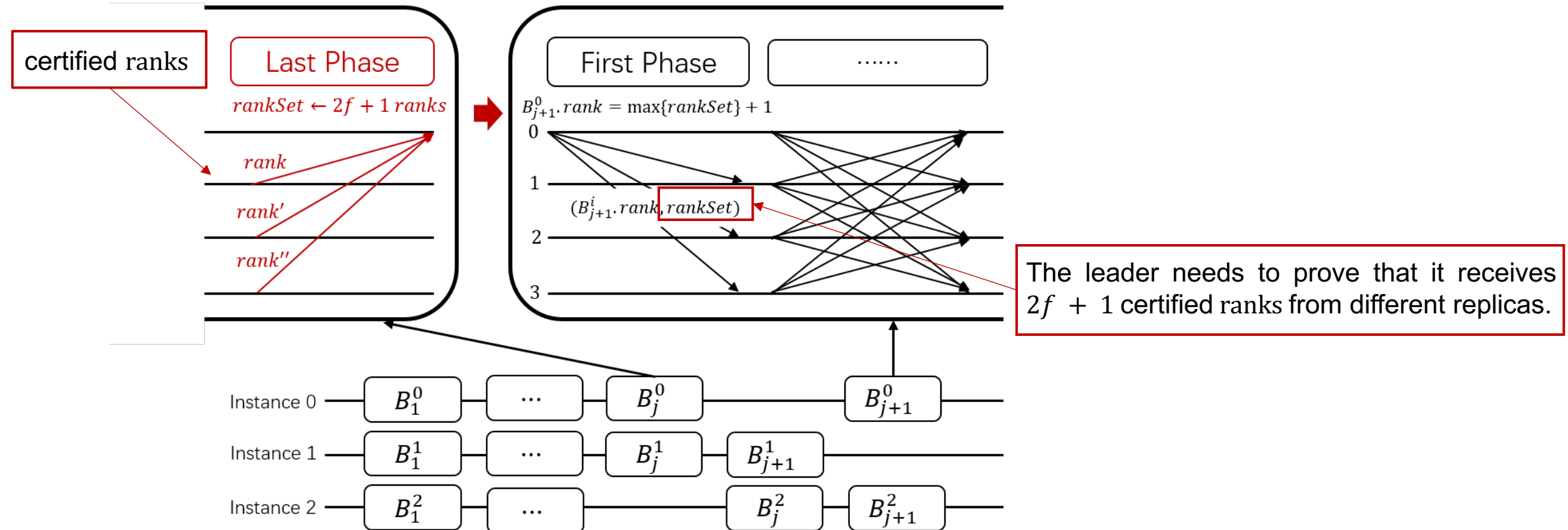
### Case 1



### Case 2



Solution: A leader collects the highest ranks from at least  $2f + 1$  replicas.



**Monotonicity:** If a block  $B'$  is generated after an intra-instance (or a partially committed inter-instance) block  $B$ , then the rank of  $B'$  is larger than the rank of  $B$ .



# Ladon Overview

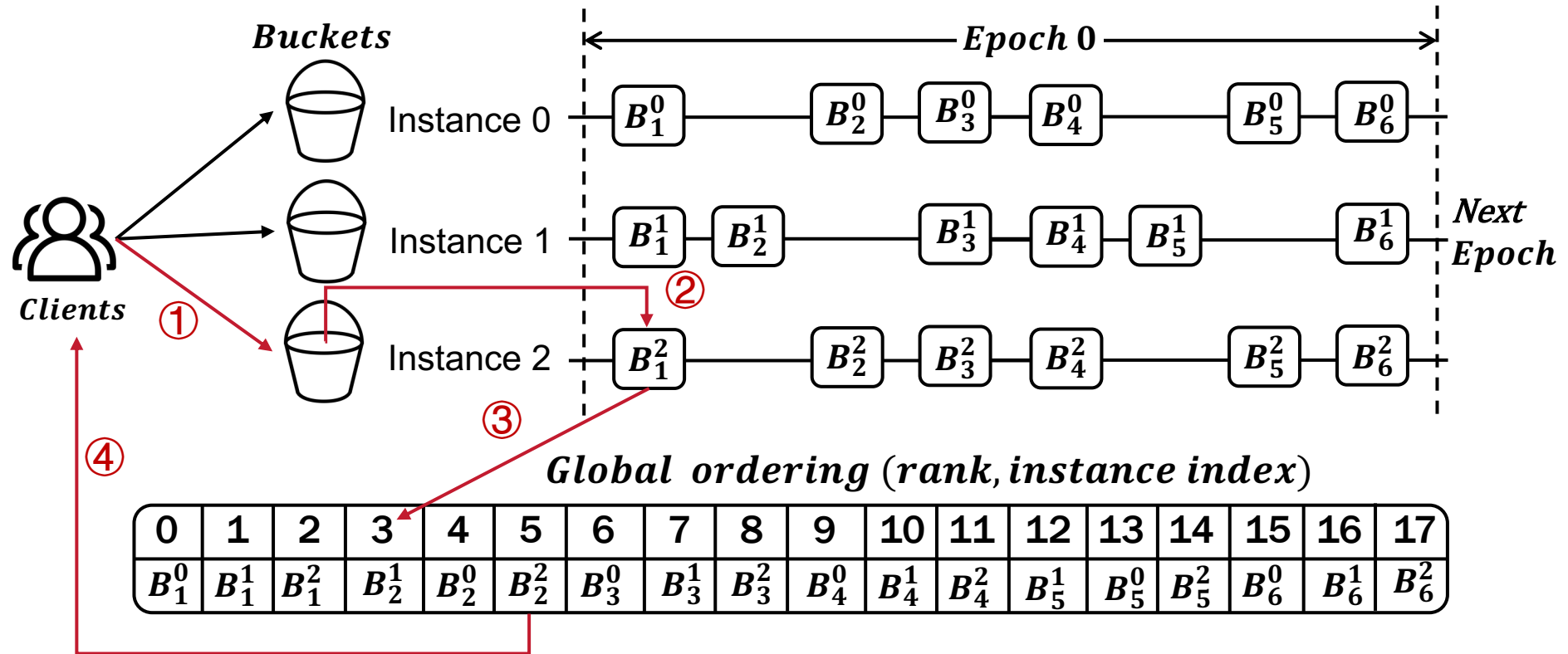


THE UNIVERSITY  
OF BRITISH COLUMBIA



SUSTech

Southern University  
of Science and  
Technology





## ***Testbeds***

- AWS EC2 c5a.2xlarge machines
- 8vCPUs and 16GB RAM Ubuntu Linux 22.04
- Deployed both in LAN and WAN
- 1Gbps bandwidths

## ***Research questions:***

- How does Ladon perform with **varying replicas** as compared to the Baseline protocols?
- How does Ladon perform under **crash and Byzantine** faults?

## ***Baseline protocols :***

- ISS<sup>[1]</sup>: PBFT/HotStuff
- Mir-BFT<sup>[2]</sup>: PBFT
- RCC<sup>[3]</sup>: PBFT
- **DQBFT<sup>[4]</sup>: PBFT**

[1] Chrysoula Stathakopoulou, Matej Pavlovic, and Marko Vukolić. State Machine Replication Scalability Made Simple, in ACM EuroSys'22.

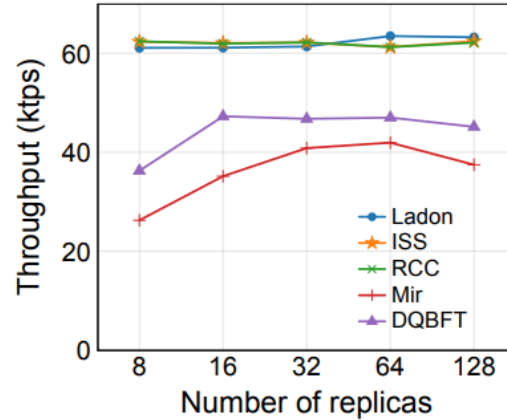
[2] Chrysoula Stathakopoulou, Tudor David, and Marko Vukolic. Mir BFT: Scalable and Robust BFT for Decentralized Networks, in Jsyst'22.

[3] Suyash Gupta, Jelle Hellings, and Mohammad Sadoghi. RCC: Resilient Concurrent Consensus for High-Throughput Secure Transaction Processing, in IEEE ICDE'21.

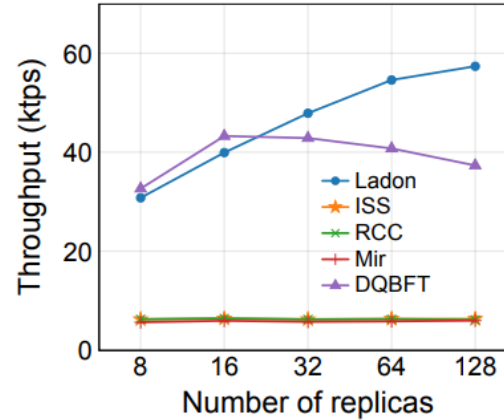
[4] Balaji Arun and Binoy Ravindran. Scalable Byzantine Fault Tolerance via Partial Decentralization, in PVLDB'22.



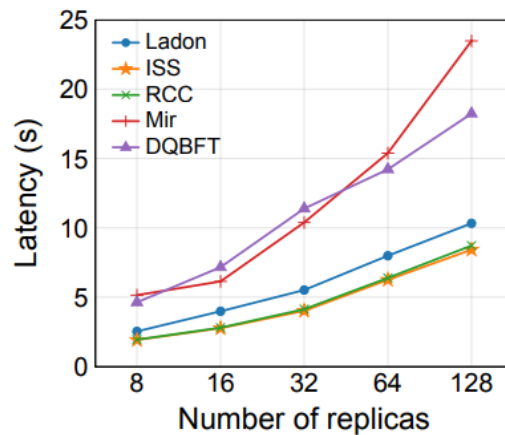
## System scalability with varying number of replicas in WAN



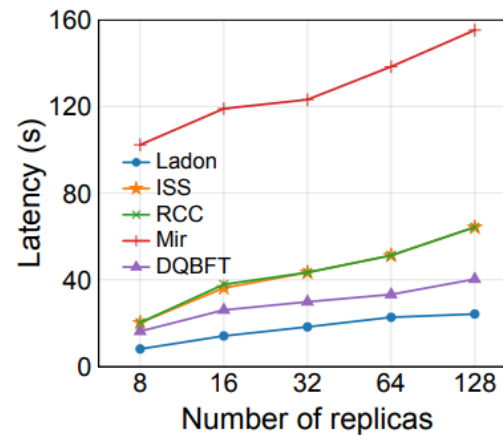
(a) #Straggler=0, WAN



(b) #Straggler=1, WAN



(c) #Straggler=0, WAN

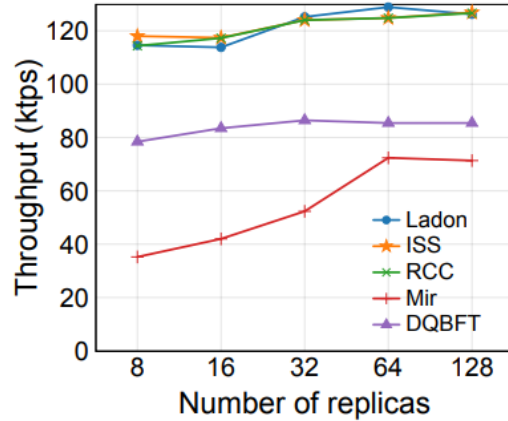


(d) #Straggler=1, WAN

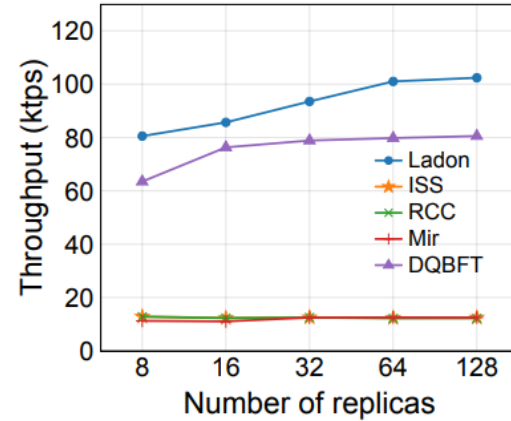
- 8/16/32/64/128 total replicas
- 0/1 stragglers
- WAN
- Without stragglers
  - Comparable throughput/latency
- With one straggler
  - Superior throughput (9X)
  - Lower latency



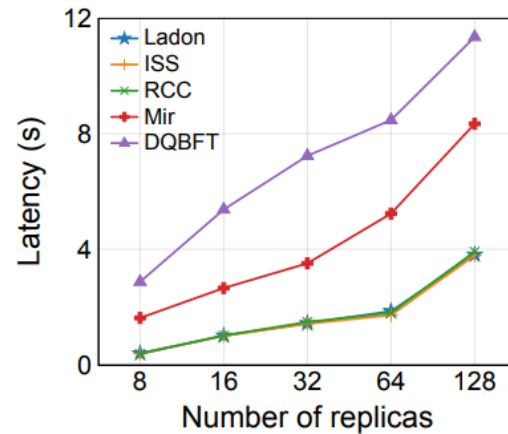
## System scalability with varying number of replicas in LAN



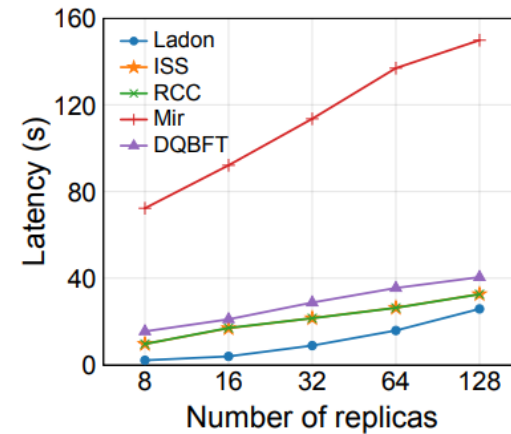
(e) #Straggler=0, LAN



(f) #Straggler=1, LAN



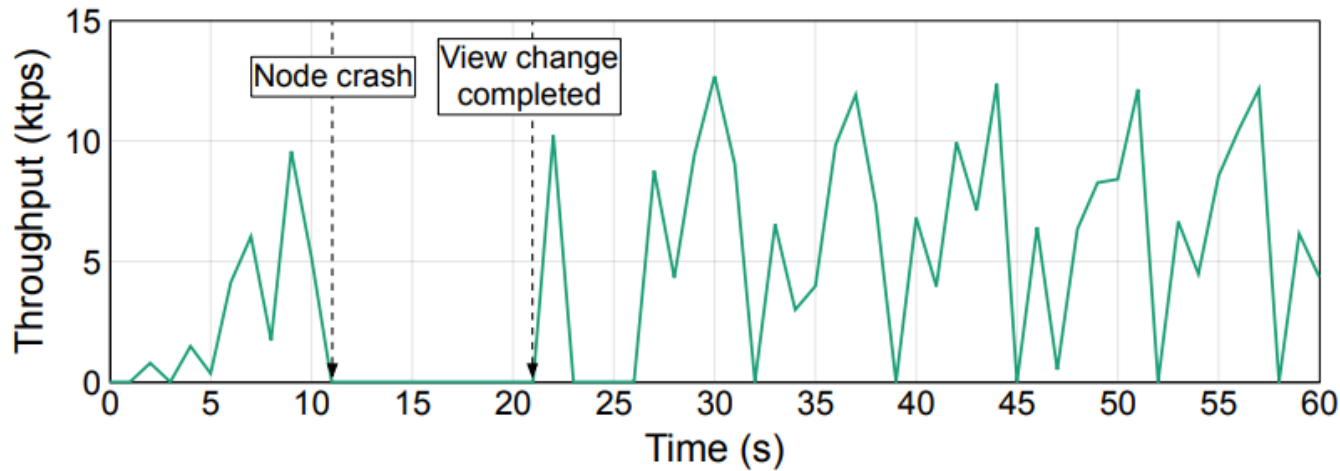
(g) #Straggler=0, LAN



(h) #Straggler=1, LAN

- 8/16/32/64/128 total replicas
- 0/1 stragglers
- LAN
- Similar performance trends to WAN

## Impact of crash faults

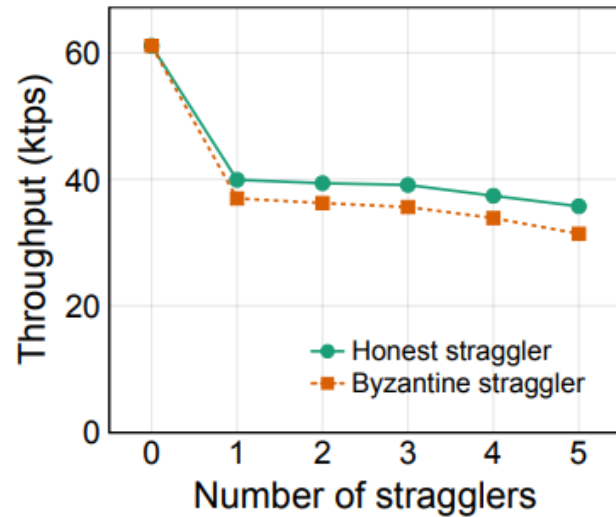


Ladon's throughput average (over 1s intervals) over time with one crash fault.

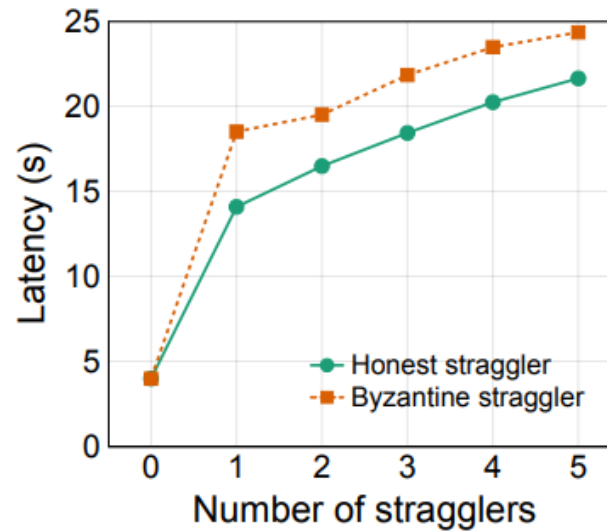
- 16 total replicas with 1 Crash leader
- View change timeout sizes: 10s
- Temporary stall when view change
- Recovers after view change



## Impact of Byzantine faults



(a) Throughput



(b) Latency

- 16 total replicas with 1-5 stragglers
- Byzantine/honest stragglers
- Slight difference
- Performance degradation is gradual



**Ladon:** A high-performance Multi-BFT protocol that mitigates the impact of stragglers

- ❑ Code: <https://github.com/eurosys2024ladon/ladon>
- ❑ Adopting **dynamic global ordering** to assign a global ordering index to blocks according to the **real-time states** of all instances.
- ❑ Using **monotonic ranks** and **pipeline** their distribution and collection with the consensus process.



THE UNIVERSITY  
OF BRITISH COLUMBIA



**SUSTech**

Southern University  
of Science and  
Technology

# Thanks for Listening!

## Hanzheng Lyu

---

BLOCKCHAIN

