# Tolerating Slowdowns in Replicated State Machines using Copilots

**Khiem Ngo**, Siddhartha Sen, Wyatt Lloyd

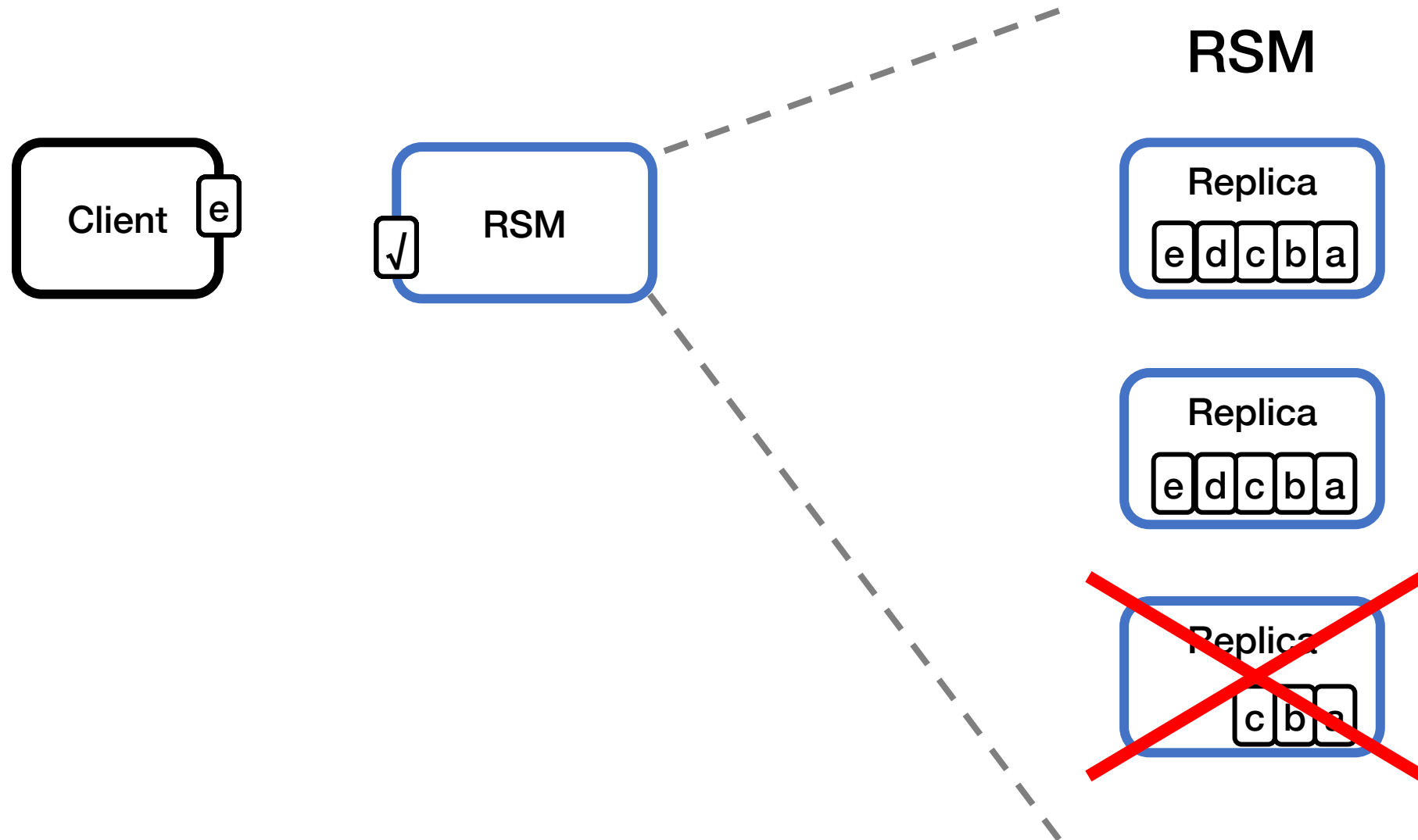Princeton University, Microsoft Research

# Replicated State Machine (RSM)

- Fault tolerant group of replicas that acts like a single machine that does not fail

- RSMs are everywhere!
  - Distributed database, cloud storage, coordination services, …

# Fault Tolerance for High Availability

Client [e]

√ RSM

RSM

Replica
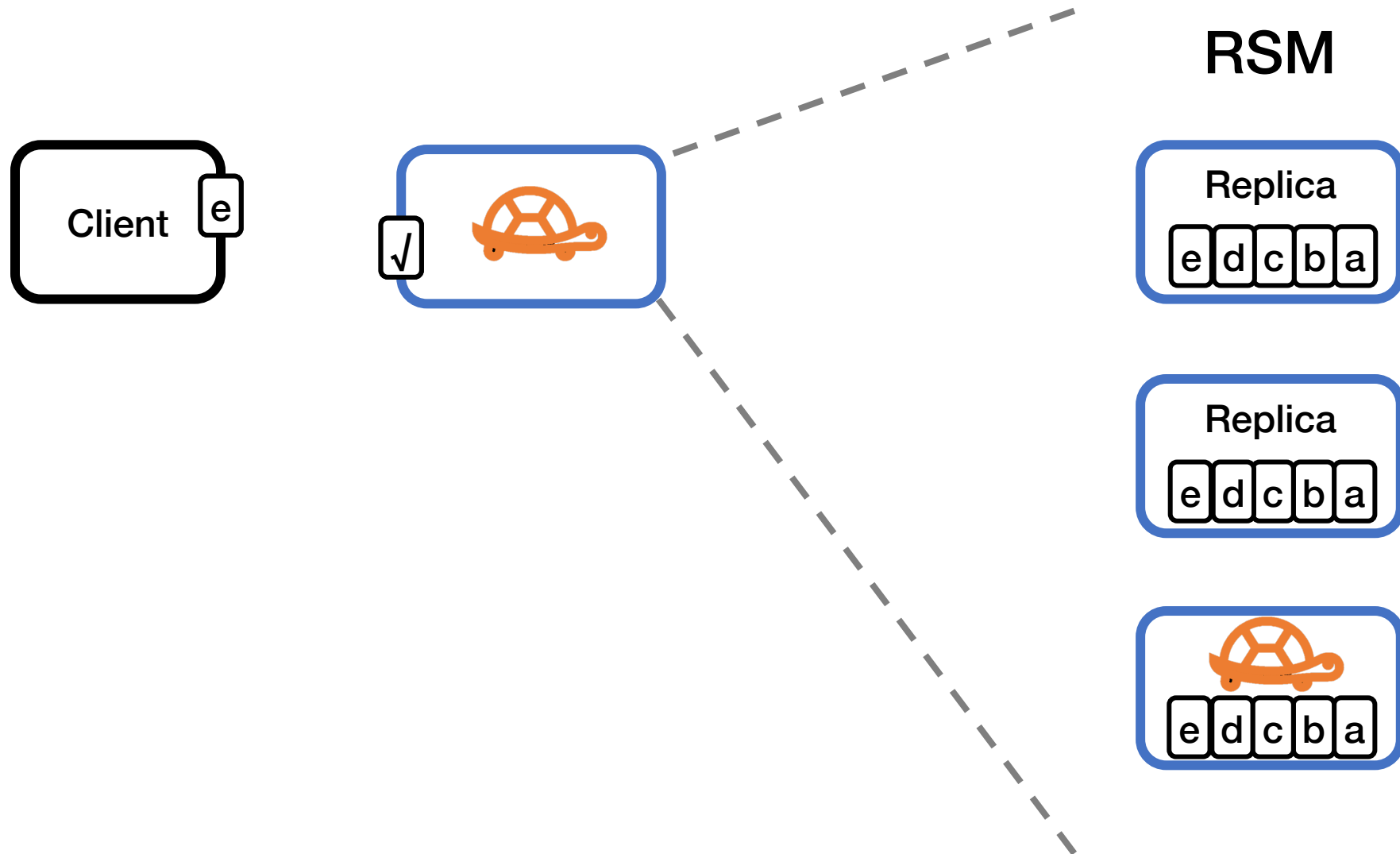[e][d][c][b][a]

Replica
[e][d][c][b][a]

Replica
[c][b][a]

# Replicas can Slowdown instead of Fail

- Many causes:
  - Misconfigurations
  - Partial hardware failures
  - Garbage collection events
  - …

- Effect: Replica takes longer than usual to send responses

# RSMs tolerate failures, not slowdowns

# Slowdowns Hurt Availability

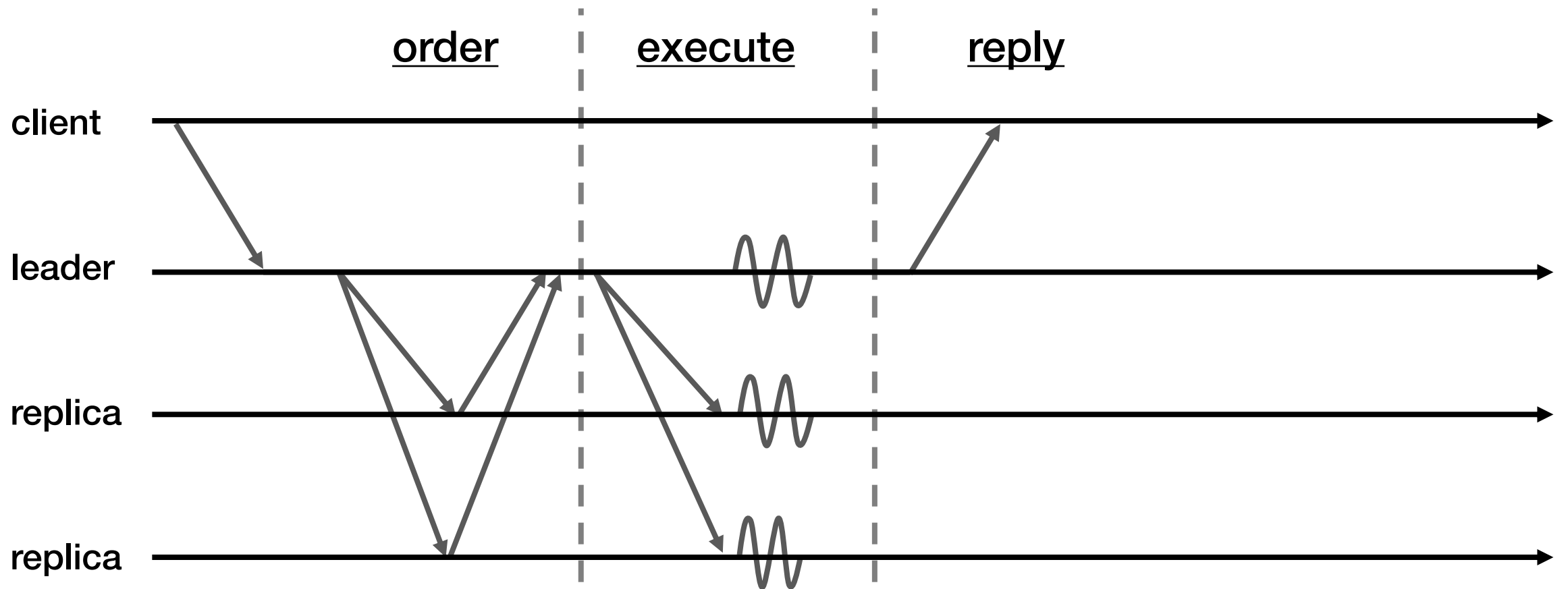# We need slowdown tolerance!

# Slowdowns Take Different Forms

- Duration
  - Transient slowdowns: not handled in general
  - Long-term slowdowns: eventually detected, but need to tolerate between onset and end of reconfiguration

- Severity
  - 10ms additional delay or 80ms?

- Scope
  - All processing paths or a subset?
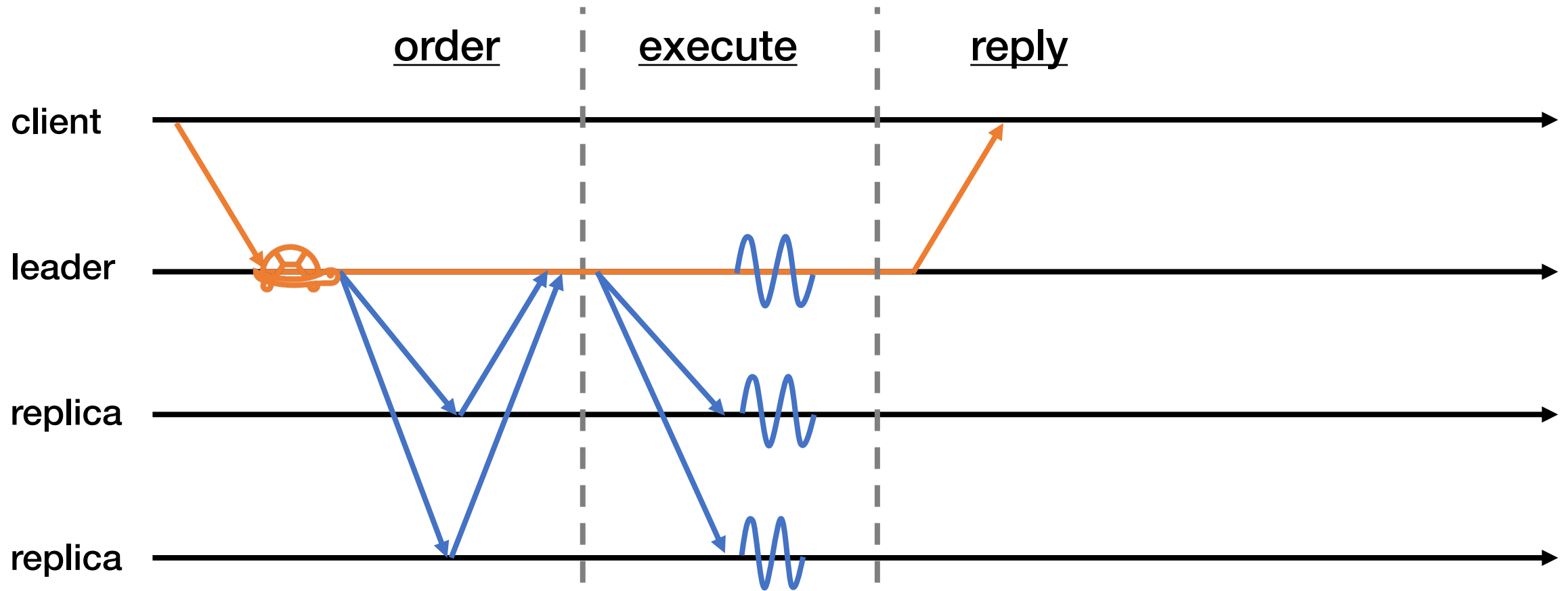
# Defining Slowdown Tolerance

- "slow" replica = responses to messages take more than threshold time $t$ over normal response time

- An RSM is $s$-slowdown-tolerant if it is not slow despite $s$ slow replicas
  - Replacing the $s$ slowest replicas with normal replicas should not change performance much

- This work's focus: 1-slowdown-tolerance
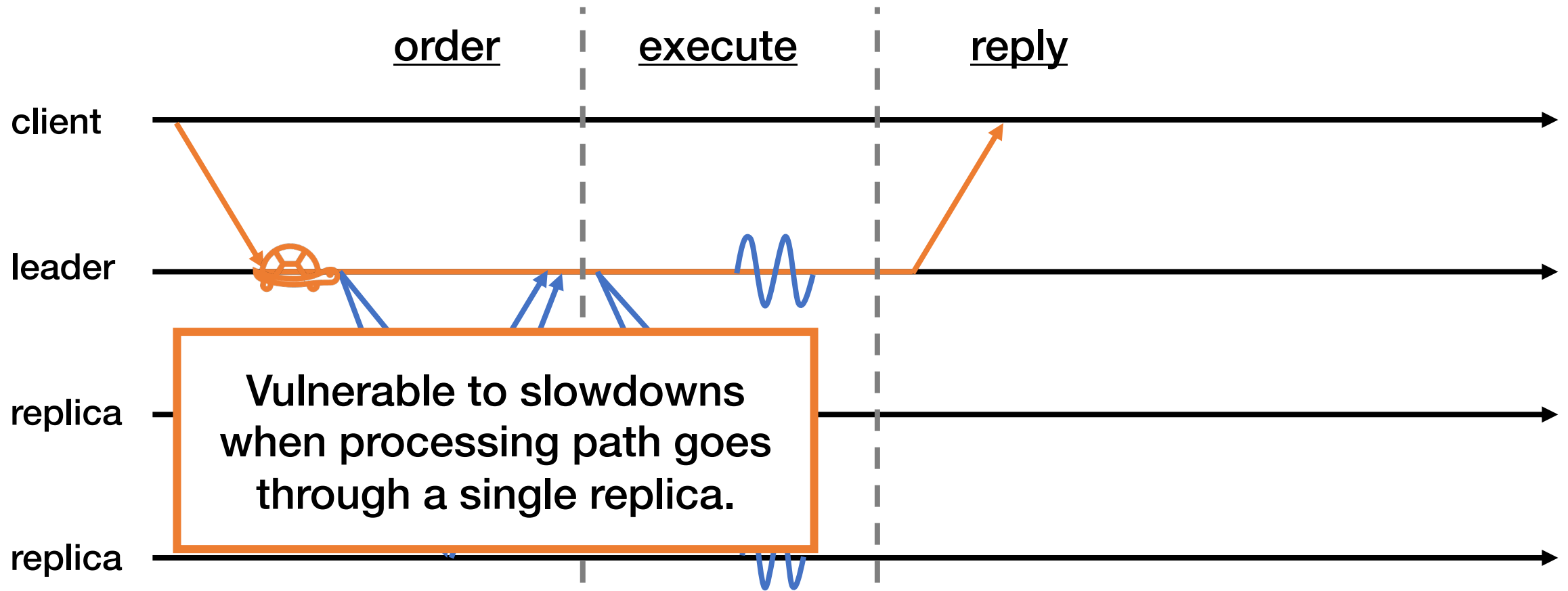
# No existing consensus protocol is 1-slowdown-tolerant
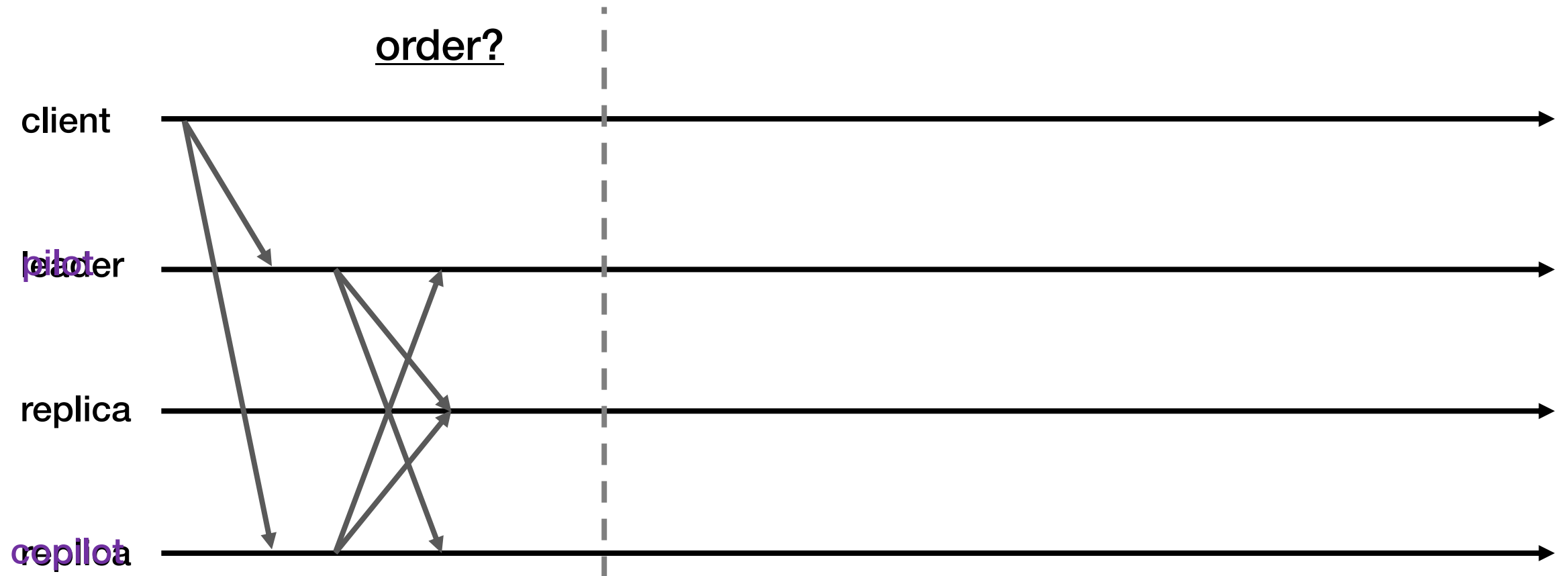
# Multi-Paxos is Not 1-Slowdown-Tolerant



order    execute    reply

client

leader

replica

replica

# Multi-Paxos is Not 1-Slowdown-Tolerant

order    execute    reply

client

leader

replica

replica

# Multi-Paxos is Not 1-Slowdown-Tolerant

order    execute    reply

client

leader

replica

replica

Vulnerable to slowdowns when processing path goes through a single replica.

# Copilot:
# First 1-Slowdown-Tolerant Protocol

order?

client

pilot leader

replica

copilot replica

# Ordering: Use Two Logs

RSM

Pilot

| e | d | c | b | a |

Copilot

| b' | e' | d' | c' | a' |

Replica

?

# Ordering: Combine Logs with Dependencies

# Ordering: Dependency Cycles

RSM

Break cycles deterministically



Pilot

e d c b a

Copilot

b' e' d' c' a'

Replica

# Ordering: A Tricky Case

RSM

Pilot

e d c b a

Copilot

b' e' d' c' a'

Possible
ordering:
a, a', b, c', c

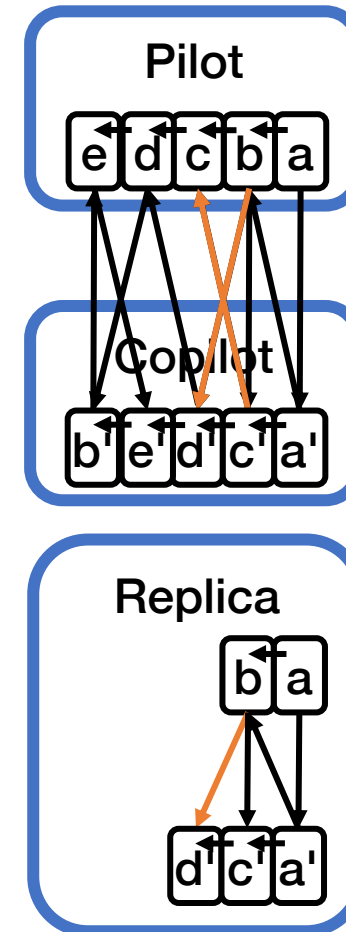Replica

c b a

c' a'

# Ordering: A Tricky Case

RSM



Possible ordering:
a, a', b, c', c
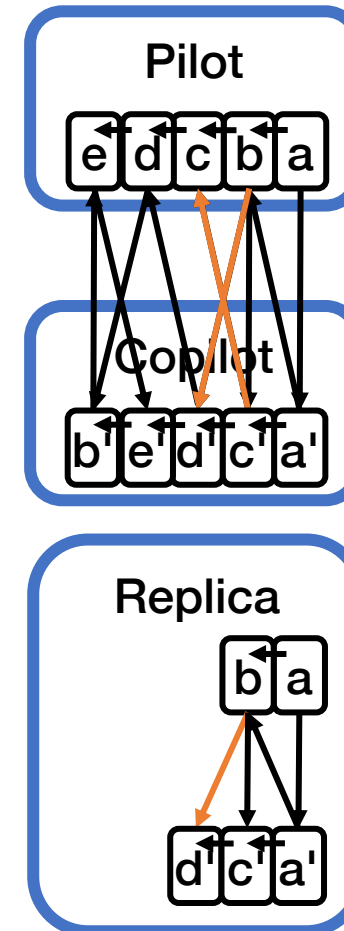
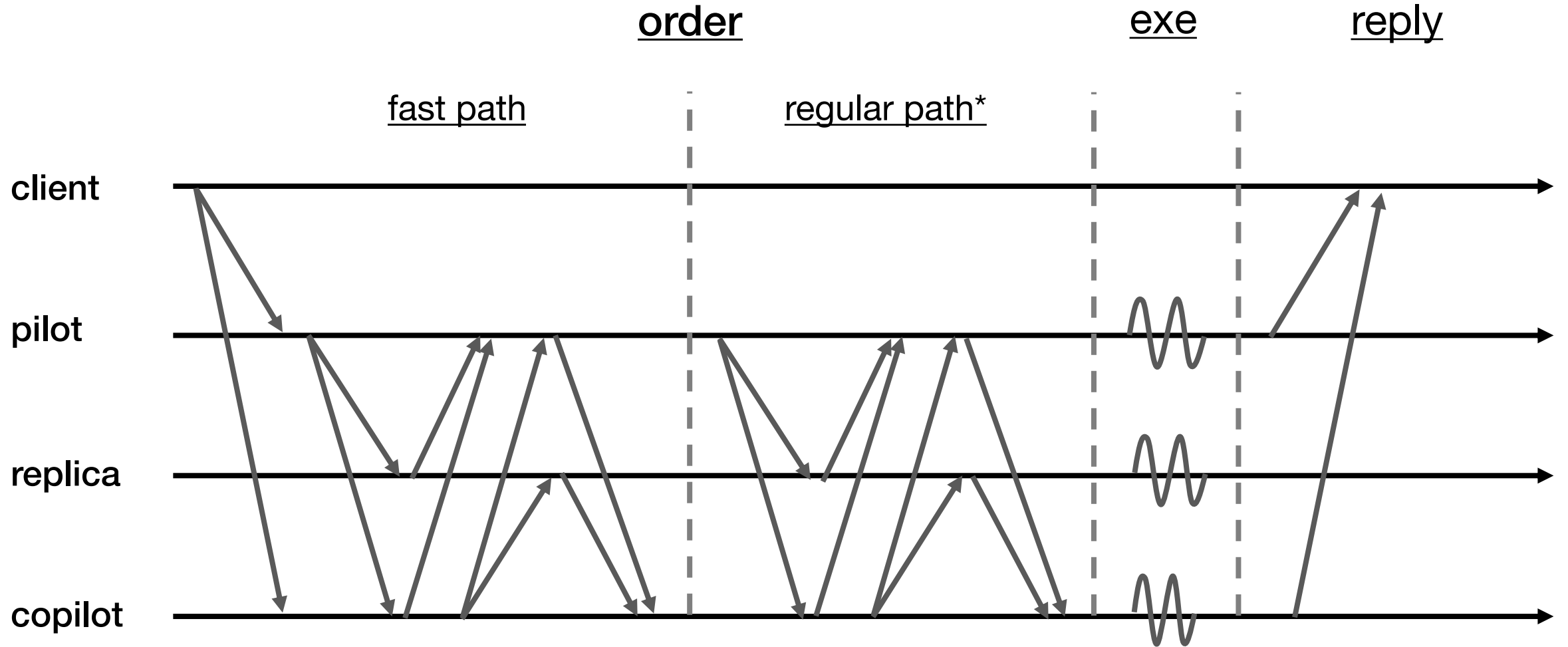Possible ordering:
a, a', b, c', d'

# Ordering: Same on All Replicas

## RSM

Compatibility check:
Only accept dependency if it
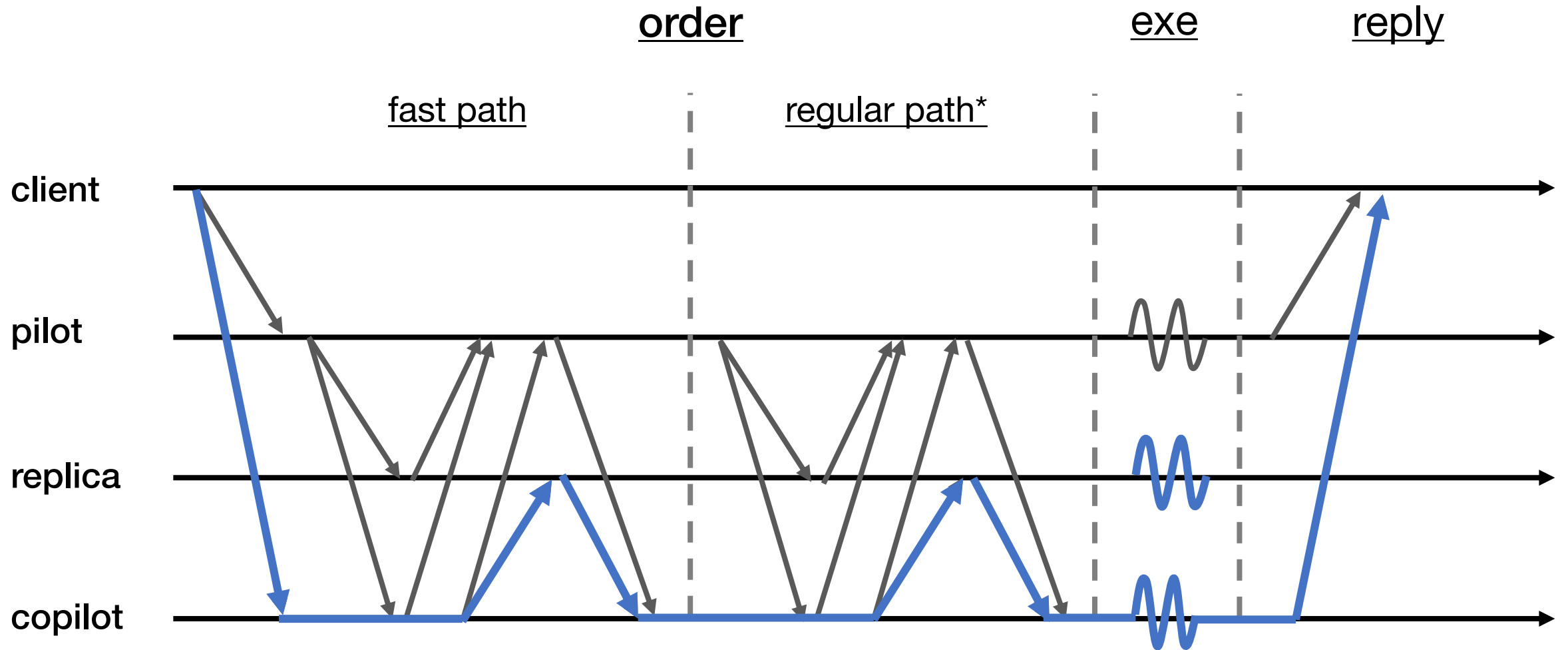cannot lead to multiple orders

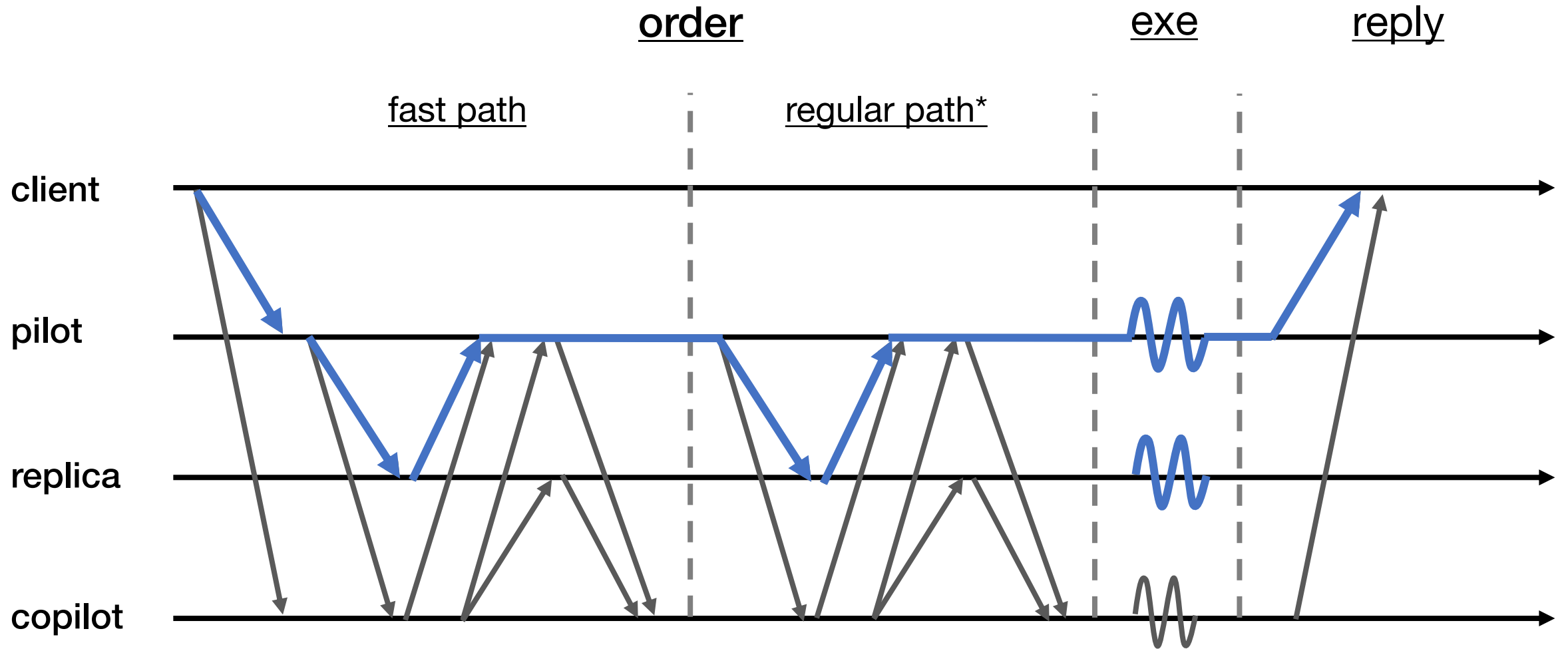Break cycles deterministically

# Copilot Protocol



order      exe      reply

fast path      regular path*

client

pilot

replica

copilot

# Copilot Protocol

# Copilot Protocol



order

exe

reply

fast path

regular path*

client

pilot

replica

copilot

# Copilot Protocol: Dependencies?

RSM

Pilot

b a

Copilot

c' a'

Solution: **fast takeover** the slow pilot's ordering work!

Replica

# Copilot Protocol: Summary

• Proactive redundancy: two pilots process all commands

• Use dependencies to combine ordering from two pilots
  • Compatibility check ensures same order on all replicas
  • Cycles broken by priority
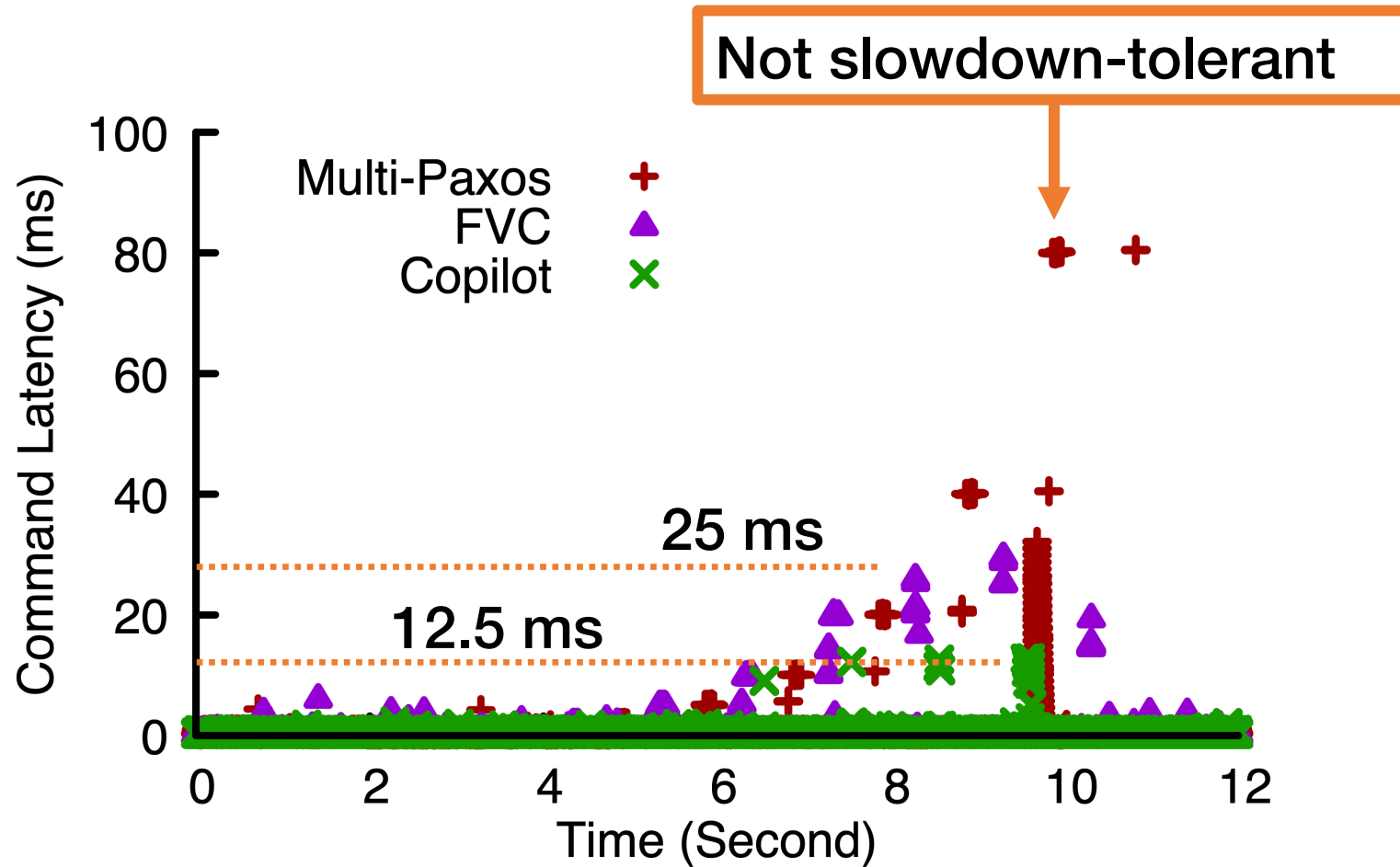  • Fast takeover to avoiding waiting on slow pilot

# Optimizations

- **Ping-Pong Batching**
  - Improve Copilot's performance when both pilots are fast
  - Pilots propose compatible orderings and commit on fast path

- **Null Dependency Elimination**
  - Improve Copilot's performance when one pilot is slow
  - Allow a fast pilot to safely avoid waiting on commits from a continually slow pilot and thus avoid fast takeover
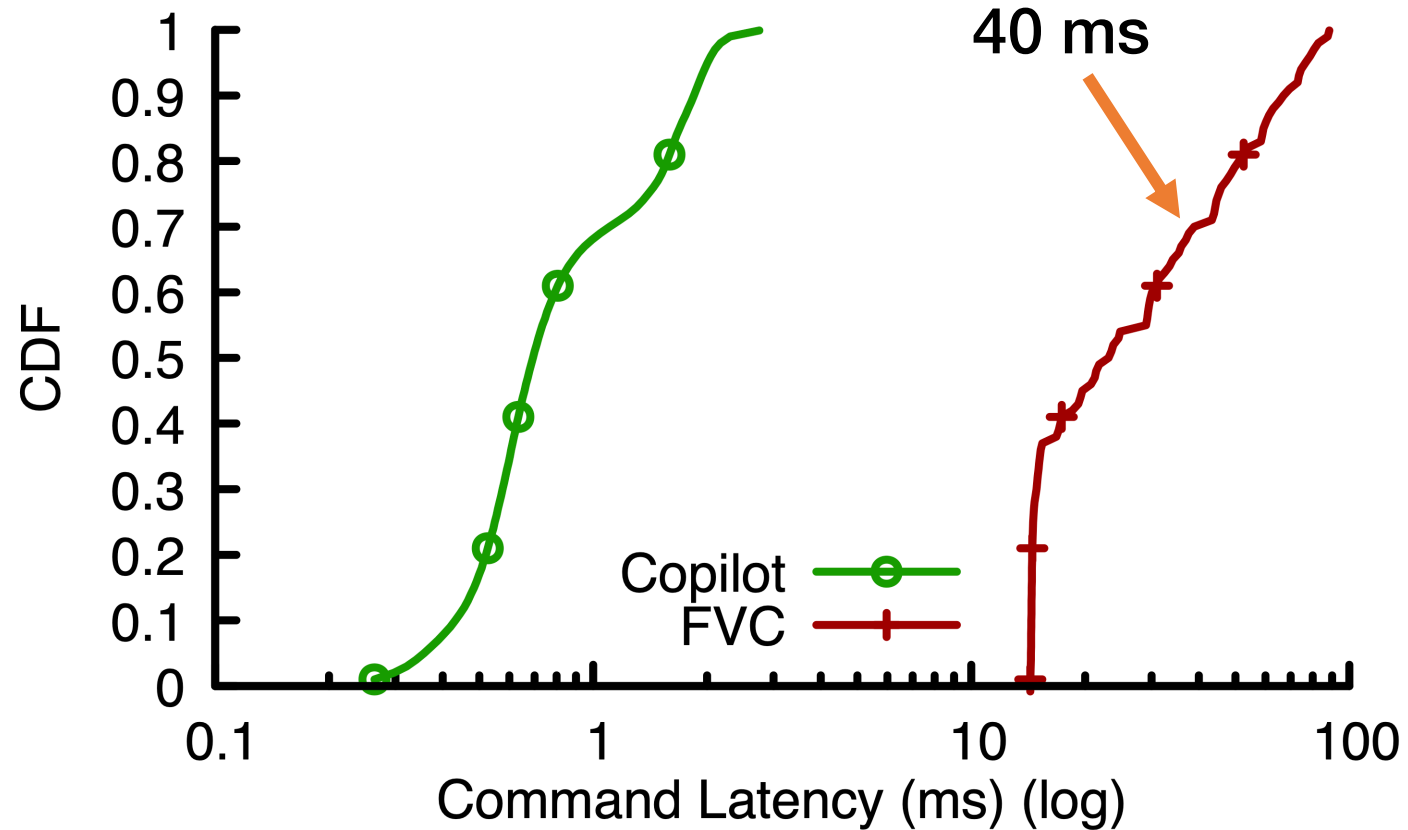
# Evaluation

- Tolerate slowdowns that are transient, have varying manifestations, have varying severity?
- How does Copilot perform without slow replicas?

- 5-replica RSM, moderate load
- Replicas and clients in the same datacenter

- Baselines:
  - EPaxos
  - Multi-Paxos
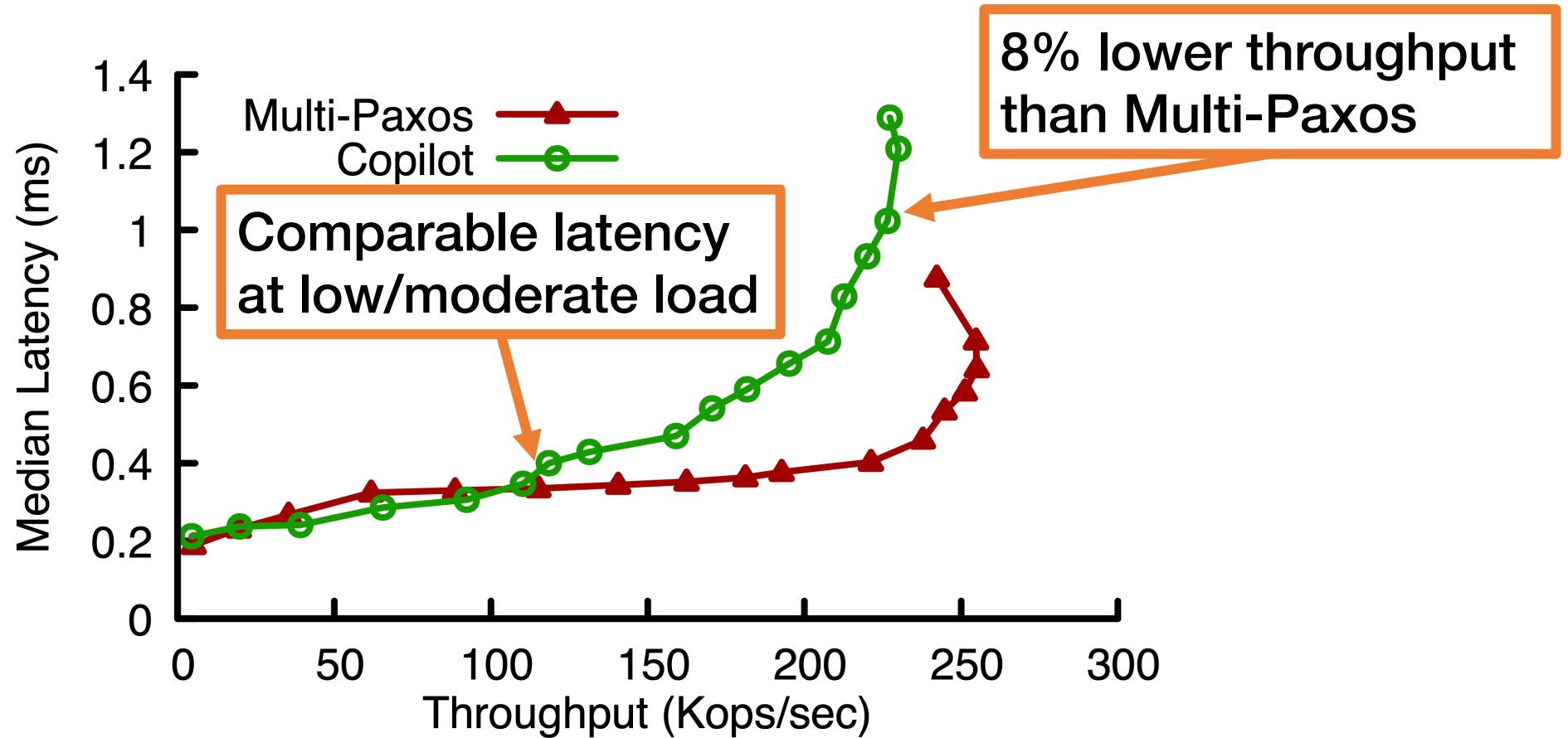  - Fast-View-Change (10 ms view-change timeout)

# Transient Slowdowns



Not slowdown-tolerant

Multi-Paxos
FVC
Copilot

25 ms

12.5 ms

Command Latency (ms)

Time (Second)

# Gradual Slowdown

# Performance Without Slow Replicas

# Conclusion

Contact Information:
Khiem Ngo
khiem@princeton.edu

- Slowdowns hurt availability, need s-slowdown-tolerant RSMs

- Copilot: first 1-slowdown-tolerant protocol
  - Slowdown tolerance: proactive redundancy and fast takeovers
  - Optimizations: ping-pong batching and null dependency elimination

- Copilot's performance without slow replicas is competitive

- Copilot is the only protocol that tolerates any one slowdown