# Avoiding Coordination with Network Ordering: NOPaxos and Eris

Ellis Michael

# Server failures are the common case in data centers

# SERVER FAILURES ARE THE COMMON CASE IN DATA CENTERS



Cloud News Daily

**Lightning Strikes Disrupt Google Data Center**

# SERVER FAILURES ARE THE COMMON CASE IN DATA CENTERS

Cloud News Daily
**Lightning Strikes Disrupt Google Data Center**

BUSINESS INSIDER
**Amazon's Cloud Crash Disaster Permanently Destroyed Many Customers' Data**

# SERVER FAILURES ARE THE COMMON CASE IN DATA CENTERS

**Cloud News Daily**

**Lightning Strikes Disrupt Google Data Center**

**Technology News**

Microsoft and Google cloud users suffer service outages

**INSIDER**

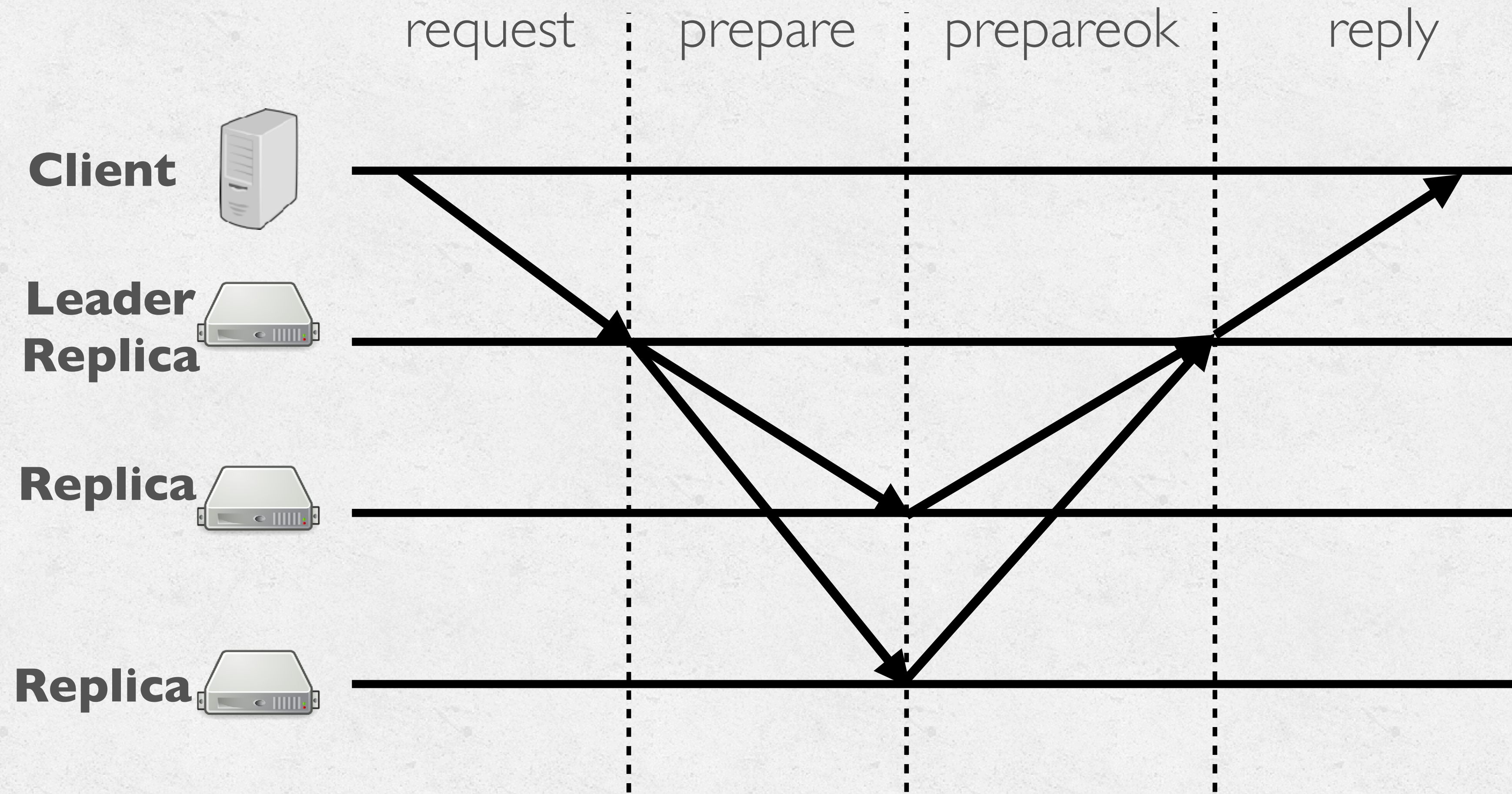**Amazon's Cloud Crash Disaster Permanently Destroyed Many Customers' Data**

# State Machine Replication

# State Machine Replication

# PAXOS FOR STATE MACHINE REPLICATION
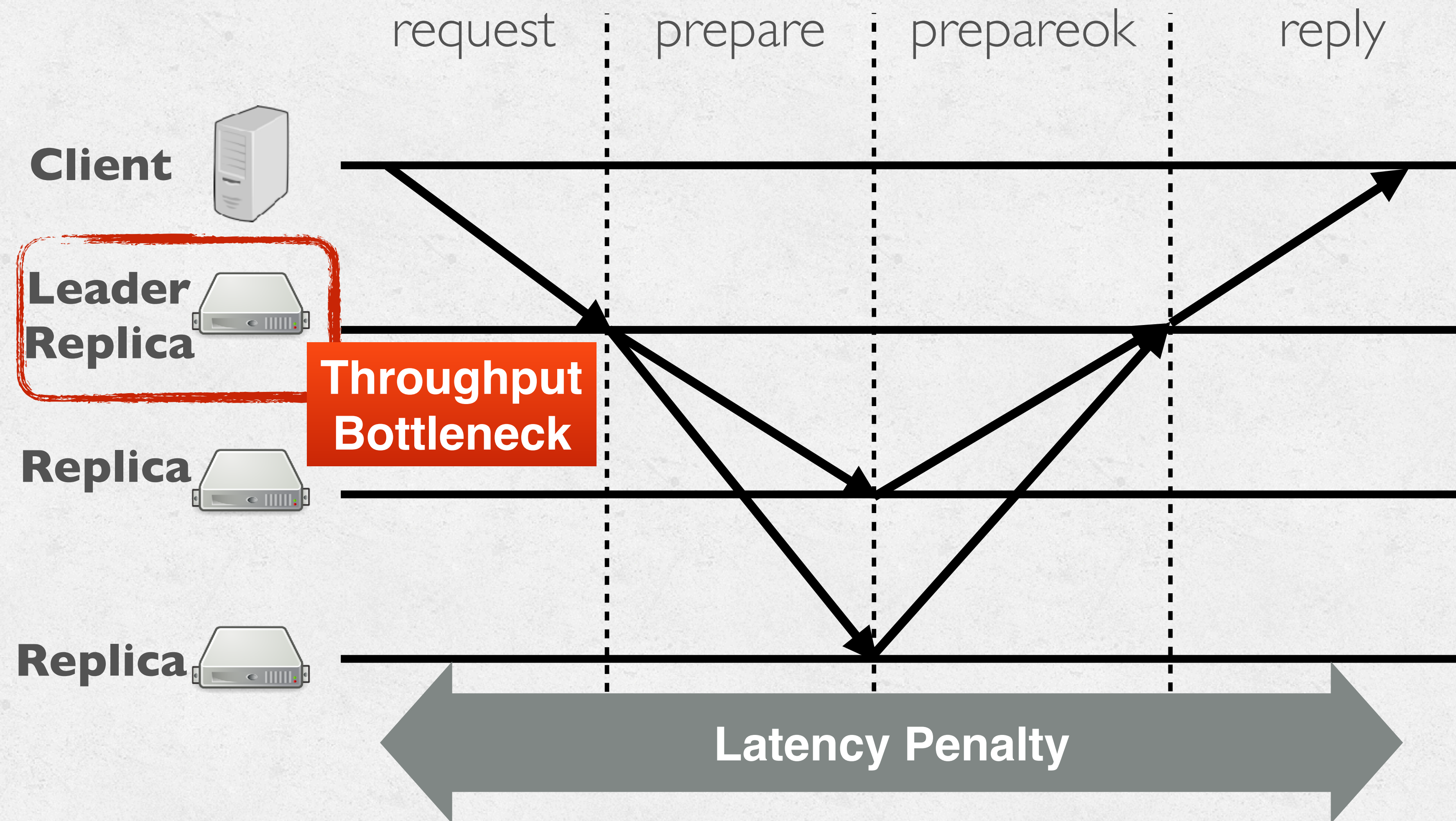
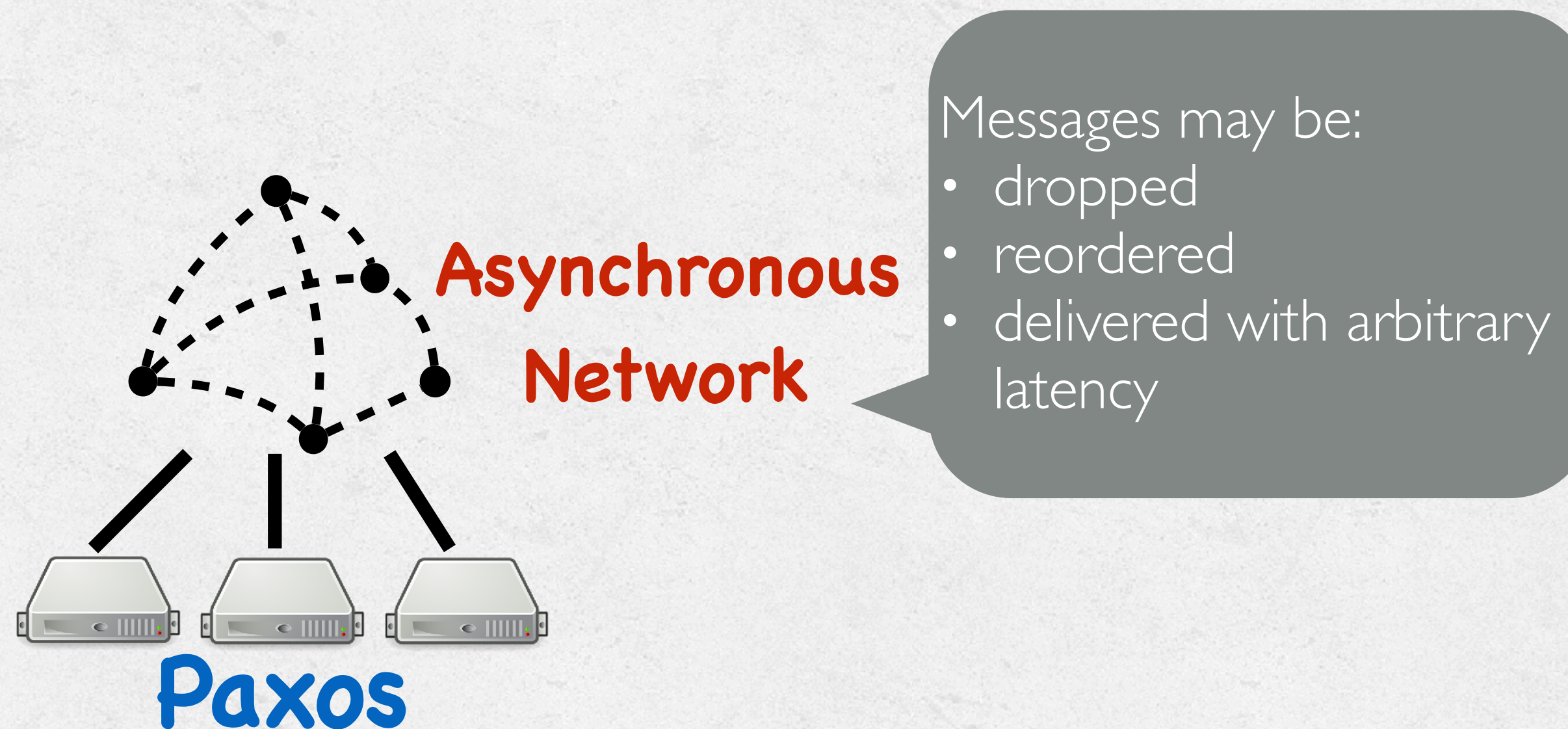request | prepare | prepareok | reply

Client

Leader
Replica

Replica

Replica

# PAXOS FOR STATE MACHINE REPLICATION

# PAXOS FOR STATE MACHINE REPLICATION
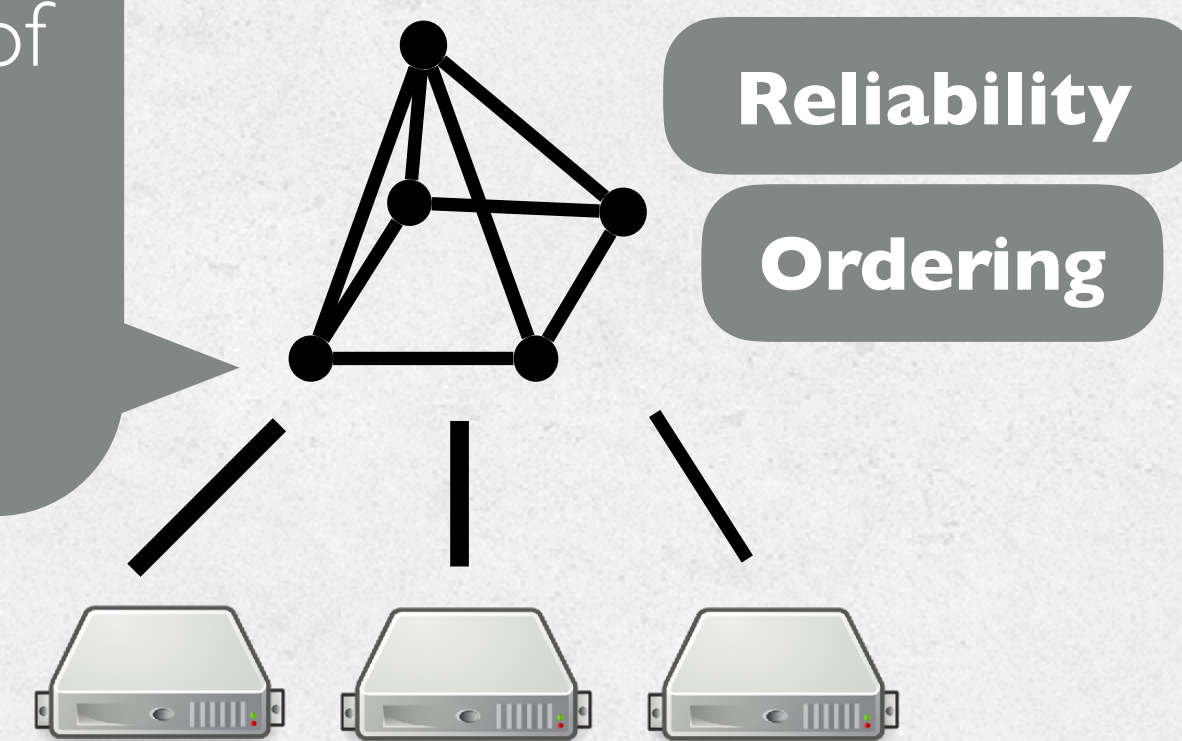
# NETWORK PROPERTIES DETERMINE REPLICATION COMPLEXITY



Asynchronous Network

Messages may be:
- dropped
- reordered
- delivered with arbitrary latency

Paxos

- Paxos protocol on every operation
- High performance cost
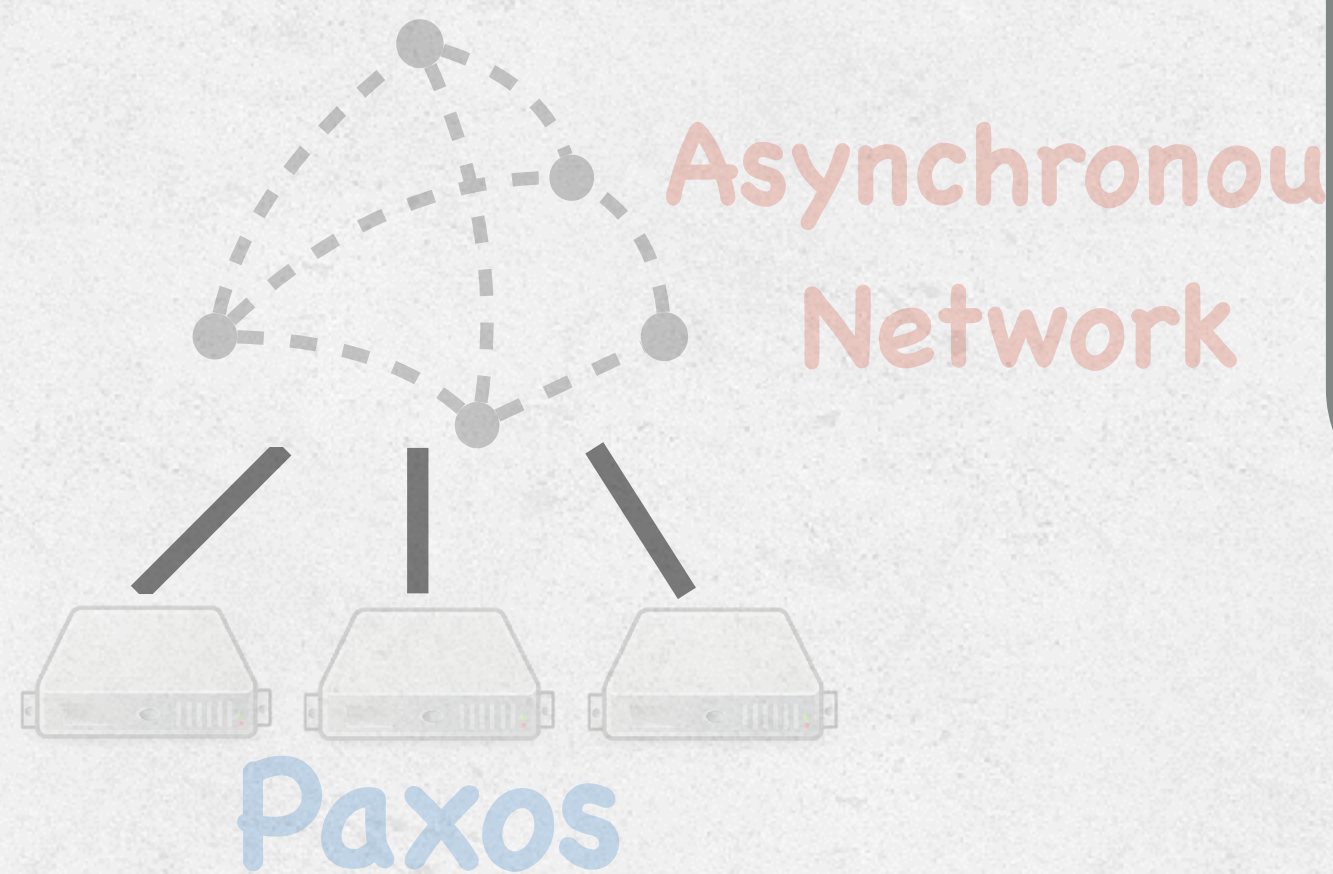
# NETWORK PROPERTIES DETERMINE REPLICATION COMPLEXITY

Asynchronous Network

Paxos

All replicas:
- receive the same set of messages
- receive them in the same order

Reliability

Ordering

- Paxos protocol on every operation
- High performance cost

# NETWORK PROPERTIES DETERMINE REPLICATION COMPLEXITY



Asynchronous Network

Paxos

All replicas:
- receive the same set of messages
- receive them in the same order

Reliability

Ordering

- Paxos protocol on every operation
- High performance cost

- Replication is trivial

# NETWORK PROPERTIES DETERMINE REPLICATION COMPLEXITY

Asynchronous Network

Paxos

All replicas:
- receive the same set of messages
- receive them in the same order

**Reliability**

**Ordering**

- Paxos protocol on every operation
- High performance cost

- Replication is trivial
- Network implementation has the same complexity as Paxos
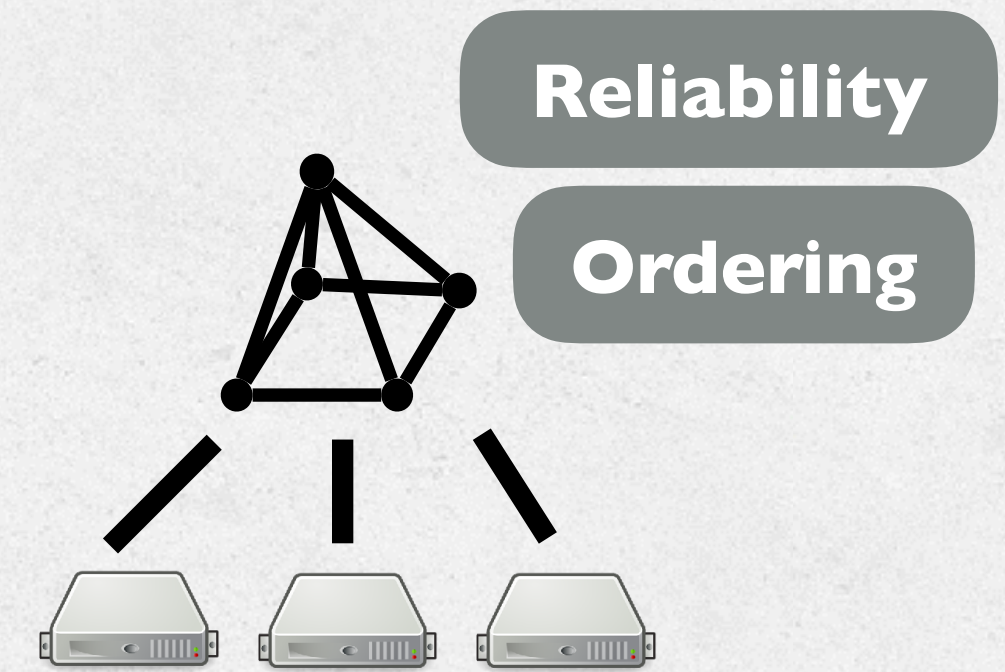
Asynchronous Network

Paxos

Reliability

Ordering

Weak

Strong

Network Guarantee

Asynchronous Network

Paxos

Reliability

Ordering

Weak

Strong

?

Network Guarantee

SPECPAXOS ASSUMED THE NETWORK WAS *MOSTLY ORDERED*

WHAT IF IT COULD PROVIDE AN *ORDERING GUARANTEE*?

# TOWARDS AN ORDERED BUT UNRELIABLE NETWORK

**Key Idea:** Separate **ordering** from **reliable delivery** in state machine replication
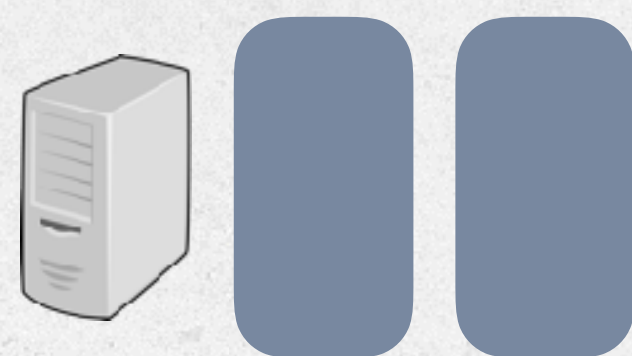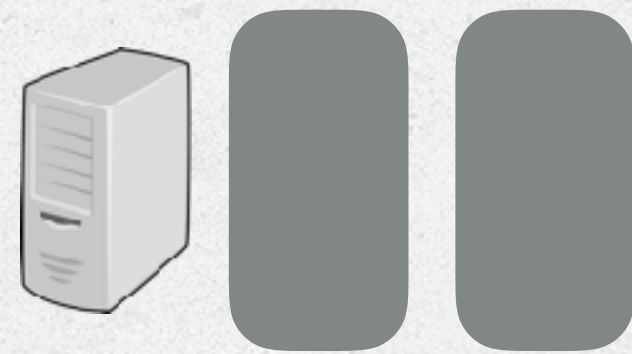
Network provides **ordering**

Replication protocol handles **reliability**

# OUM APPROACH

- Designate one **sequencer** in the network

- Sequencer maintains a counter for each OUM group

  1. Forward OUM messages to the sequencer

  2. Sequencer increments counter and writes counter value into packet headers

  3. Receivers use sequence numbers to detect **reordering** and **message drops**

Ordered Unreliable Multicast

Counter: 2

Senders

Receivers

Ordered Unreliable Multicast

Counter: 4

Senders

Receivers

Ordered Unreliable Multicast

Counter: 4

Senders

Receivers

**Ordered Unreliable**

**Ordered Multicast:**
no coordination required to determine order of messages

Counter: 4

1 2 DROP 4

1 2 3 4

**Senders**

**Receivers**

## Ordered Unreliable

**Ordered Multicast:**
no coordination required to determine order of messages

Counter:

1 2 DROP 4

**Drop Detection:**
coordination only required when messages are dropped

**Senders**                    **Receivers**

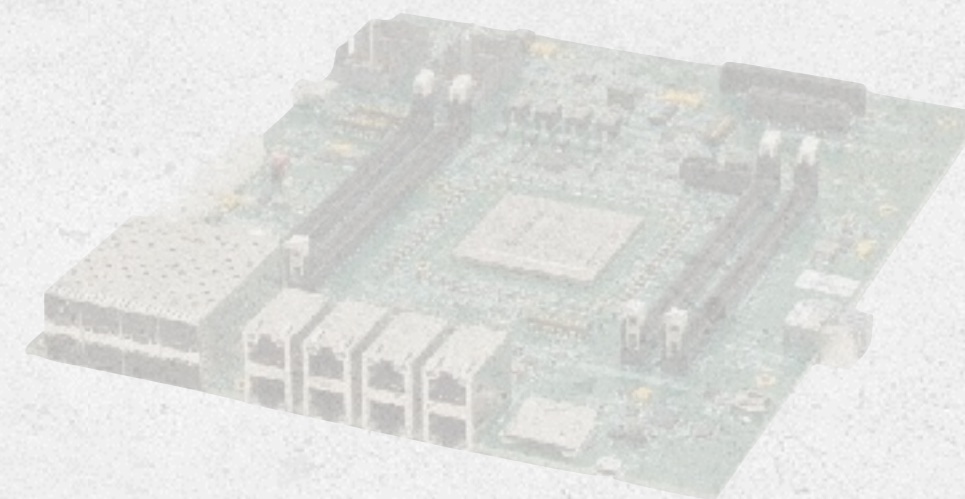# SEQUENCER IMPLEMENTATIONS

## In-switch sequencing

- next generation programmable switches
- implemented in P4
- nearly **zero cost**

## Middlebox prototype

- Cavium Octeon network processor
- connects to root switches
- adds 8 us latency

## End-host sequencing

- no specialized hardware required
- incurs higher latency penalties
- similar throughput benefits

# SEQUENCER IMPLEMENTATIONS

## In-switch sequencing

- next generation programmable switches
- implemented in P4
- nearly **zero cost**

## Middlebox prototype

- Cavium Octeon network processor
- connects to root switches
- adds 8 us latency

## End-host sequencing

- no specialized hardware required
- incurs higher latency penalties
- similar throughput benefits

# SEQUENCER IMPLEMENTATIONS

## In-switch sequencing
- next generation programmable switches
- implemented in P4
- nearly **zero cost**

## Middlebox prototype
- Cavium Octeon network processor
- connects to root switches
- adds 8 us latency

## End-host sequencing
- no specialized hardware required
- incurs higher latency penalties
- similar throughput benefits

# NOPAXOS OVERVIEW

- Built on top of the guarantees of OUM

- Client requests are **totally ordered** but can be dropped

- **No coordination** in the common case

- Replicas run agreement on drop detection

- View change protocol for leader or sequencer failure

# Normal Operation

# Normal Operation

# Normal Operation

# Normal Operation

request | reply

waits for replies from majority including leader's

**Client**

OUM

**Replica (leader)**

Execute

**Replica**

**Replica**

# GAP AGREEMENT

Replicas detect message drops.

- **Non-leader replicas:** recover the missing message from the leader

- **Leader replica:** coordinates to commit a NO-OP (Paxos)

- Efficient recovery from network anomalies

# WHY DO FOLLOWERS NOT EXECUTE?

- Request logs in NOPaxos are **non-authoritative**. The followers might not be involved in the quorum to commit a no-op. The leader might get replaced.

- Followers simply log operations. Operations are permanently committed with periodic **synchronization**.

- If a leader gets replaced and discovers that some of its commands weren't actually committed, it can roll-back or get a state transfer.

# VIEW CHANGE

- Handles leader or sequencer failure

- Ensures that all replicas are in a **consistent state** and agree on all of the commands and no-ops committed in the previous view.

- Runs a view change protocol similar to VR

- *view-number* is a tuple of *<leader-number, session-number>*

# NOPAXOS ACHIEVES BETTER THROUGHPUT AND LATENCY

Latency (us)

better ↓

better →

Throughput (ops/sec)

# NOPAXOS ACHIEVES BETTER THROUGHPUT AND LATENCY

# NOP**AXOS** **ACHIEVES BETTER THROUGHPUT AND LATENCY**

Latency (us)

better ↓

1000

750

500

250

0

Paxos

Fast Paxos

NOPaxos

4.7X throughput and more than 40% reduction in latency

0          65,000       130,000      195,000      260,000

better →

Throughput (ops/sec)

# NOPAXOS ACHIEVES BETTER THROUGHPUT AND LATENCY

# NOPAXOS ACHIEVES BETTER THROUGHPUT AND LATENCY

# NOPAXOS IS RESILIENT TO NETWORK ANOMALIES



Legend: ● NOPaxos  ◆ Paxos  ■ SpecPaxos

Y-axis: Throughput (ops/sec) — 0, 65,000, 130,000, 195,000, 260,000

X-axis: Packet Drop Rate — 0.001%, 0.01%, 0.1%, 1%

# NOPAXOS IS RESILIENT TO NETWORK ANOMALIES

# NOPₐxₒₛ ATTAINS THROUGHPUT WITHIN 2% OF AN UNREPLICATED SYSTEM

# NOPAXOS ATTAINS THROUGHPUT WITHIN 2% OF AN UNREPLICATED SYSTEM

# NOPAXOS ATTAINS THROUGHPUT WITHIN 2% OF AN UNREPLICATED SYSTEM

# NOPAXOS ATTAINS THROUGHPUT WITHIN 2% OF AN UNREPLICATED SYSTEM



Latency (us)

500

375

250    Paxos

better ↓

125    within 2% throughput and 16us latency of an unreplicated system    NOPaxos

Unreplicated

0

0    65,000    130,000    195,000    260,000

better →

Throughput (ops/sec)

# NOPAXOS ATTAINS THROUGHPUT WITHIN 2% OF AN UNREPLICATED SYSTEM

Latency (us)

500

375

250 — Paxos

better ↓

NOPaxos using end-host sequencer

within 2% throughput and 16us latency of an unreplicated system

125

NOPaxos

Unreplicated

0

0          65,000          130,000          195,000          260,000

better →

Throughput (ops/sec)

# NOPAXOS ATTAINS THROUGHPUT WITHIN 2% OF AN UNREPLICATED SYSTEM

# SUMMARY

- Separate ordering from reliable delivery in state machine replication

- A network model OUM that provides ordered but unreliable message delivery

- A more efficient replication protocol NOPaxos that ensures reliable delivery

- The combined system achieves performance equivalent to an unreplicated system

# THE ERIS TRANSACTION PROTOCOL

# Existing transactional systems: extensive coordination

# EXISTING TRANSACTIONAL SYSTEMS: EXTENSIVE COORDINATION

# EXISTING TRANSACTIONAL SYSTEMS: EXTENSIVE COORDINATION

# EXISTING TRANSACTIONAL SYSTEMS: EXTENSIVE COORDINATION

# ERIS

- Processes independent transactions **without coordination** in the normal case

- Performance within **3%** of a nontransactional, unreplicated system on TPC-C

- Strongly consistent, fault tolerant transactions with **minimal performance penalties**

# Key Contributions

A **new architecture** that divides the responsibility for transactional guarantees by

...leveraging the **datacenter network** to order messages within and across shards

...and a co-designed **transaction protocol** with minimal coordination.

# TRADITIONAL LAYERED APPROACH

# TRADITIONAL LAYERED APPROACH

Atomic Commitment (2PC)

Concurrency Control (2PL)

Concurrency Control (2PL)

Replication (Paxos)

Replica

Replica

Replica

Replication (Paxos)

Replica

Replica

Replica

Reliability (within shard)

Ordering (within shard)

# TRADITIONAL LAYERED APPROACH

# Traditional Layered Approach

# A New Way to Divide Responsibilities

Eris

| Isolation | General Transaction Protocol |

| Reliability (across shards) | Reliability (within shard) | Independent Transaction Protocol |

| Ordering (across shard) | Ordering (within shard) | Multi-sequencing |

# A NEW WAY TO DIVIDE RESPONSIBILITIES

Eris

Isolation

General Transaction Protocol

Reliability (across shards)

Reliability (within shard)

Independent Transaction Protocol

Application

- - - - - - - - - - - - - - - -

Network

Ordering (across shard)

Ordering (within shard)

Multi-sequencing

# GOAL

Client

Sequencer

# In-Network Concurrency Control Goals

- **Globally consistent ordering** across messages delivered to multiple destination shards

- No reliable delivery guarantee

- Recipients can **detect dropped messages**

A

**T1**
(ABC)

**T2**
(AB)

B

**T1**
(ABC)

**T2**
(AB)

C

**T1**
(ABC)

Receivers

# Multi-Sequenced Groupcast

- Groupcast: message header specifies a **set** of destination multicast groups

- Multi-sequenced groupcast: messages are sequenced **atomically** across all recipient groups

- Sequencer keeps a counter for each group

- Extends OUM in NOPaxos

Sequencer

T1
(ABC)

Counter:
A0  B0  C0

A

B

C

Receivers

Sequencer

T1
(ABC)

Counter:
A1  B1  C1

A

B

C

Receivers

Sequencer

T1
(ABC)

A1
B1
C1

Counter:
A1  B1  C1

A

B

C

Receivers

Sequencer

**T2**
(AB)

Counter:
A1  B1  C1

A

**T1**
(ABC)

**A1**
B1
C1

B

**T1**
(ABC)

A1
**B1**
C1

C

**T1**
(ABC)

A1
B1
**C1**

Receivers

Sequencer

**T2**
(AB)

Counter:
A2  B2  C1

A

**T1**
(ABC)

**A1**
B1
C1

B

**T1**
(ABC)

A1
**B1**
C1

C

**T1**
(ABC)

A1
B1
**C1**

Receivers

Sequencer

T2
(AB)
A2
B2

Counter:
A2  B2  C1

A

T1
(ABC)
**A1**
B1
C1

B

T1
(ABC)
A1
**B1**
C1

C

T1
(ABC)
A1
B1
**C1**

Receivers

Sequencer

**T3**
(A)

Counter:
A2  B2  C1

A

**T1**
(ABC)

**A1**
B1
C1

B

**T1**
(ABC)

A1
**B1**
C1

**T2**
(AB)

A2
**B2**

C

**T1**
(ABC)
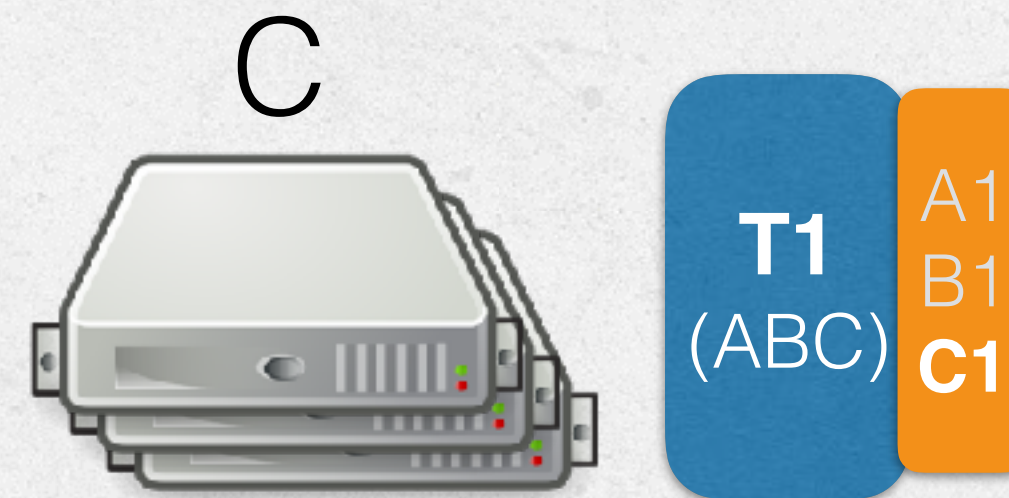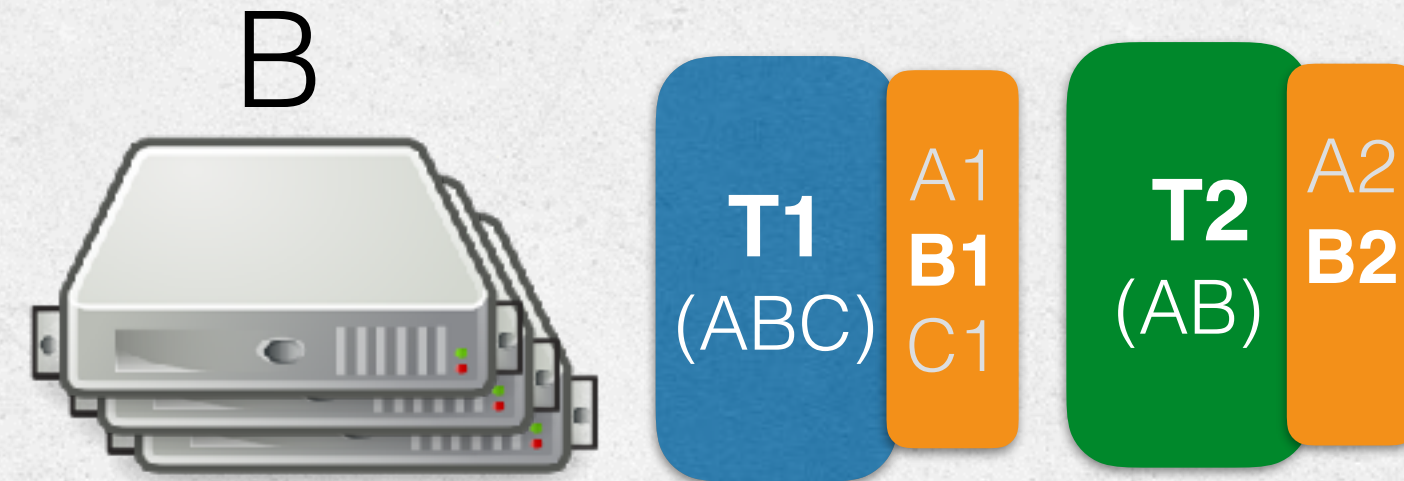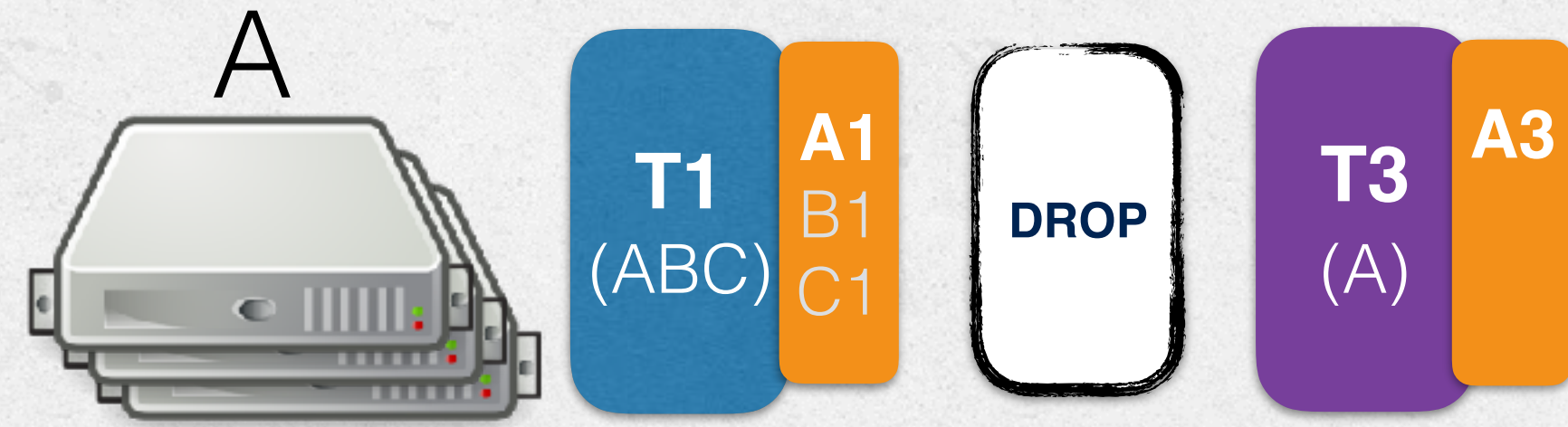
A1
B1
**C1**

Receivers

# WHAT HAVE WE ACCOMPLISHED SO FAR?

- Consistently ordered groupcast primitive with drop detection

- How do we go from multi-sequenced groupcast to transactions?

# TRANSACTION MODEL

Eris supports two types of transactions

- **Independent transactions**:

  ✤ One-shot (stored procedures)

  ✤ No cross-shard dependencies

  ✤ Proposed by H-Store [VLDB '07] and Granola [ATC '12]

- Fully general transactions

# INDEPENDENT TRANSACTION

```
START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT
```

| Name | Salary |
|------|--------|
| Alice | 600 |



| Name | Salary |
|------|--------|
| Bob | 350 |



| Name | Salary |
|------|--------|
| Charlie | 400 |

# INDEPENDENT TRANSACTION

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Alice | 600 |

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Bob | 350 |

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Charlie | 400 |

# INDEPENDENT TRANSACTION

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Alice | 600 |

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Bob | 450 |

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Charlie | 500 |

# INDEPENDENT TRANSACTION

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE 500 < (SELECT AVG(t2.Salary) FROM tb t2)
COMMIT

Not Independent!

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Alice | 600 |

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Bob | 450 |

START TRANSACTION
UPDATE tb t1
SET t1.Salary = t1.Salary + 100
WHERE t1.Salary < 500
COMMIT

| Name | Salary |
|------|--------|
| Charlie | 500 |

# Independent Transaction

| | |
|---|---|
| START TRANSACTION | |
| UPDATE tb t1 | |
| SET t1.Salary = t1.Salary + 100 | |
| WHERE t1.Salary < 500 | |
| COMMIT | |

| Name | Salary |
|---|---|
| Alice | 600 |

| | |
|---|---|
| START TRANSACTION | |
| UPDATE tb t1 | |
| SET t1.Salary = t1.Salary + 100 | |
| WHERE t1.Salary < 500 | |
| COMMIT | |

| Name | Salary |
|---|---|
| Bob | 450 |

| | |
|---|---|
| START TRANSACTION | |
| UPDATE tb t1 | |
| SET t1.Salary = t1.Salary + 100 | |
| WHERE t1.Salary < 500 | |
| COMMIT | |

| Name | Salary |
|---|---|
| Charlie | 500 |

# WHY INDEPENDENT TRANSACTIONS?

- **No** coordination/communication across shards

- Executing them serially at each shard in a **consistent order** guarantees **serializability**

- Multi-sequenced groupcast establishes such an order

- How to handle **message drops** and sequencer/ server **failures**?

# NORMAL CASE

Client

Sequencer

Shard 1   Learner

      Replica

      Replica

Shard 2   Learner

      Replica

      Replica

Shard 3   Learner

      Replica

      Replica

# Normal Case

# Normal Case

Client

Sequencer

Shard 1 Learner

Replica

Replica

Shard 2 Learner

Replica

Replica

Shard 3 Learner

Replica

Replica

# NORMAL CASE

# Normal Case

Client

Sequencer

Shard 1 Learner

Replica

Replica

Shard 2 Learner

Replica

Replica

Shard 3 Learner

Replica

Replica

NORMAL CASE

1 round trip

Client
Sequencer
Shard 1 Learner
Replica
Replica
Shard 2 Learner
Replica
Replica
Shard 3 Learner
Replica
Replica

# How to handle dropped messages?

# HOW TO HANDLE DROPPED MESSAGES?

# HOW TO HANDLE DROPPED MESSAGES?

# How to handle dropped messages?

# HOW TO HANDLE DROPPED MESSAGES?

# HOW TO HANDLE DROPPED MESSAGES?

Global coordination problem

# THE FAILURE COORDINATOR

# The Failure Coordinator

# The Failure Coordinator

# The Failure Coordinator

# The Failure Coordinator

# The Failure Coordinator

# The Failure Coordinator

# DESIGNATED LEARNER AND SEQUENCER FAILURES

Designated learner (DL) failure:

- View change based protocol

- Ensures new DL learns all **committed transactions** from previous views

Sequencer failure:

- Higher epoch number from the new sequencer

- Epoch change ensures all replicas across all shards start the new epoch in **consistent states**. They should all agree on the exact set of transactions completed in the previous epoch.

# CAN WE PROCESS NON-INDEPENDENT TRANSACTIONS EFFICIENTLY?

# APPROACH: DIVIDE INTO INDEPENDENT TRANSACTIONS

- Relies on the **linearizable execution** of independent transactions

- This means that we have the abstraction of a single, correct machine that processes independent transactions only.

- Uses **locks** to provide strong isolation

- Two phases:

  - Independent transaction 1: execute reads and acquire locks

  - Independent transaction 2: commit/abort changes and release locks

# BENEFITS OF OUR LAYERED ARCHITECTURE

- Simple solution to handle client failures: if the client fails, any server can **unilaterally** send the abort command for its general transactions *as an independent transaction.*

- No deadlocks/deadlock detection. Locks are acquired in a single step.

- Furthermore, we don't even need aborts! Wait queues are easy.

- Takes advantage of the efficient independent transaction processing layer. General transactions are processed in two round trips in the normal case.

# Evaluation Comparison Systems

- Lock-Store (2PC + 2PL + Paxos)

- TAPIR [SOSP '15]

- Granola [ATC '12]

- Non-transactional, unreplicated (NT-UR)

# ERIS PERFORMS WELL ON INDEPENDENT TRANSACTIONS

Distributed independent transactions

# Eris performs well on independent transactions

Distributed independent transactions

Throughput (txns/sec)

1,200K

900K

600K

300K

0K

Eris outperforms Lock-Store, TAPIR and Granola by more than 3X

Lock-Store    TAPIR    Granola    **Eris**    NT-UR

# ERIS PERFORMS WELL ON INDEPENDENT TRANSACTIONS

Distributed independent transactions

Throughput (txns/sec)

1,200K

900K

600K

300K

0K

Lock-Store    TAPIR    Granola    **Eris**    NT-UR

Eris outperforms Lock-Store, TAPIR and Granola by more than 3X

Eris achieves throughput within 10% of NT-UR

# ERIS PERFORMS WELL ON INDEPENDENT TRANSACTIONS

Distributed independent transactions

More than **70% reduction** in latency compared to Lock-Store, and **within 10%** latency of NT-UR

**Granola by more than 3X**

Throug

300K

0K

Lock-Store          TAPIR          Granola          **Eris**          NT-UR

# ERIS ALSO PERFORMS WELL ON GENERAL TRANSACTIONS

Distributed general transactions

Eris maintains throughput within 10% of NT-UR

Throughput (txns/sec)

1,200K

900K

600K

300K

0K

Lock-Store    TAPIR    Granola    **Eris**    NT-UR

# ERIS EXCELS AT COMPLEX TRANSACTIONAL APPLICATIONS



TPC-C benchmark

# ERIS EXCELS AT COMPLEX TRANSACTIONAL APPLICATIONS

TPC-C benchmark

7.6X and 6.4X higher throughput than Lock-Store and Tapir

# ERIS IS RESILIENT TO NETWORK ANOMALIES

# ERIS IS RESILIENT TO NETWORK ANOMALIES

# ERIS RECAP

- A new division of responsibility for transaction processing

  - ✤ An in-network concurrency control mechanism that establishes a **consistent order** of transactions across shards

  - ✤ An efficient protocol that ensures **reliable delivery** of independent transactions

  - ✤ A **general transaction** layer atop independent transaction processing

- Result: strongly consistent, fault-tolerant transactions with **minimal performance overhead**

# ERIS AND NOPAXOS DISCUSSION

- Can we use an end-host sequencer for Eris? In NOPaxos, it's not a problem.

- What properties are important to NOPaxos's "scalability"?

- How deployable are these approaches?

- How scalable is Eris compared to two-phase commit?