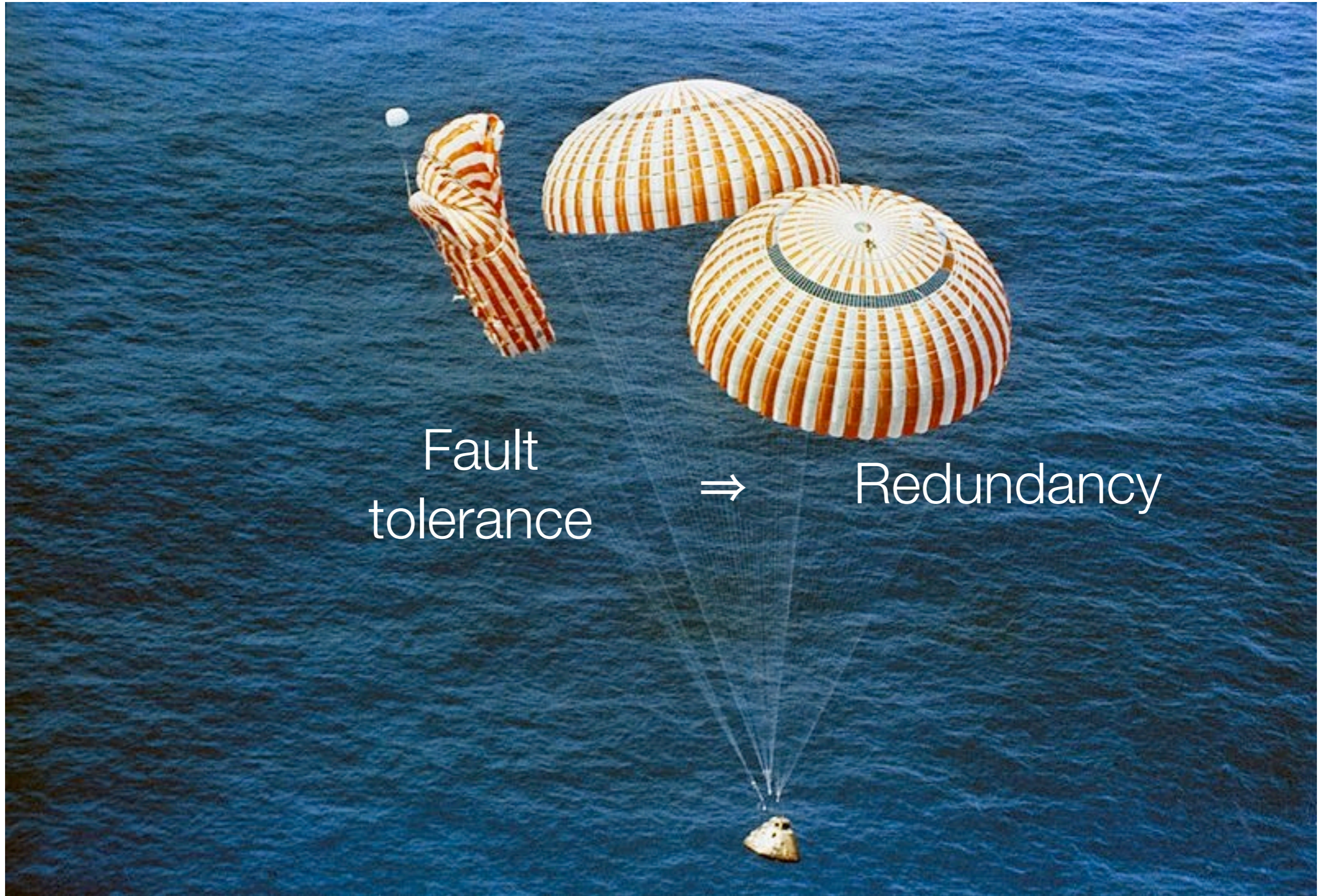


There Is More Consensus in Egalitarian Parliaments

Iulian Moraru, David Andersen, Michael Kaminsky

Carnegie Mellon University

Intel Labs



Fault
tolerance



Redundancy

State Machine Replication

State Machine Replication

Execute the same commands in the same order

State Machine Replication

Execute the same commands in the same order



1,9 km
Turn right onto Lewis Run Rd

1,5 km
Turn left onto Clairton Blvd

54,1 km
Turn right to merge onto US-119 S

40 5,7 km
Continue onto Uniontown Byp S

21,6 km
The destination is on your right



1,9 km
Turn right onto Lewis Run Rd

1,5 km
Turn left onto Clairton Blvd

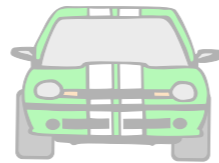
54,1 km
Turn right to merge onto US-119 S

40 5,7 km
Continue onto Uniontown Byp S

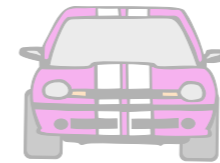
21,6 km
The destination is on your right

State Machine Replication

Execute the same commands in the same order



- 1,9 km
Turn right onto Lewis Run Rd
- 1,5 km
Turn left onto Clairton Blvd
- 54,1 km
Turn right to merge onto US-119 S
- 5,7 km
Continue onto Uniontown Byp S
- 21,6 km
The destination is on your right

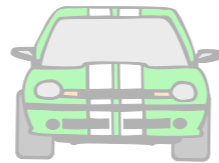


- 1,9 km
Turn right onto Lewis Run Rd
- 1,5 km
Turn left onto Clairton Blvd
- 54,1 km
Turn right to merge onto US-119 S
- 5,7 km
Continue onto Uniontown Byp S
- 21,6 km
The destination is on your right

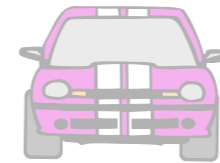
Paxos

State Machine Replication

Execute the same commands in the same order



→	1,9 km
	Turn right onto Lewis Run Rd
←	1,5 km
	Turn left onto Clairton Blvd
→	54,1 km
	Turn right to merge onto US-119 S
40	5,7 km
	Continue onto Uniontown Byp S
→	21,6 km
	The destination is on your right



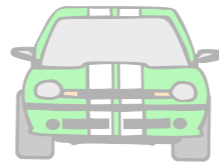
→	1,9 km
	Turn right onto Lewis Run Rd
←	1,5 km
	Turn left onto Clairton Blvd
→	54,1 km
	Turn right to merge onto US-119 S
40	5,7 km
	Continue onto Uniontown Byp S
→	21,6 km
	The destination is on your right

Paxos

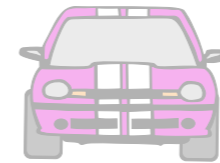
- No external failure detector required

State Machine Replication

Execute the same commands in the same order



→	1,9 km
	Turn right onto Lewis Run Rd
←	1,5 km
	Turn left onto Clairton Blvd
→	54,1 km
	Turn right to merge onto US-119 S
40	5,7 km
	Continue onto Uniontown Byp S
→	21,6 km
	The destination is on your right



→	1,9 km
	Turn right onto Lewis Run Rd
←	1,5 km
	Turn left onto Clairton Blvd
→	54,1 km
	Turn right to merge onto US-119 S
40	5,7 km
	Continue onto Uniontown Byp S
→	21,6 km
	The destination is on your right

Paxos

- No external failure detector required
- Fast fail-over (high availability)

Paxos is important in clusters

Chubby, Boxwood, SMARTER, ZooKeeper

- Synchronization
- Resource discovery
- Data replication

High throughput
High availability



Paxos is important in the wide-area

Spanner, Megastore

- Bring data closer to clients
- Tolerate datacenter outages

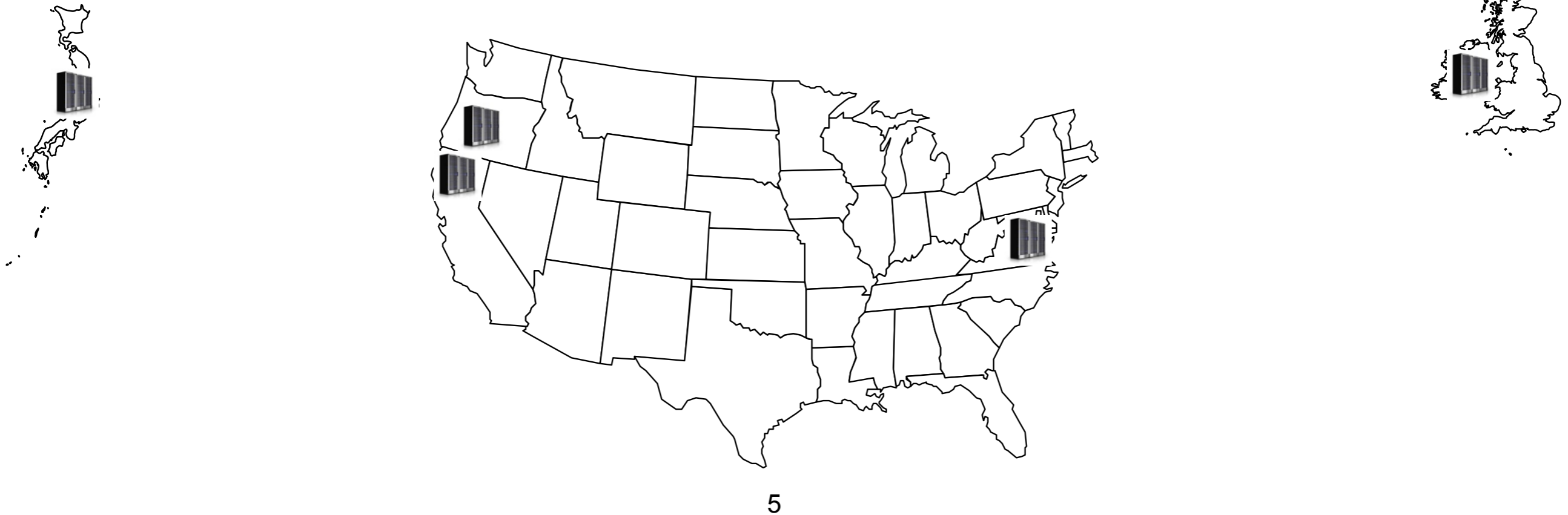
Low latency



Paxos is important in the wide-area

Spanner, Megastore

- Bring data closer to clients **Low latency**
- Tolerate datacenter outages

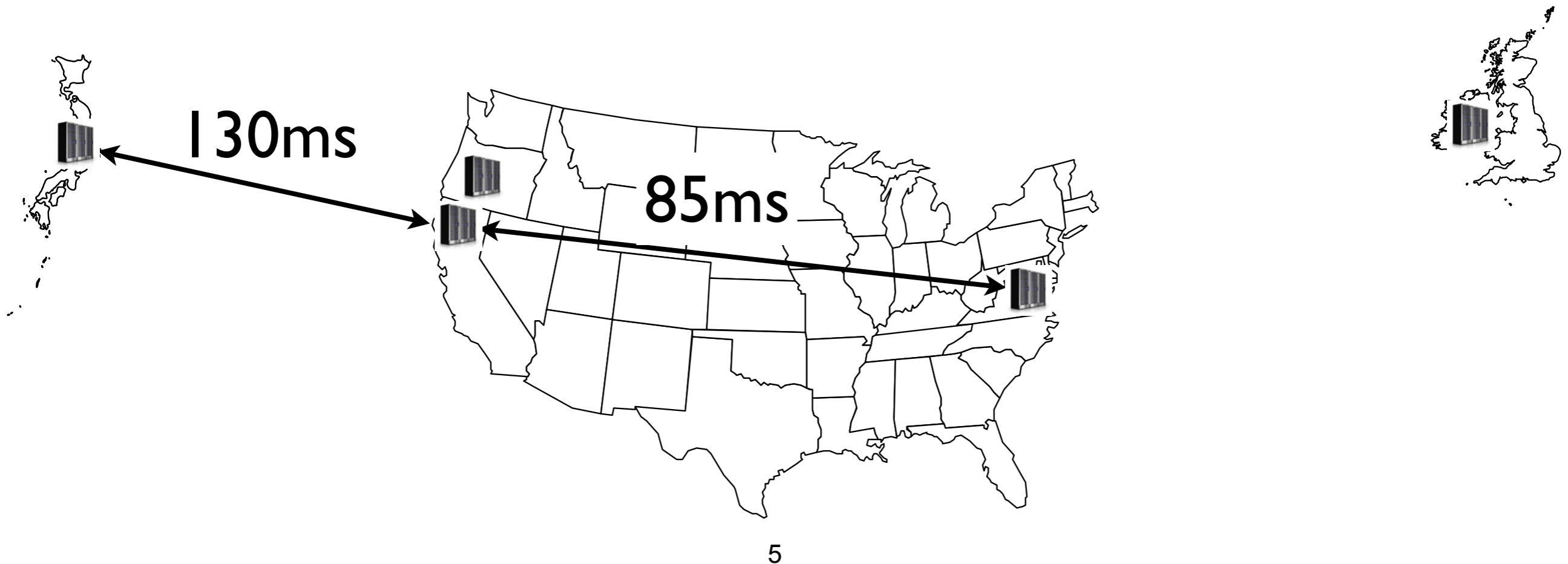


Paxos is important in the wide-area

Spanner, Megastore

- Bring data closer to clients
- Tolerate datacenter outages

Low latency



Paxos overview

Paxos overview

- Agreement protocol

Paxos overview

- Agreement protocol
- Tolerates F failures with $2F+1$ replicas (optimal)
 - No external failure detector required

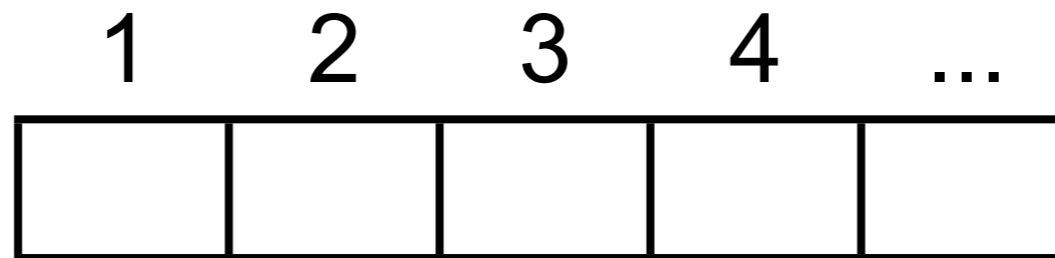
Paxos overview

- Agreement protocol
- Tolerates F failures with $2F+1$ replicas (optimal)
 - No external failure detector required
- Replicas can fail by crashing (non-Byzantine)

Paxos overview

- Agreement protocol
- Tolerates F failures with $2F+1$ replicas (optimal)
 - No external failure detector required
- Replicas can fail by crashing (non-Byzantine)
- Asynchronous communication

Using Paxos to order commands



Using Paxos to order commands

C



A



B



1

2

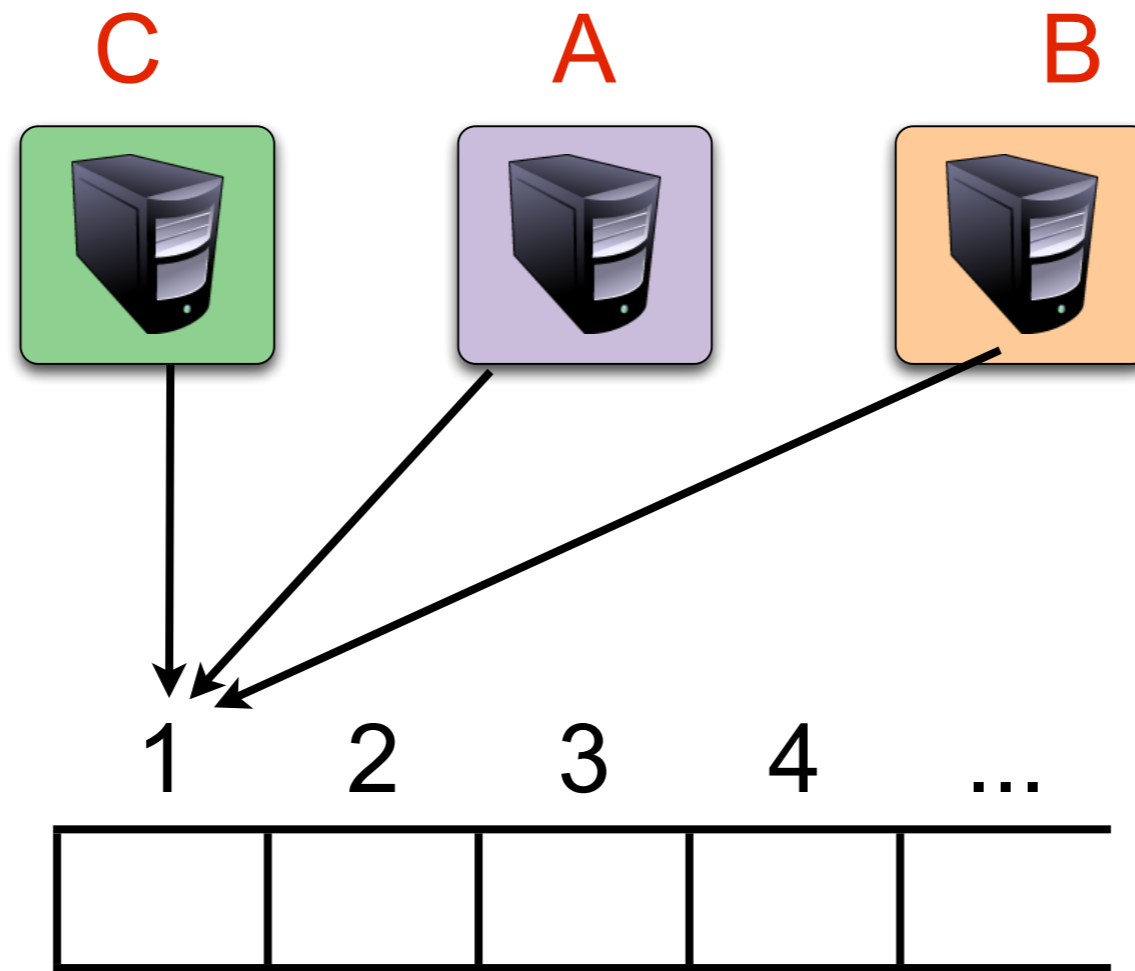
3

4

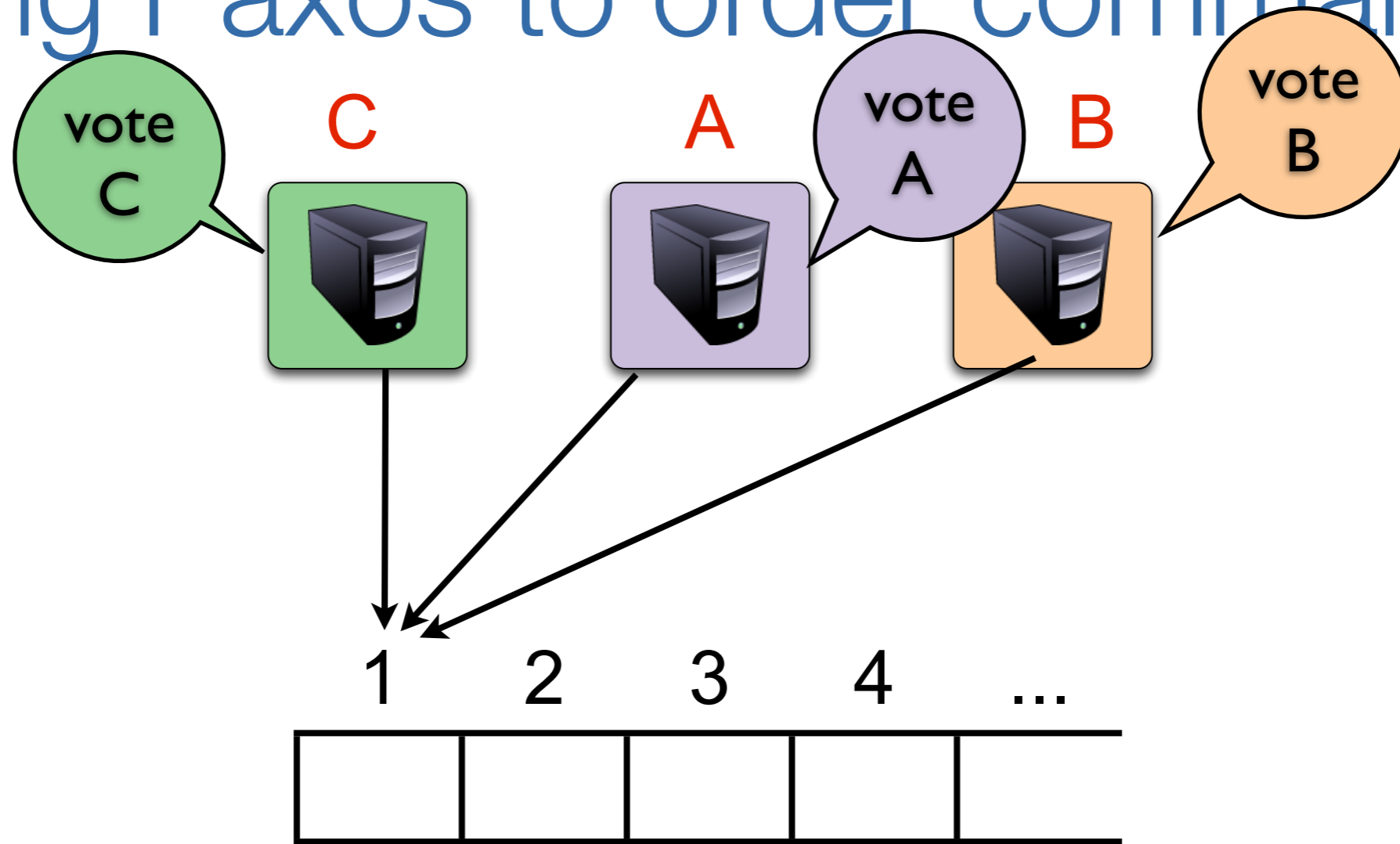
...



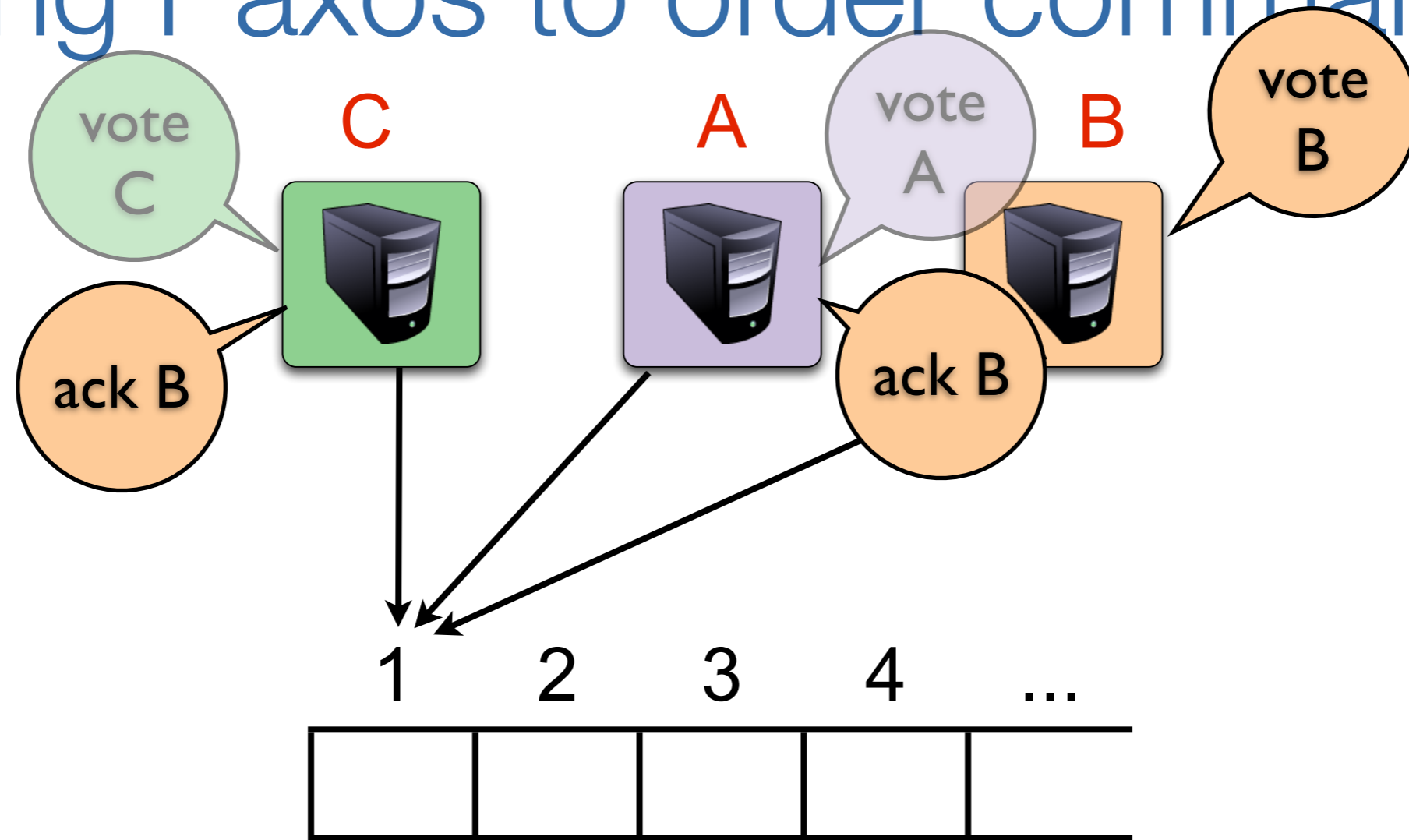
Using Paxos to order commands



Using Paxos to order commands



Using Paxos to order commands



Using Paxos to order commands

C



A



1

2

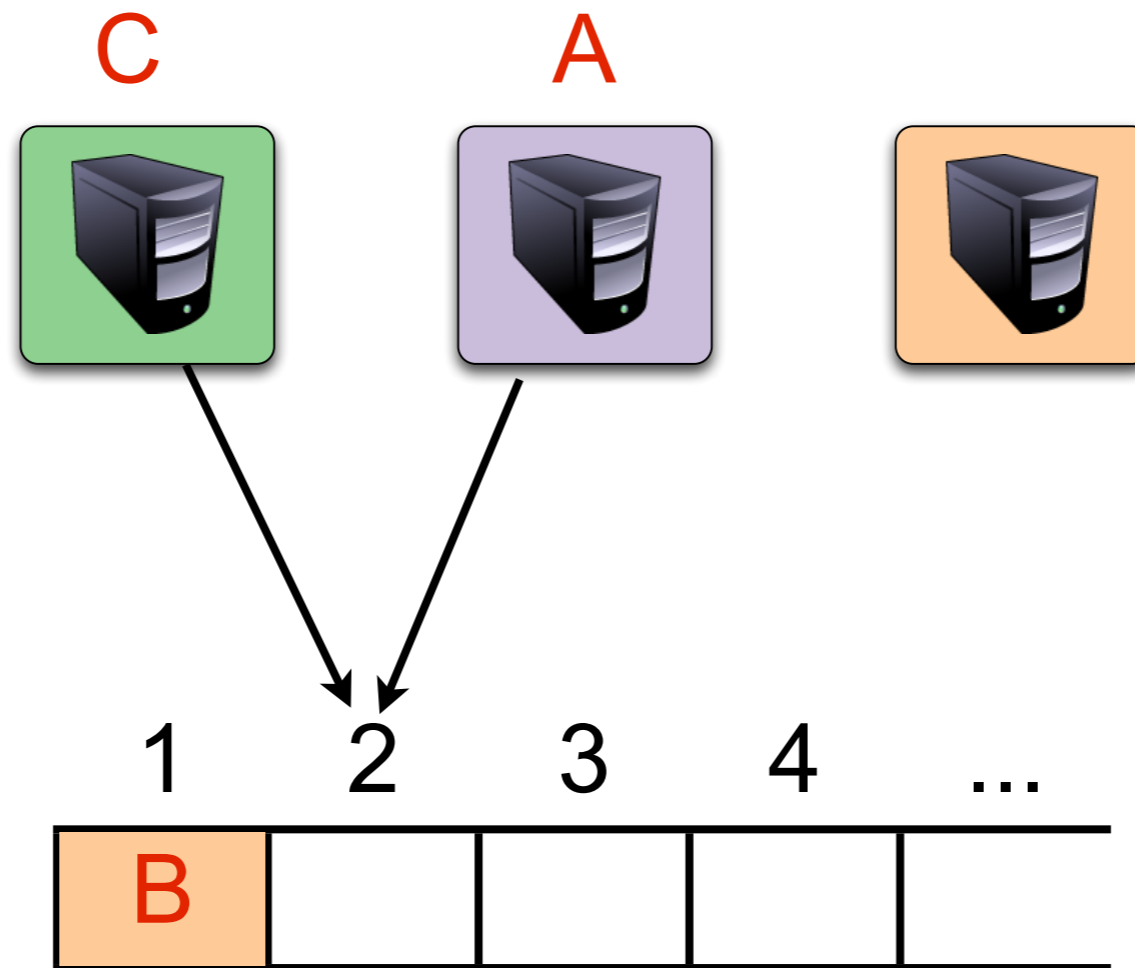
3

4

...



Using Paxos to order commands



Using Paxos to order commands

C



1

2

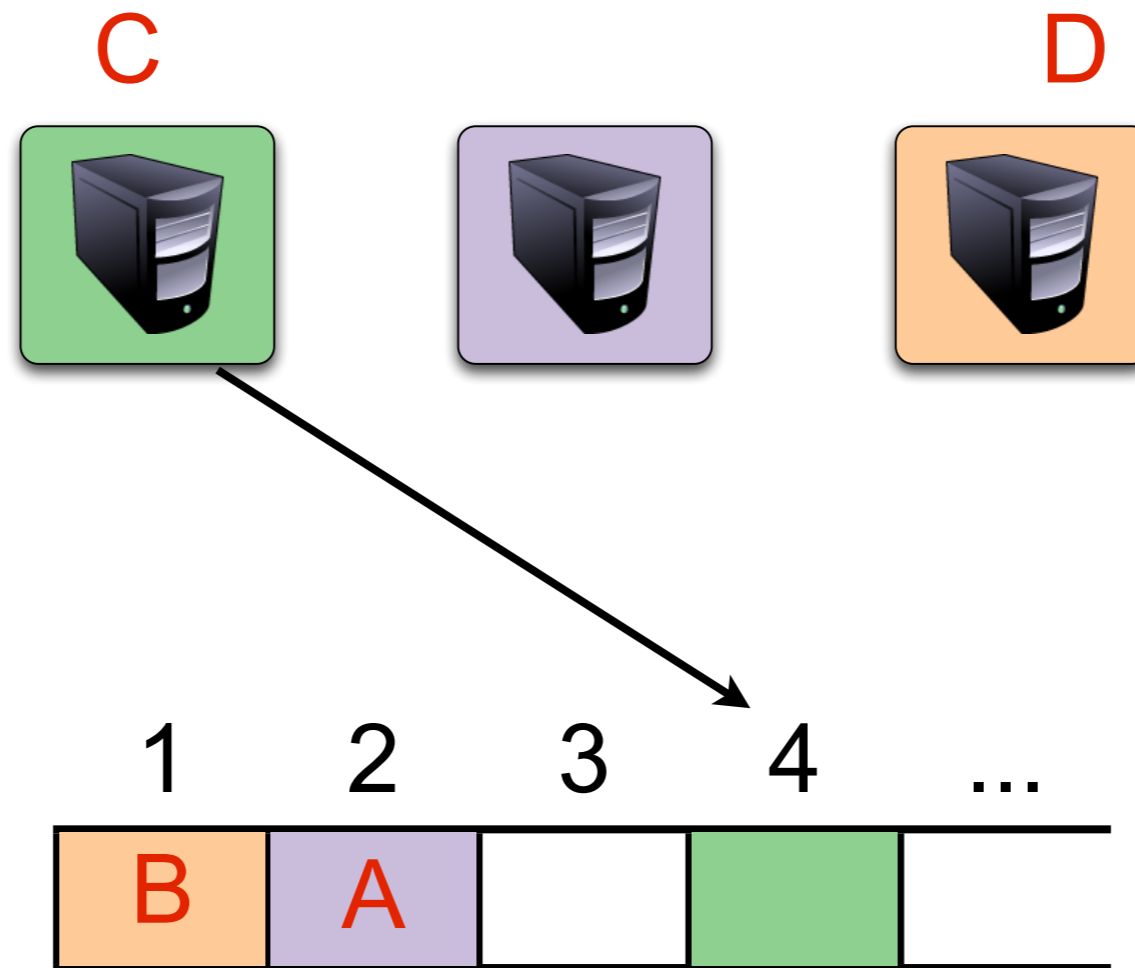
3

4

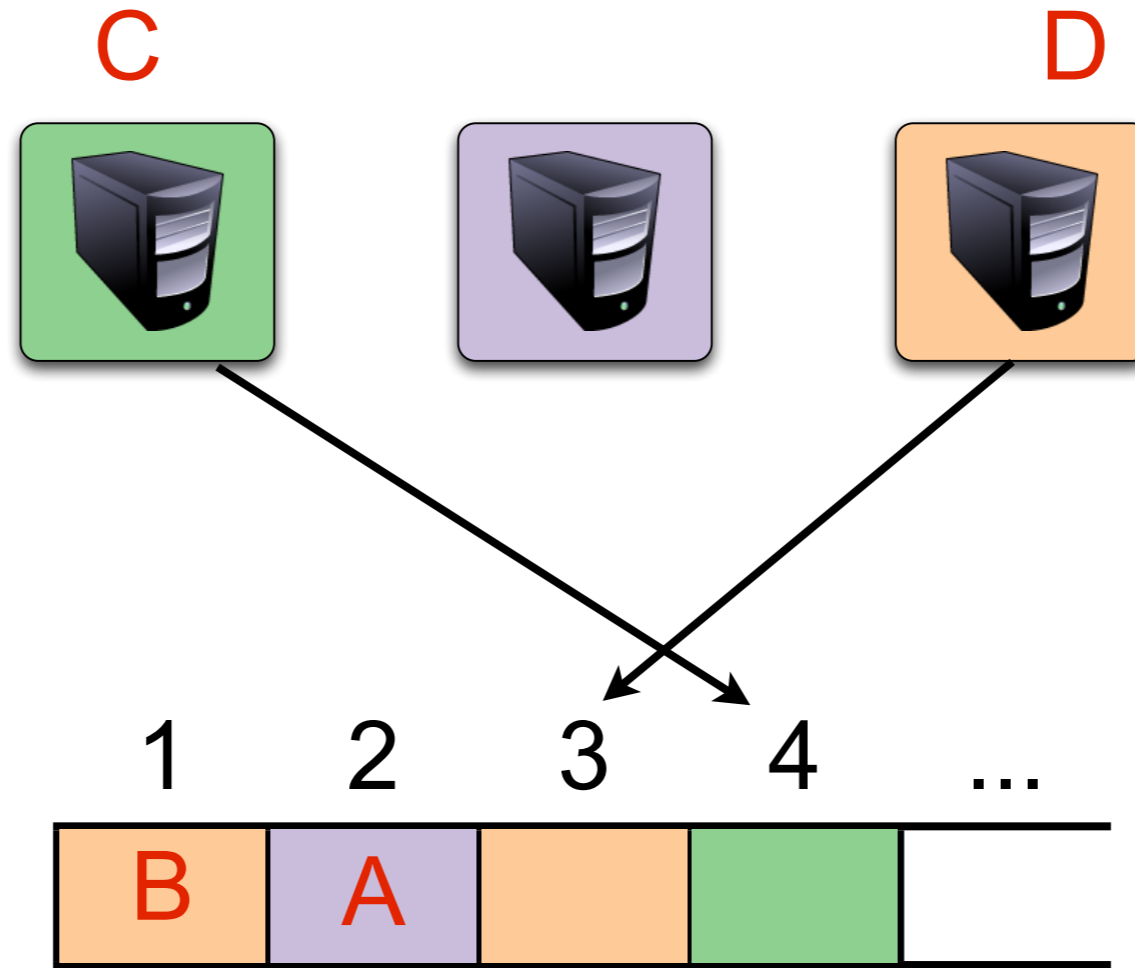
...



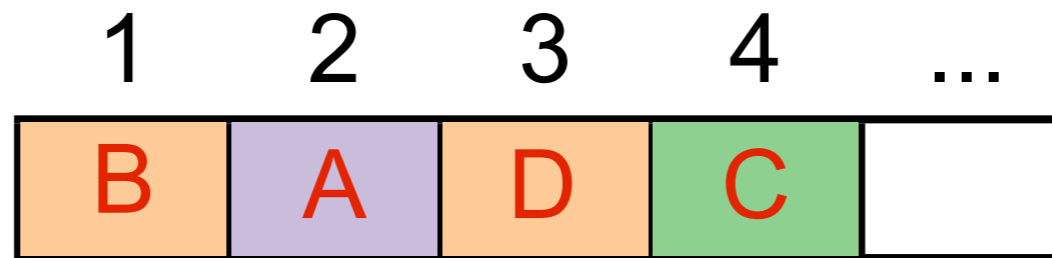
Using Paxos to order commands



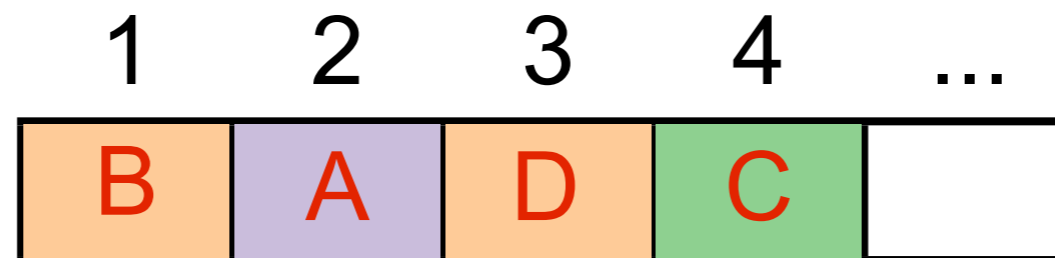
Using Paxos to order commands



Using Paxos to order commands

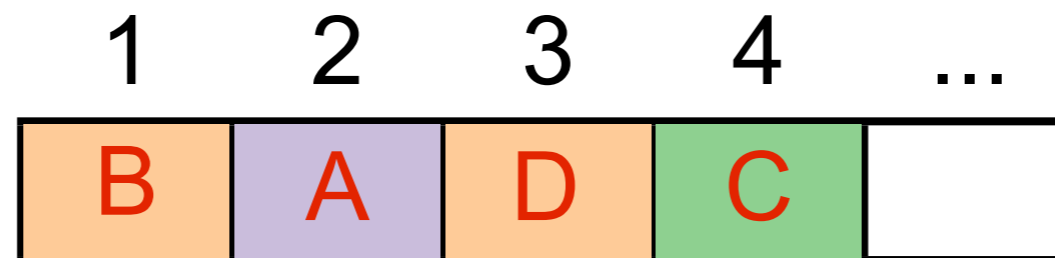


Using Paxos to order commands



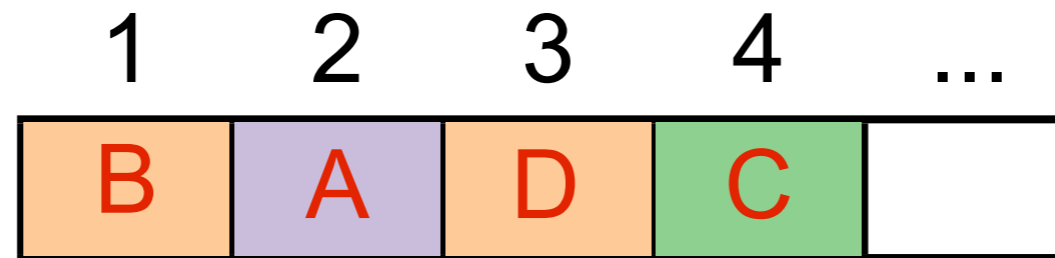
- Choose commands *independently* for each slot

Using Paxos to order commands



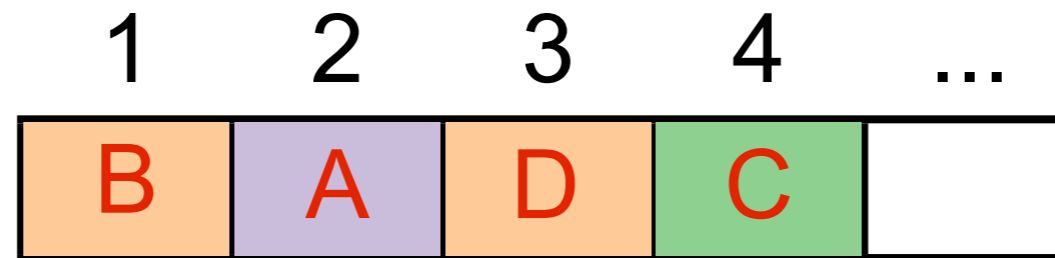
- Choose commands *independently* for each slot
- At least 2 RTTs per slot:

Using Paxos to order commands



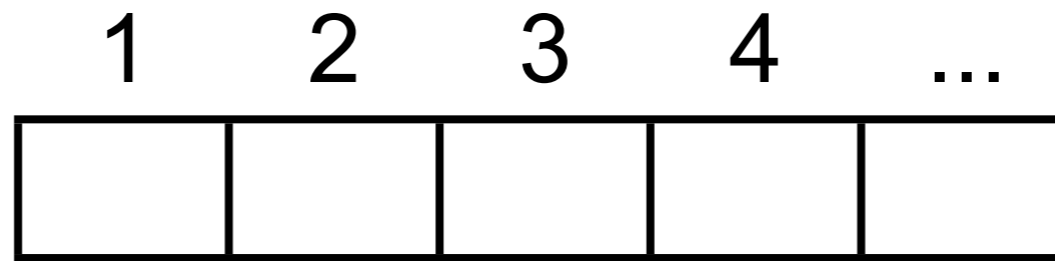
- Choose commands *independently* for each slot
- At least 2 RTTs per slot:
 1. Take ownership of a slot

Using Paxos to order commands

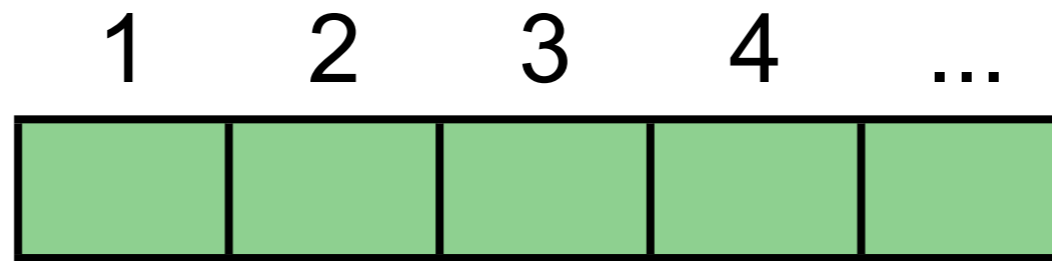


- Choose commands *independently* for each slot
- At least 2 RTTs per slot:
 1. Take ownership of a slot
 2. Propose command

Multi-Paxos



Multi-Paxos



Multi-Paxos

ABCD



1

2

3

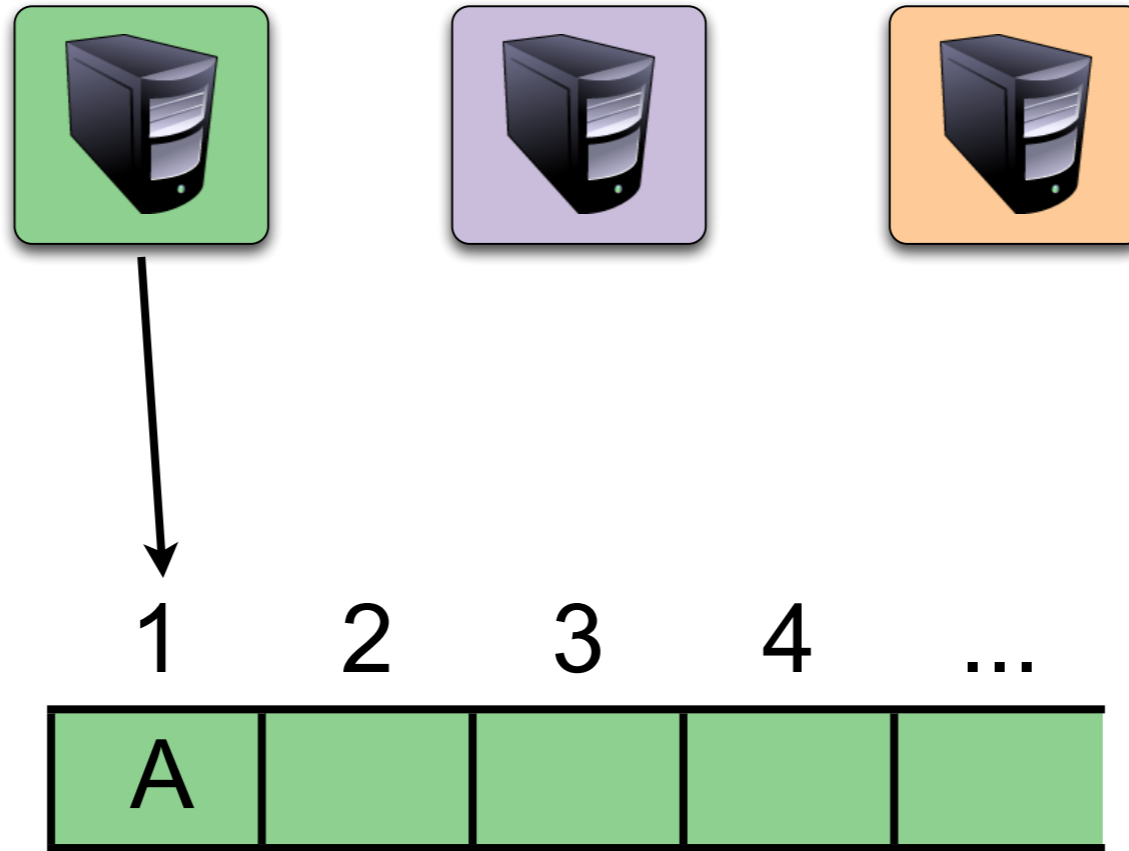
4

...

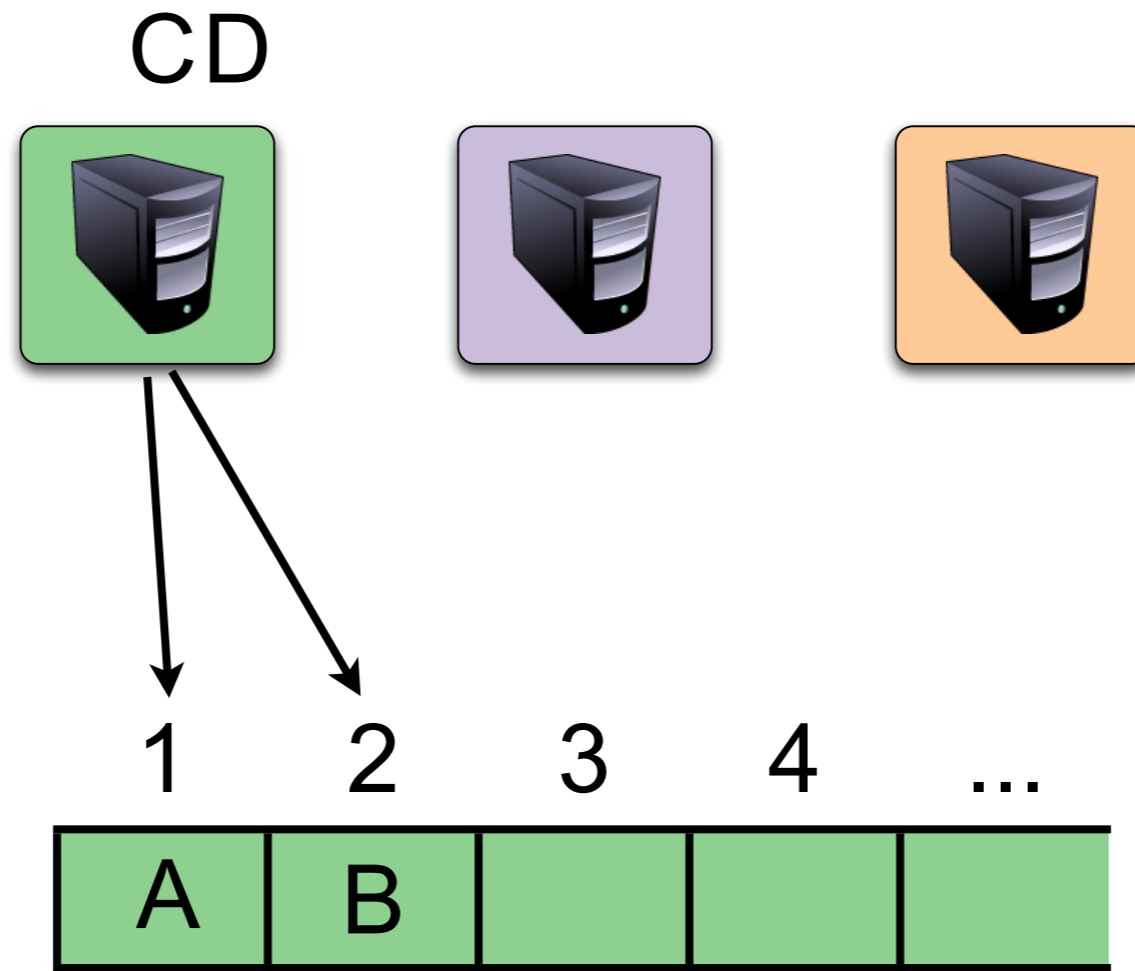


Multi-Paxos

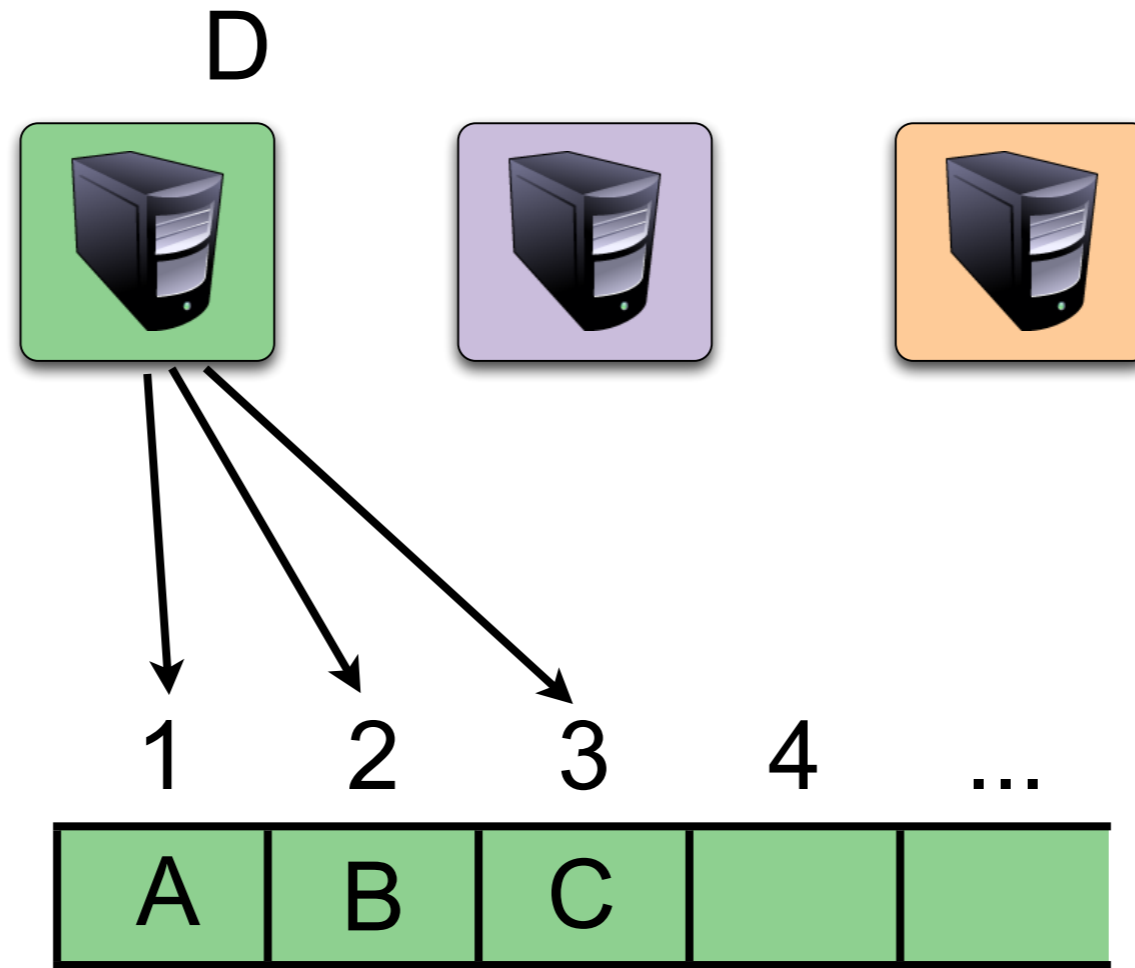
BCD



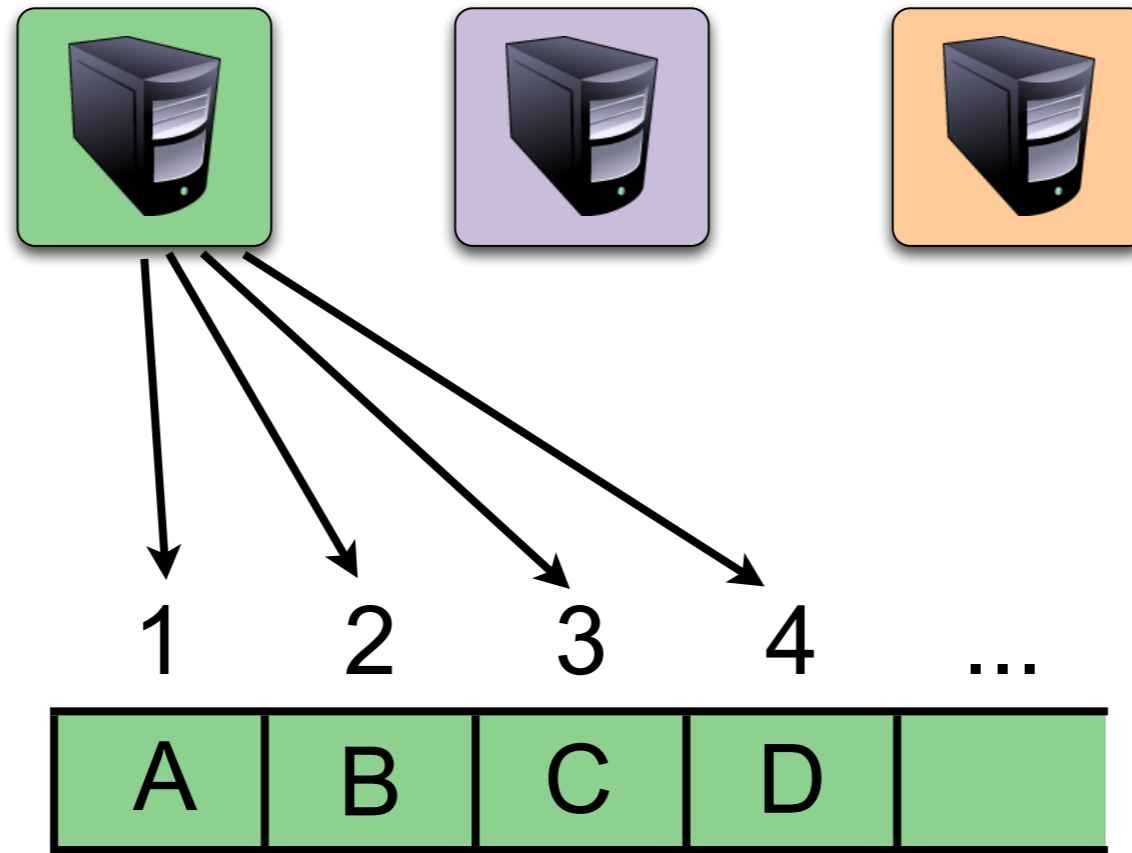
Multi-Paxos



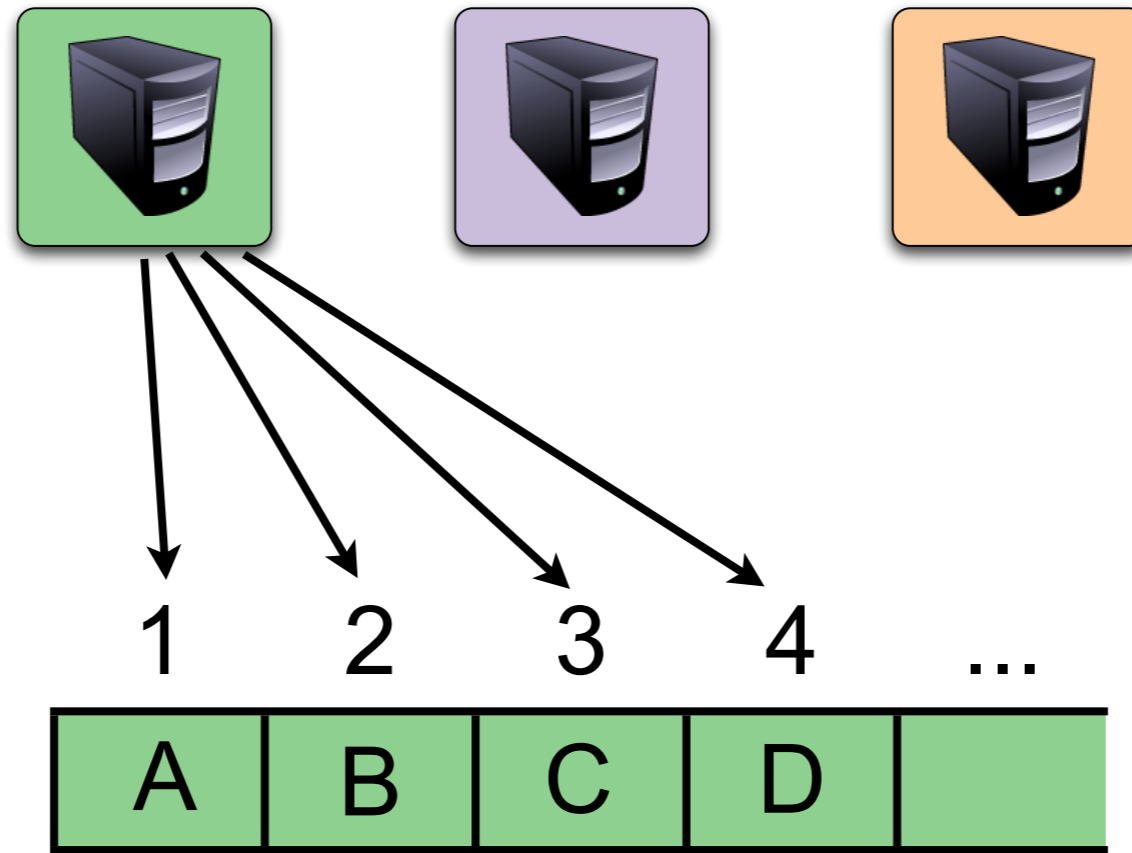
Multi-Paxos



Multi-Paxos

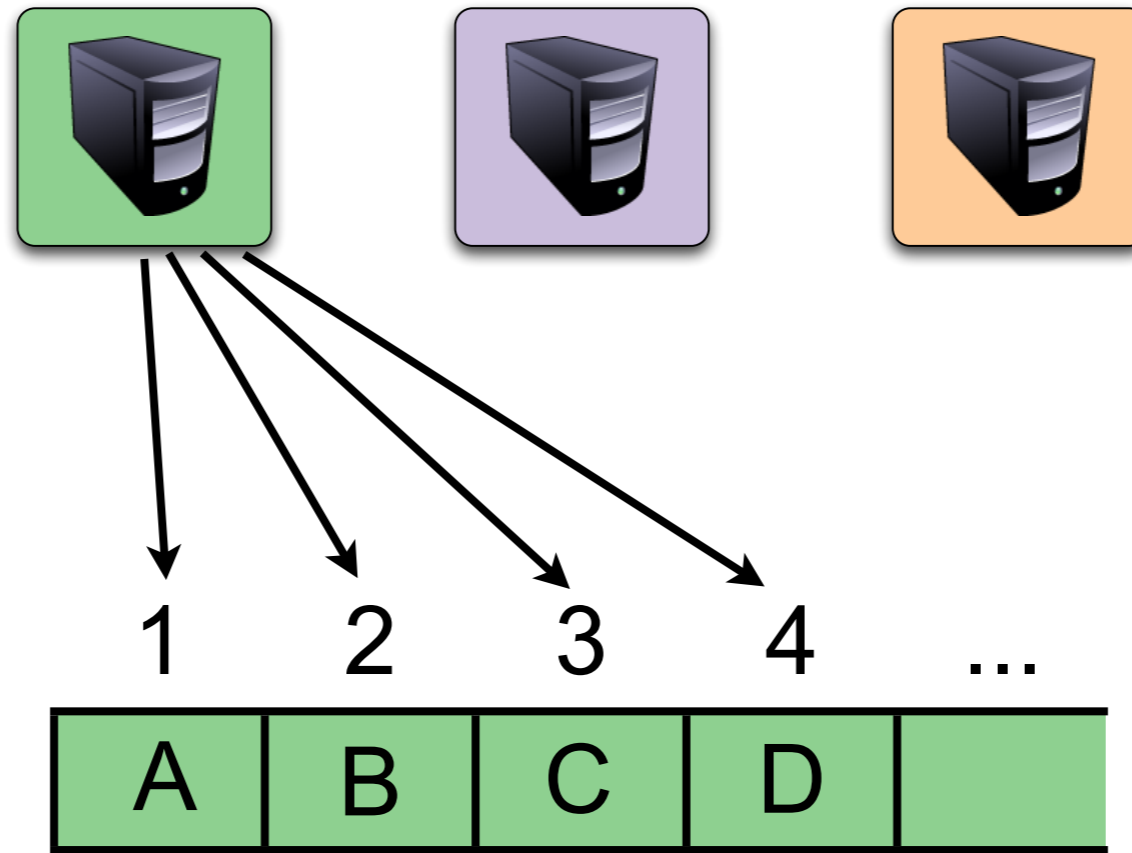


Multi-Paxos



- 1 RTT to commit

Multi-Paxos



- 1 RTT to commit
- Bottleneck for performance and availability

Can we have it all?

Can we have it all?

- High throughput, low latency

Can we have it all?

- High throughput, low latency
- Constant availability

Can we have it all?

- High throughput, low latency
- Constant availability
- Distribute load evenly across all replicas

Can we have it all?

- High throughput, low latency
- Constant availability
- Distribute load evenly across all replicas
- Use fastest replicas

Can we have it all?

- High throughput, low latency
- Constant availability
- Distribute load evenly across all replicas
- Use fastest replicas
- Use closest (lowest latency) replicas

Can we have it all?

- High throughput, low latency
- Constant availability
- Distribute load evenly across all replicas
- Use fastest replicas
- Use closest (lowest latency) replicas

Paxos

Can we have it all?

- High throughput, low latency
- Constant availability
- Distribute load evenly across all replicas
- Use fastest replicas
- Use closest (lowest latency) replicas

Multi-
Paxos

Can we have it all?

- High throughput, low latency
- Constant availability
- Distribute load evenly across all replicas
- Use fastest replicas
- Use closest (lowest latency) replicas

Multi-
Paxos

Egalitarian Paxos
(EPaxos)

EPaxos is all about ordering

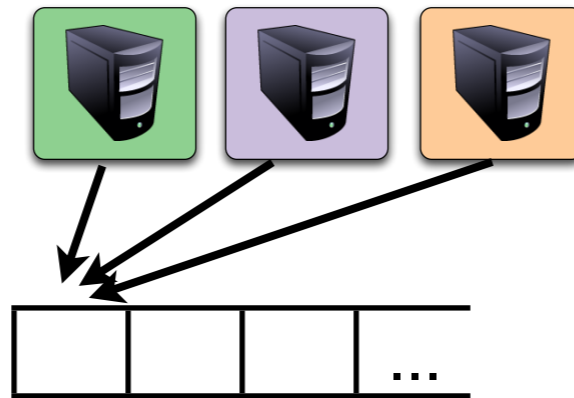
EPaxos is all about ordering

Previous strategies:

EPaxos is all about ordering

Previous strategies:

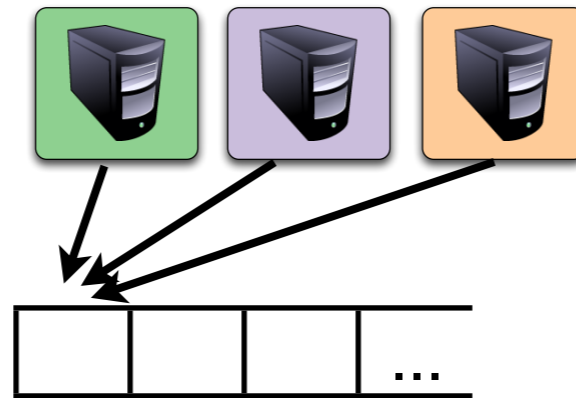
Contend
for slots



EPaxos is all about ordering

Previous strategies:

Contend
for slots

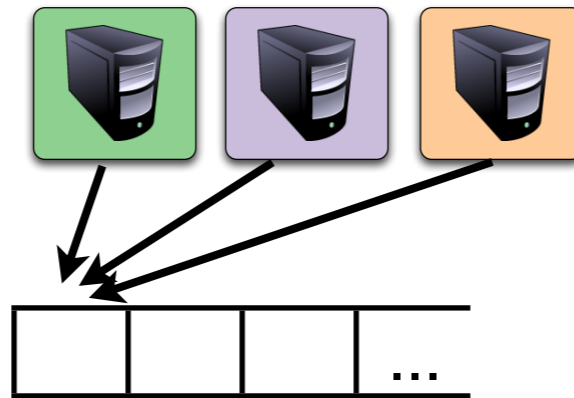


Paxos

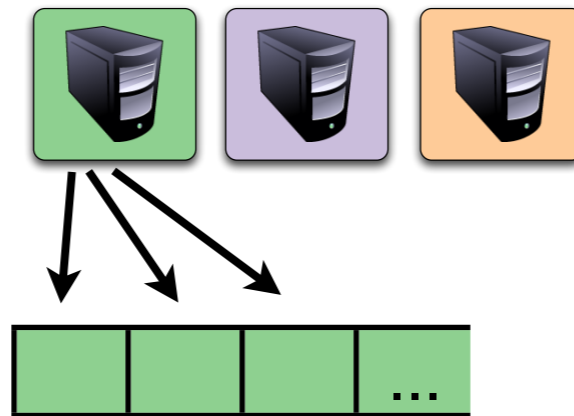
EPaxos is all about ordering

Previous strategies:

Contend
for slots



One replica
decides

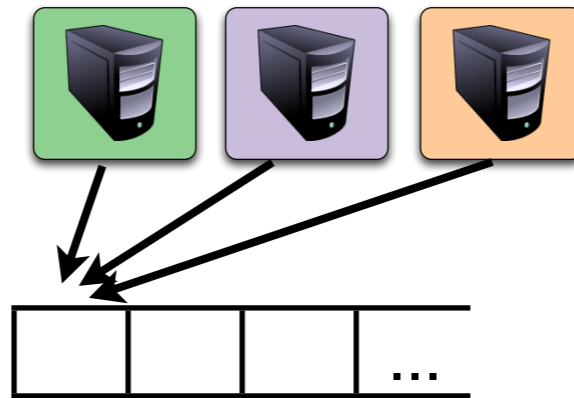


Paxos

EPaxos is all about ordering

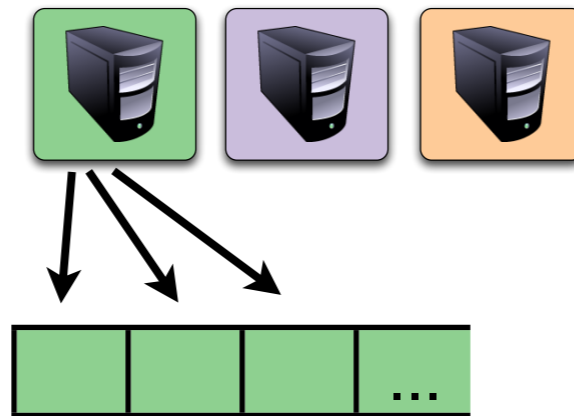
Previous strategies:

Contend
for slots



Paxos

One replica
decides

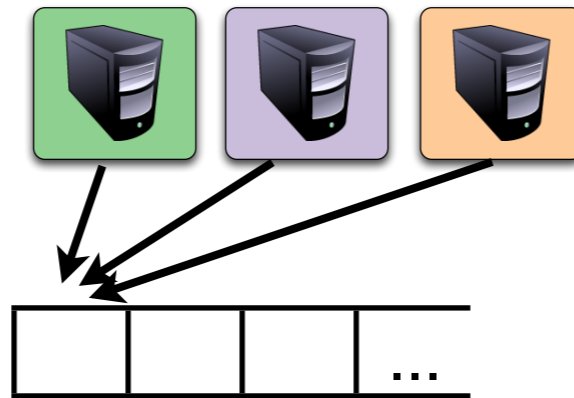


Multi-Paxos, Fast Paxos,
Generalized Paxos

EPaxos is all about ordering

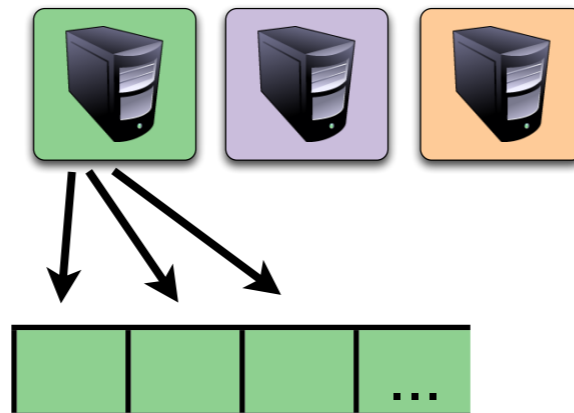
Previous strategies:

Contend for slots



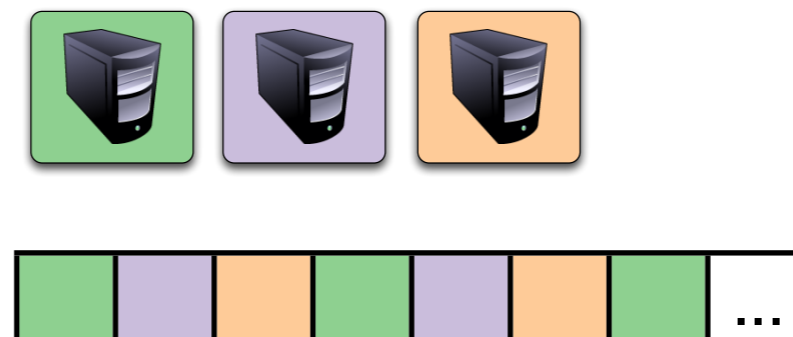
Paxos

One replica decides



Multi-Paxos, Fast Paxos, Generalized Paxos

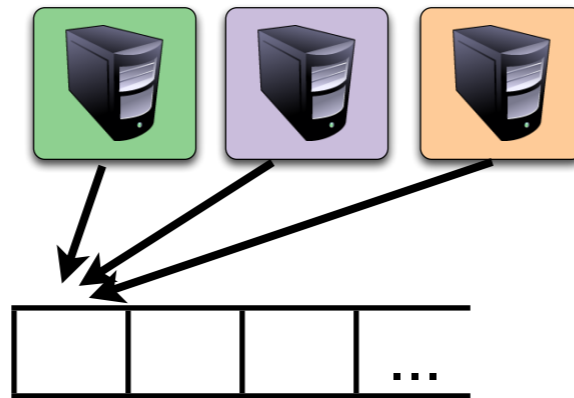
Take turns round-robin



EPaxos is all about ordering

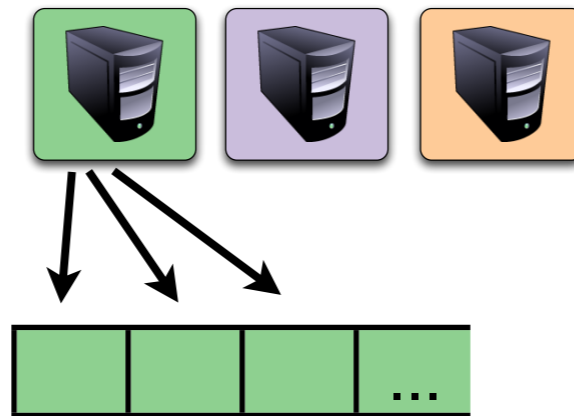
Previous strategies:

Contend for slots



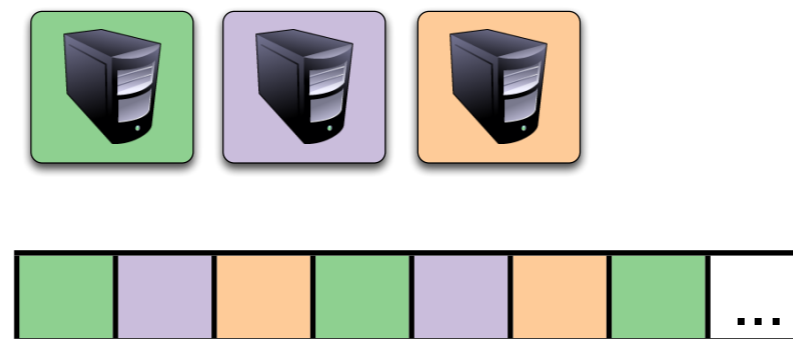
Paxos

One replica decides



Multi-Paxos, Fast Paxos, Generalized Paxos

Take turns round-robin



Mencius

EPaxos intuition



EPaxos intuition



1

2

3

4

...

EPaxos intuition

A



1

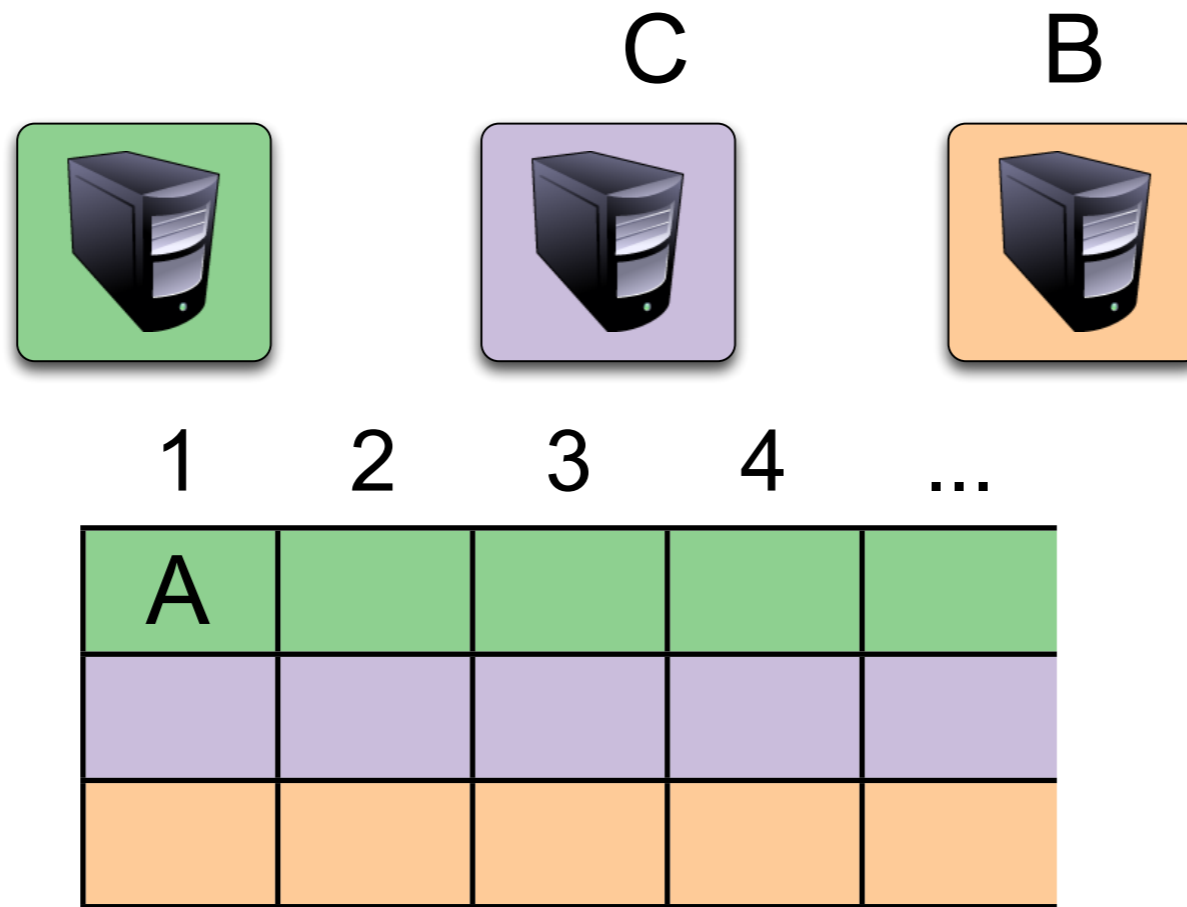
2

3

4

...

EPaxos intuition



EPaxos intuition

D



1

2

3

4

...

A				
C				
B				

EPaxos intuition



1

2

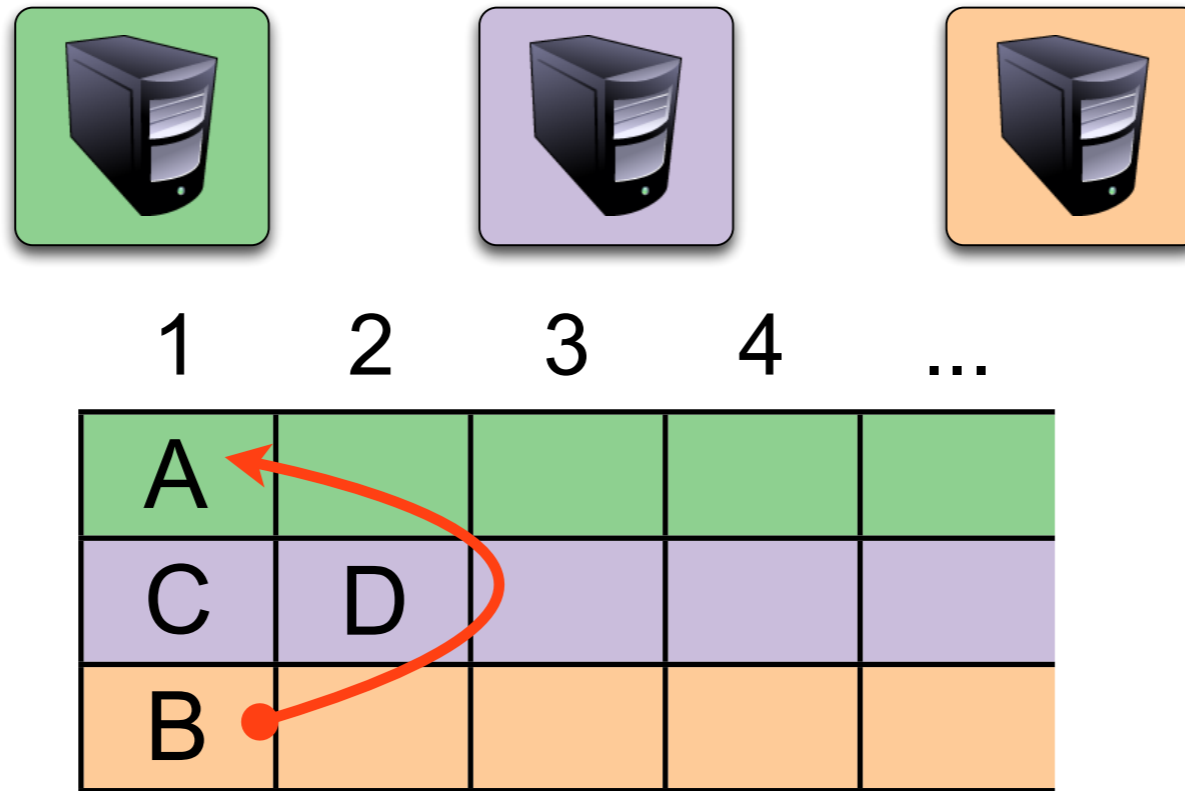
3

4

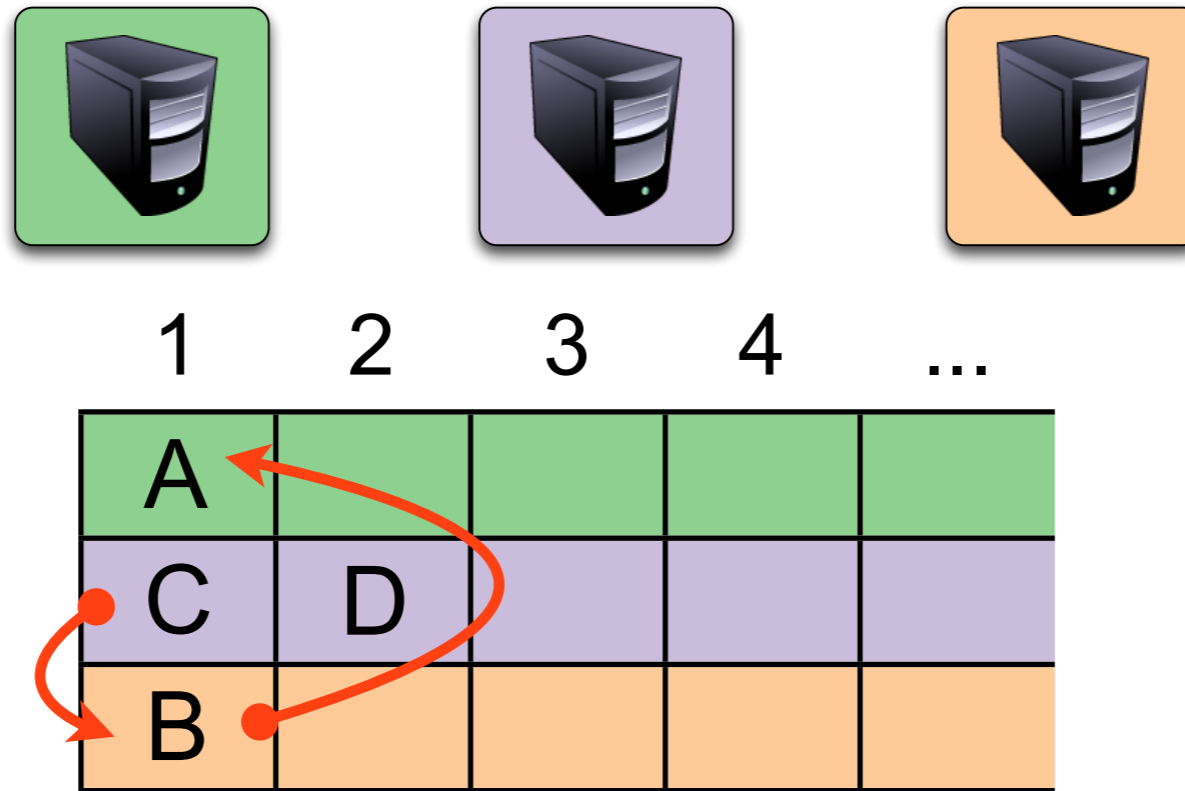
...

A				
C	D			
B				

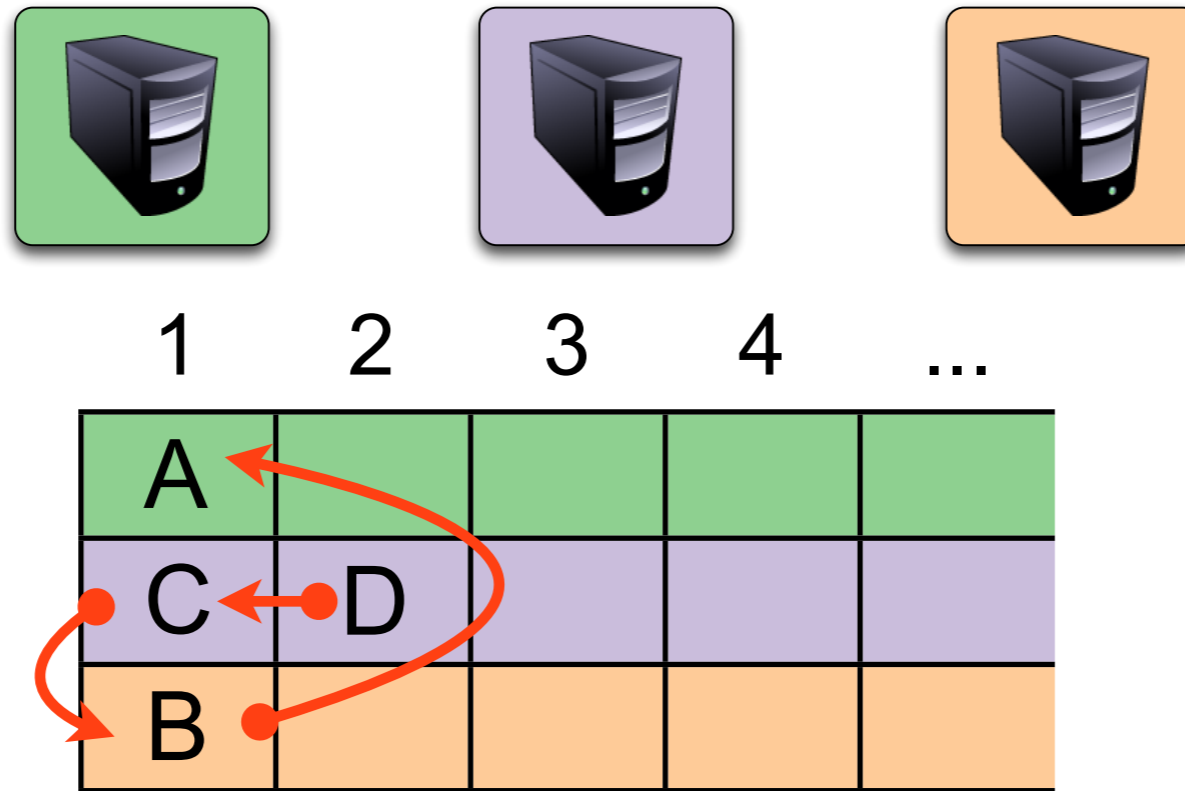
EPaxos intuition



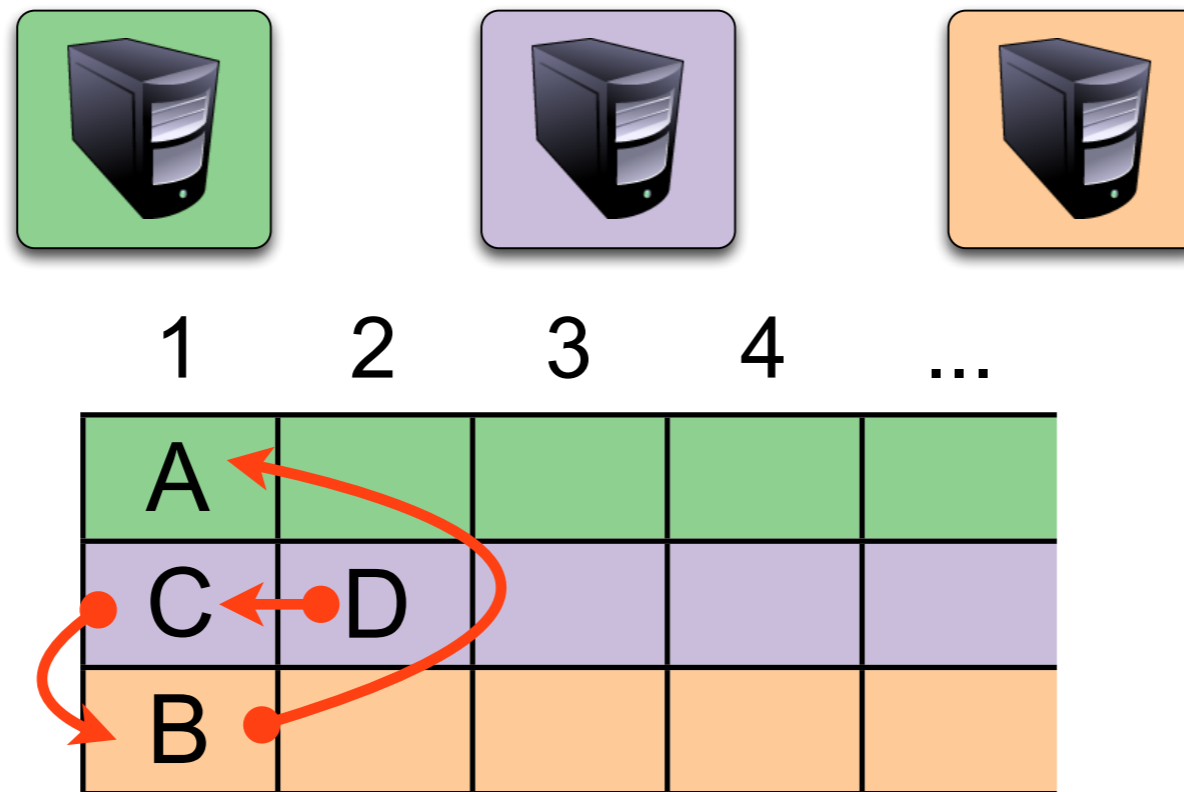
EPaxos intuition



EPaxos intuition

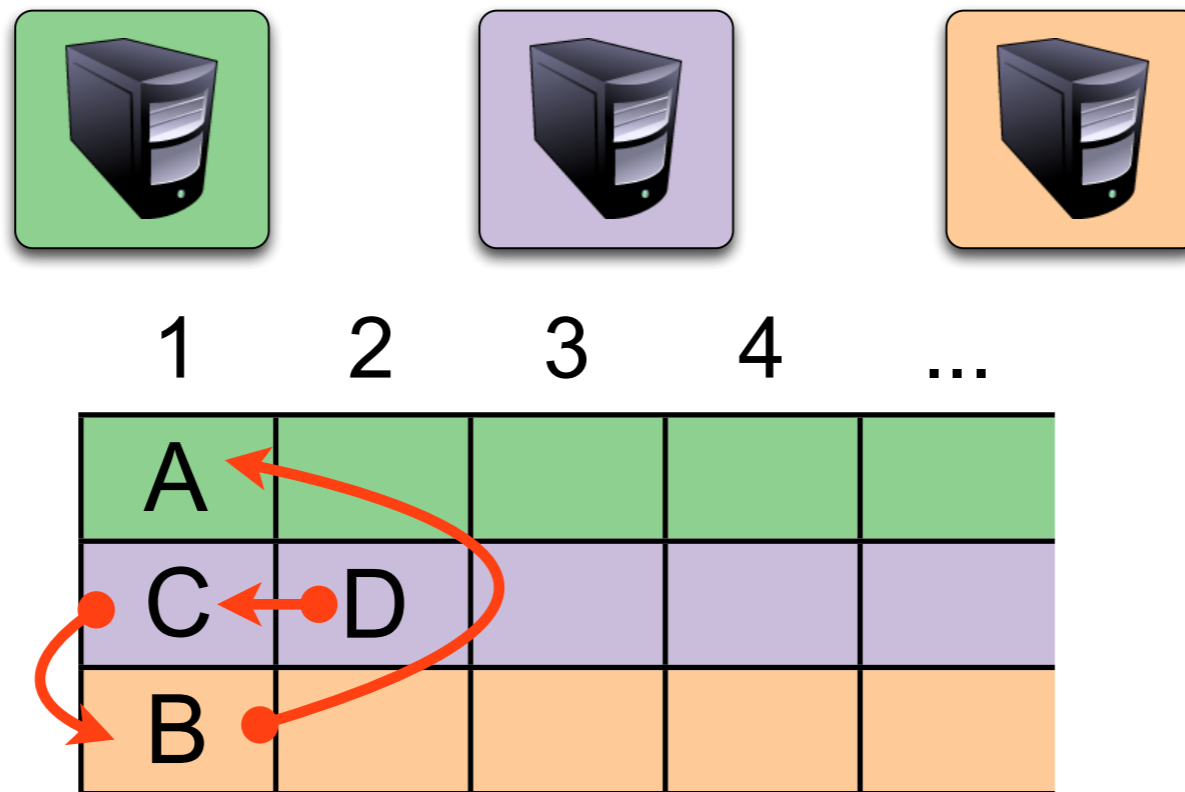


EPaxos intuition



After commit
@ each replica

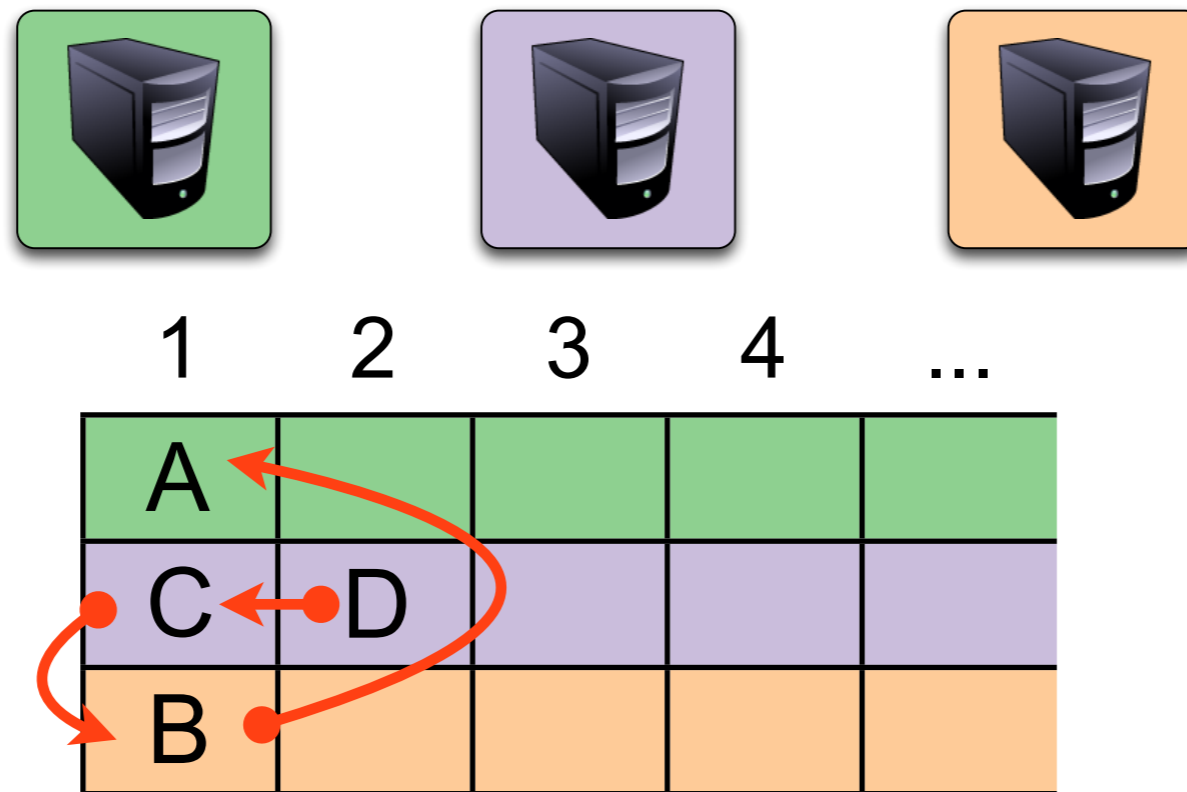
EPaxos intuition



After commit
@ each replica



EPaxos intuition

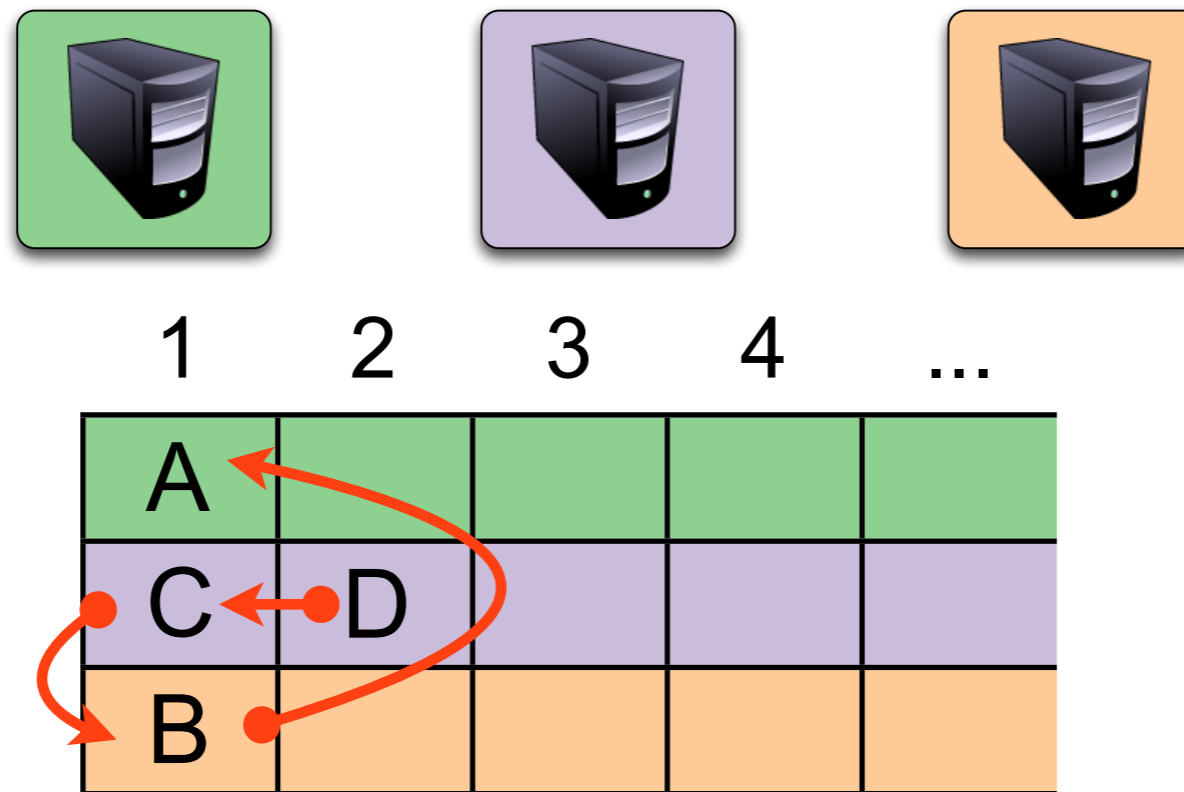


After commit
@ each replica



- Load balance (every replica is a leader)

EPaxos intuition



After commit
@ each replica



- Load balance (every replica is a leader)
- EPaxos can choose *any* quorum for each command

EPaxos commit protocol

R1 _____

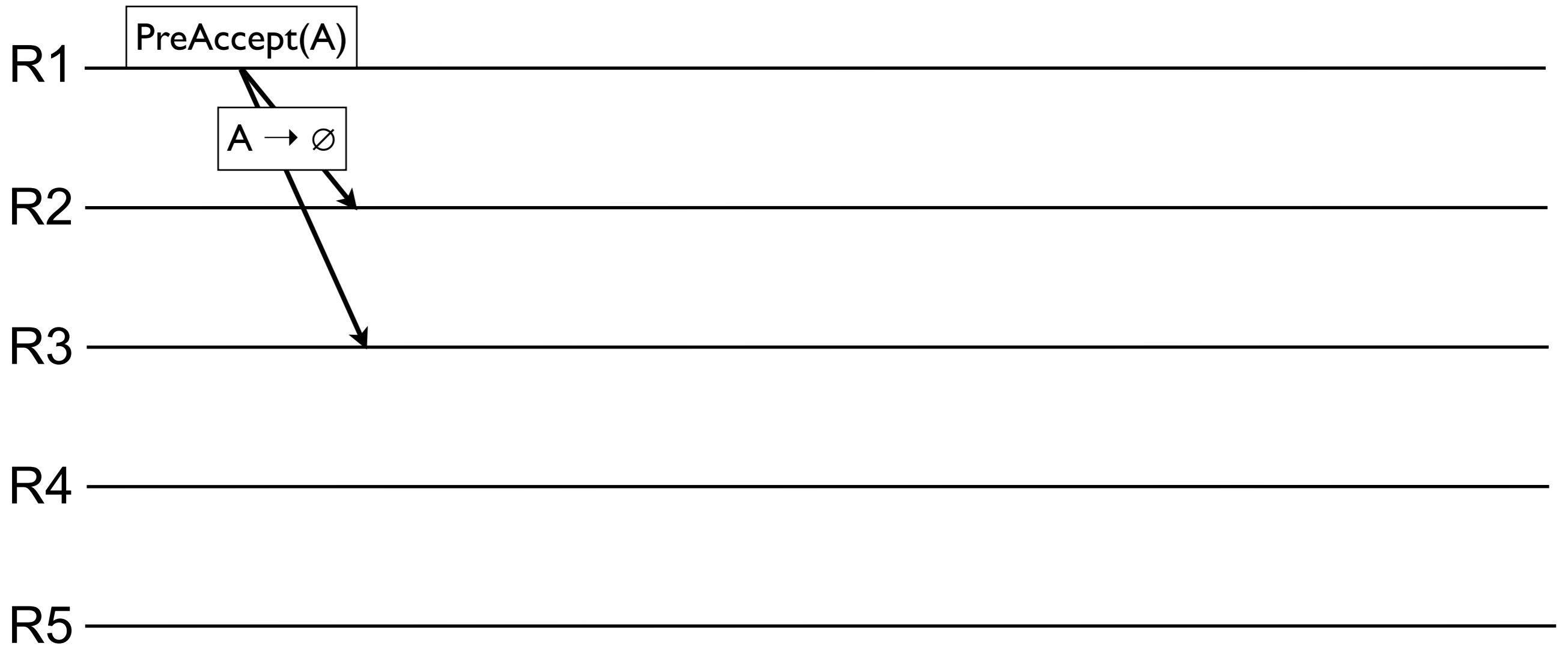
R2 _____

R3 _____

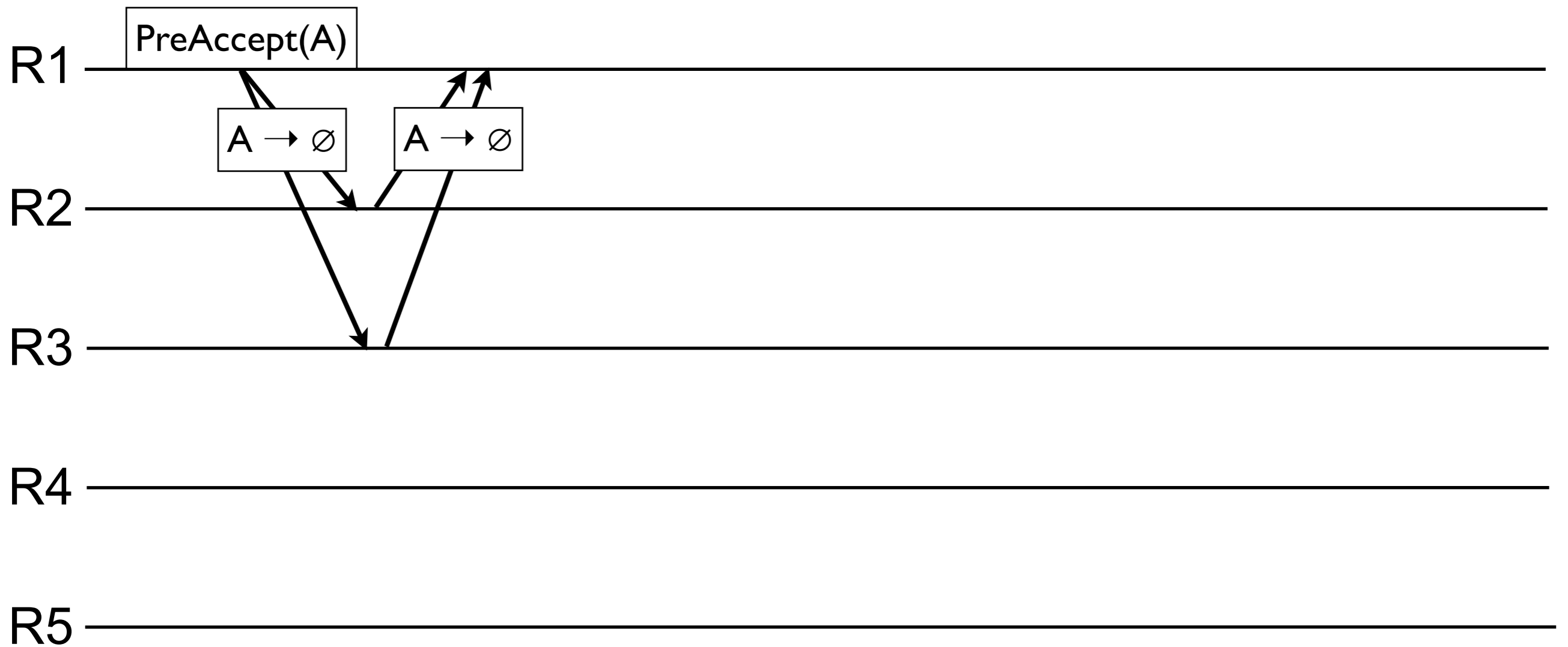
R4 _____

R5 _____

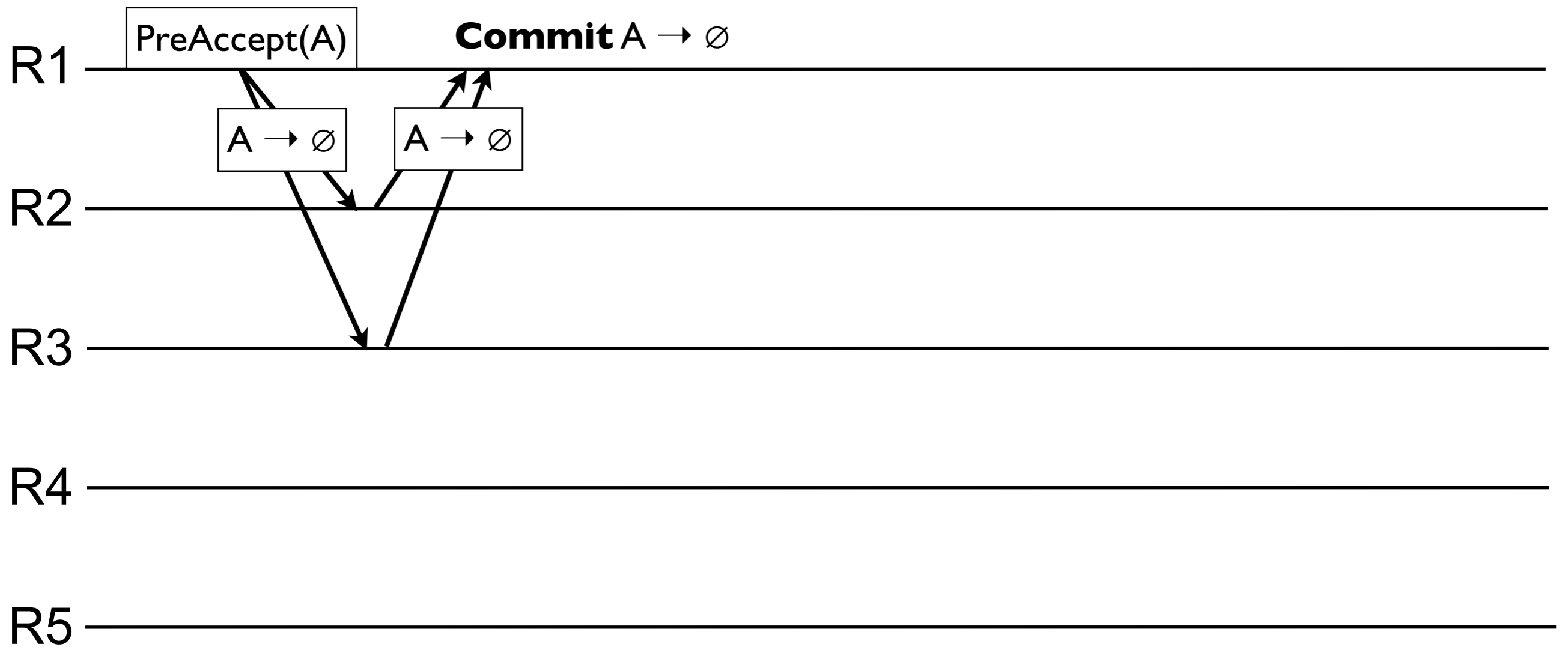
EPaxos commit protocol



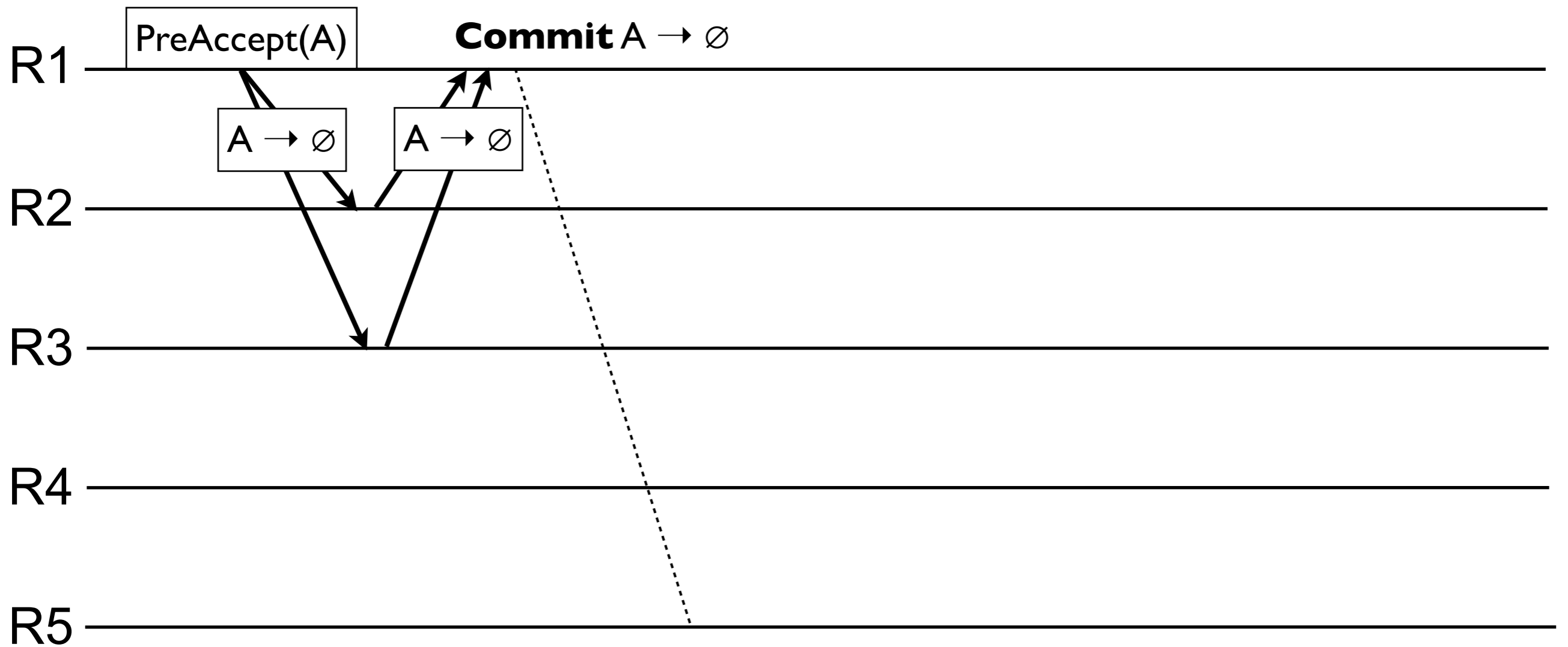
EPaxos commit protocol



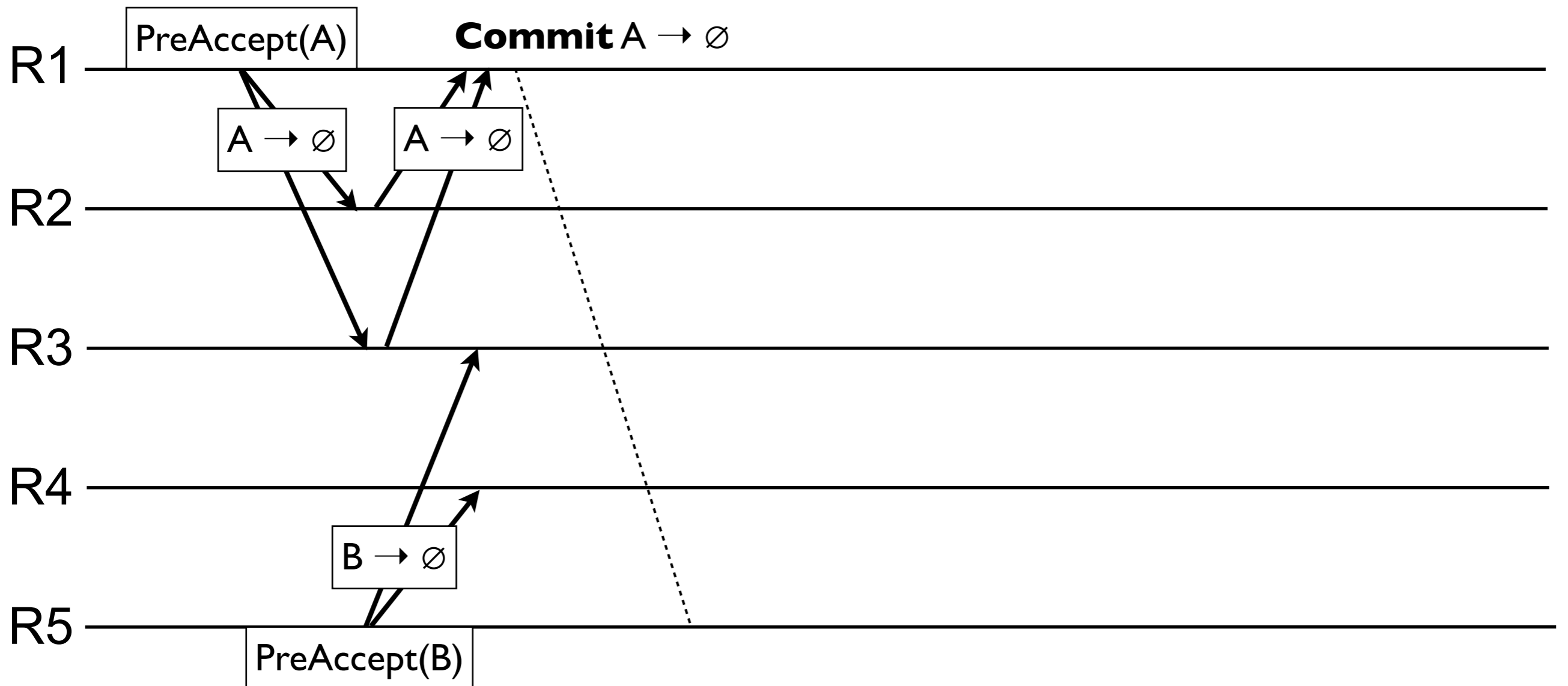
EPaxos commit protocol



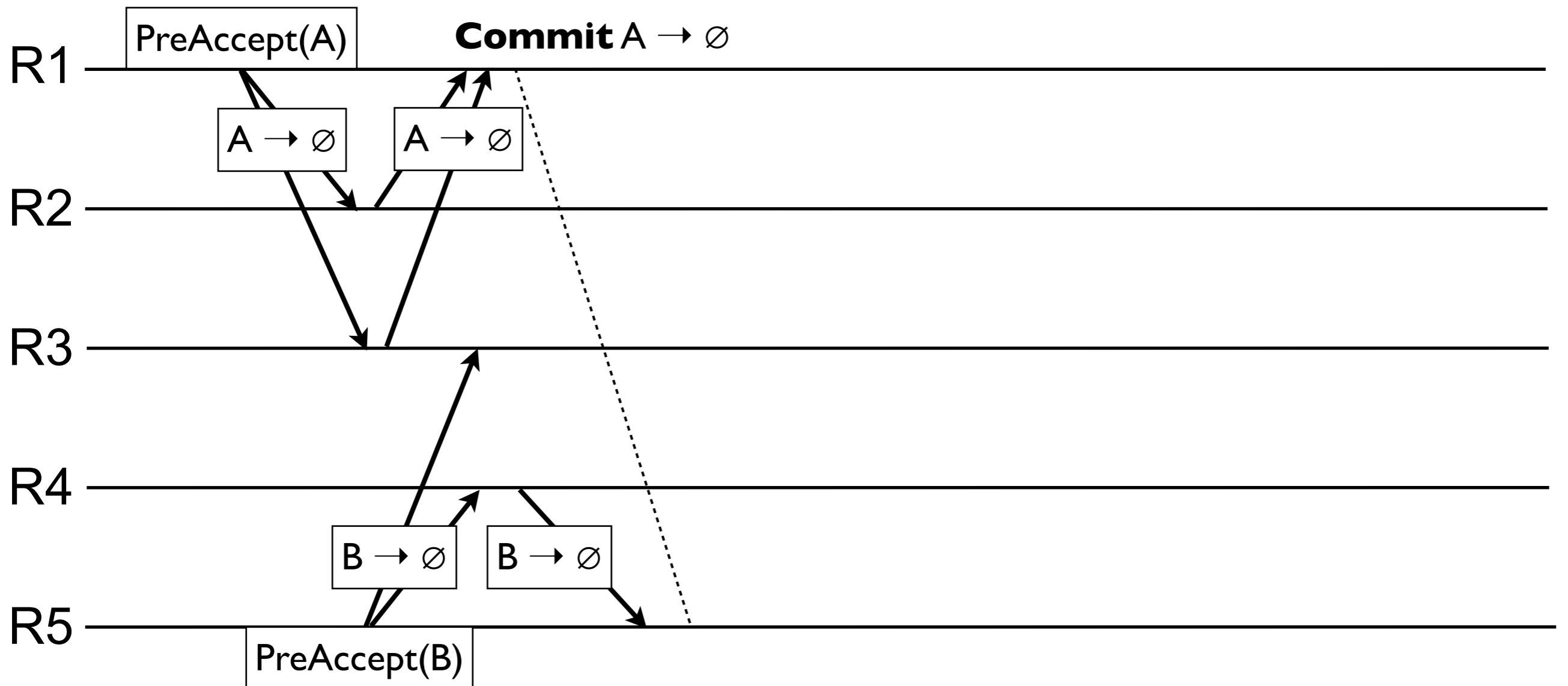
EPaxos commit protocol



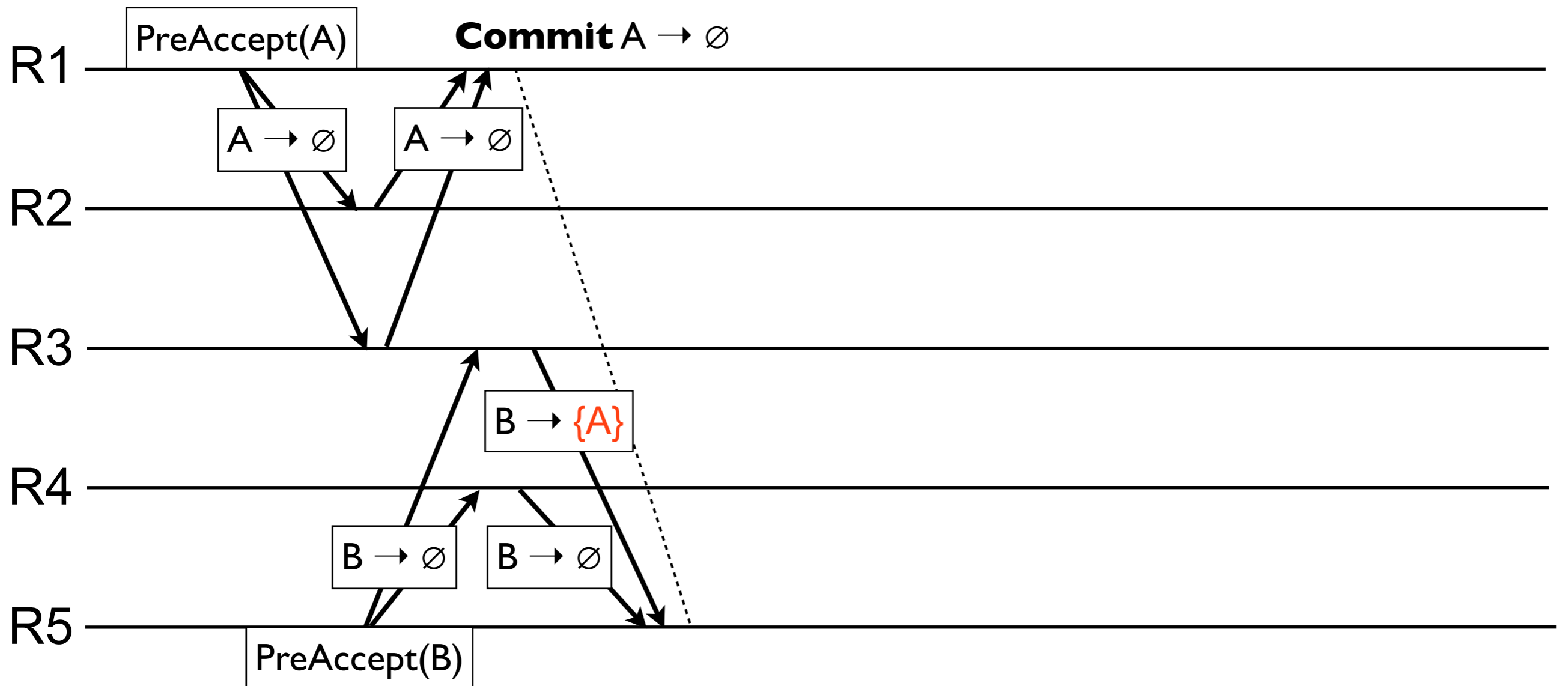
EPaxos commit protocol



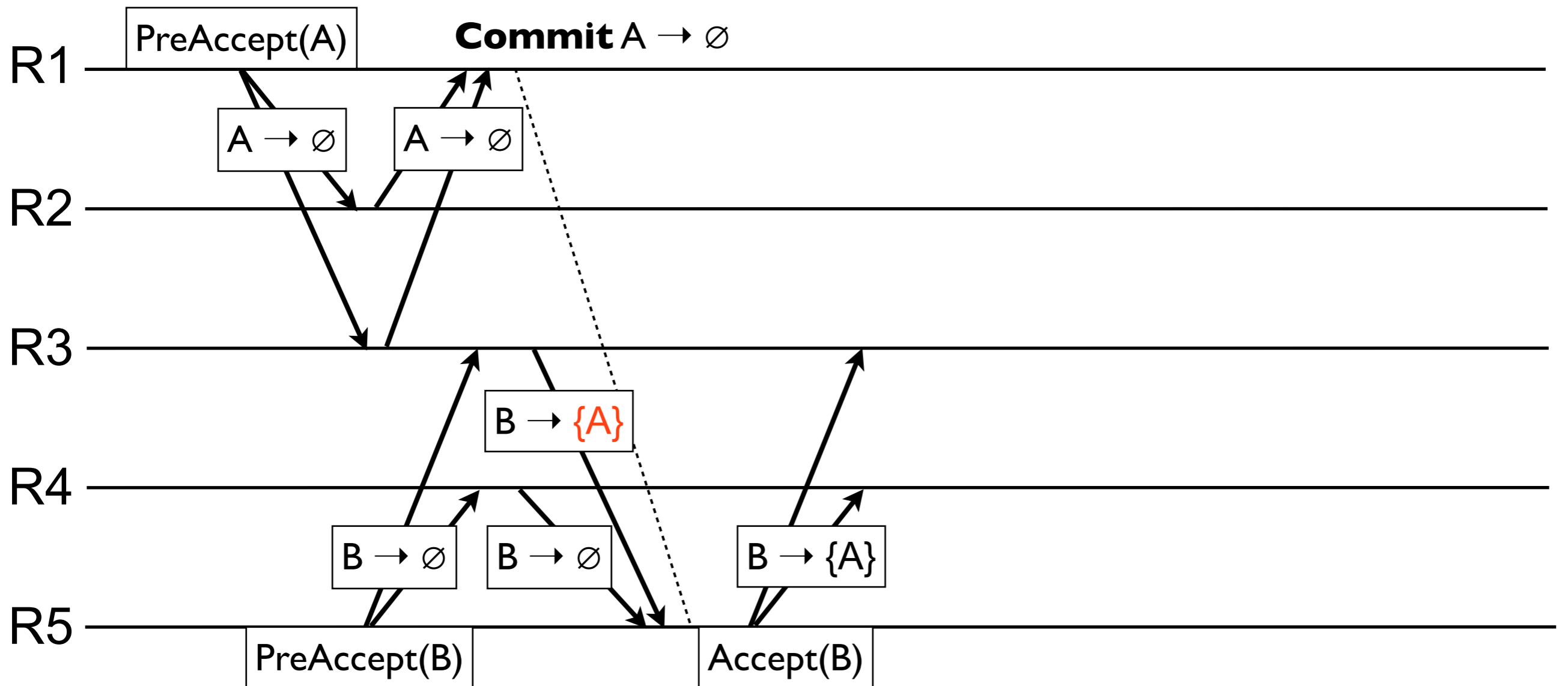
EPaxos commit protocol



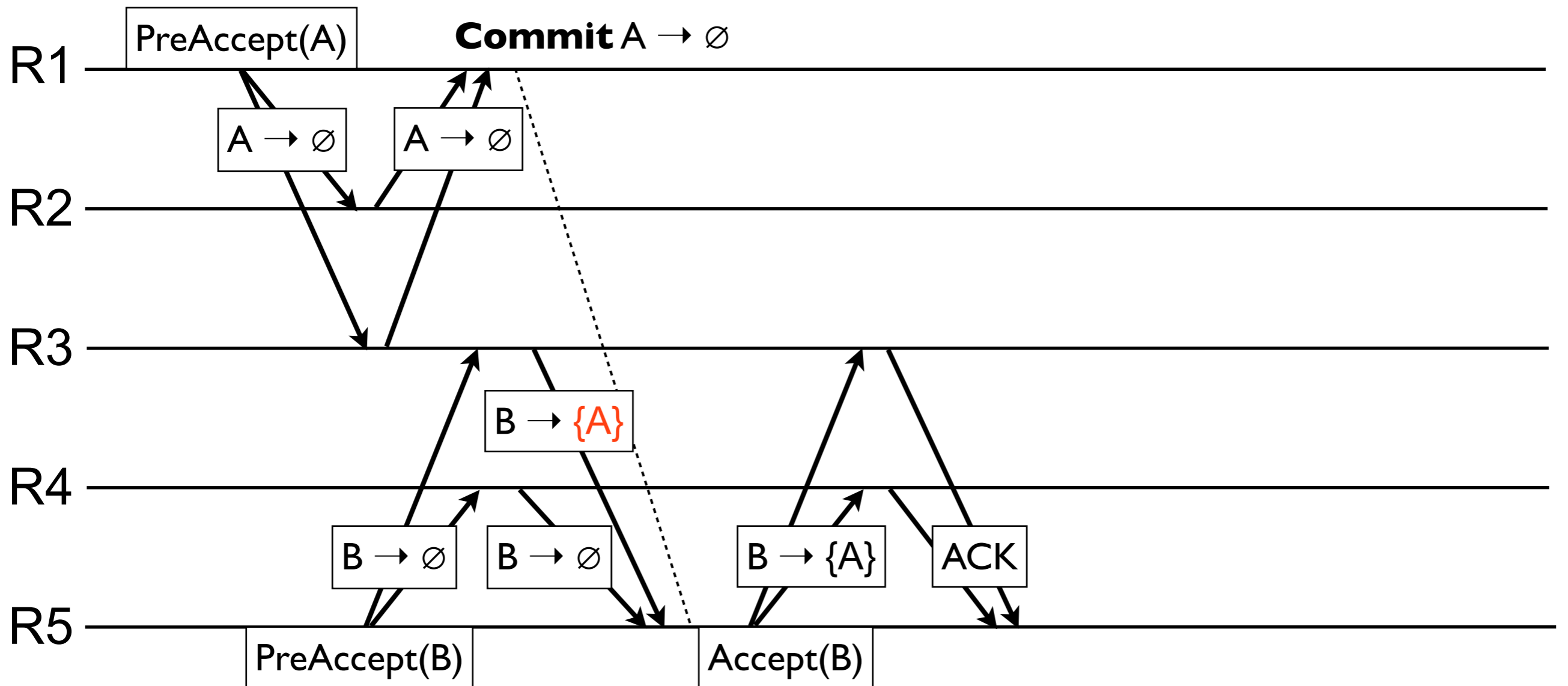
EPaxos commit protocol



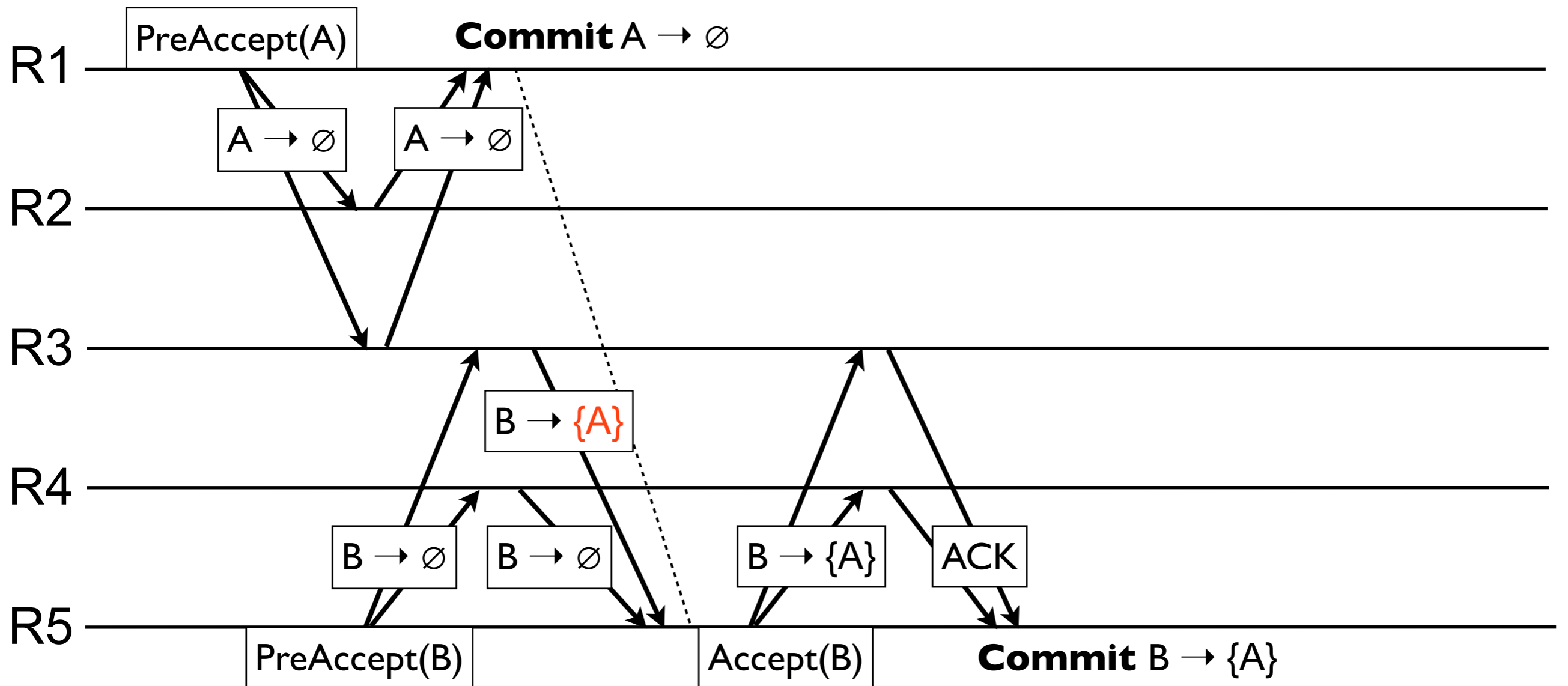
EPaxos commit protocol



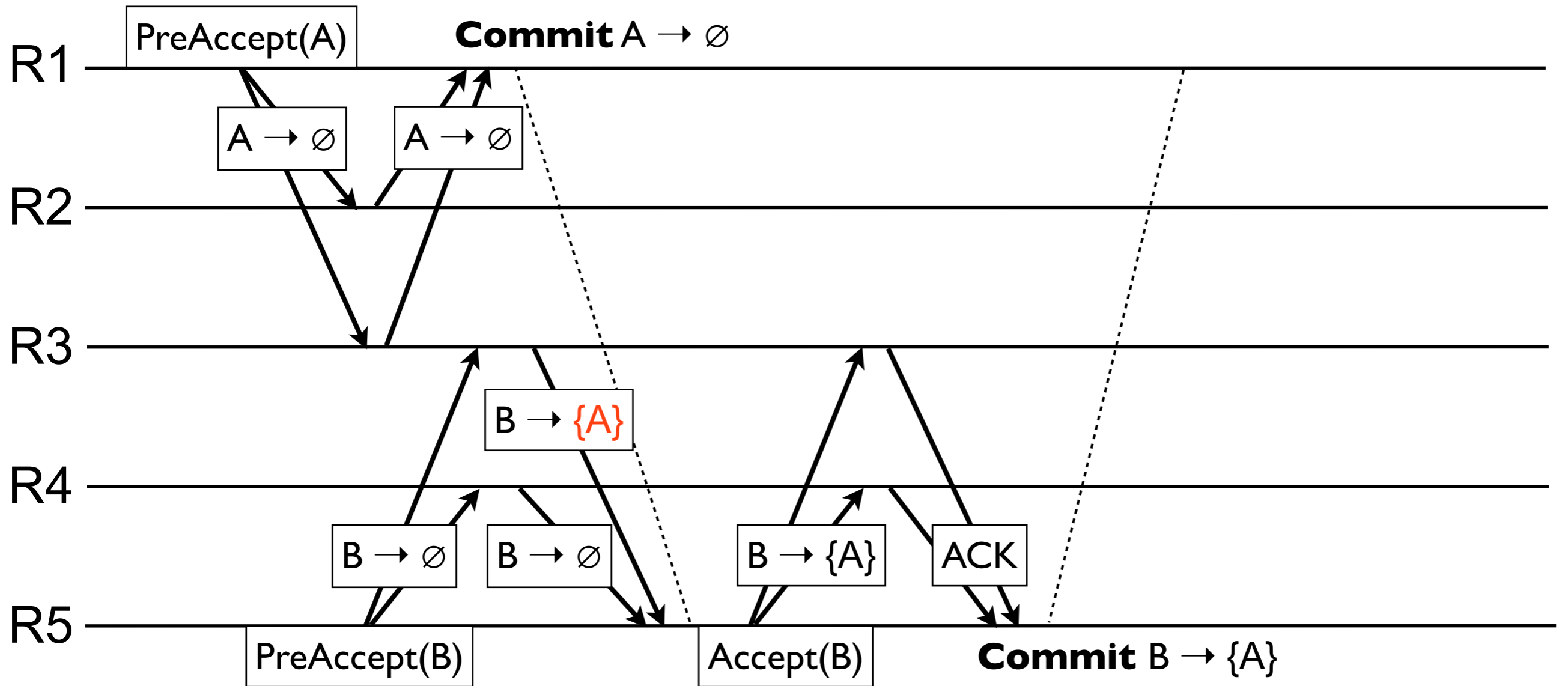
EPaxos commit protocol



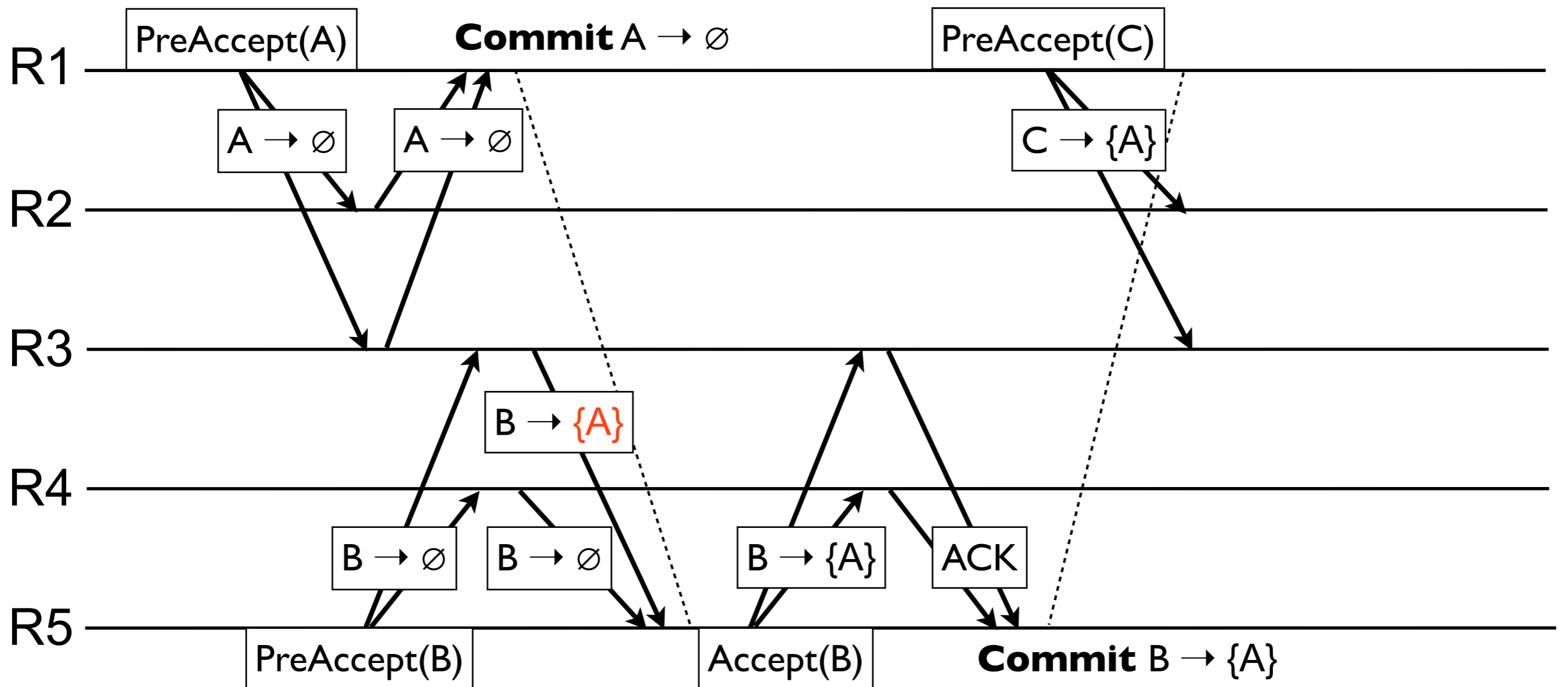
EPaxos commit protocol



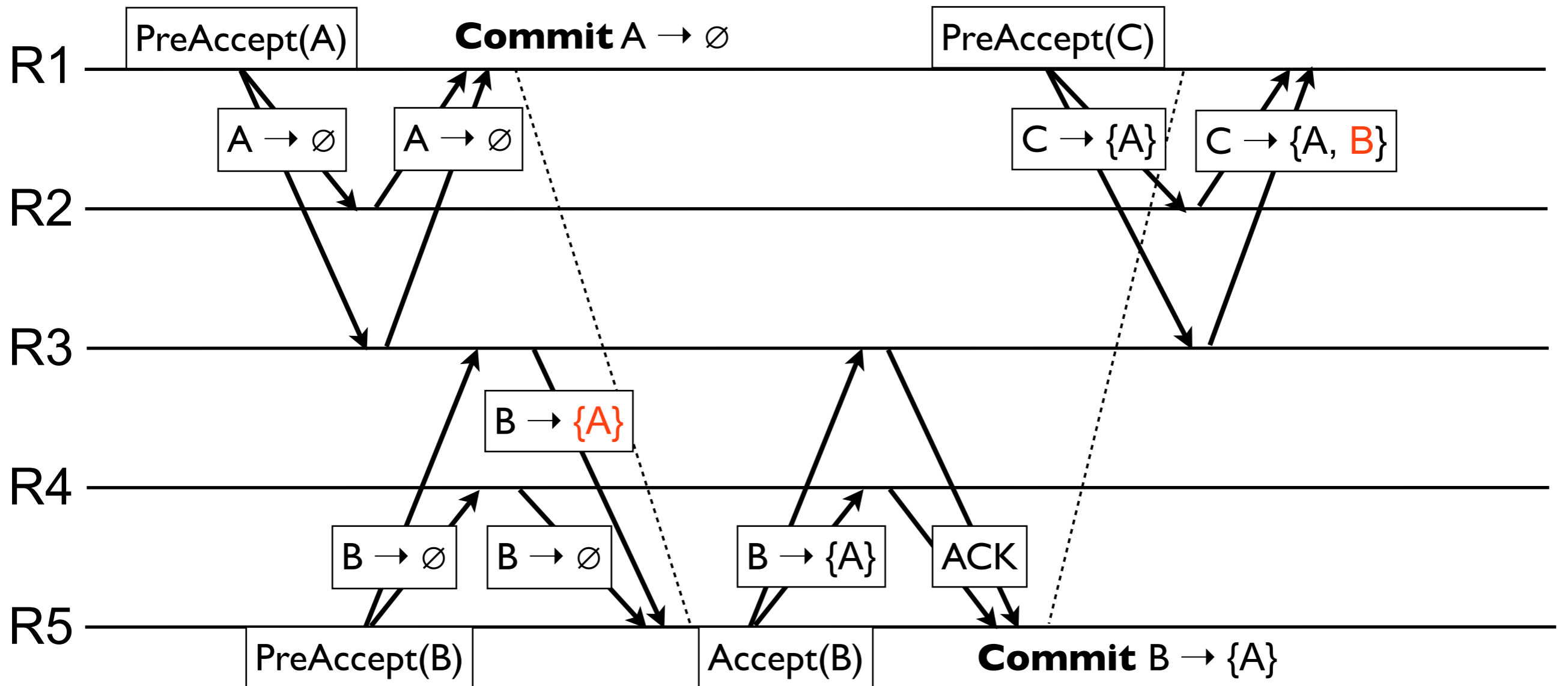
EPaxos commit protocol



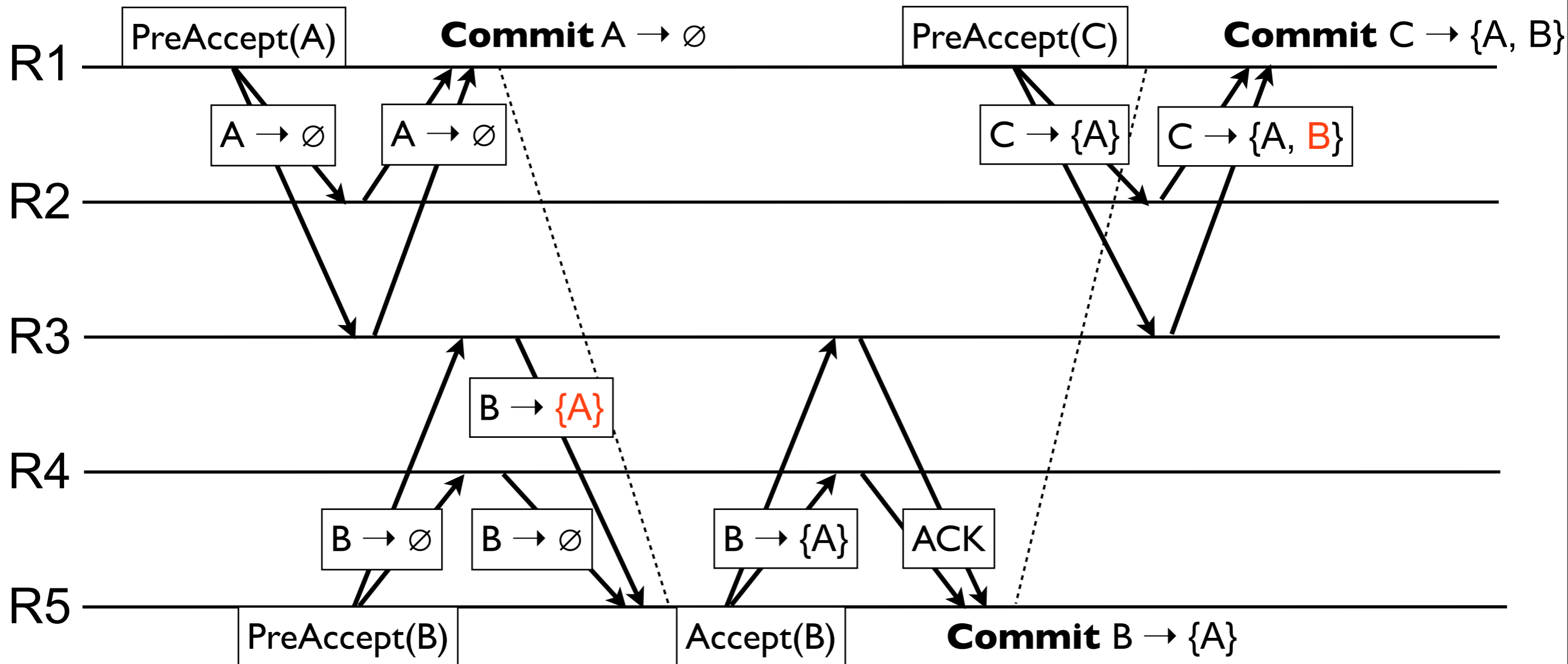
EPaxos commit protocol



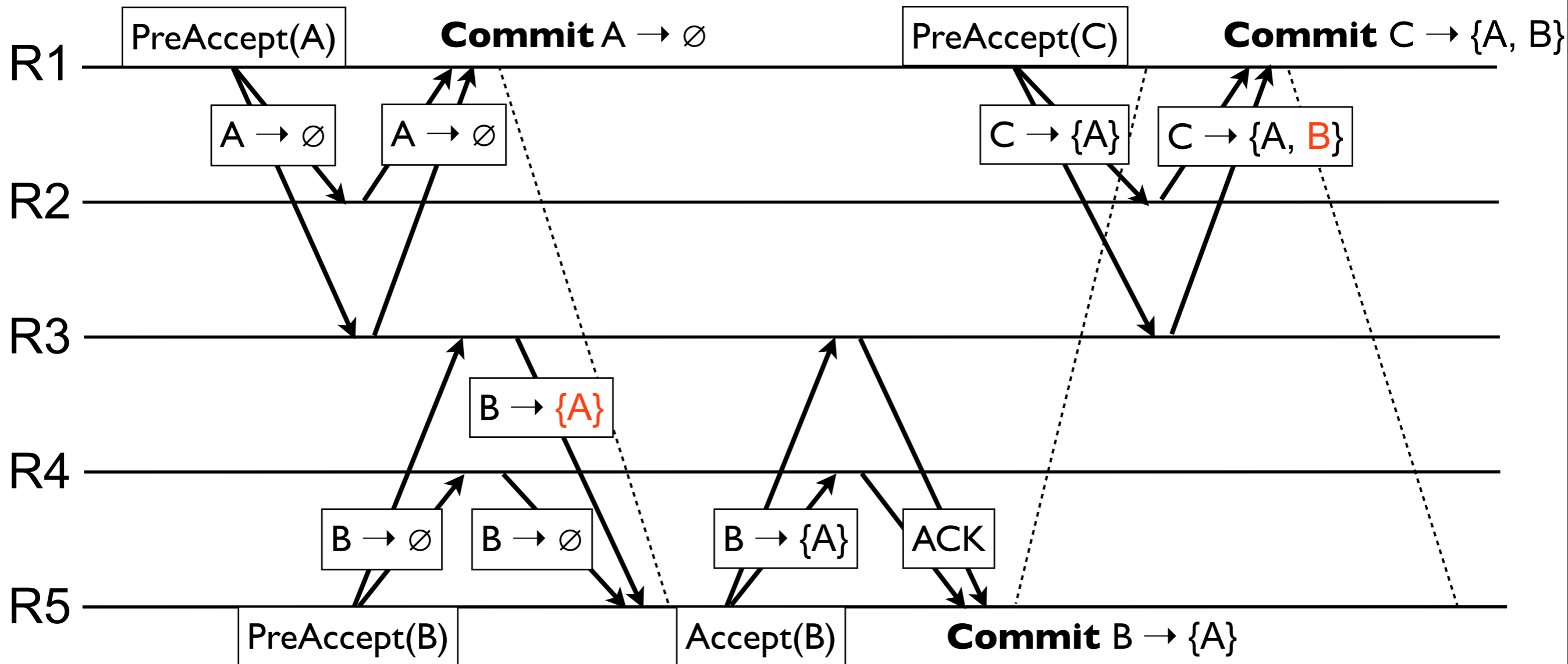
EPaxos commit protocol



EPaxos commit protocol



EPaxos commit protocol



Order only interfering commands

- 1 RTT
 - Non-concurrent commands
 - **OR non-interfering commands**
- 2 RTTs
 - Concurrent **AND interfering**

Interference is application-specific

Interference is application-specific

KV store

- Infer from operation key

Interference is application-specific

KV store

- Infer from operation key

Google App Engine

- Programmer-specified

Interference is application-specific

KV store

- Infer from operation key

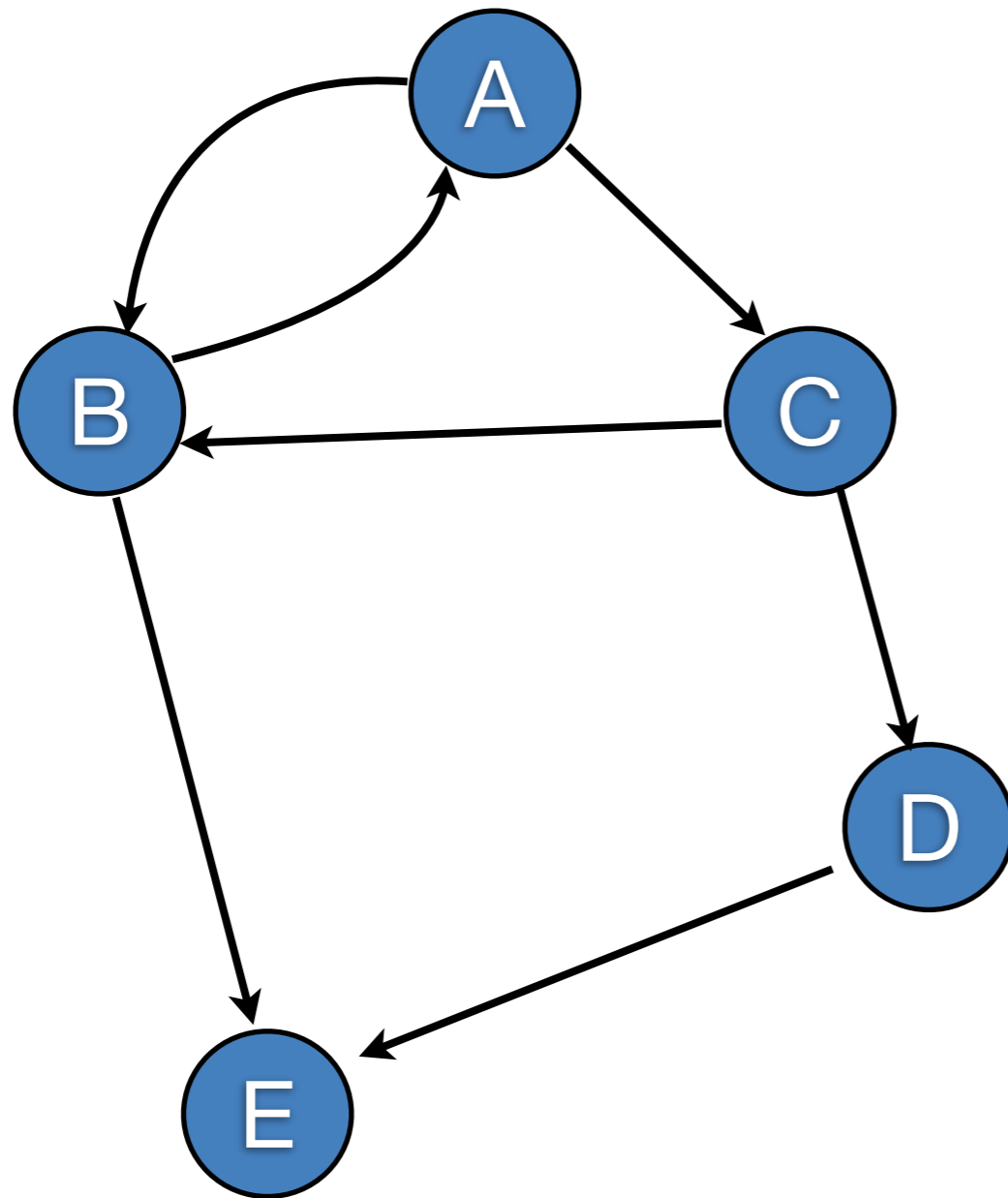
Google App Engine

- Programmer-specified

Relational databases

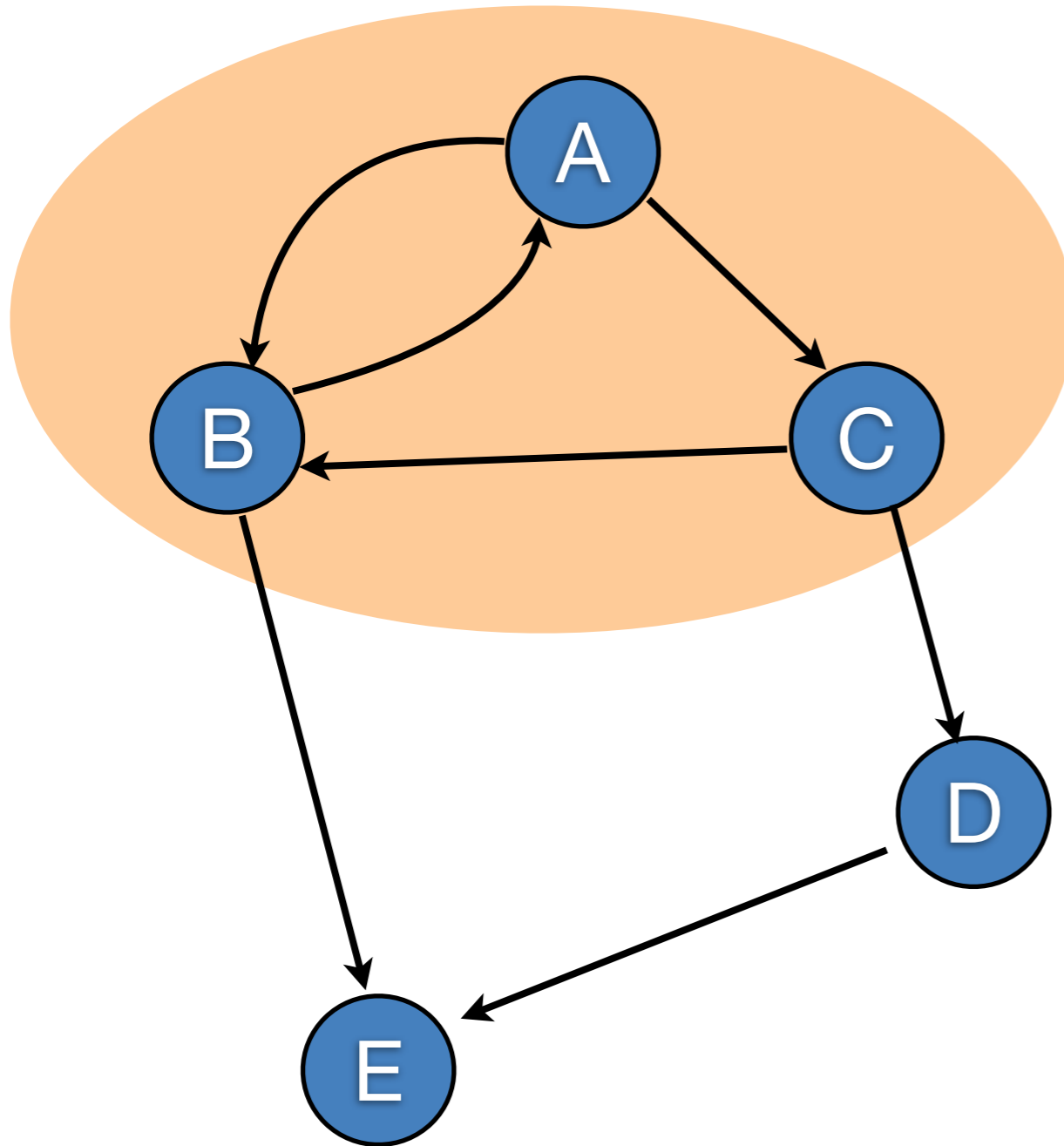
- Most transactions are simple, can be analyzed
- Few remaining transactions interfere w/ everything

Execution



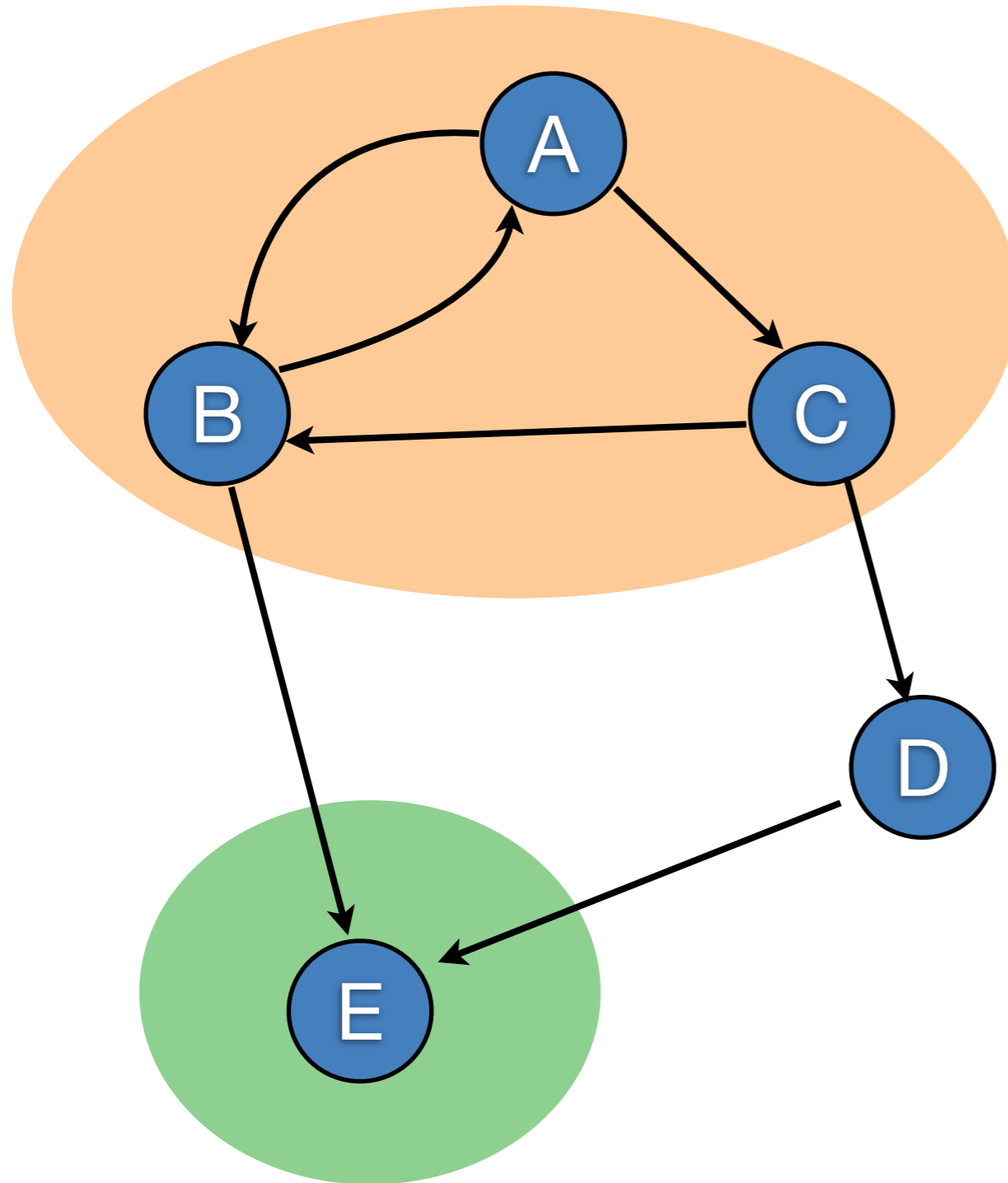
Execution

strongly-connected component



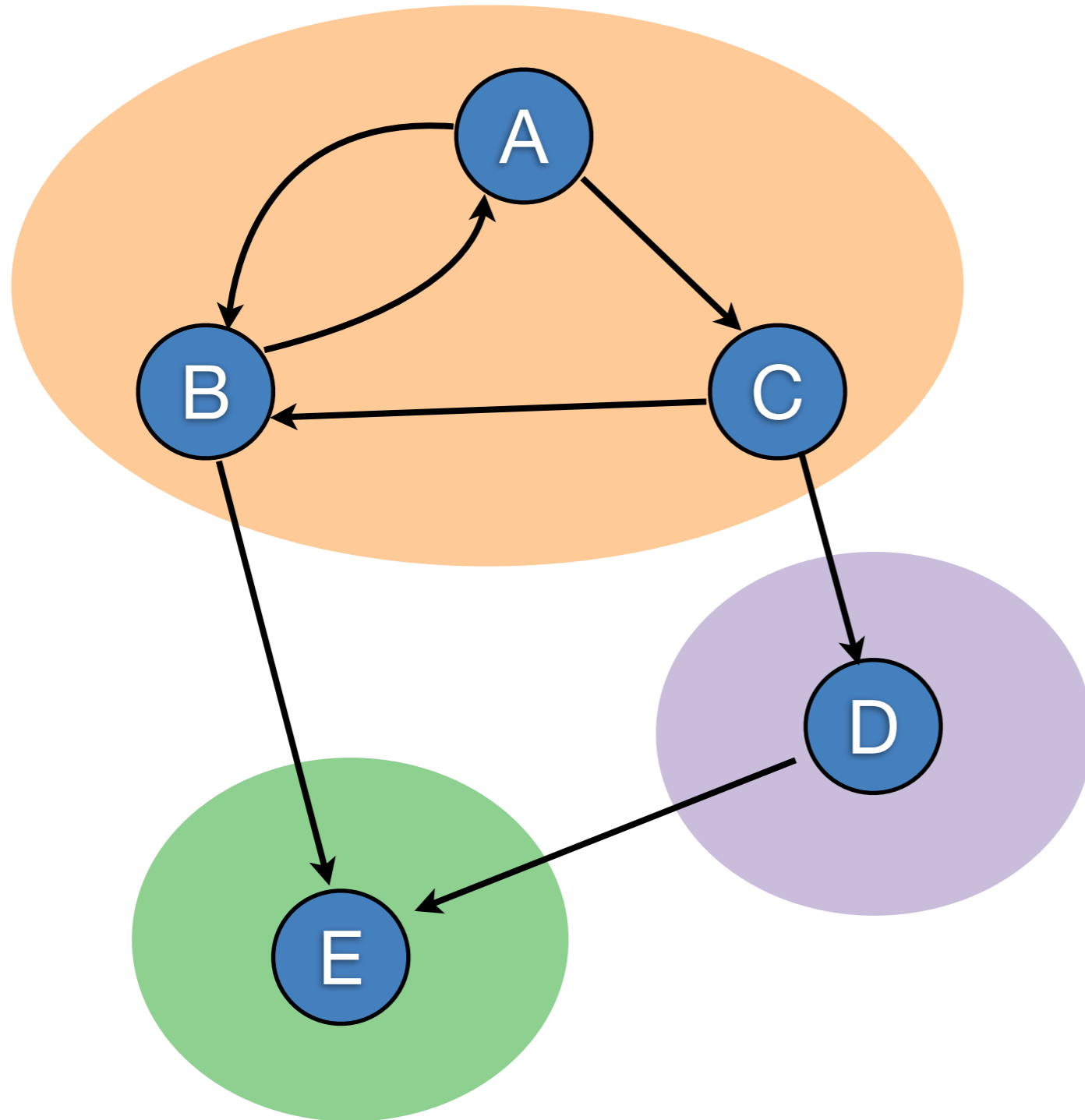
Execution

strongly-connected component



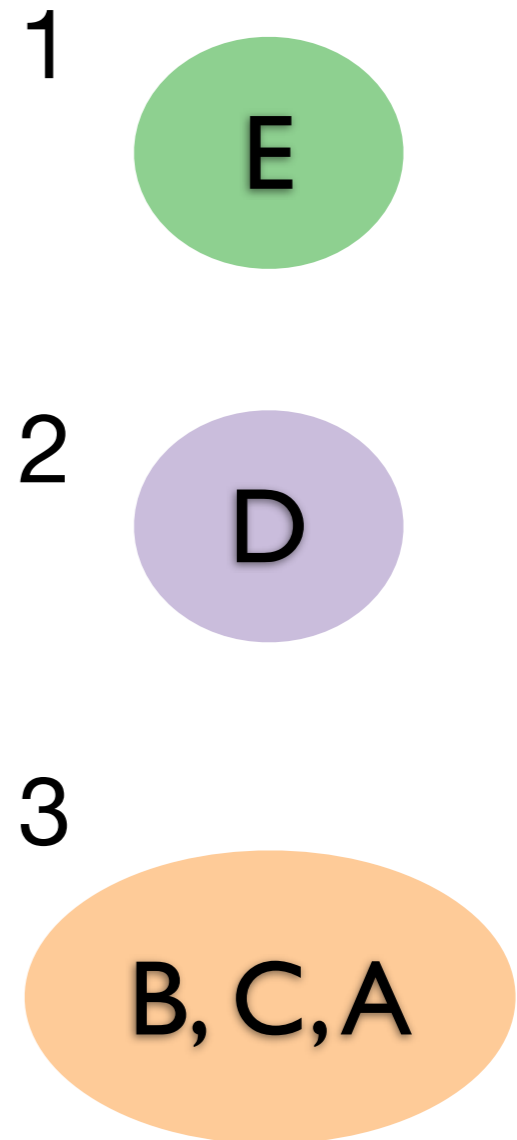
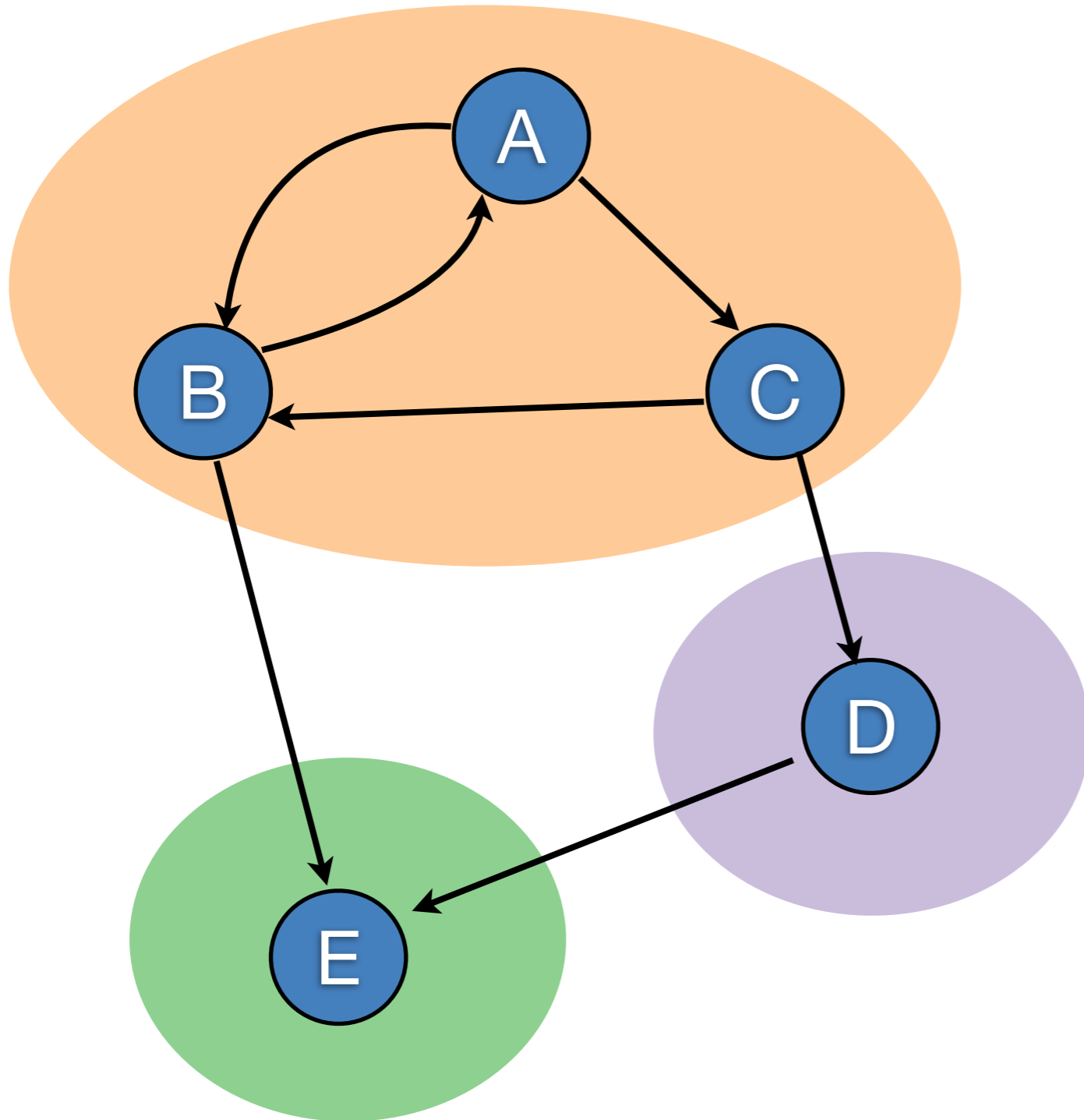
Execution

strongly-connected component



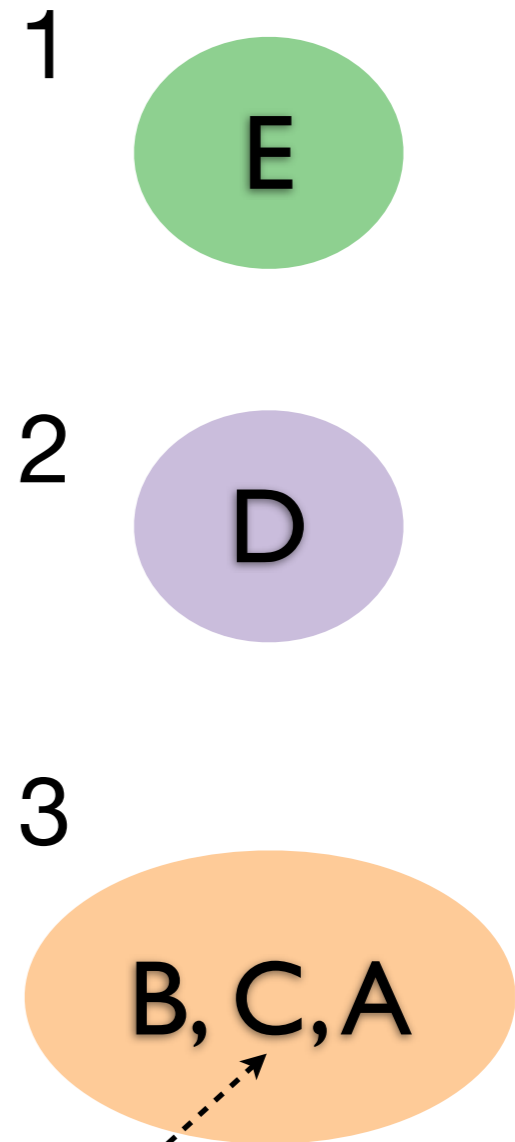
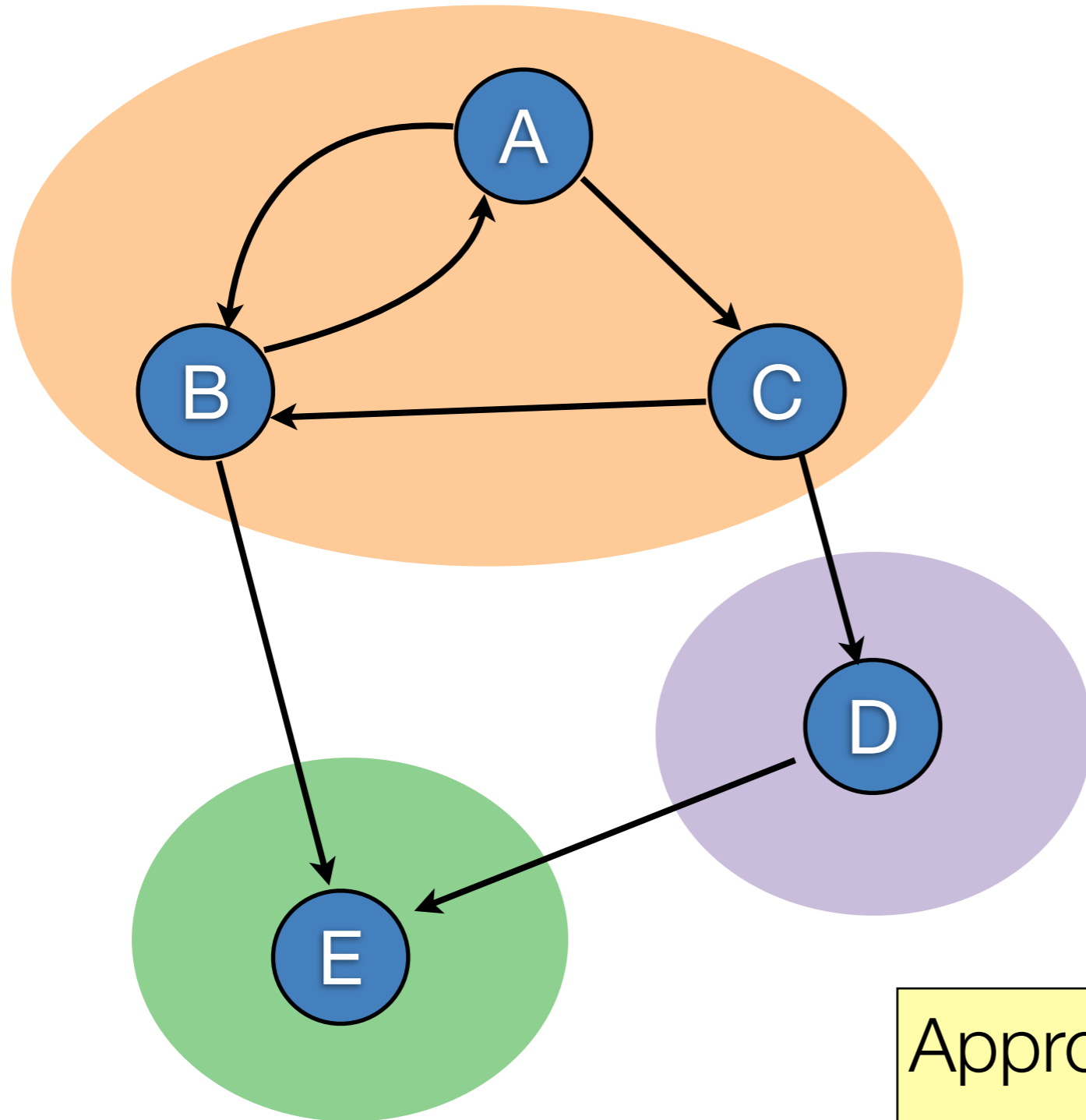
Execution

strongly-connected component



Execution

strongly-connected component



Approximate sequence # order
(Lamport clock)

EPaxos properties

EPaxos properties

Linearizability: If $A \sim B$, and A committed before B proposed
then A will be executed before B.

EPaxos properties

Linearizability: If $A \sim B$, and A committed before B proposed
then A will be executed before B.

Fast-path quorum: $F + \lceil F / 2 \rceil$

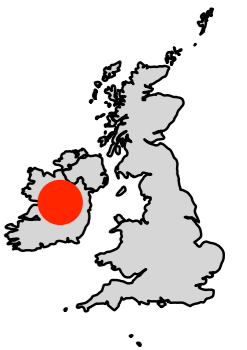
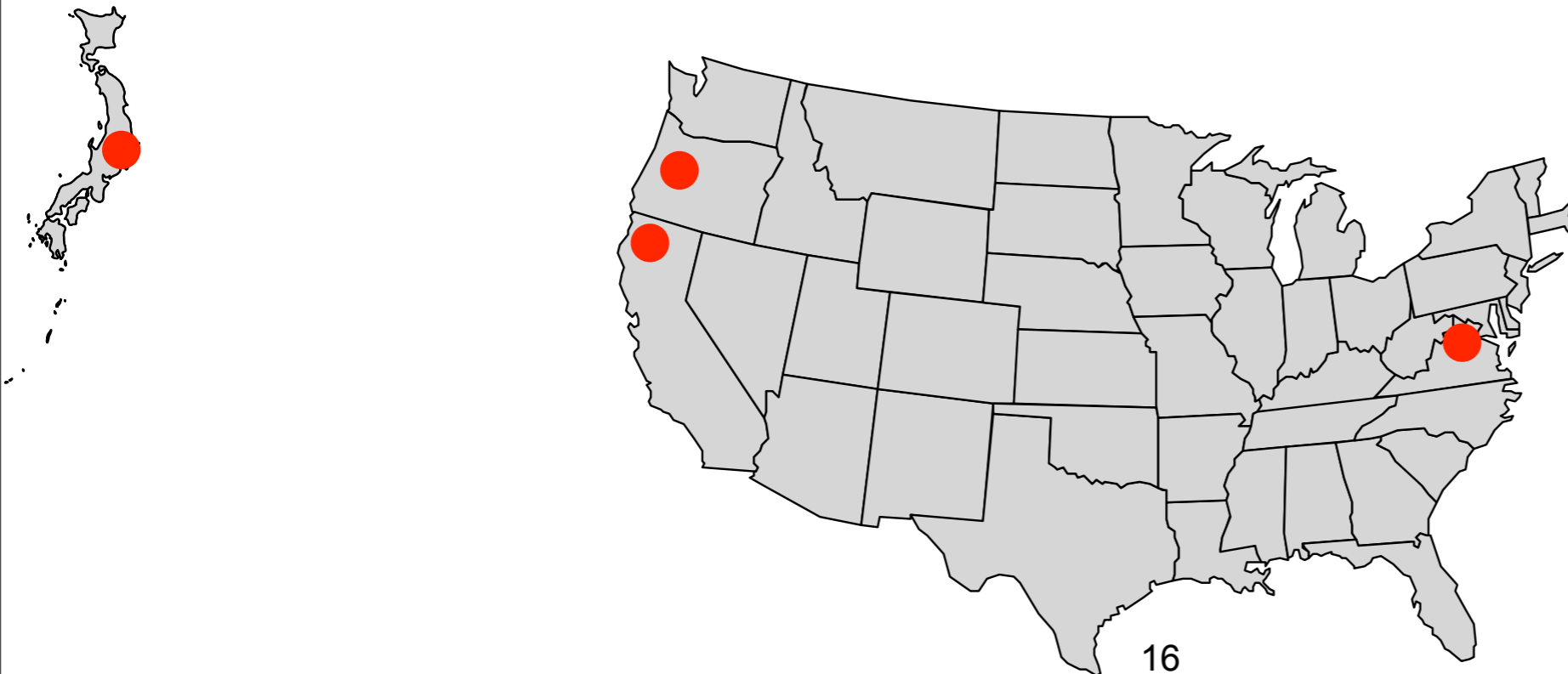
- Optimal for 3 and 5 replicas
- Better than Fast / Generalized Paxos by 1

EPaxos properties

Linearizability: If $A \sim B$, and A committed before B proposed then A will be executed before B.

Fast-path quorum: $F + \lceil F / 2 \rceil$

- Optimal for 3 and 5 replicas
- Better than Fast / Generalized Paxos by 1

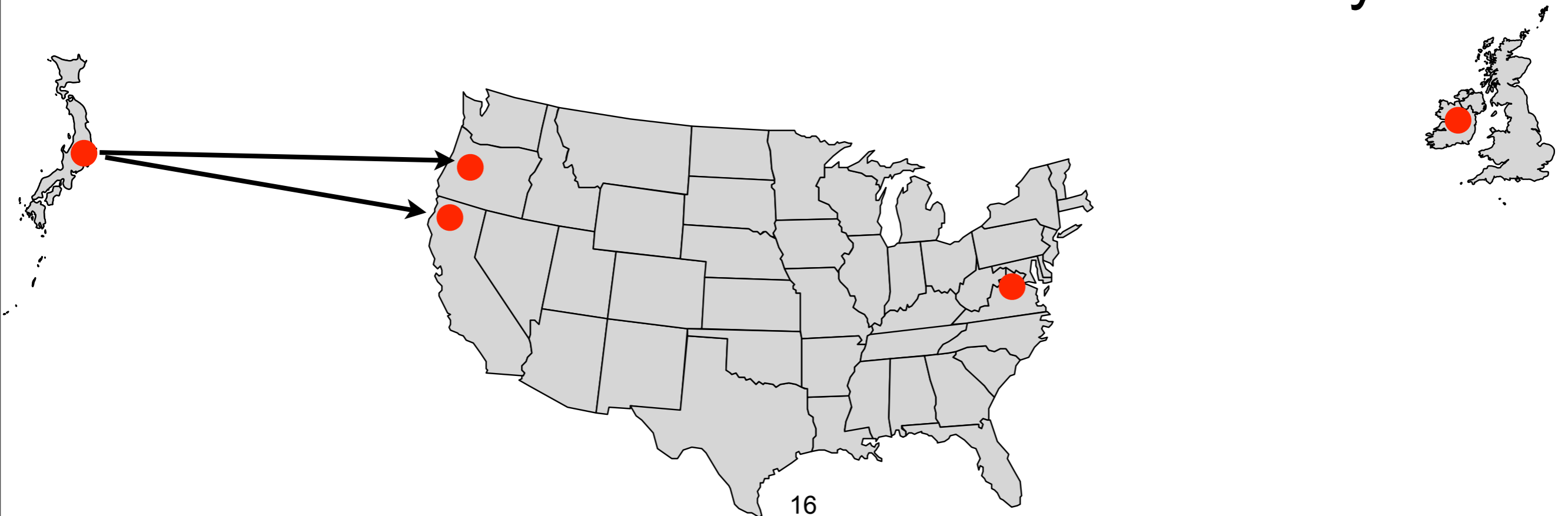


EPaxos properties

Linearizability: If $A \sim B$, and A committed before B proposed then A will be executed before B.

Fast-path quorum: $F + \lceil F / 2 \rceil$

- Optimal for 3 and 5 replicas
- Better than Fast / Generalized Paxos by 1

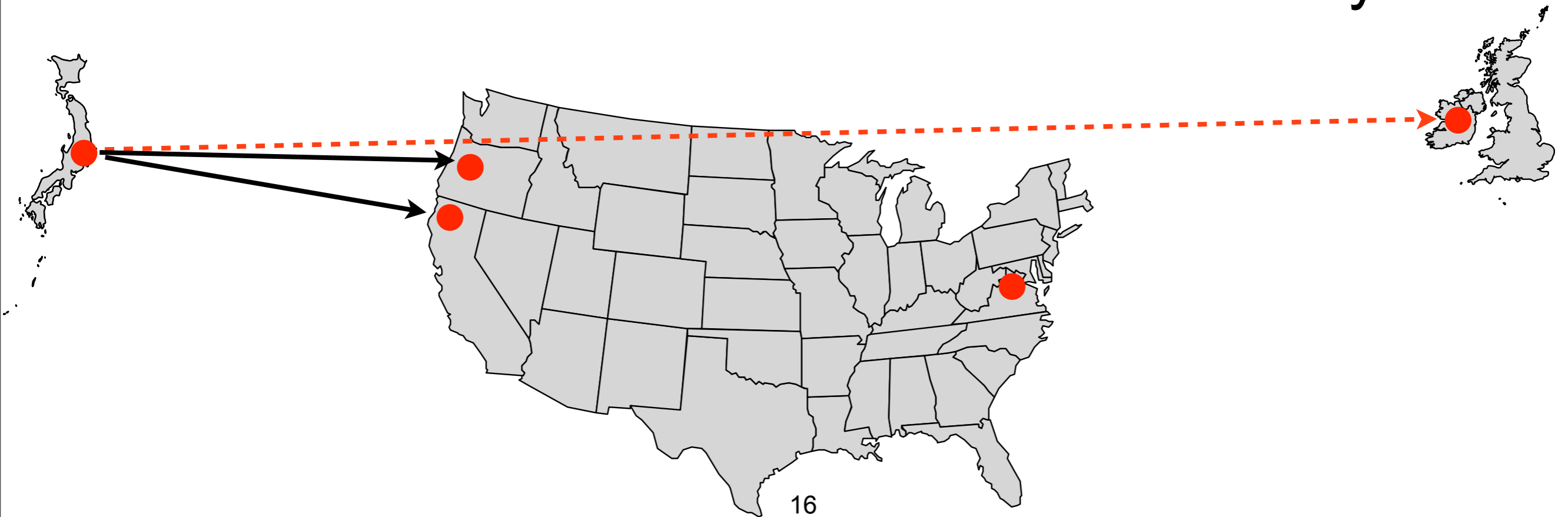


EPaxos properties

Linearizability: If $A \sim B$, and A committed before B proposed then A will be executed before B.

Fast-path quorum: $F + \lceil F / 2 \rceil$

- Optimal for 3 and 5 replicas
- Better than Fast / Generalized Paxos by 1

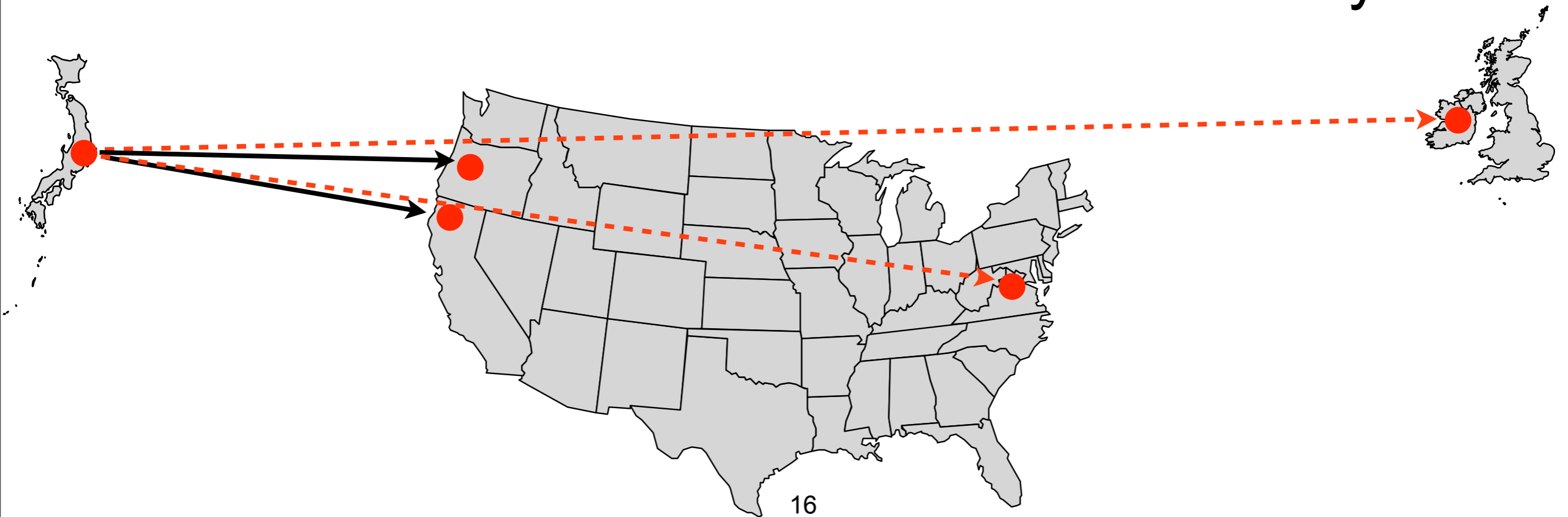


EPaxos properties

Linearizability: If $A \sim B$, and A committed before B proposed then A will be executed before B.

Fast-path quorum: $F + \lceil F / 2 \rceil$

- Optimal for 3 and 5 replicas
- Better than Fast / Generalized Paxos by 1



Results

Optimal wide-area commit latency



Optimal wide-area commit latency



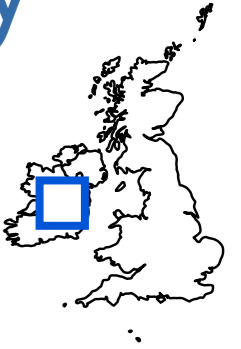
Optimal wide-area commit latency



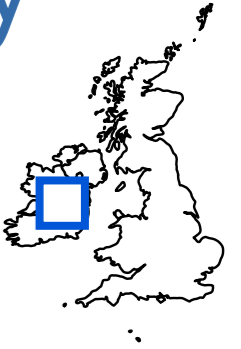
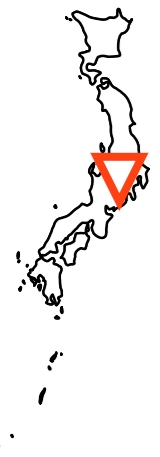
Optimal wide-area commit latency



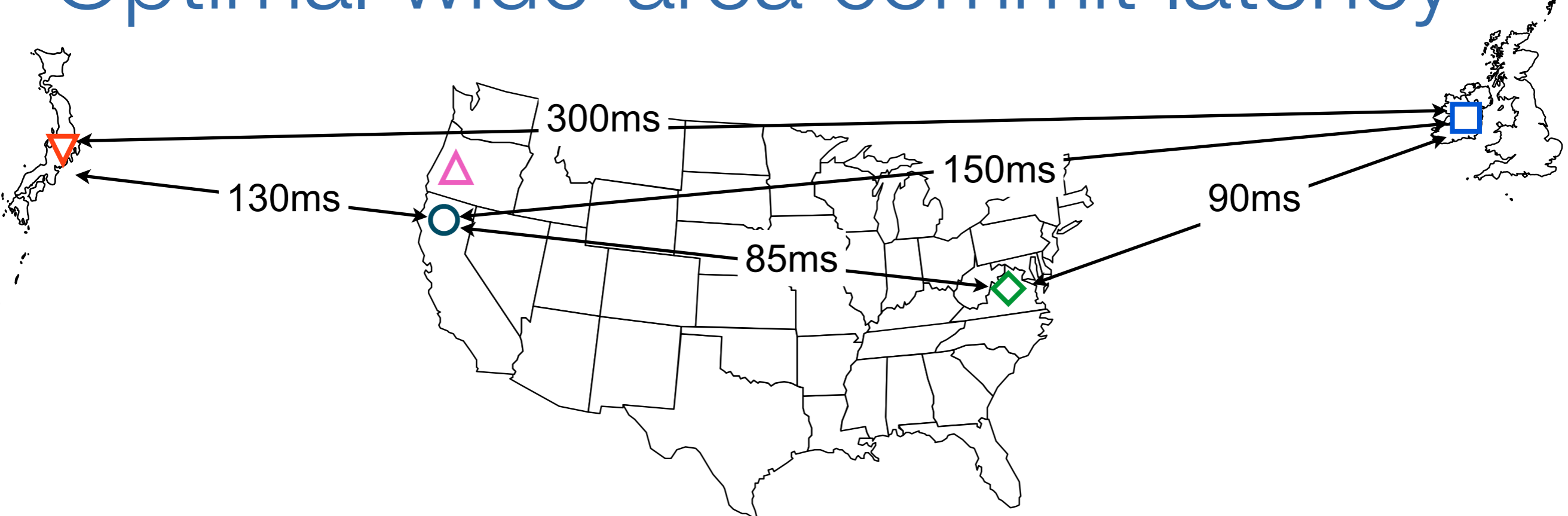
Optimal wide-area commit latency



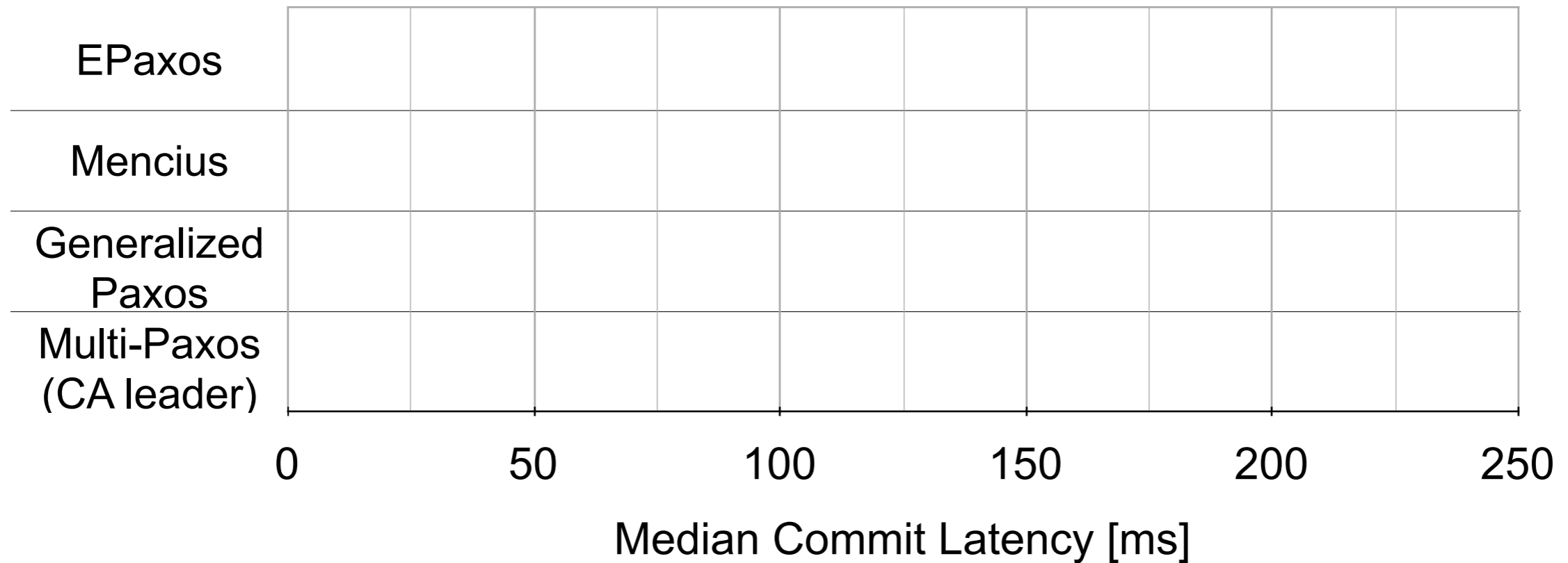
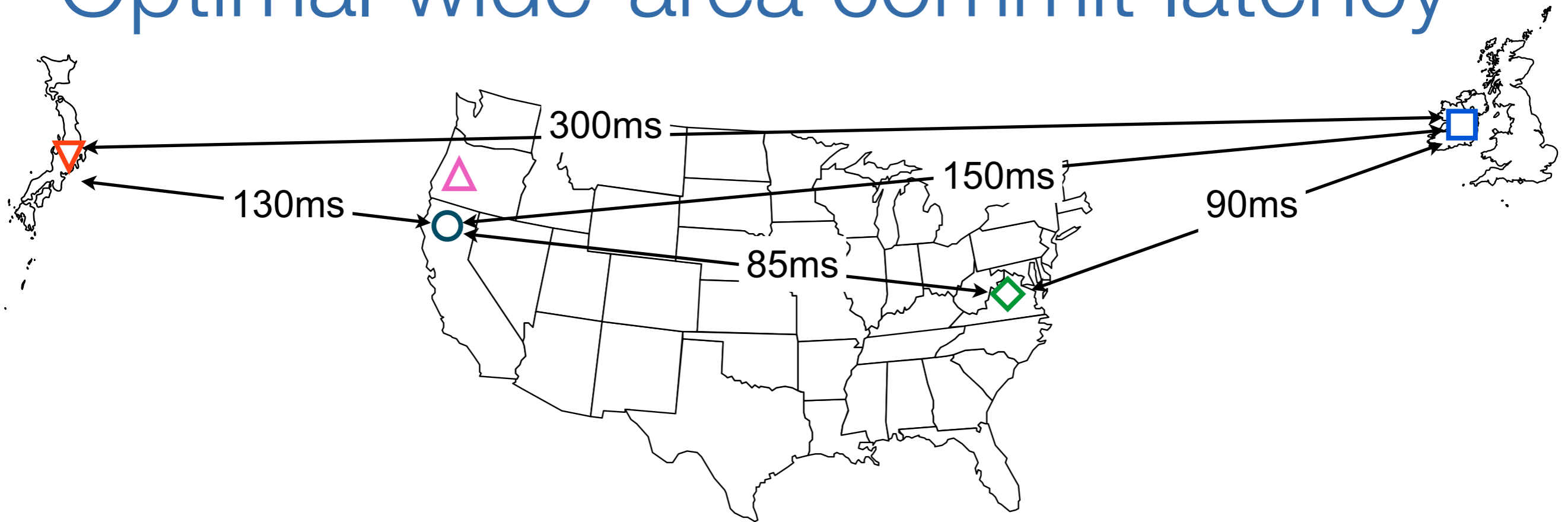
Optimal wide-area commit latency



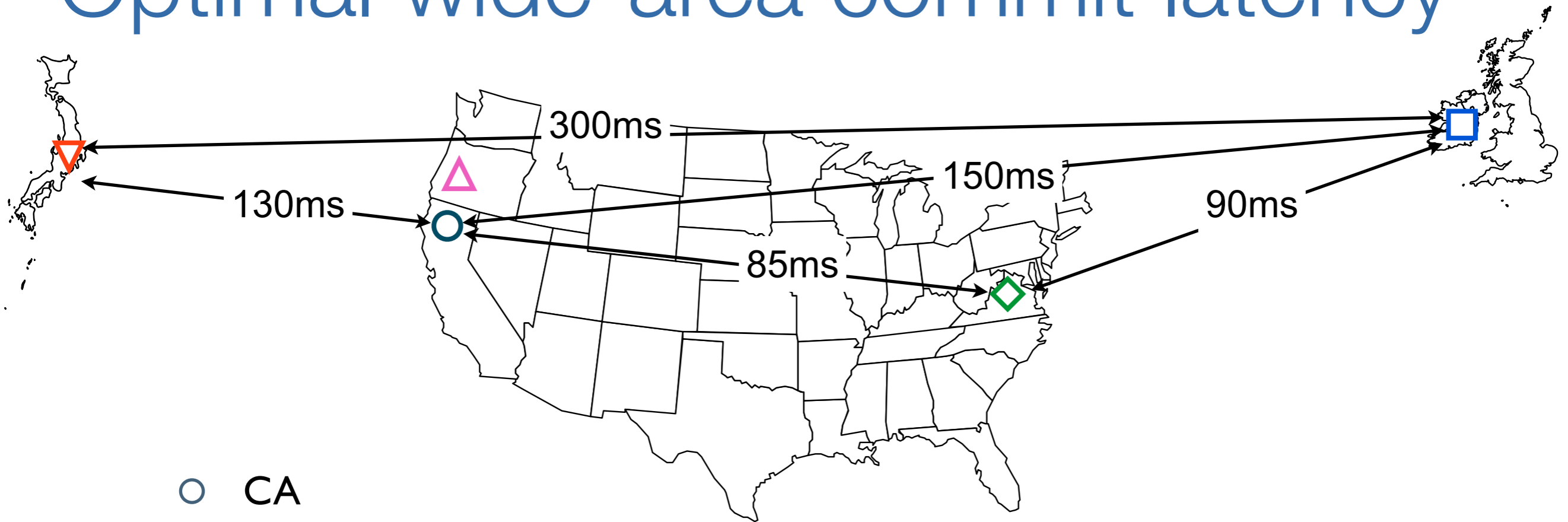
Optimal wide-area commit latency



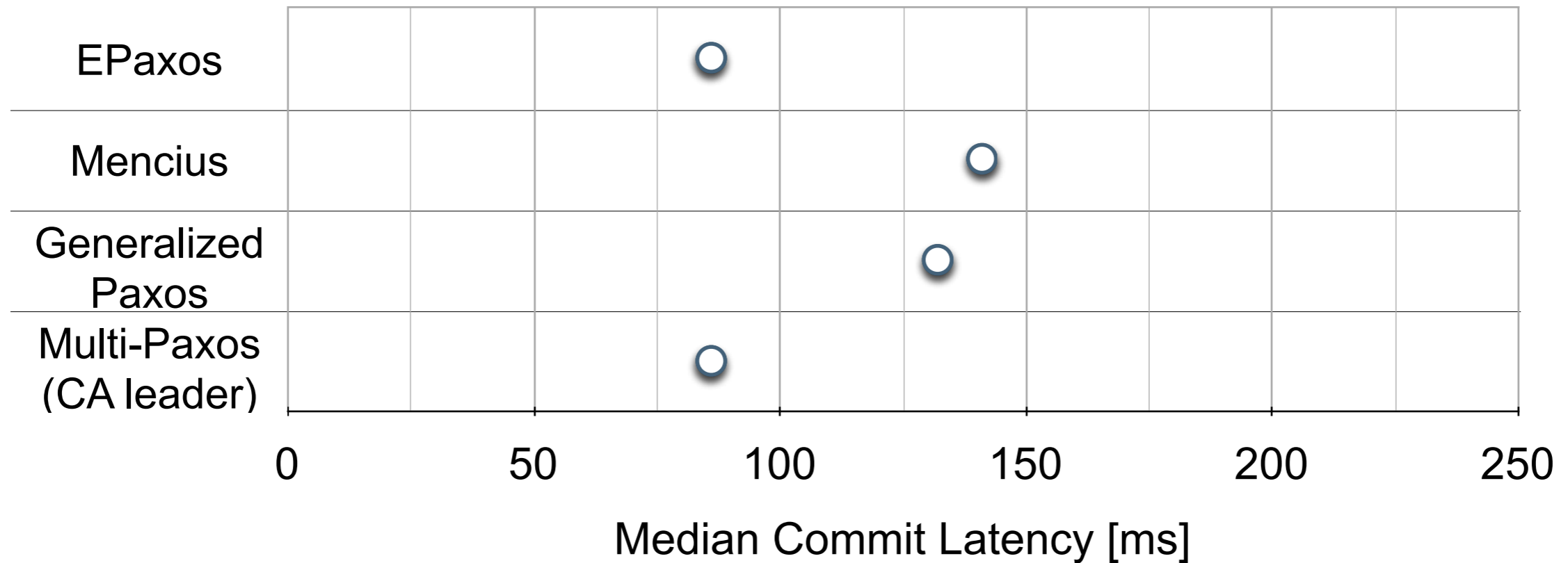
Optimal wide-area commit latency



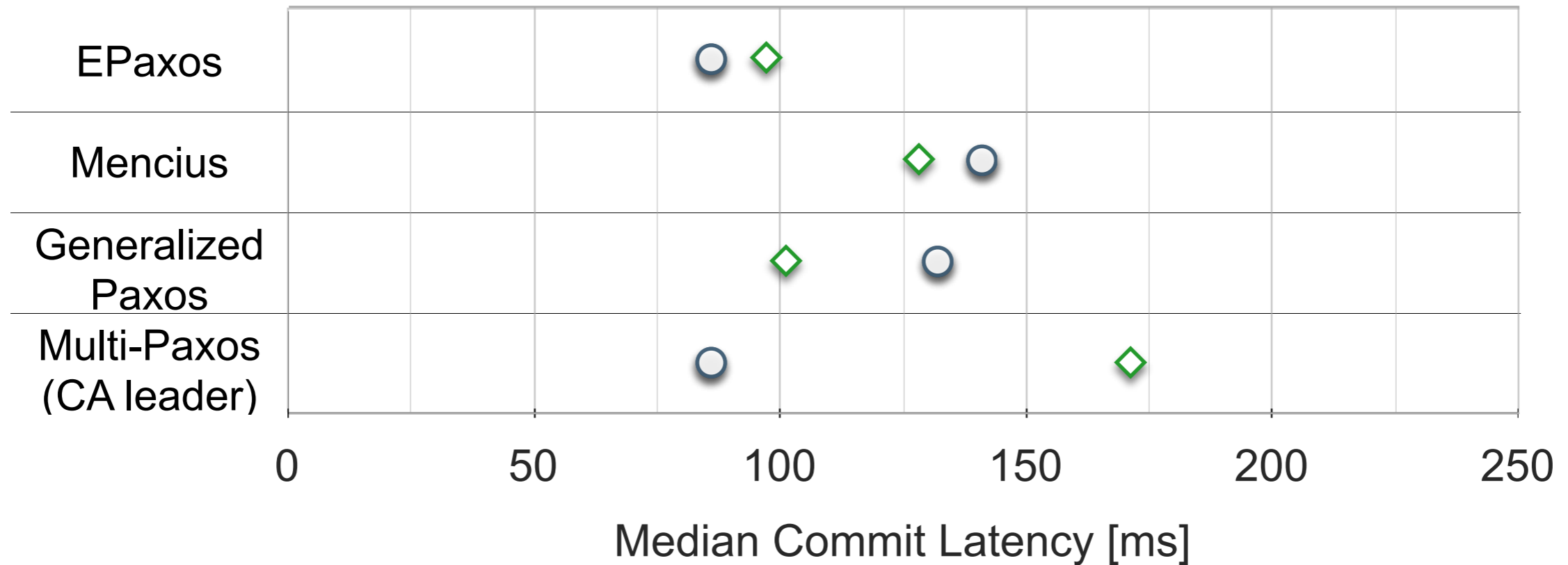
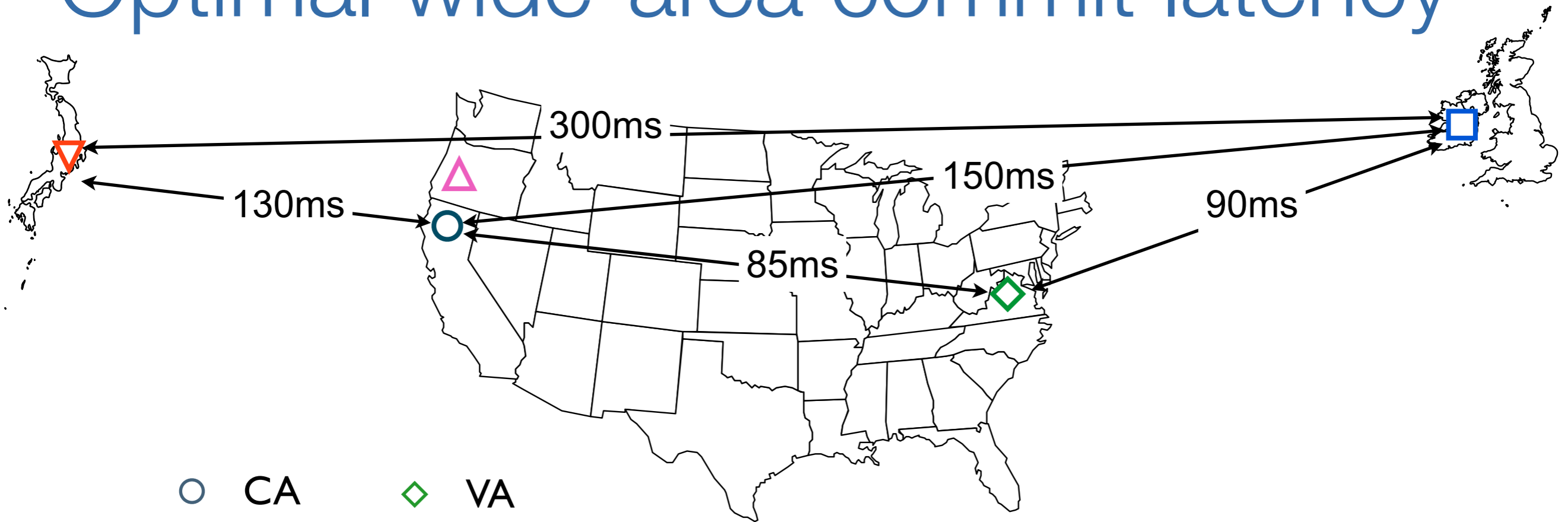
Optimal wide-area commit latency



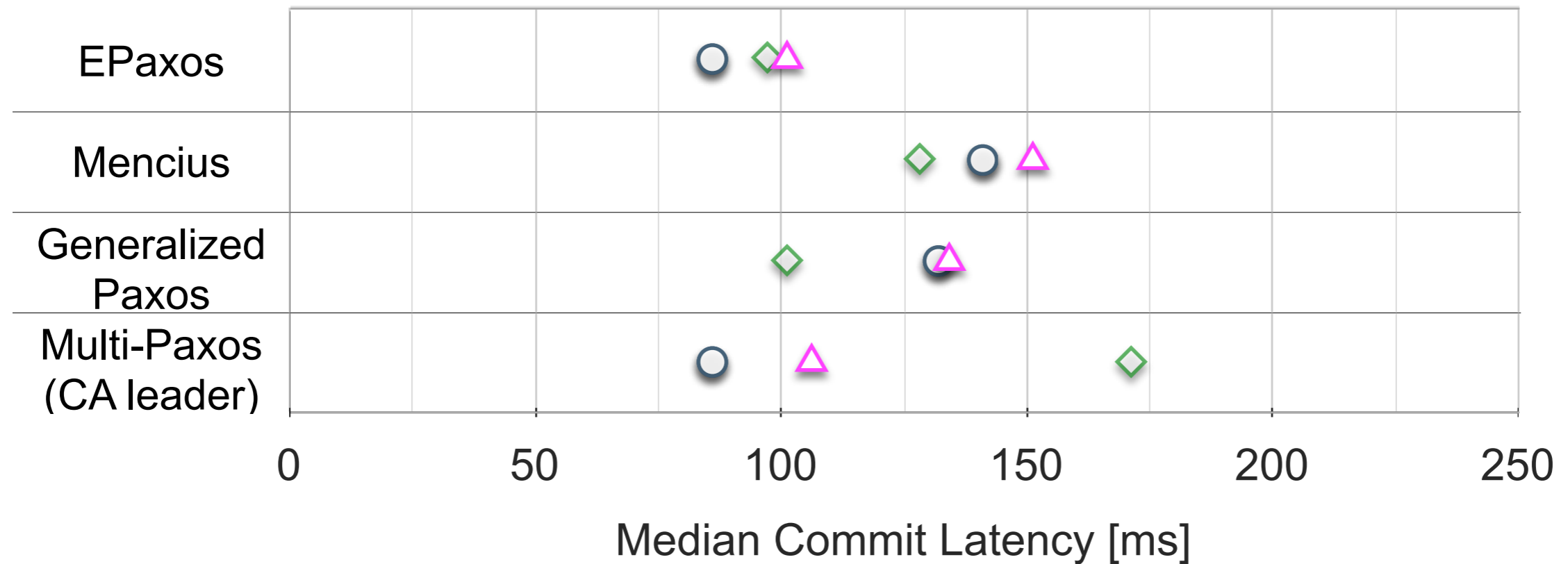
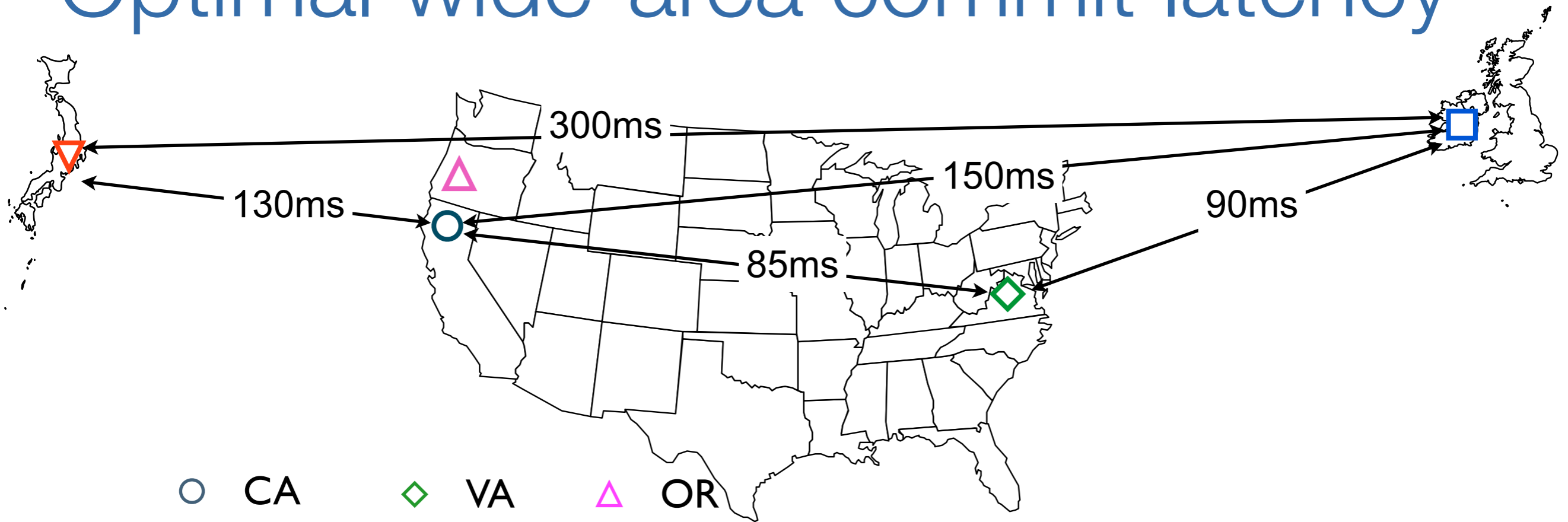
○ CA



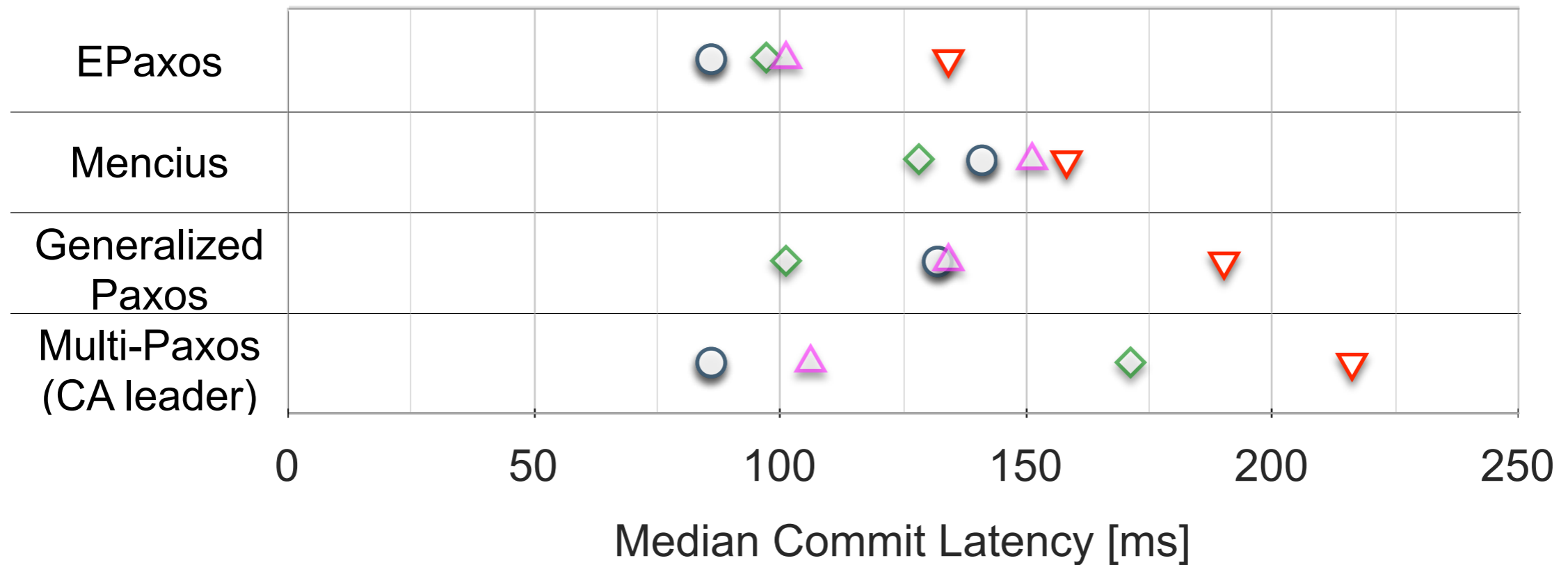
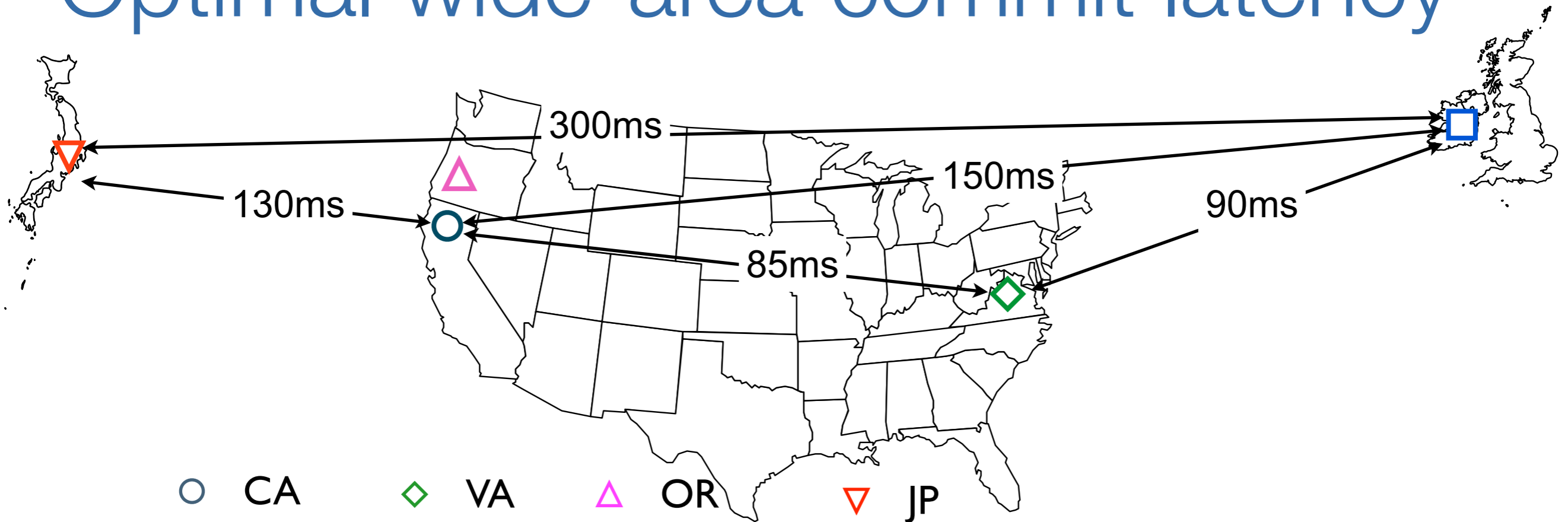
Optimal wide-area commit latency



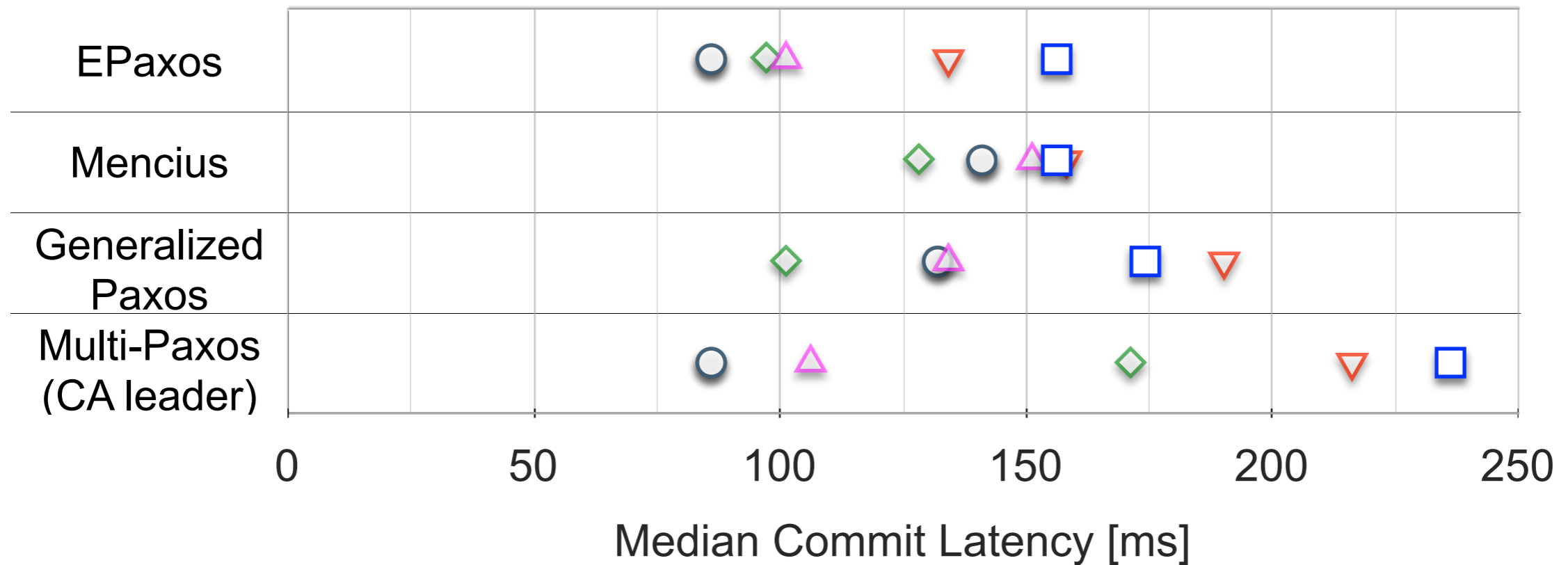
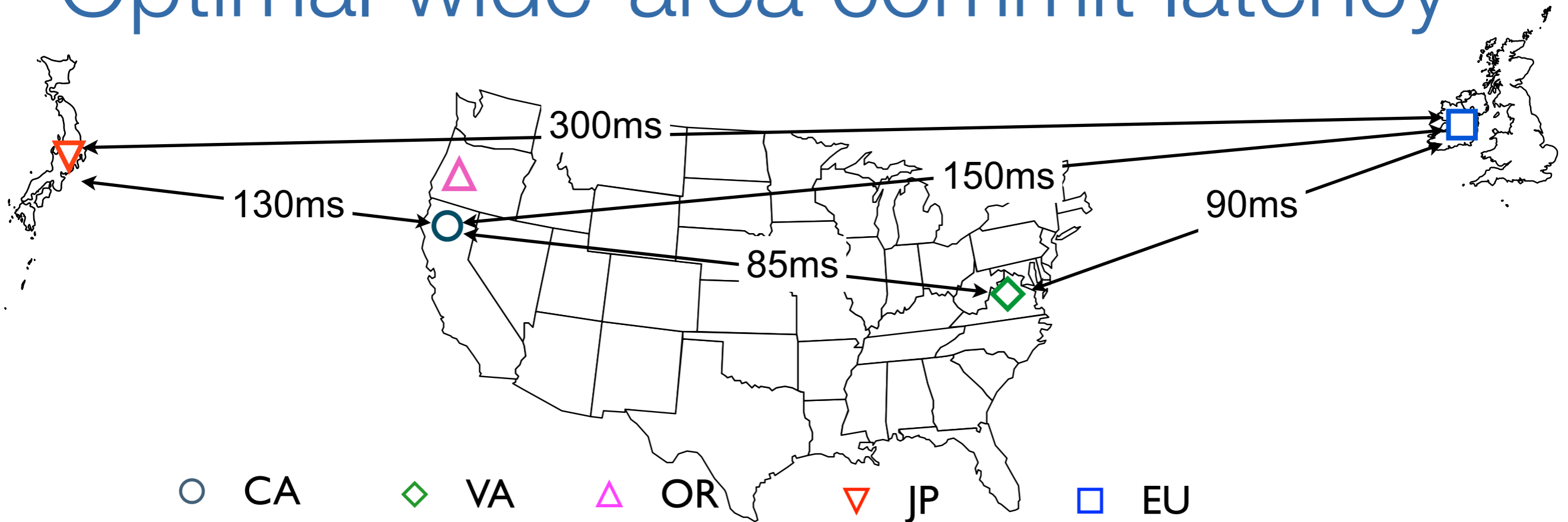
Optimal wide-area commit latency



Optimal wide-area commit latency



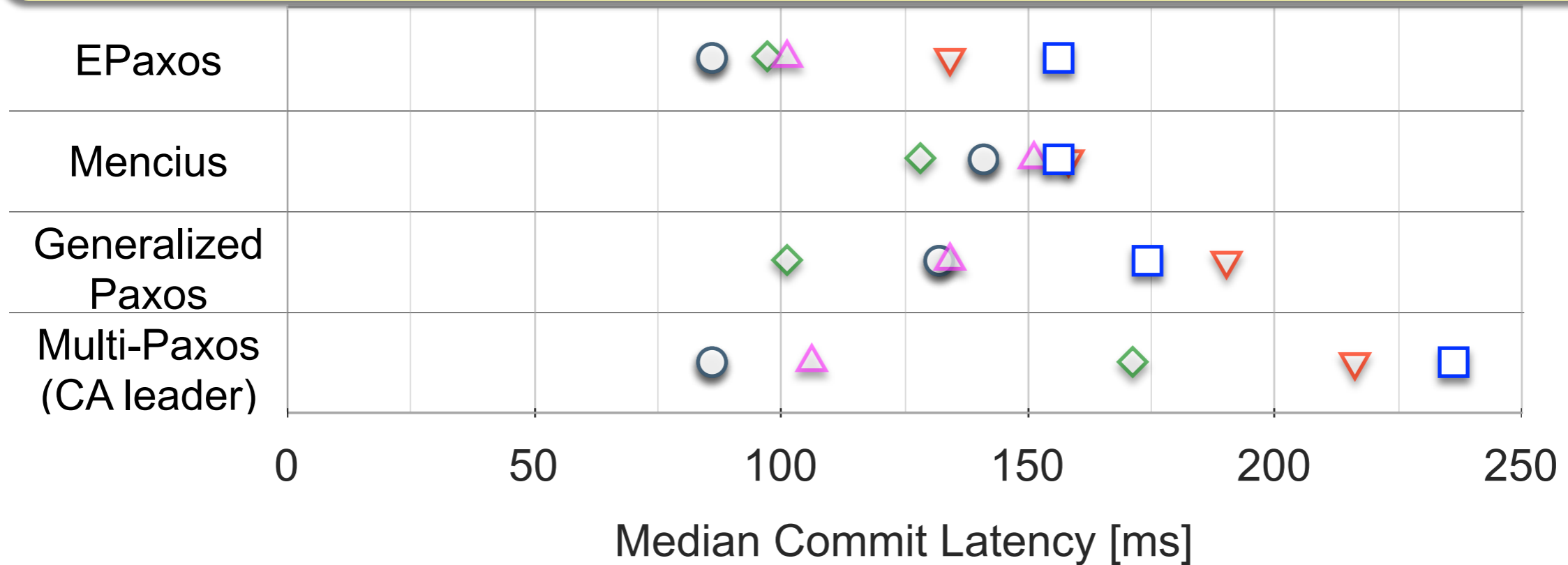
Optimal wide-area commit latency



Optimal wide-area commit latency



EPaxos: Optimal commit latency in wide-area for 3 and 5 replicas



Higher + more stable throughput

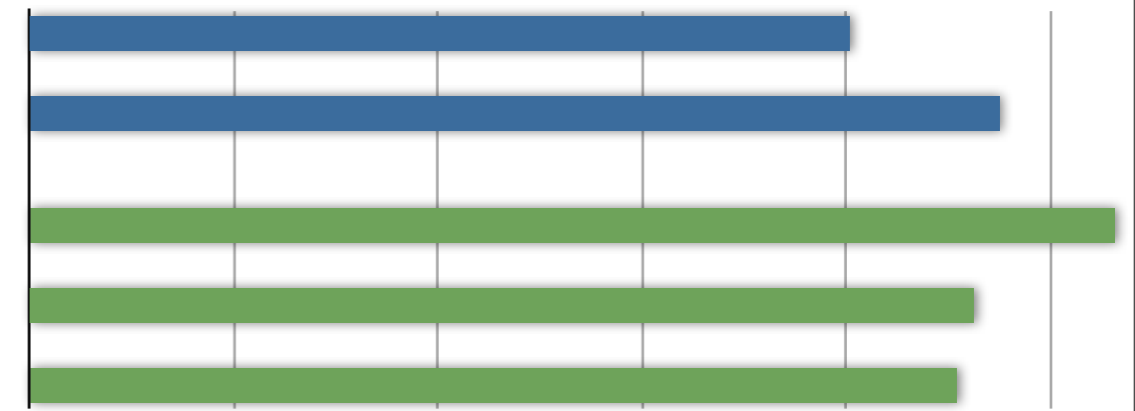
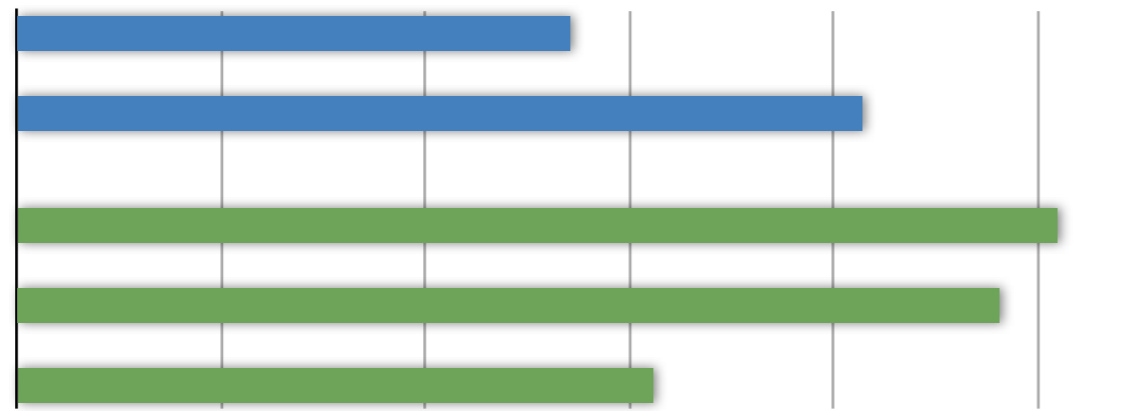
5 replicas

3 replicas

Multi-Paxos
Mencius
EPaxos 0%
EPaxos 2%
EPaxos 100%

0 10000 20000 30000 40000 50000
Operations / sec

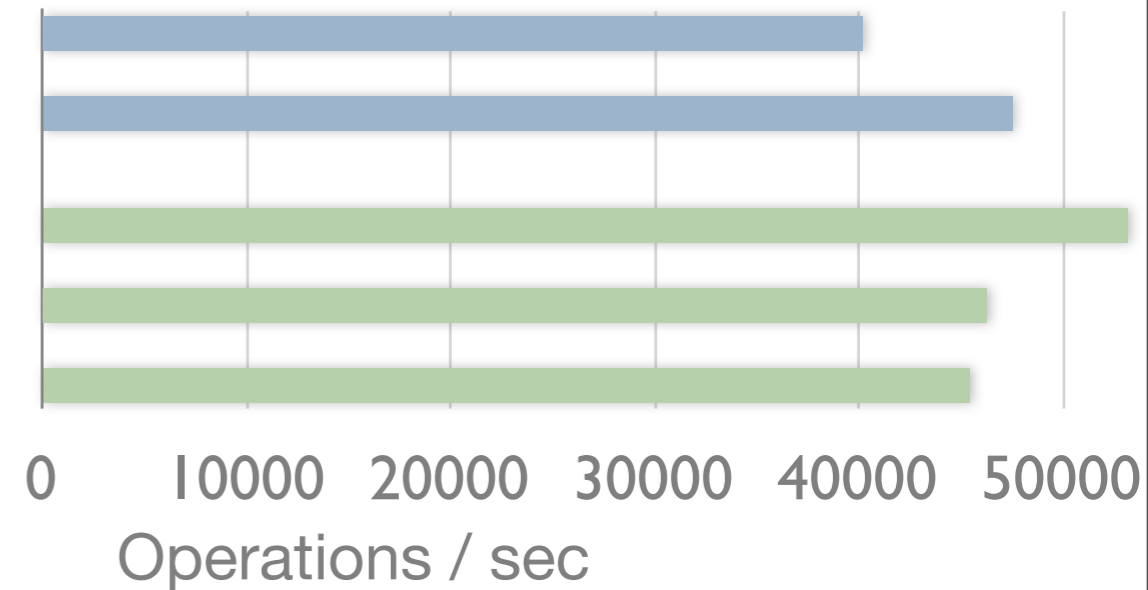
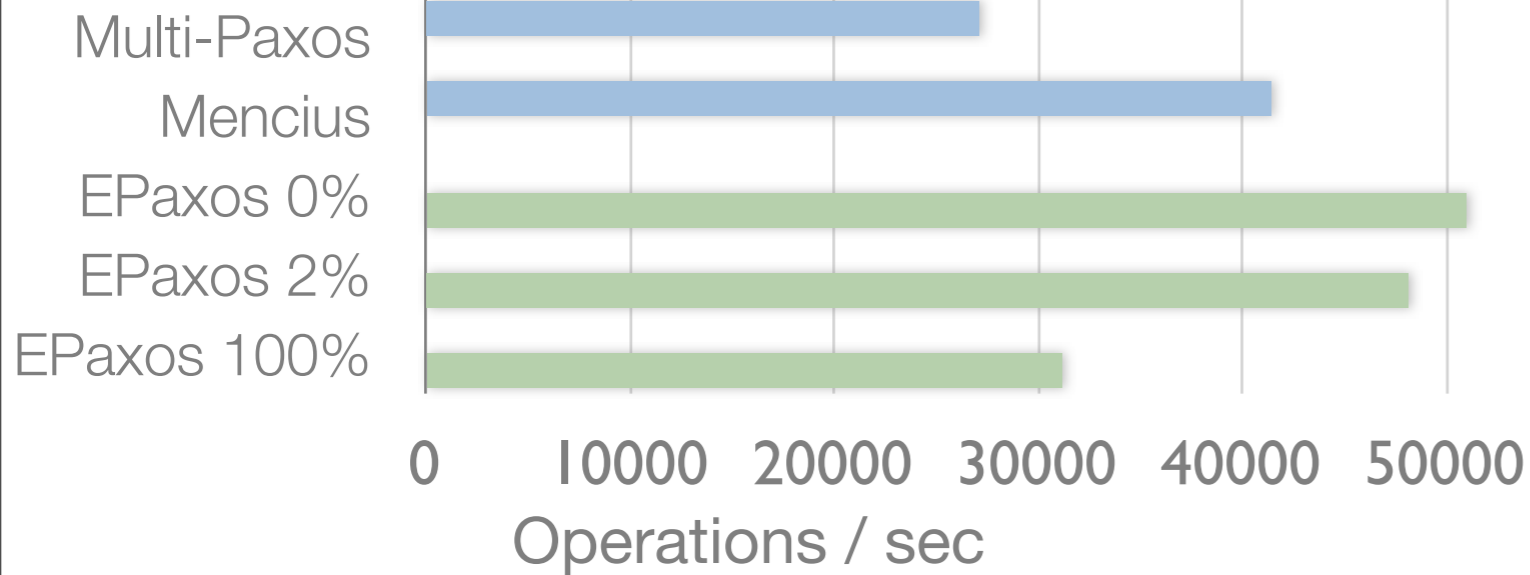
0 10000 20000 30000 40000 50000
Operations / sec



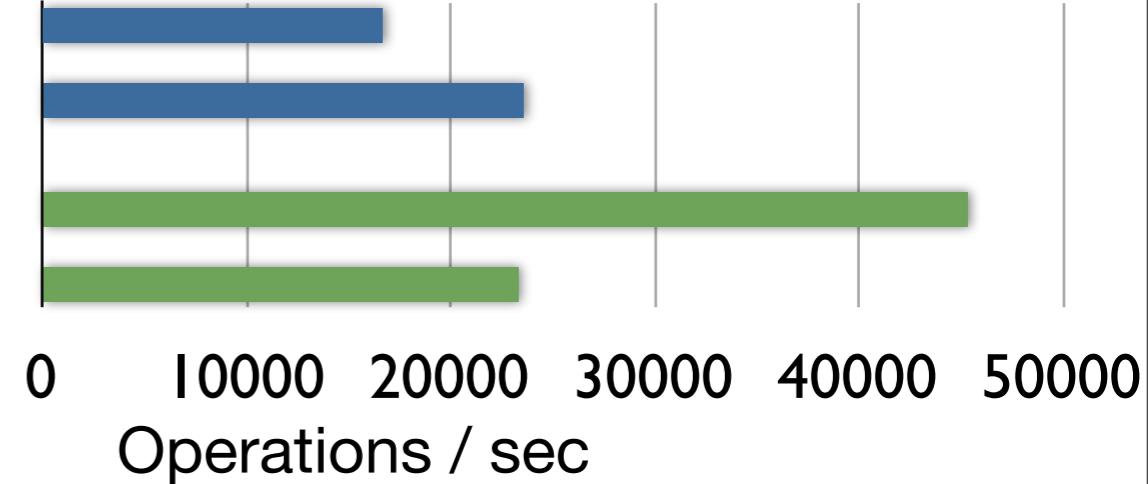
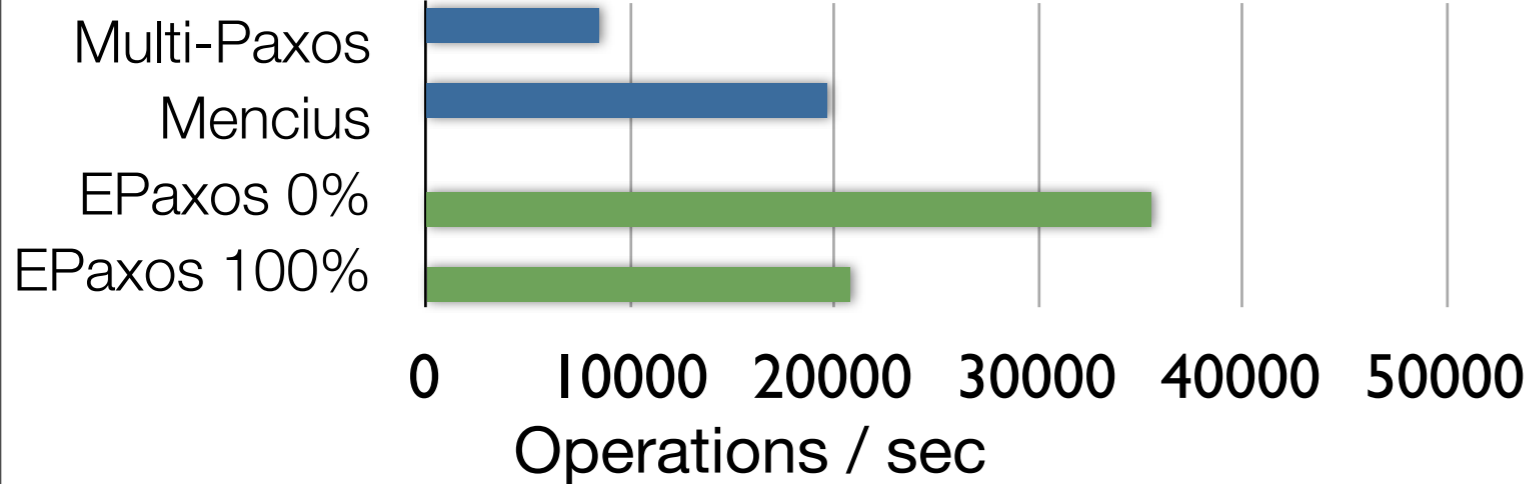
Higher + more stable throughput

5 replicas

3 replicas



When one replica is slow



EPaxos: higher throughput w/ batching

5 ms batching, local area

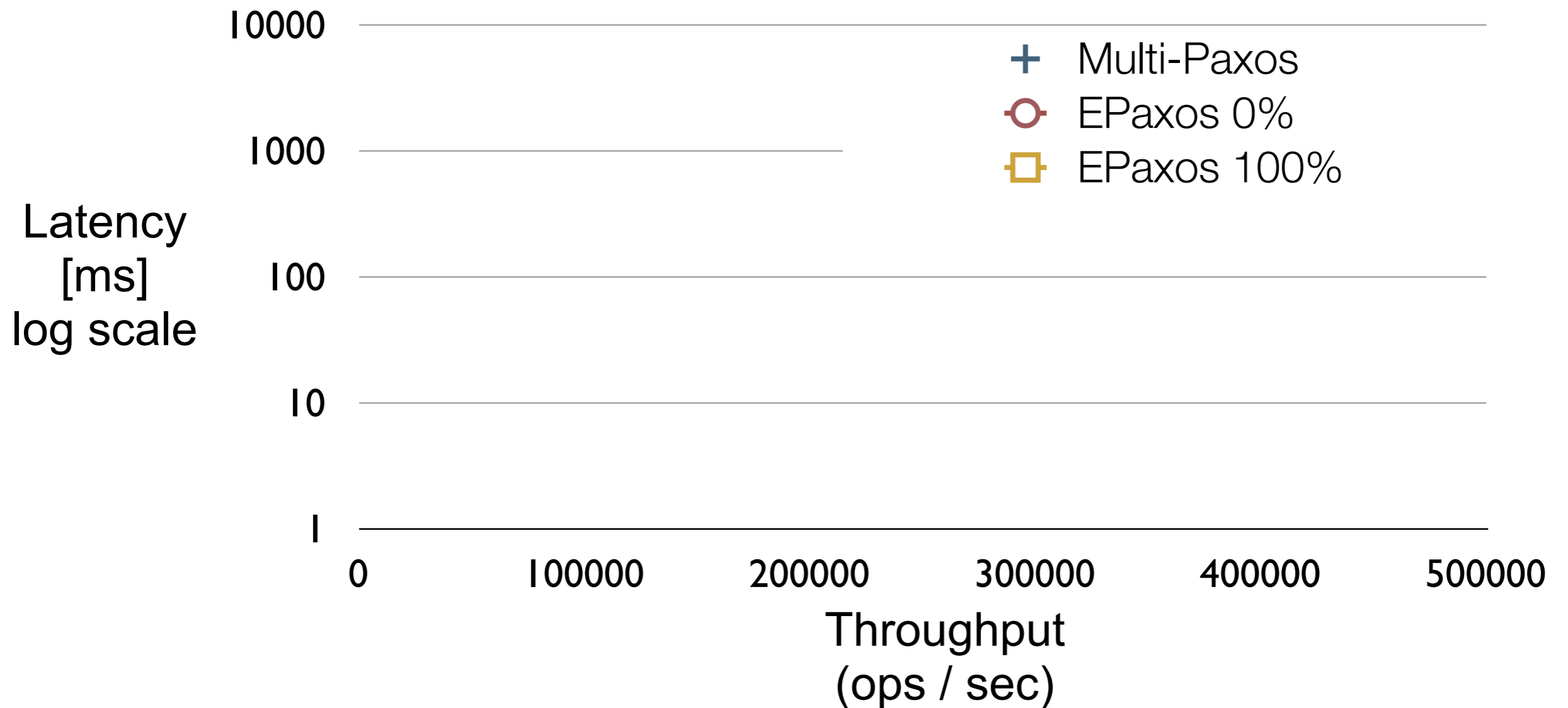
Latency
[ms]
log scale

- + Multi-Paxos
- EPaxos 0%
- EPaxos 100%

Throughput
(ops / sec)

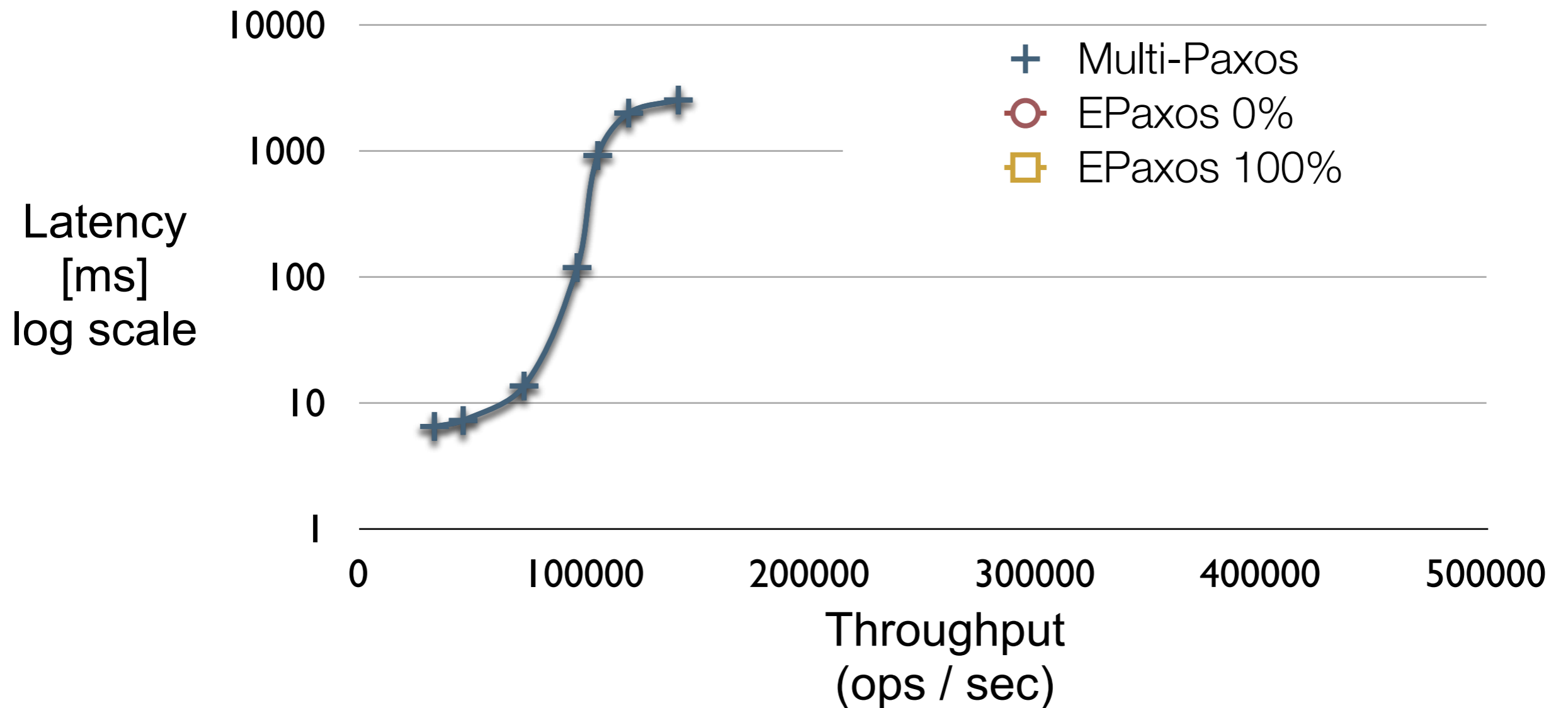
EPaxos: higher throughput w/ batching

5 ms batching, local area



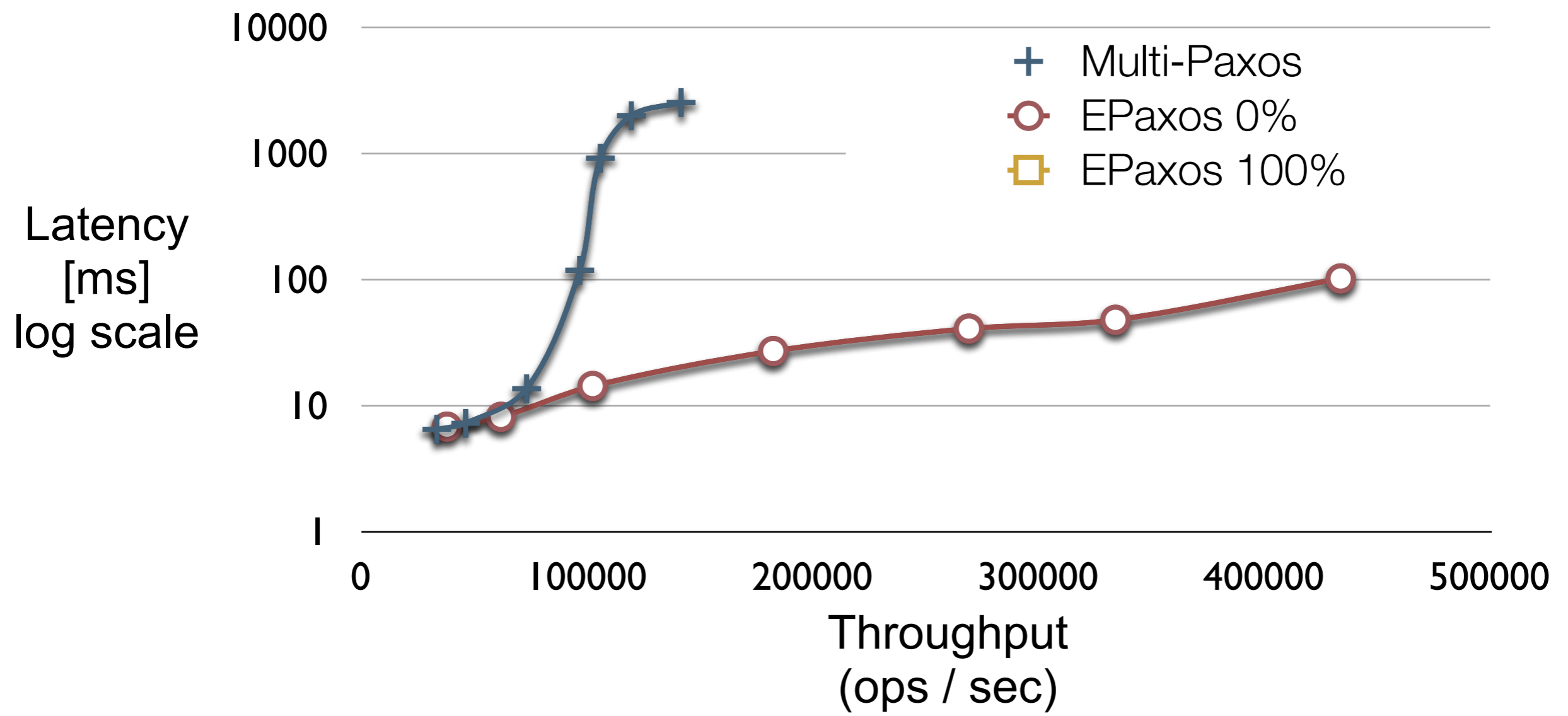
EPaxos: higher throughput w/ batching

5 ms batching, local area



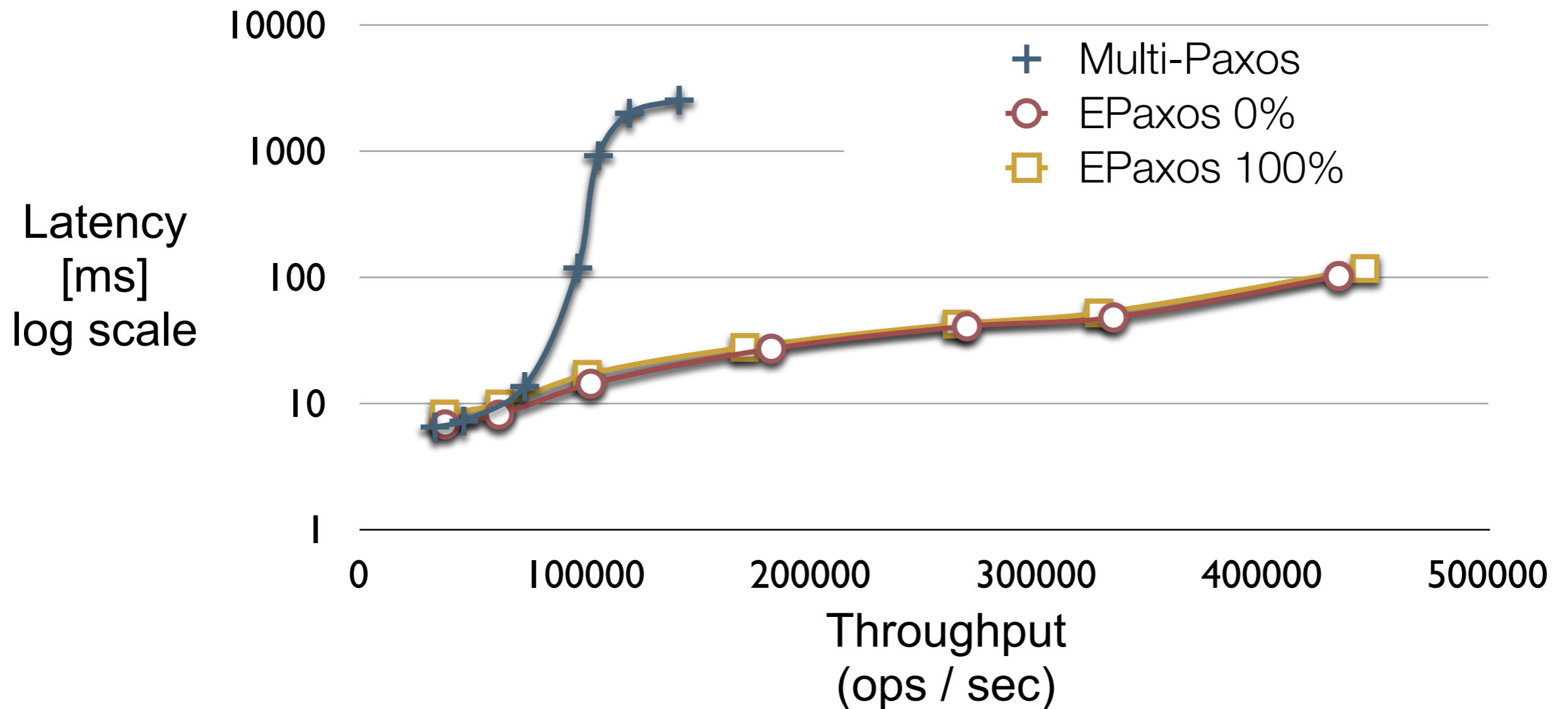
EPaxos: higher throughput w/ batching

5 ms batching, local area

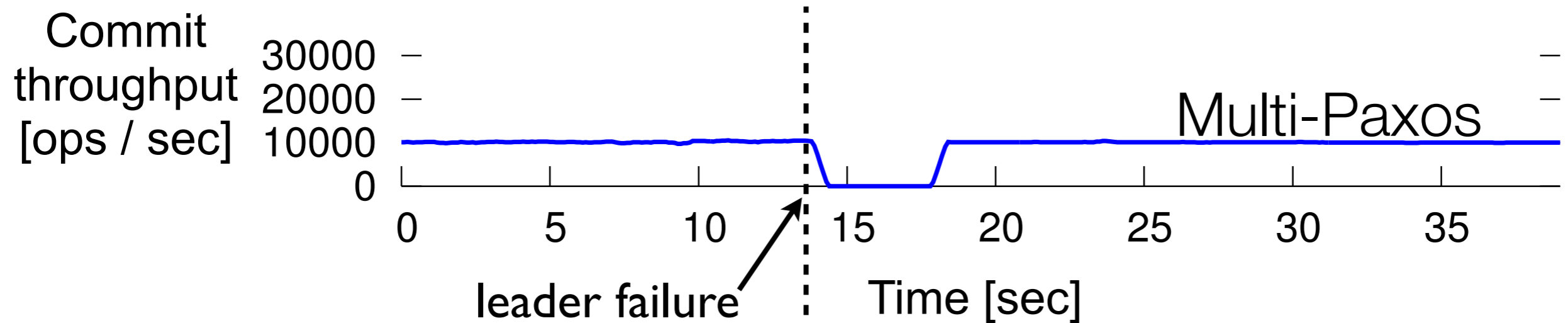


EPaxos: higher throughput w/ batching

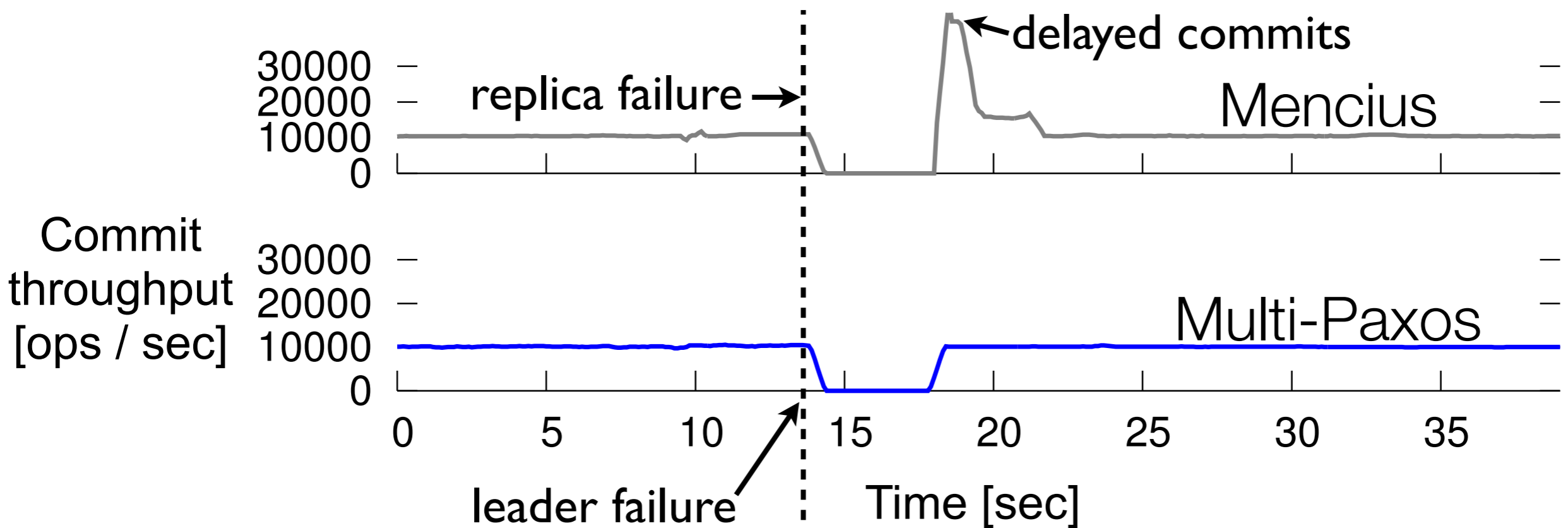
5 ms batching, local area



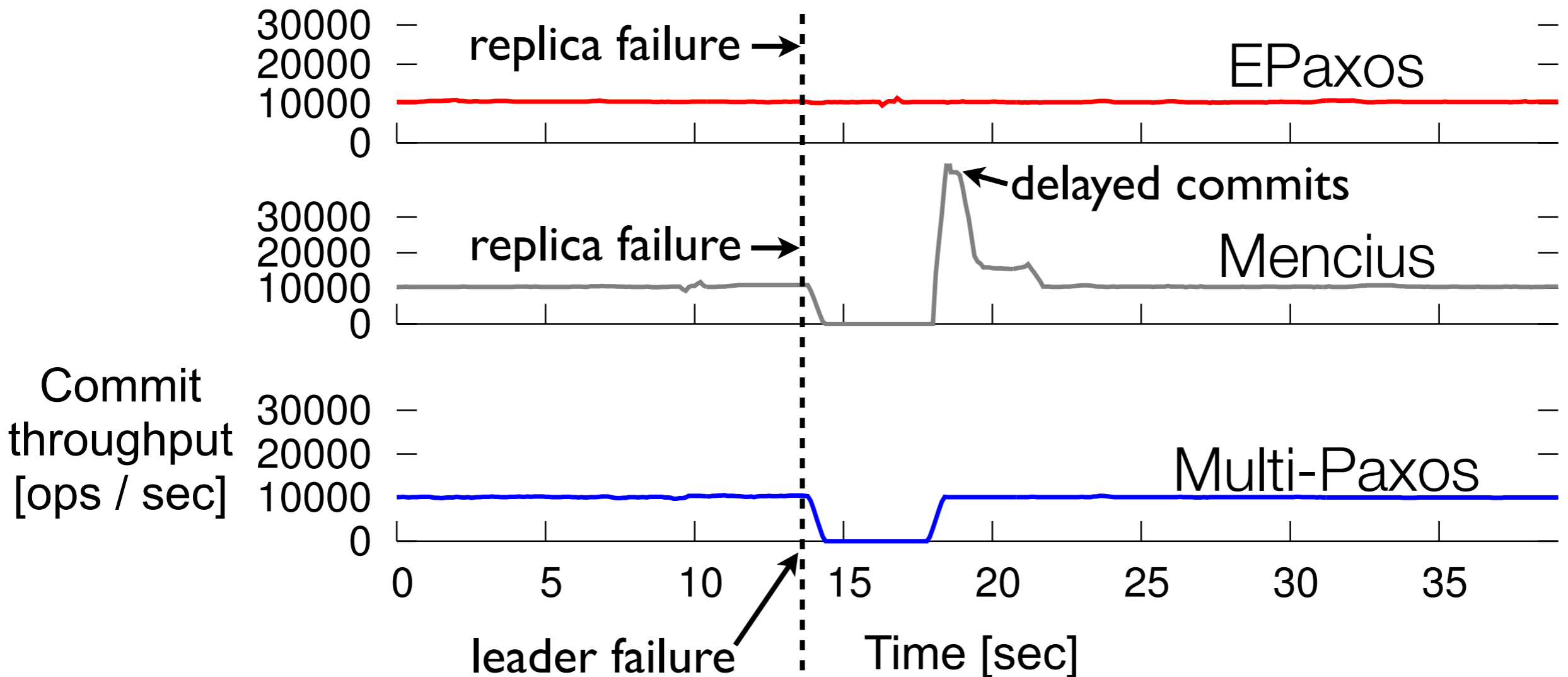
Constant availability



Constant availability



Constant availability



EPaxos insights

EPaxos insights

Order commands explicitly

EPaxos insights

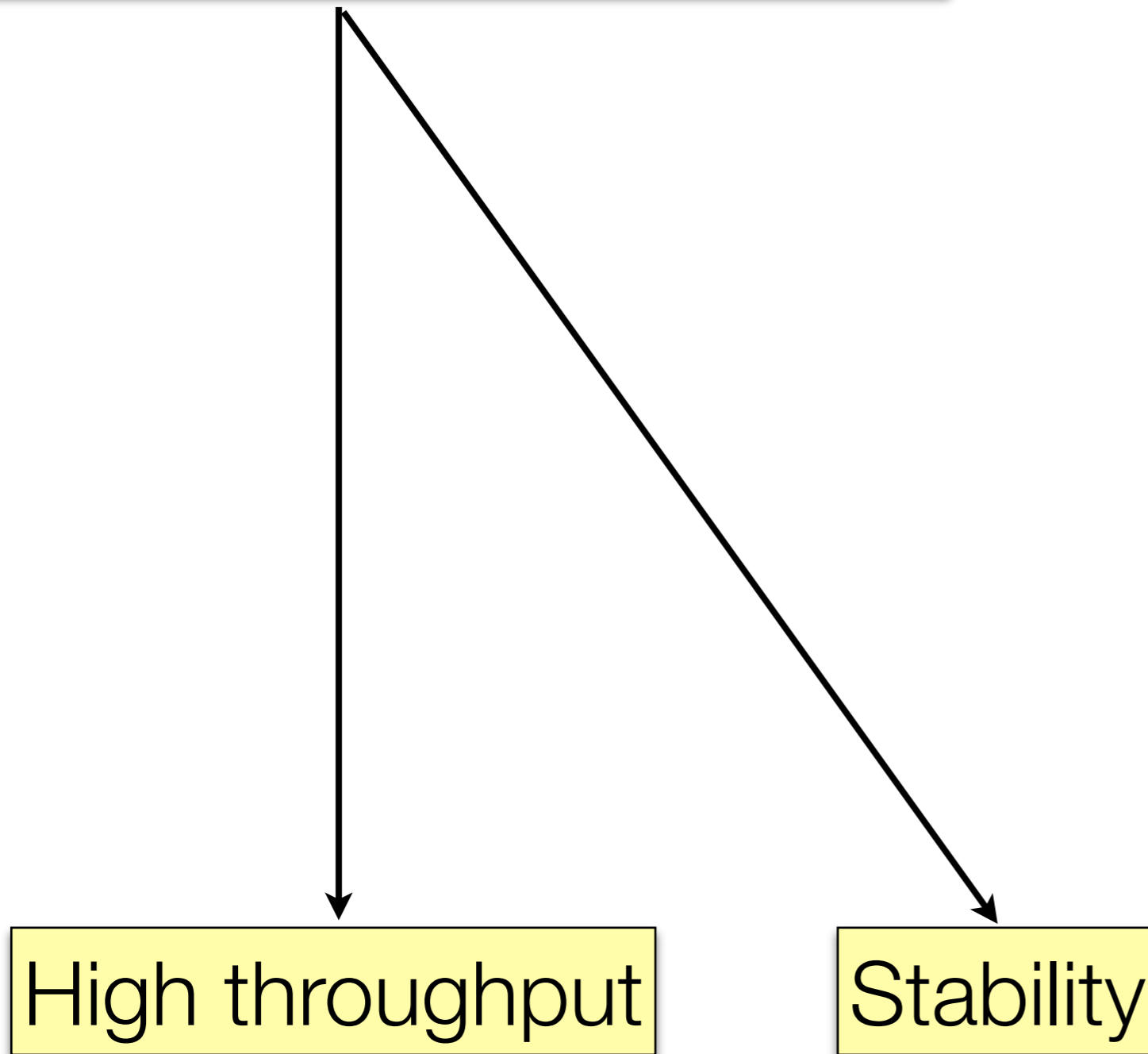
Order commands explicitly



High throughput

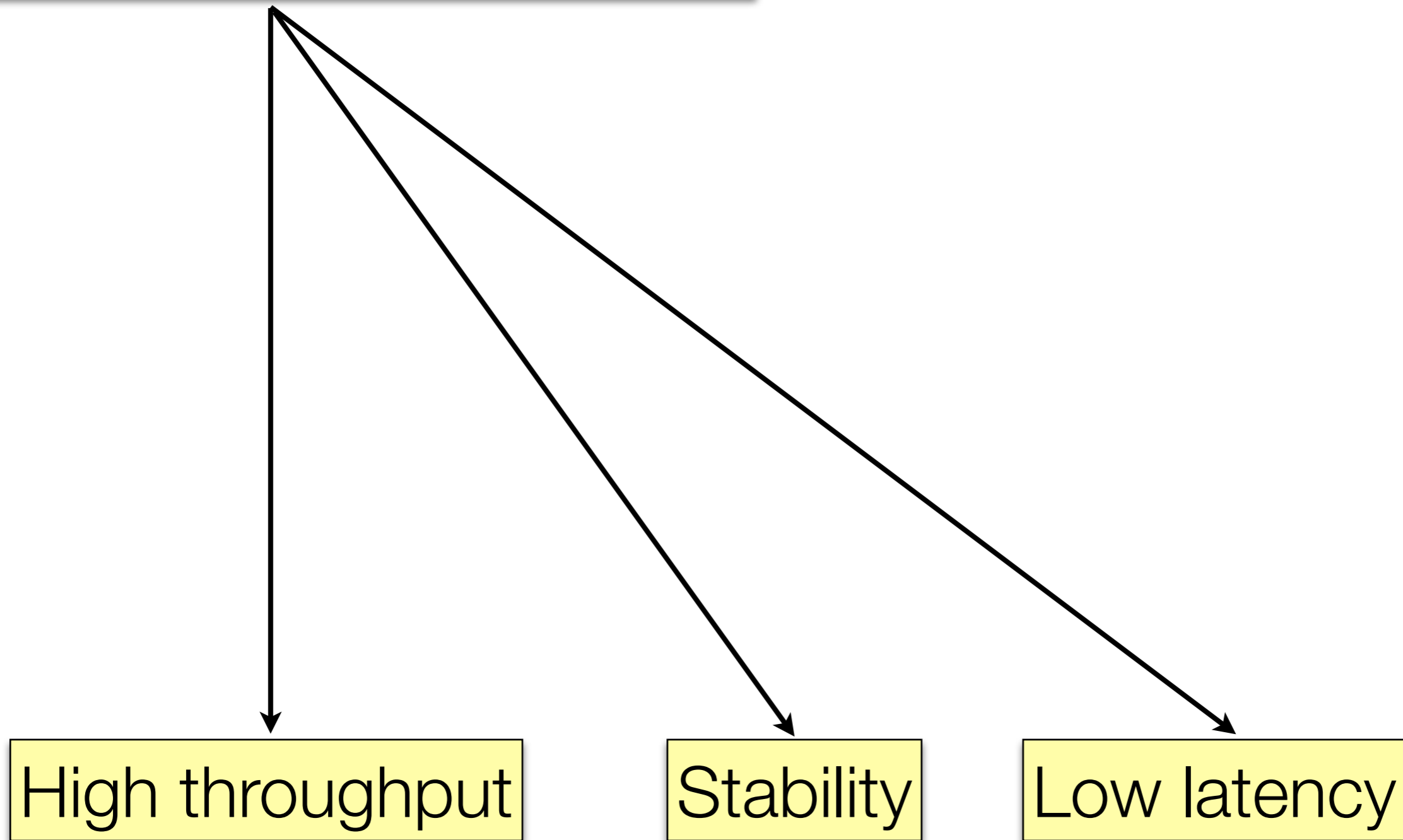
EPaxos insights

Order commands explicitly



EPaxos insights

Order commands explicitly



EPaxos insights

Order commands explicitly

Optimize only delays that matter
(clients co-located w/ closest replica)

High throughput

Stability

Low latency

EPaxos insights

Order commands explicitly

Optimize only delays that matter
(clients co-located w/ closest replica)

Smaller quorums

High throughput

Stability

Low latency

EPaxos insights

Order commands explicitly

Optimize only delays that matter
(clients co-located w/ closest replica)

Smaller quorums

High throughput

Stability

Low latency

Formal Proof

TLA+ Spec

<http://cs.cmu.edu/~imoraru/epaxos/tr.pdf>

Open Source Release



<http://github.com/efficient/epaxos>