# Leveraging Network Hardware in Distributed Systems Design

Jialin Li

NUS | Computing
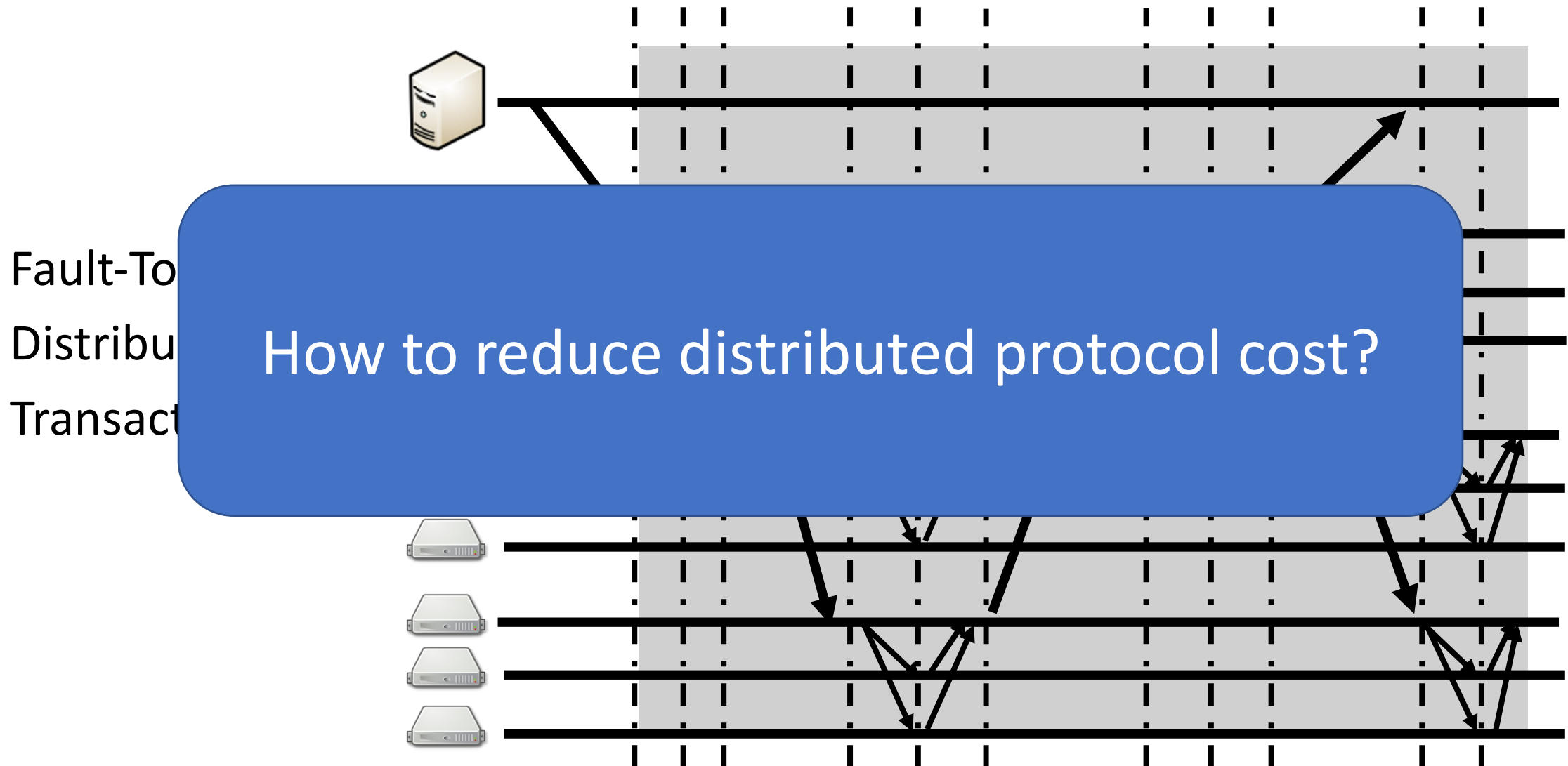National University of Singapore

# Distributed Systems Challenges in Data Centers

- Increasing network speed

- Stalling CPU speedups

- Strict user Service-Level Agreement
  - *us-scale* tail latency requirement
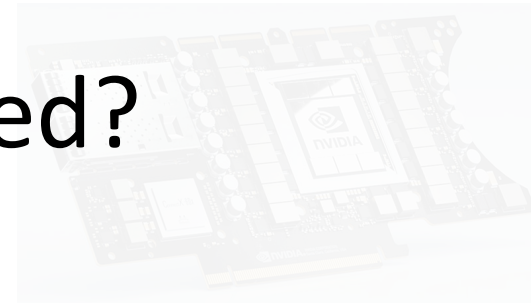
# Complex protocol communication patterns!

Fault-To

Distribu

Transact

How to reduce distributed protocol cost?

# Protocol offloading to the **network**?

In-network computation at line rate

Resource constrained devices

What should be offloaded?

software defined networking

Barefoot Tofino 2-powered 12.8 Tbps 32xQSFP56-DD System

reconfigurable switches

DPUs, SmartNICs

# Theme: *partial protocol offloading to the network* for distributed systems

- Use programmability in the network to offload **simple protocol primitives**

- **Efficient** network implementation

- **Co-design** distributed protocols and the network

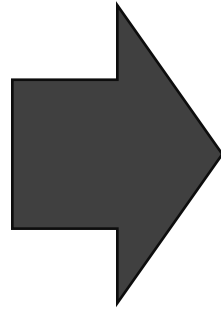- Result: practical distributed systems with both **strong guarantees** and **high performance**

# Serialization-Free Network Ordering for Strongly Consistent Distributed Applications
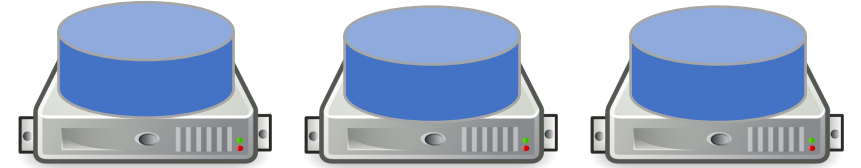
Scalability by *Sharding*
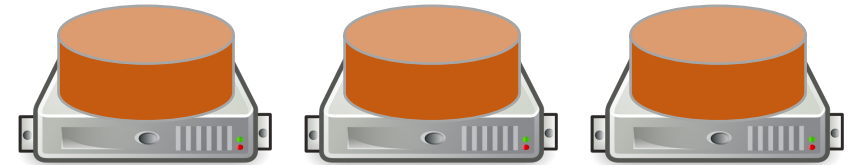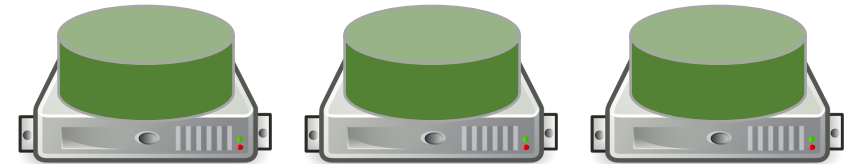Fault Tolerance by *Replication*
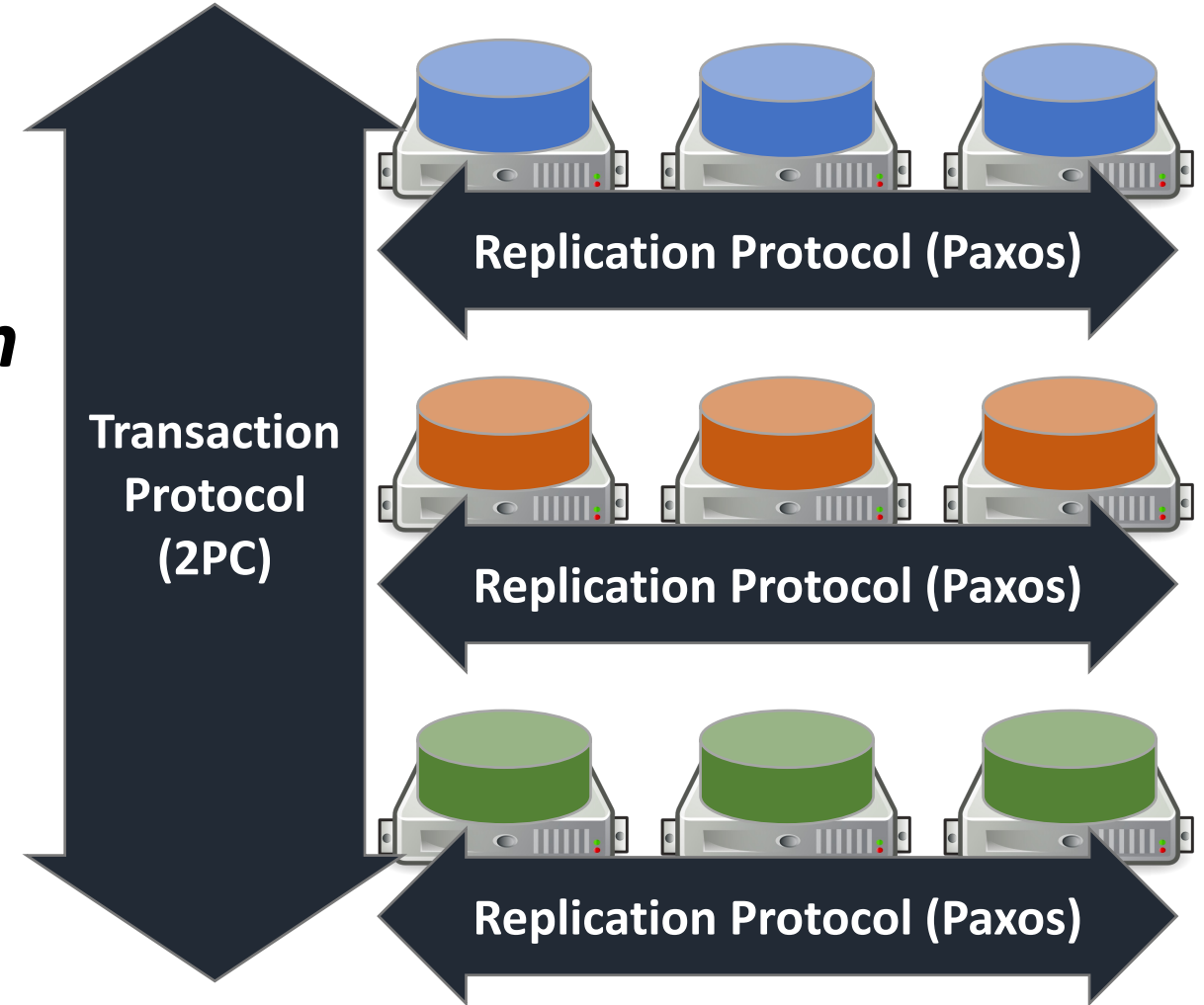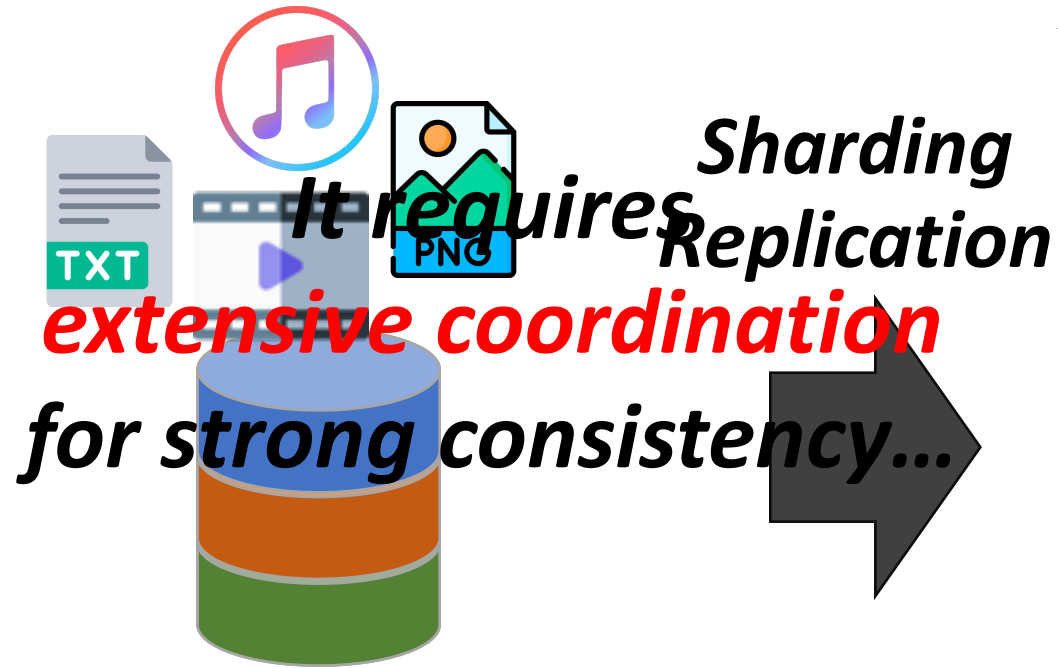
*Sharding*
*Replication*

Shard 1

Shard 2

Shard 3

# Scalability by *Sharding*
# Fault Tolerance by *Replication*



*It requires*
**extensive coordination**
*for strong consistency...*

*Sharding*
*Replication*

Transaction Protocol (2PC)

Replication Protocol (Paxos)

Replication Protocol (Paxos)

Replication Protocol (Paxos)

# *Network Ordering*
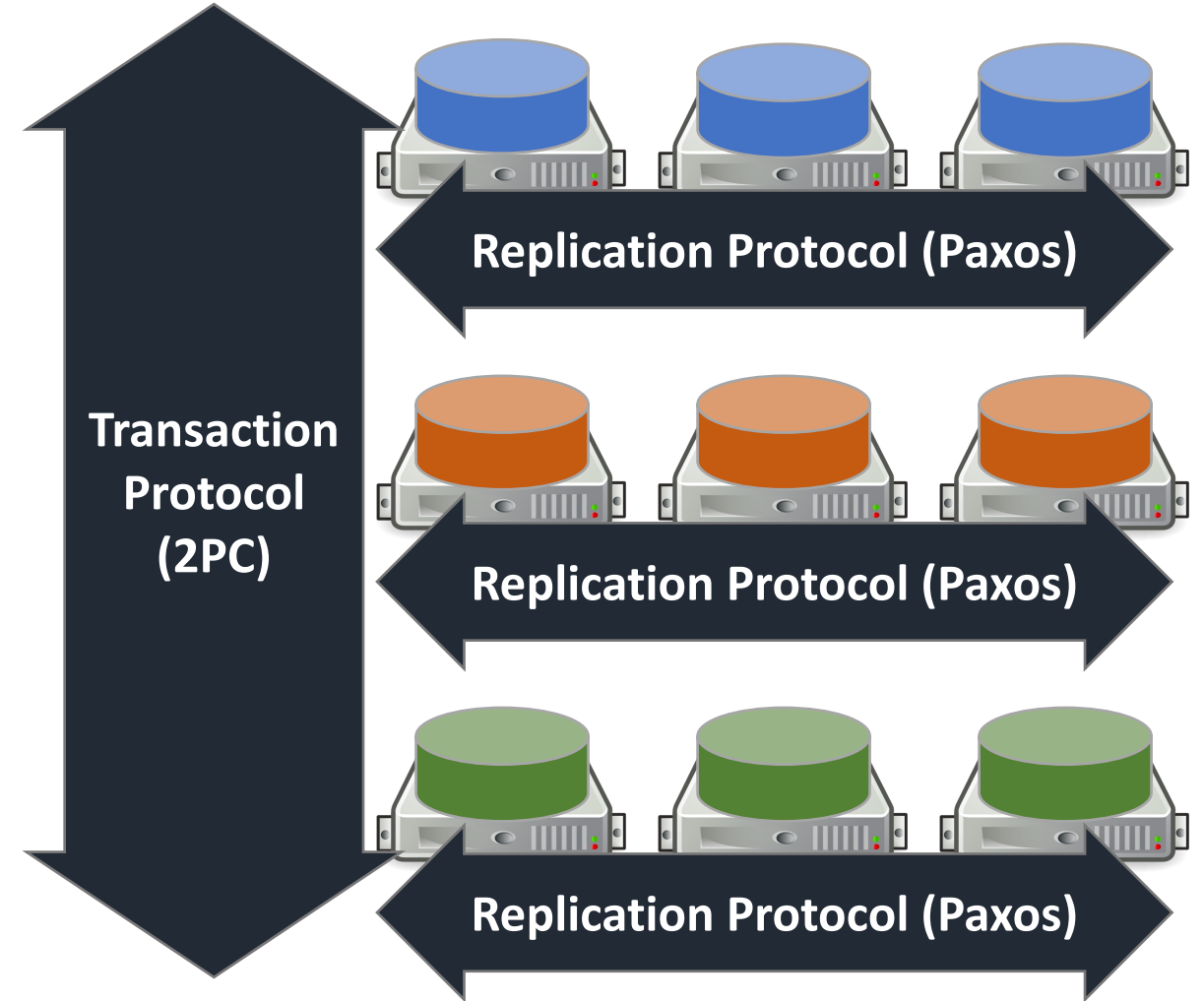## to Eliminate Coordination

**Network**

**Sequencer**
*(e.g., Programmable Switch)*



**Consistent Ordering**

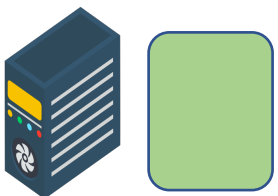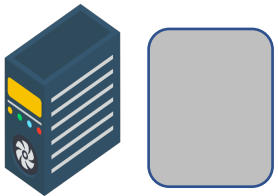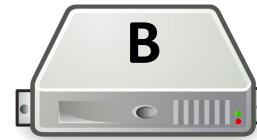Related Works:
**NOPaxos** [OSDI '16], **Eris** [SOSP '17]

**Transaction Protocol (2PC)**

**Replication Protocol (Paxos)**

**Replication Protocol (Paxos)**

**Replication Protocol (Paxos)**

# *Network Ordering*

# Network Ordering

**Senders**

**Receivers**



**Sequencer**

Counter
0

A

B

C

# Network Ordering

**Senders**

**Receivers**

**Sequencer**

Counter
1

**1**

A

B

C

# Network Ordering

**Senders**

**Receivers**

**Sequencer**

*Multicast!*

**1**

Counter
1

A

B

C

# Network Ordering

**Senders**

**Receivers**

**Sequencer**

A  1

2

Counter
2

B  1

C  1

# Network Ordering

**Senders**

**Receivers**

**Sequencer**

*Multicast!*

Counter
2

A  1

B  1

C  1

2

# Network Ordering

**Senders**

**Receivers**

**Sequencer**



Counter
3

A  1 2

B  1 2

C  1 2

# Network Ordering

**Senders**

**Receivers**



**Guarantees**

**Consistent Ordering**

- **Partial ordering across shards**
- **Total ordering across replicas**

**Sequencer**

**Multicast!**

**3**

Counter
3

# Network Ordering

**Guarantees**

Consistent Ordering

Drop Detection

**Sequencer**

**Receivers**

A  1 2 3

B  1 2 3

C  1 2 3

# Drawbacks due to a Single Sequencer



**Sequencer**

*All request traffic*
**must go through the *single sequencer***
**(*Network Serialization*)**

☒ ***Network load imbalance ⇒ high latency***
☒ ***Sequencer scalability bottleneck***
☒ ***Prolonged sequencer failover***

# Drawbacks due to a Single Sequencer

**Sequencer**

**All request traffic
must go through the *single sequencer*
(Network Serialization)**

⬇

☒ **Network load imbalance ⇒ high latency**
☒ **Sequencer scalability bottleneck**
☒ **Prolonged sequencer failover**

# Drawbacks due to a Single Sequencer

**Sequencer**

*All request traffic*
**must go through the** *single sequencer*
**(Network Serialization)**

☒ *Network load imbalance ⇒ high latency*
☒ *Sequencer scalability bottleneck*
☒ *Prolonged sequencer failover*

Well, let's try multiple sequencers

# Multi-Sequencer Challenge:

**Guarantees**

**Consistent Ordering**

**Drop Detection**

*Sequencer* ①

Counter
0

*Sequencer* ②

Counter
0

A

B

C

# Multi-Sequencer Challenge:

**Guarantees**

**Consistent Ordering**

**Drop Detection**

*Sequencer* ①

*Multicast!*

**1**

Counter
1

*Sequencer* ②

Counter
0

**A**    **1**

**B**    **1**

**C**    **1**

# Multi-Sequencer Challenge:

**Guarantees**

Consistent Ordering

Drop Detection

**Sequencer ①**

Counter 1

**Sequencer ②**

Counter 0

A 1

B 1

C 1

# Multi-Sequencer Challenge:

**Guarantees**

~~**Consistent Ordering**~~

**Drop Detection**

*Sequencer* ① 

Counter 1

*Sequencer* ② 

*Multicast!*

**1**

Counter 1

A **1** **1** **?**

B **1** **1** **?**

C **1** **1** **?**

# Multi-Sequencer Challenge:

**Sequencer ①**

*Multicast!*

A    **1**  **1**  **2**

*Guarantees*

**2**

Counter
1

Naively using multiple sequencers does not work!

C    **1**  **1**  **2**

# *Solution*: Combine sequence number with *physical clock*



| Sequence Number | → | Drop Detection |

| Physical Clock Timestamp | → | Consistent Ordering |

- Consistent across receivers
- Strictly monotonically increasing
- Physical clocks can be (loosely) synchronized across sequencers

# *Solution*: Combine sequence number with *physical clock*

# Hydra Network Primitive - Sequencers

# Hydra Network Primitive - Sequencers

# Hydra Network Primitive - Sequencers

# Hydra Network Primitive - Sequencers

# Hydra Network Primitive - Sequencers

# Hydra Network Primitive - Sequencers

# Hydra Network Primitive - Receivers

**Message Buffer**

Received

**Q: When can a message be delivered?**

Ordering

⏱ *12*
1

⏱ *19*
1

- **Delivered to the application layer**

# A: Once **ALL** messages with lower timestamp have been received

# Deliver messages up to
# *the minimum barrier*

# Other Hydra designs

- Flush messages to ensure progress
    - Receiver-side solicitation
    - In-network aggregation
- Adding/removing sequencers
- Congestion-aware routing

# Scales beyond a single sequencer

**Throughput (txn/s)**

Serialization approach (Eris) is limited by network capacity of the single sequencer

Higher throughput by reducing server coordination

| | | |
|---|---|---|
| 4M | | |
| 3M | | |
| 2M | | |
| 1M | | |
| 0M | | |

Lock-Store (Spanner)    Granola    Eris

# Applying Network Programmability to BFT Protocols

# BFT SMR protocols

- Systems today face sophisticated failures
  - Adversaries, malicious participants

- Byzantine fault tolerance protocols

- Permissioned deployment in data centers
  - High-performance applications
  - Low latency requirement

# Challenges of Applying In-Network Ordering to BFT

- Adversaries can generate **conflicting** message order

# Challenges of Applying In-Network Ordering to BFT

- Adversaries can generate **conflicting** message order

# Challenges of Applying In-Network Ordering to BFT

- Adversaries can generate **conflicting** message order



Sequencer

# Challenges of Applying In-Network Ordering to BFT

- Adversaries can generate **conflicting** message order



Sequencer

# Challenges of Applying In-Network Ordering to BFT

- Adversaries can generate **conflicting** message order

# New Approach: Network Ordering with Authentication

- Can replicas determine the "correct" message order?
  - Yes, if messages from network primitives can be **authenticated**
- Solution: sequencer switch **signs** messages
- Replicas independently **verify** message signatures
  - Messages are indeed generated by sequencer

# Authenticated Ordered Multicast Illustration

# Authenticated Ordered Multicast Illustration

- Ordered: receivers receive AOM messages in the same order

# Authenticated Ordered Multicast Illustration

- Ordered: receivers receive AOM messages in the same order

# Authenticated Ordered Multicast Illustration

- Authentication: receiver can verify that messages are correctly processed by AOM sequencer

# Authenticated Ordered Multicast Illustration

- Authentication: receiver can verify that messages are correctly processed by AOM sequencer

# Authenticated Ordered Multicast Illustration

- Authentication: receiver can verify that messages are correctly processed by AOM sequencer

# Authenticated Ordered Multicast Illustration

- Authentication: receiver can verify that messages are correctly processed by AOM sequencer

# Implementing In-Network Authentication Is Challenging

**Cryptographic Algorithms**

- Complex computation

- Involves large prime numbers

- Unbounded loops and large vectors

**Programmable Data Plane**

- Limited computation support

- Restricted to small, fixed-width numbers

- Small number of pipeline stages

- Highly resource constrained

# Implementing In-Network Authentication Is Challenging

**Cryptographic Algorithms**

- Complex computation

- Involves large prime numbers

- Unbounded loops and large vectors

**Programmable Data Plane**

- Limited computation support

- Restricted to small, fixed-width numbers

- Small number of pipeline stages

- Highly resource constrained

Lightweight Message Authentication Code (MAC)

# MAC on a Switch

- Dedicated pipeline for computing MAC vector

- Unrolled HalfSipHash implementation

- 4 parallel MAC generation instances

# MAC on a Switch

- Dedicated pipeline for computing MAC vector

- Unrolled HalfSipHash implementation

- 4 parallel MAC generation instances

# MAC on a Switch

- Dedicated pipeline for computing MAC vector
- Unrolled HalfSipHash implementation
- 4 parallel MAC generation instances

# MAC on a Switch

- Dedicated pipeline for computing MAC vector
- Unrolled HalfSipHash implementation
- 4 parallel MAC generation instances

# MAC on a Switch

- Dedicated pipeline for computing MAC vector

- Unrolled HalfSipHash implementation

- 4 parallel MAC generation instances

# MAC on a Switch

- Dedicated pipeline for computing MAC vector

- Unrolled HalfSipHash implementation

- 4 parallel MAC generation instances

# MAC on a Switch

- Dedicated pipeline for computing MAC vector

- Unrolled HalfSipHash implementation

- 4 parallel MAC generation instances

# Another Authentication Implementation: Public-Key Cryptography

- Problem with MAC vector: poor scalability
  - One MAC per receiver
- Public-key cryptography offers "*infinite*" scalability
  - Too complex to implement directly in the switch data plane
- New switch architecture that integrates a **cryptographic coprocessor**
- FPGA-based implementation in our prototype

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA



Switch

FPGA

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA



Switch

FPGA

# FPGA-based Public-Key Cryptography Prototype

- SECP256K1 signature implemented using FPGA



Switch

FPGA

- Message Authentication Code
  - Up to 76.24Mpps
- Public-Key Cryptography
  - 1.11Mpps

# NeoBFT: Normal Operation

- Client commit in 1 RTT
- No coordination among replicas
- Can tolerate $f$ faulty replicas



Speculative commit and reply

# NeoBFT Achieves both **Lower Latency** and **Higher Throughput**

# NeoBFT Achieves both **Lower Latency** and **Higher Throughput**

# NeoBFT Achieves both **Lower Latency** and **Higher Throughput**



Legend:
- NeoBFT (MAC)
- NeoBFT (public key)
- NeoBFT (untrusted network)
- Zyzzyva
- Zyzzyva (with faulty node)
- PBFT
- HotStuff

1.8x - 3.1x throughput improvement

# NeoBFT Achieves both **Lower Latency** and **Higher Throughput**



8.5x – 42x latency improvement

Legend:
- NeoBFT (MAC)
- NeoBFT (public key)
- NeoBFT (untrusted network)
- Zyzzyva
- Zyzzyva (with faulty node)
- PBFT
- HotStuff

Latency (us) — vertical axis: 0, 500, 1000, 1500, 2000

Throughput (Kop/sec) — horizontal axis: 0, 50, 100, 150, 200, 250, 300

# NeoBFT Achieves both **Lower Latency** and **Higher Throughput**



Still better performance when the network is adversarial

- NeoBFT (MAC)
- NeoBFT (public key)
- NeoBFT (untrusted network)
- Zyzzyva
- Zyzzyva (with faulty node)
- PBFT
- HotStuff

Latency (us)

Throughput (Kop/sec)

# Leveraging Network Programmability for Load Balancing

# Challenge: Unpredictable Server Overloads



**Bursty Traffic**

**Unpredictability**

**Load Imbalance**

# Existing Load Balancing Approaches

**L7 per-request frontend proxy**

- NGINX

☒ *Scalability*

Clients

···

Backend          Frontend          Backend

# Existing Load Balancing Approaches



**L7 per-request frontend proxy**

- NGINX

☒ *Scalability*

**Load balancing by switch**

**Per-packet distribution**

- SwitchKV [NSDI '16]
- NetCache [SOSP '17]
- Pegasus [OSDI '20]

☒ *TCP support*

Clients

Switch

Backend  Frontend  Backend

# Existing Load Balancing Approaches

L7 per-request frontend proxy

What if servers can migrate
live TCP connections?

states for
consistency (PCC)

Clients
...

Switch

☒ *TCP support*

☒ *Balancing skewed workloads*

...

# Existing Load Balancing Approaches



L7 per-request frontend proxy

What if servers can migrate
live TCP connections?

→ Strong load balancing systems with
✅ *Scalability*
✅ *TCP support*
✅ *Balancing skewed workloads*

☒ *TCP support*        ☒ *Balancing skewed workloads*

Clients
...

Switch

...

# Disruptive or slow migration can make things even worse



**Disruptive (Dropping) Migration**

**Slow Migration**

# *Capybara*

**Design goal: $\mu$s-scale-fast and <u>client-transparent</u> TCP migration**

*"without disconnection or packet drops"*

**User Space**

**Customized protocols**

| Application | ⟷ | Kernel-bypass Library OS |

**No context switch**

**Kernel Space**

**NIC**

RX/TX Queues

Kernel-bypassing OS

**Two-phase protocol**
- Phase 1: Migration handshake to buffer the connection
- Phase 2: Connection state transfer

-aware rwarding

ding ion

Programmable Switch

✔ **Transparently migrate a live TCP connection within 15 $\mu$s**

# Server Architecture

# Server Architecture

Implements TCP migration protocol in Demikernel LibOS [SOSP '21]

**TCP**

- Maintains TCP stats (e.g., per-connection request rate)
- TCP state management

**TCPMig (Layer 4)**

- Migration protocol implementation
- Process migration msgs (PREPARE_MIG, etc.)
- Transient packet buffering



User Space

Application

LibOS    TCP    TCPMig

DPDK (rte_mempool & PMD)

HW    NIC

**Capybara Server**

**TCP**

**TCPMig**

# Switch Architecture

1. **Migration-aware packet forwarding**

| Migration Directory | | | Workload | |
|---|---|---|---|---|
| Client | Origin | Target | Total | 0 |
| | | | S0 | 0 |
| | | | S1 | 0 |

**TCP** **TCPMig** S0 S1 **TCP** **TCPMig**

**Workflow:**
1. Monitor workload

C0 ... Cn

| Migration Directory | | |
|---|---|---|
| Client | Origin | Target |
| | | |
| | | |

| Workload | |
|---|---|
| Total | 0 |
| S0 | 0 |
| S1 | 0 |

**TCP**

| Workload | |
|---|---|
| C0 | 0 |
| C2 | 0 |
| ... | ... |

**TCPMig**

S0

S1

**TCP**

| Workload | |
|---|---|
| C1 | 0 |
| C3 | 0 |
| ... | ... |

**TCPMig**

**Workflow:**
1. Monitor workload

**C0** ... **Cn**

| Migration Directory | | | Workload | |
|---|---|---|---|---|
| Client | Origin | Target | Total | 50 |
| | | | S0 | 25 |
| | | | S1 | 25 |

**TCP**

| Workload | |
|---|---|
| C0 | 10 |
| C2 | 5 |
| ... | ... |

**TCPMig**

Threshold: 60%
↓
Ok!

**WORKLOAD**
Total: 50
S0: 25

**S0**

**WORKLOAD**
Total: 50
S1: 25

**S1**

**TCP**

| Workload | |
|---|---|
| C1 | 5 |
| C3 | 15 |
| ... | ... |

**TCPMig**

**Workflow:**
1. Monitor workload

C0 ... Cn

| Migration Directory | | |
|---|---|---|
| Client | Origin | Target |
| | | |
| | | |

| Workload | |
|---|---|
| Total | 100 |
| S0 | 70 |
| S1 | 30 |

**TCP**

| Workload | |
|---|---|
| C0 | 30 |
| C2 | 10 |
| ... | ... |

**TCPMig**

Threshold: 60%
↓
Migrate 10

**WORKLOAD**
Total: 100
S0: 70

S0

**WORKLOAD**
Total: 100
S1: 30

S1

**TCP**

| Workload | |
|---|---|
| C1 | 10 |
| C3 | 15 |
| ... | ... |

**TCPMig**

**Workflow:**
1. Monitor workload
2. Prepare migration

C0 ... Cn

**Migration Directory**

| Client | Origin | Target |
|--------|--------|--------|
|        |        |        |
|        |        |        |

**Workload**

| Total | 100 |
|-------|-----|
| S0    | 70  |
| S1    | 30  |

**TCP**

**Workload**

| C0 | 30 |
|----|----|
| C2 | 10 |
| ... | ... |

**TCPMig**

Threshold: 60%

C2

PREPARE_MIG
Origin: S0
Conn: C2

S0

PREPARE_MIG
Origin: S0
Conn: C2
Target: S1

S1

**TCP**

**Workload**

| C1 | 10 |
|----|----|
| C3 | 15 |
| ... | ... |

**TCPMig**

C2

**Workflow:**
1. Monitor workload
2. Prepare migration

C0 ... Cn

| Migration Directory | | |
|---|---|---|
| Client | Origin | Target |
|  |  |  |
|  |  |  |

| Workload | |
|---|---|
| Total | 100 |
| S0 | 70 |
| S1 | 30 |

**TCP**

| Workload | |
|---|---|
| C0 | 30 |
| C2 | 10 |
| ... | ... |

**TCPMig**

Threshold: 60%

C2

S0

**PREPARE_MIG_ACK**

Origin: S0
Conn: C2
Target: S1

S1

**TCP**

| Workload | |
|---|---|
| C1 | 10 |
| C3 | 15 |
| ... | ... |

**TCPMig**

C2

**Workflow:**
1. Monitor workload
2. Prepare migration

C0 ... Cn

| Migration Directory | | |
|---|---|---|
| Client | Origin | Target |
| C2 | S0 | S1 |
| | | |

| Workload | |
|---|---|
| Total | 100 |
| S0 | 70 |
| S1 | 30 |

**PREPARE_MIG_ACK**

**Origin: S0**
**Conn: C2**
**Target: S1**

**PREPARE_MIG_ACK**

**Origin: S0**
**Conn: C2**
**Target: S1**

**TCP**

| Workload | |
|---|---|
| C0 | 30 |
| C2 | 10 |
| ... | ... |

**TCPMig**

**Threshold: 60%**

**C2**

**S0**

**S1**

**TCP**

| Workload | |
|---|---|
| C1 | 10 |
| C3 | 15 |
| ... | ... |

**TCPMig**

**C2**

**Workflow:**
1. Monitor workload
2. Prepare migration
3. State transfer

C0  ...  Cn

| Migration Directory | | |
|---|---|---|
| Client | Origin | Target |
| C2 | S0 | S1 |
| | | |

| Workload | |
|---|---|
| Total | 100 |
| S0 | 70 |
| S1 | 30 |

**TCP**

| Workload | |
|---|---|
| C0 | 30 |
| C2 | 10 |
| ... | ... |

**TCPMig**

Threshold: 60%

**Serialize** ▶ C2

PREPARE_MIG_ACK
Origin: S0
-Conn: C2
Target: S1

S0

CONN_STATE
Origin: S0
Target: S1
C2

S1

**TCP**

| Workload | |
|---|---|
| C1 | 10 |
| C3 | 15 |
| ... | ◀ **Merge** |

**TCPMig**

*1*

C2

**Workflow:**
1. Monitor workload
2. Prepare migration
3. State transfer
4. Migration complete

C0 ... Cn

**Src: S0**
**Dst: C2**

| Migration Directory | | |
|---|---|---|
| Client | Origin | Target |
| C2 | S0 | S1 |
| | | |

| Workload | |
|---|---|
| Total | 100 |
| S0 | 60 |
| S1 | 40 |

**Src: S1**
**Dst: C2**

**TCP**

| Workload | |
|---|---|
| C0 | 30 |
| ... | ... |
| | |

**TCPMig**

Threshold: 60%

S0

S1

**TCP**

| Workload | |
|---|---|
| C1 | 10 |
| C3 | 15 |
| C2 | 10 |

**TCPMig**

# Load Balancing Benefit of Capybara



Static L4

# Conclusion

- New approach to designing distributed systems

  - Leverage programmability in data center networks for partial protocol offloading

- Co-designed distributed systems with both **strong guarantees** and **high performance**

  - Serialization-free network ordering for strongly consistent dapps

  - Authenticated network ordering for BFT protocols

  - us-scale live TCP migration for load balancing