# SwiftPaxos: Fast Geo-Replicated State Machines

Fedor Ryabinin[1][2], Alexey Gotsman[1], Pierre Sutra[3]
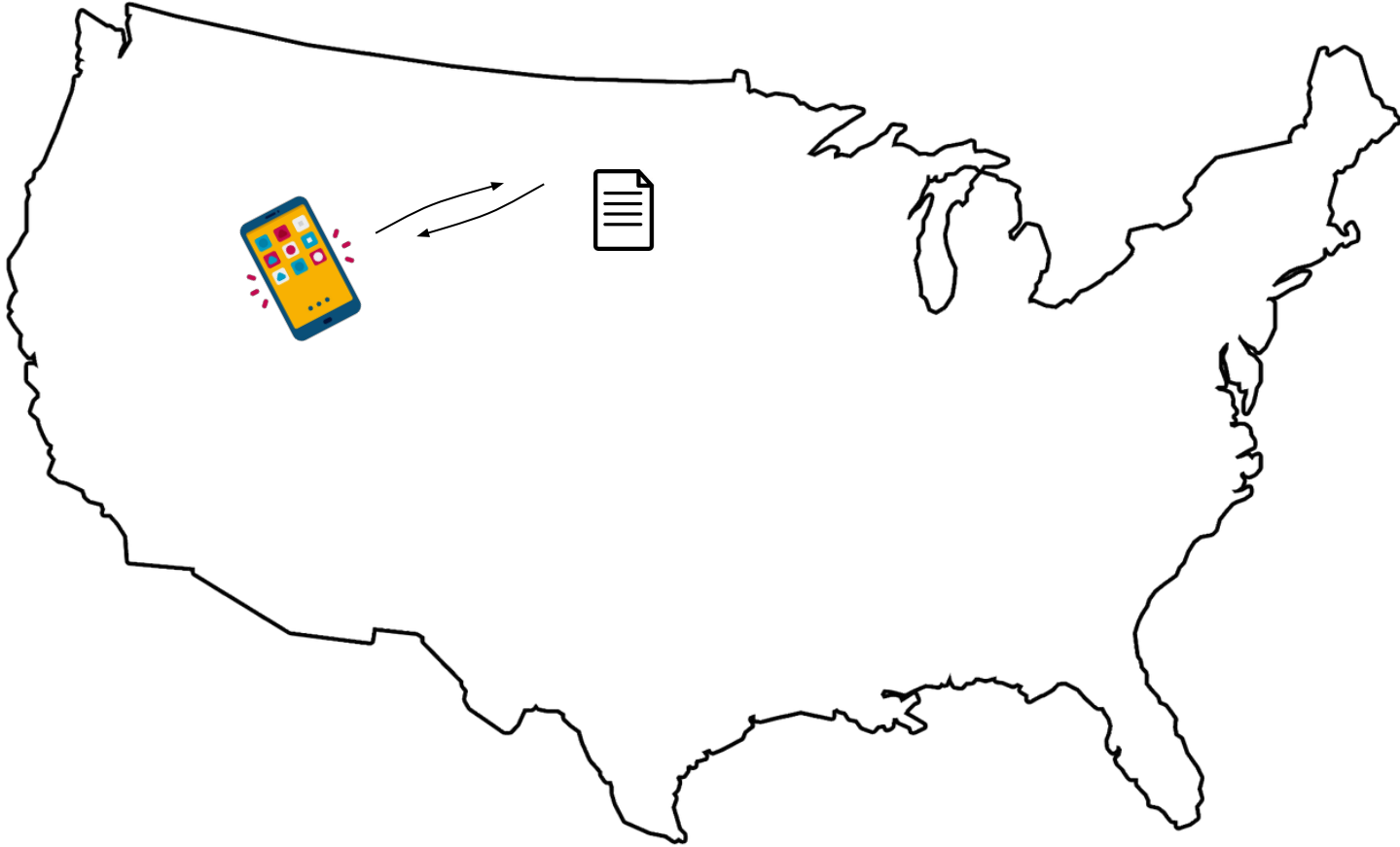
[1] *IMDEA Software*

[2] *Universidad Politécnica de Madrid*
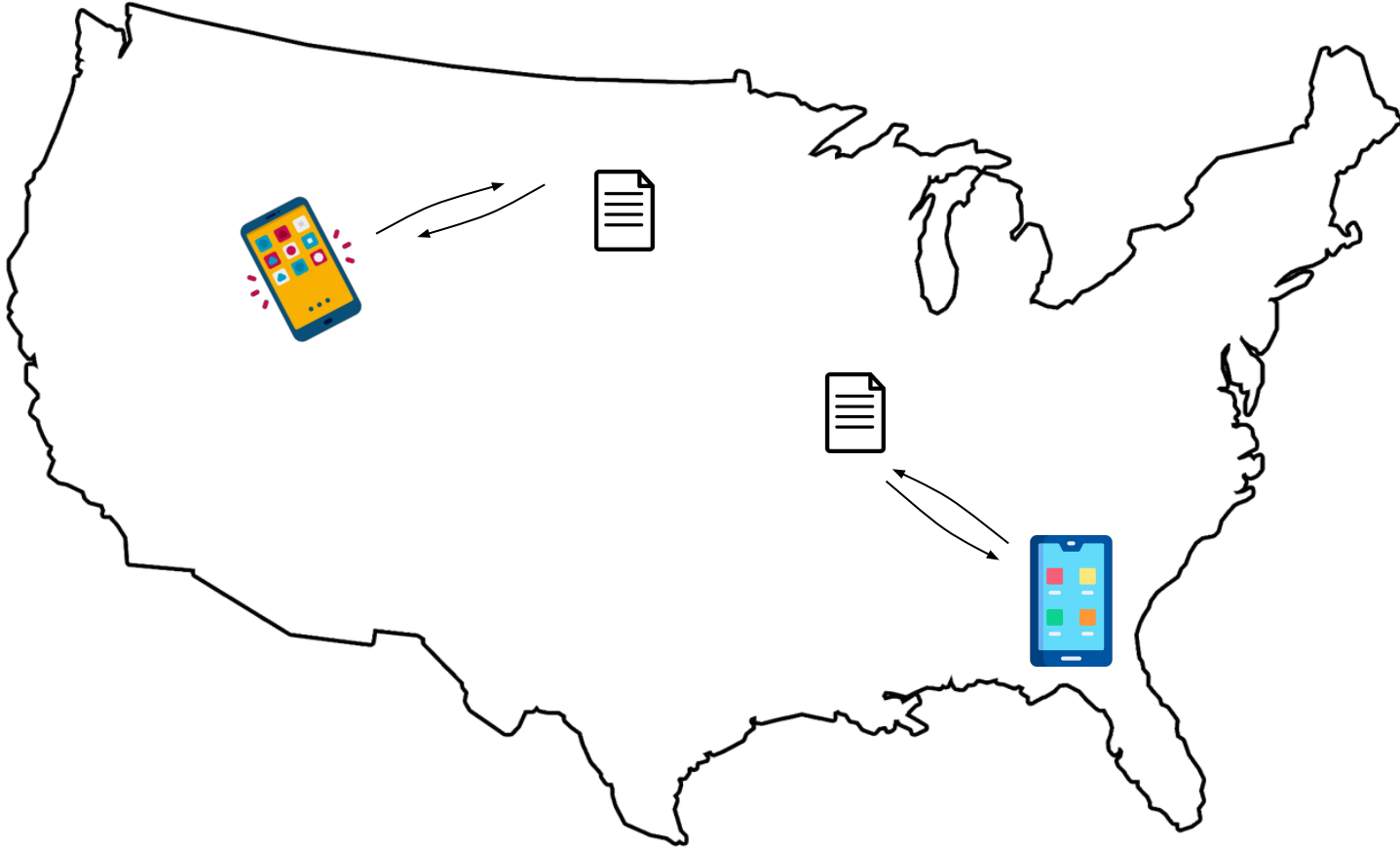
[3] *Telecom SudParis & INRIA*

nsdi'24

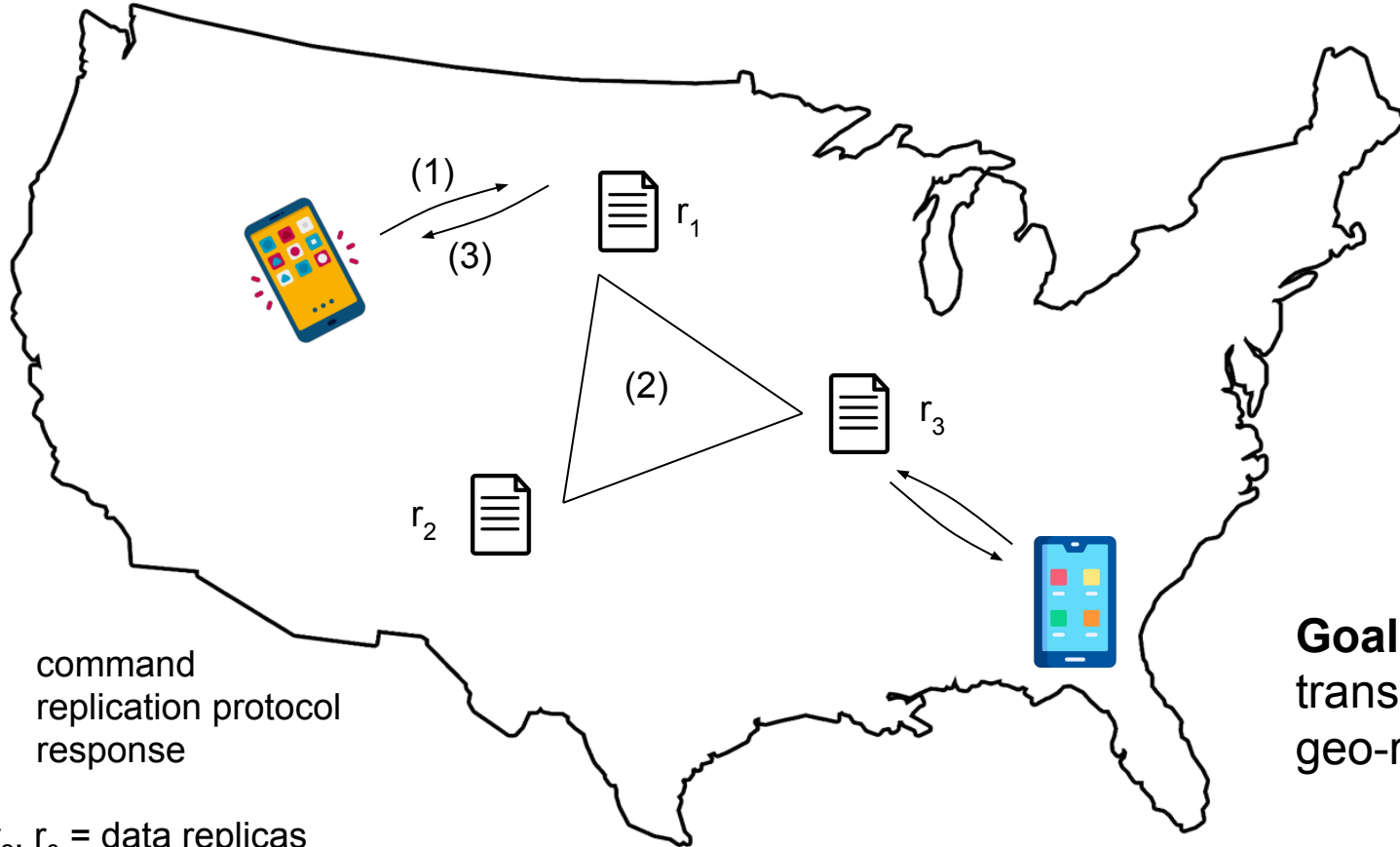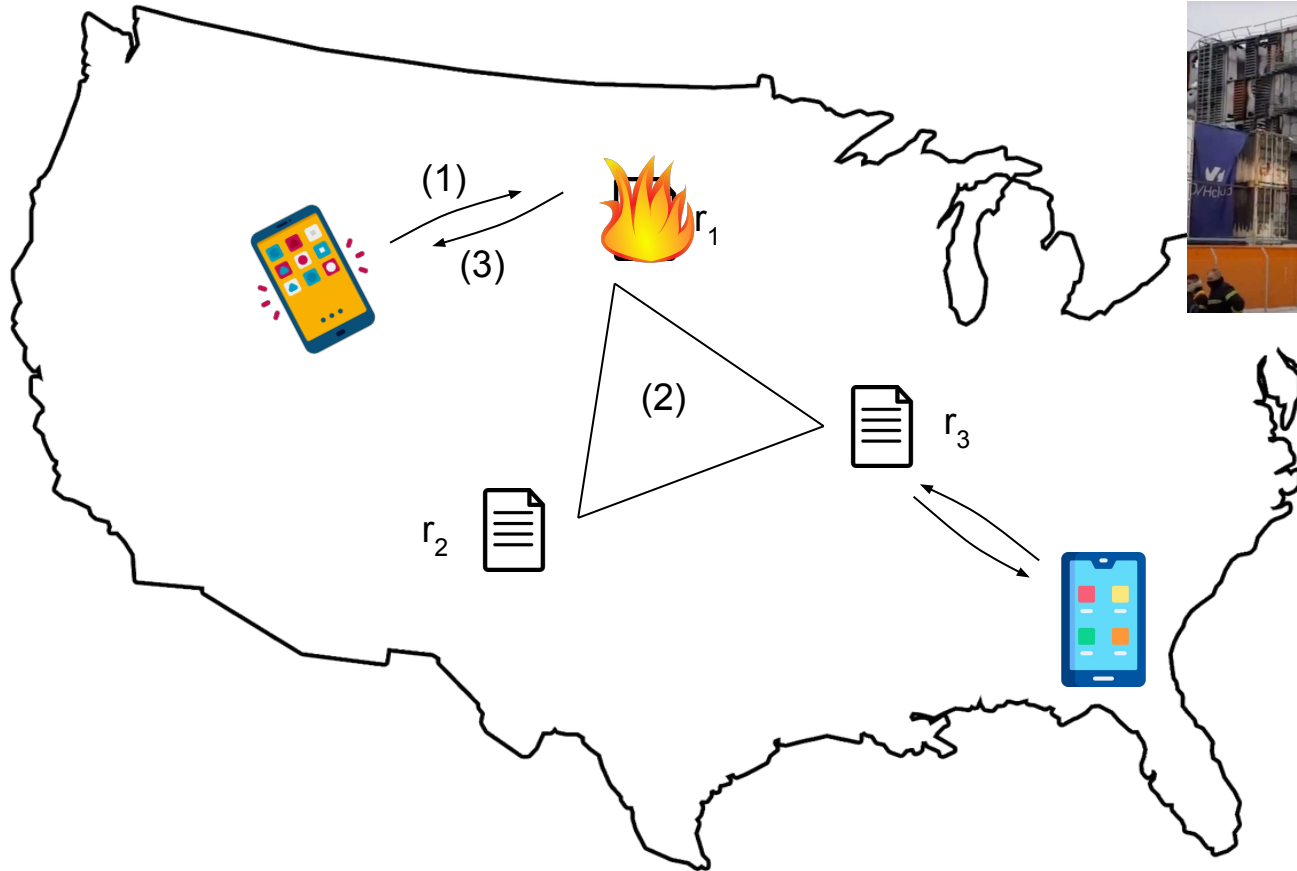# Geo-replicated data

# Geo-replicated data



(1)    command
(2)    replication protocol
(3)    response

$r_1$, $r_2$, $r_3$ = data replicas

**Goal:**
transparent efficient
geo-replication

1

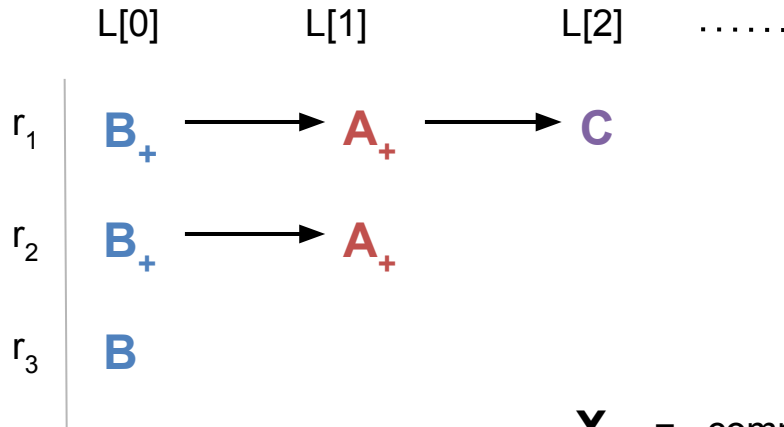**Goal:**
transparent efficient
and *robust*
geo-replication

# State-machine Replication

Each replica holds a log L

Execute commands in log order

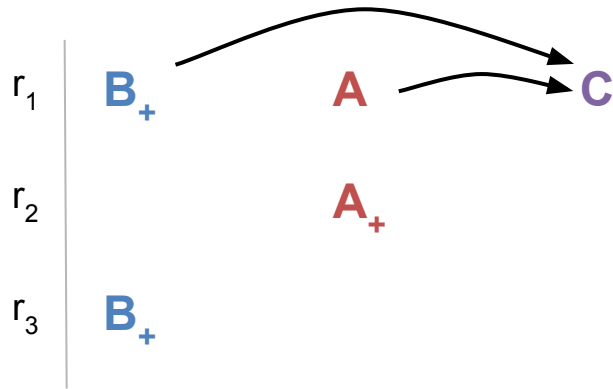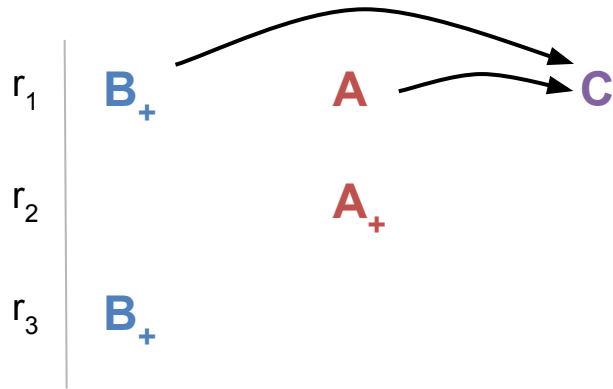To decide a command at position L[i]

- run i-th consensus



L[0]          L[1]          L[2]      · · · · · ·

$r_1$   **B$_+$**  →  **A$_+$**  →  **C**

$r_2$   **B$_+$**  →  **A$_+$**

$r_3$   **B**

**X$_+$**  =  command
             is executed

Execute *conflicting* commands in the same order



$r_1$   B$_+$     A     C

$r_2$     A$_+$

$r_3$   B$_+$

A $= x \leftarrow 42$

B $= y \leftarrow 7$

C $= z \leftarrow x + y$
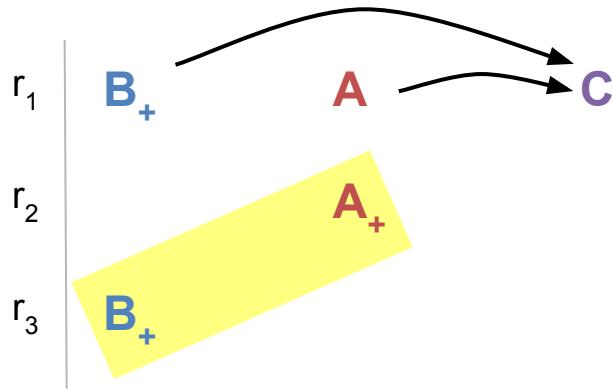
Execute *conflicting* commands in the same order



**How**:
- For each command, compute its *dependencies.*
- Execute commands wrt. dependencies.

A = x ← 42

B = y ← 7

C = z ← x + y

Execute *conflicting* commands in the same order

$r_1$ | $B_+$     $A$     $C$

$r_2$ | $A_+$

$r_3$ | $B_+$

In this example,

- $A$ can execute before or after $B$.

$A$ = x ← 42

$B$ = y ← 7

$C$ = z ← x + y

Execute *conflicting* commands in the same order



In this example,

- A can execute before or after B.

- C depends on both A and B.

A = x ← 42

B = y ← 7

C = z ← x + y

3

Execute *conflicting* commands in the same order



**Invariants:**

- At each replica, dependencies are acyclic.

A $= x \leftarrow 42$

B $= y \leftarrow 7$

C $= z \leftarrow x + y$
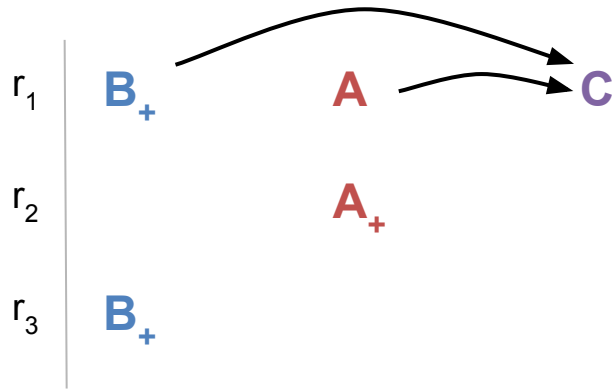
Execute *conflicting* commands in the same order



**Invariants**:

- At each replica, dependencies are acyclic.
- Replicas agree on dependencies.

$A$ = x ← 42

$B$ = y ← 7

$C$ = z ← x + y

Execute *conflicting* commands in the same order



$r_1$   **B**$_+$       **A**       **C**

$r_2$           **A**$_+$

$r_3$   **B**$_+$

**Invariants**:

- At each replica, dependencies are acyclic.
- Replicas agree on dependencies.
- For two conflicting commands **X** and **Y**, either **X** is a dependency of **Y**, or the converse is true.

**A** $= x \leftarrow 42$

**B** $= y \leftarrow 7$

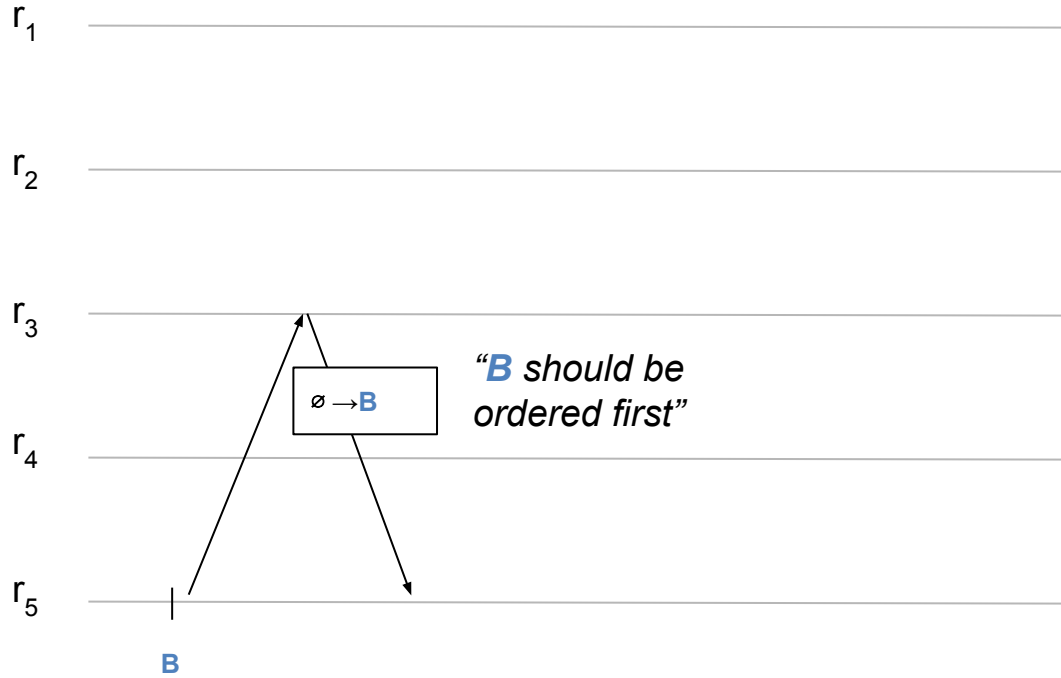**C** $= z \leftarrow x + y$

4

## Summary

- A *new* strongly-consistent replication protocol
- Maintains *at least* Paxos latency
- Executes commands in *optimal time*:
    - 1 RTT when no contention (conflicts already solved by the network)
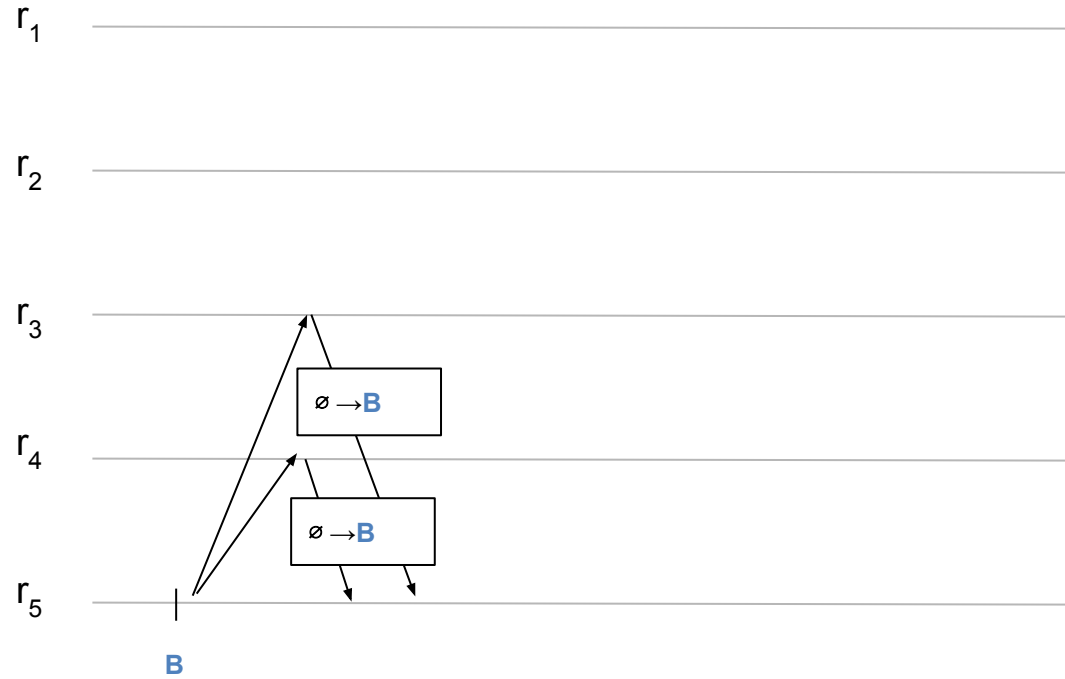    - 1.5 RTT otherwise

## Key novelty

*(double-voting)* in consensus, a replica can vote *twice*, once for its own proposal then for the leader's.
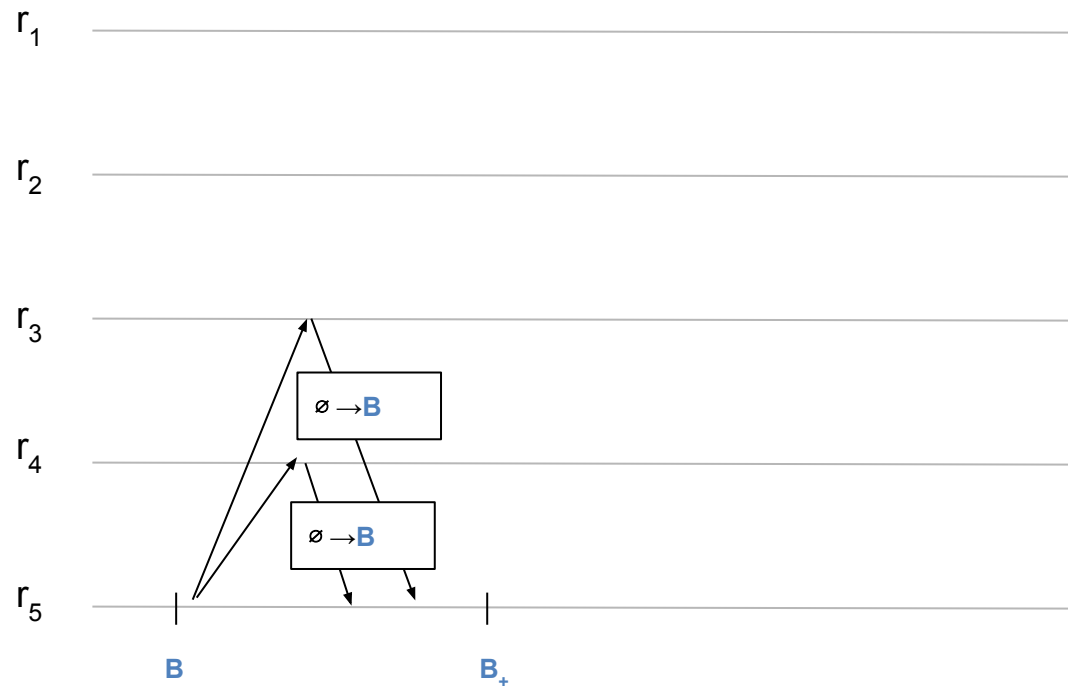
"*B should be ordered first*"

$X_+$ = command is executed

# SwiftPaxos / *best-case latency (2δ)*



$r_1$

$r_2$

$r_3$

$δ$

$r_4$

$δ$

$r_5$

**B**     **B₊**

**Features**
- Optimal best-case latency

$X_+$ = command is executed

$B \rightarrow A$

*"B should be ordered before A"*

$r_1$

**A**

$r_2$

$\emptyset \rightarrow$ **A**

*"nothing precedes A"*

$r_3$

$r_4$

$r_5$

**B**

**A** $= x \leftarrow 42$

**B** $= y \leftarrow 7$

Commands **A** and **B** can execute in any order

# SwiftPaxos / *dependencies tracking*



Features
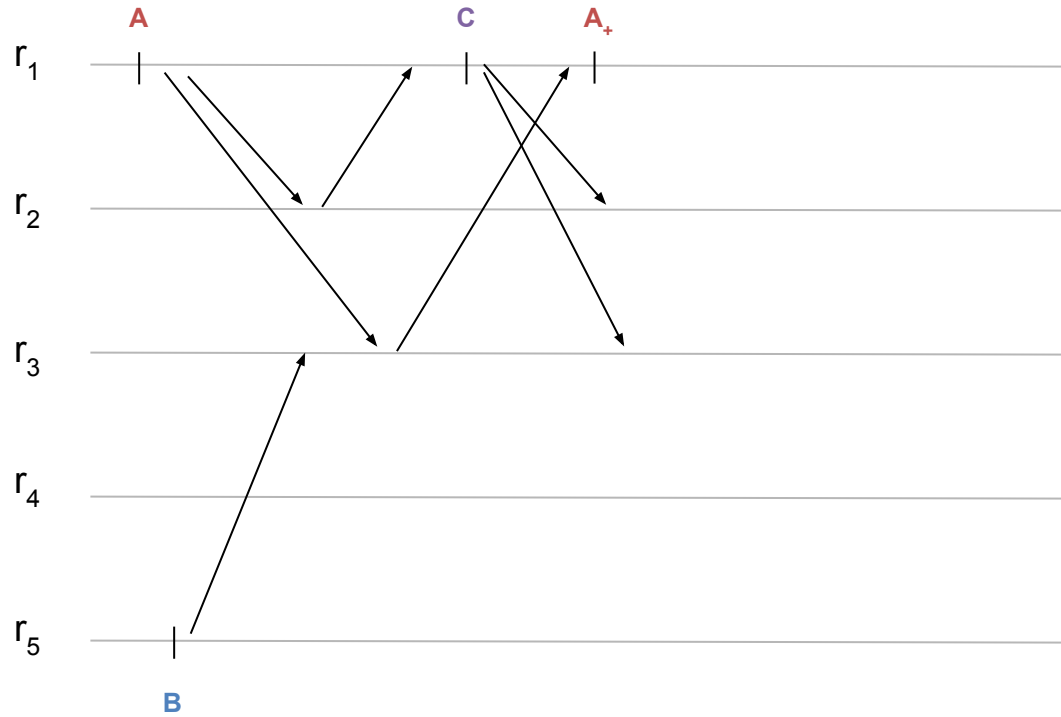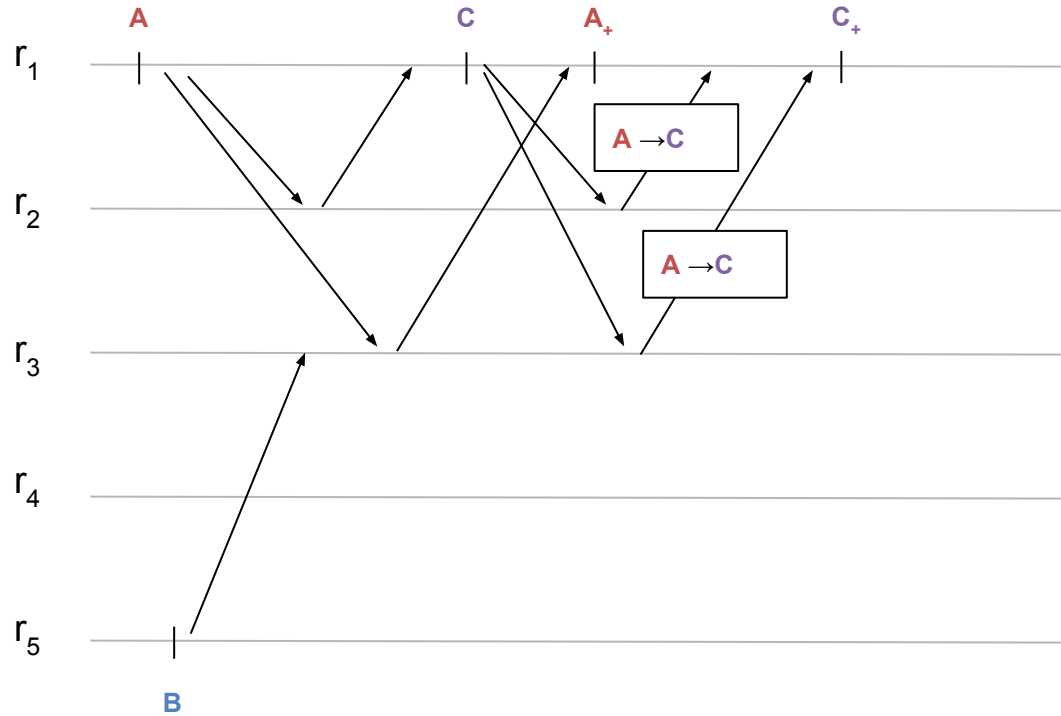- Optimal best-case latency
- Leverage commutativity

A = x ← 42

B = y ← 7

A = x ← 42

B = y ← 7

C = x ← 2x

# SwiftPaxos / *dependencies tracking*



A = x ← 42

B = y ← 7

C = x ← 2x

# SwiftPaxos / *dependencies tracking*



**Features**
- Optimal best-case latency
- Leverage commutativity
- Use *spontaneous order* in the network

A →C

A →C

A = x ← 42

B = y ← 7

C = x ← 2x

# SwiftPaxos / *double voting*



$\mathcal{L}$ = leader replica

B = y ← 7

D = return x + y

$(\mathcal{L})$

B = y ← 7

D = return x + y

$r_1$

$r_2$

D

$r_3$

D →B

(𝓛)   $r_4$

ø →B

$r_5$

B

B = y ← 7

D = return x + y

9

# SwiftPaxos / *double voting*



$(\mathcal{L})$

**B** = y ← 7

**D** = return x + y

# SwiftPaxos / *double voting*



$(\mathcal{L})$

vote twice!

$\varnothing \rightarrow$ **B**

**B** $= y \leftarrow 7$

**D** $=$ return x + y

vote twice!

$\varnothing \rightarrow$ **B**

($\mathcal{L}$)

**B**

*hint*: the leader is part of all the quorums
of the ongoing ballot

**B** $= y \leftarrow 7$

**D** $=$ `return x + y`

9

# SwiftPaxos / *double voting*



**Features**
- Optimal best-case latency
- Leverage commutativity
- Use *spontaneous order* in the network
- At least Paxos speed

$B$ = y ← 7

$D$ = return x + y

$r_1$

$r_2$

$r_3$

$r_4$

$r_5$

D

B $= y \leftarrow 7$

D $= \texttt{return x + y}$

10

$r_1$

$r_2$

$r_3$

$r_4$

$r_5$

D+

D

0

B = y ← 7

D = return x + y

10

($\mathcal{L}$)

$r_1$

$r_2$

$r_3$

$r_4$

$r_5$

$D_+$

$D$

**0**

B = y ← 7

D = return x + y

$r_1$

$r_2$

$(\mathcal{L})$  $r_3$

$D_+$

**0**

$r_4$

$r_5$

$\varnothing \rightarrow$ **D**

**D**

**B** = $y \leftarrow 7$

**D** = `return x + y`

10

$(\mathcal{L})$

B = y ← 7

D = return x + y

$r_1$

B

$r_2$

$(\mathcal{L})$   $r_3$    $B_+$    $D_+$

7    $B \rightarrow D$

$r_4$

$r_5$

D

B = $y \leftarrow 7$

D = return x + y

10

**Features**
- Optimal best-case latency
- Leverage commutativity
- Use *spontaneous order* in the network
- At least Paxos speed
- **Quick reply to non-collocated clients**

B = y ← 7

D = return x + y

10

# SwiftPaxos / *related work*

| | *sequential* | *conflict-free* | *contention-free* | *general* |
|---|---|---|---|---|
| Paxos | 4δ | | | |
| FastPaxos+ | 2δ+1 | 3δ+1 | | |
| Generalized Paxos | 2δ+1 | | | 6δ+1 |
| Egalitarian Paxos | 2δ+1 | | | $O(nδ)$ |
| CURP | 2δ | 3δ+1 | | |
| SwiftPaxos | 2δ | | | 3δ |

*(sequential)* no concurrent commands
*(conflict-free)* concurrent commands do not conflict
*(contention-free)* concurrent conflicting commands are received in the same order everywhere
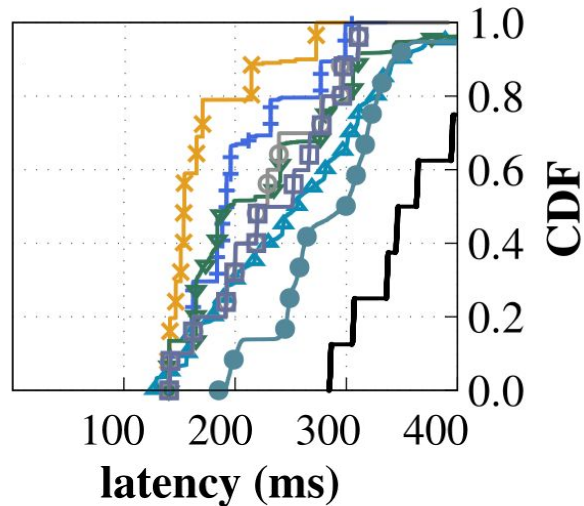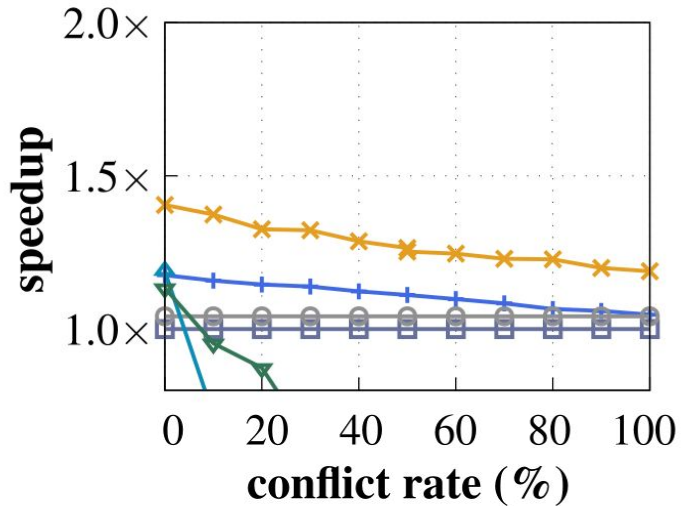
# SwiftPaxos / *related work*

| | *fast path* | *dependencies tracking* | *optimistic execution* | *double-voting* |
|---|---|---|---|---|
| Paxos | ✗ | ✗ | ✗ | ✗ |
| FastPaxos+ | ✓ | ✗ | ✗ | ✗ |
| Generalized Paxos | ✓ | ✓ | ✗ | ✗ |
| Egalitarian Paxos | ✓ | ✓ | ✗ | ✗ |
| CURP | ✓ | ✓ | ✓ | ✗ |
| SwiftPaxos | ✓ | ✓ | ✓ | ✓ |

✓ *supported*
✓ *partially supported*
✗ *not supported*

SwiftPaxos ✳   EPaxos ▲   FastPaxos+ ●   Paxos ▱

CURP+ +   GPaxos ▼   N²Paxos ⊝   Mencius ─

varying conflict rate
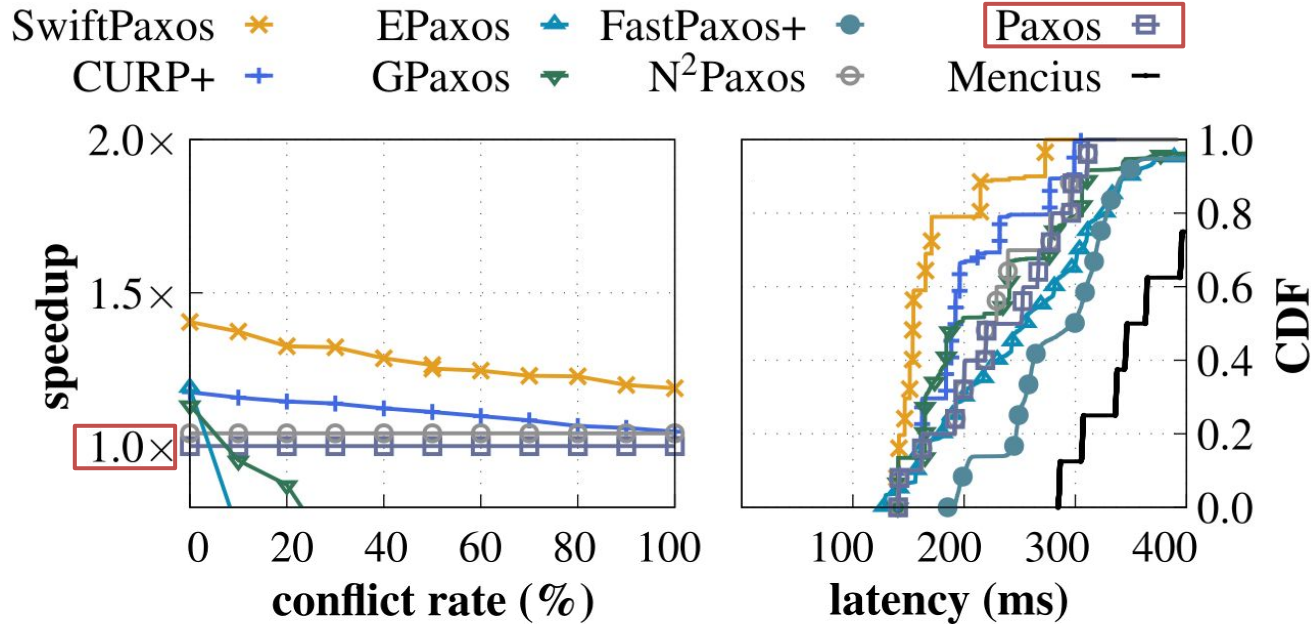
fixed 2% conflict rate

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (1KB)

13

SwiftPaxos  ✳    EPaxos  ▲    FastPaxos+  ●    Paxos  ⊟
CURP+  +    GPaxos  ▼    N$^2$Paxos  ⊖    Mencius  ―



13 AWS EC2 sites
- 5 replica sites
- 10 client sites
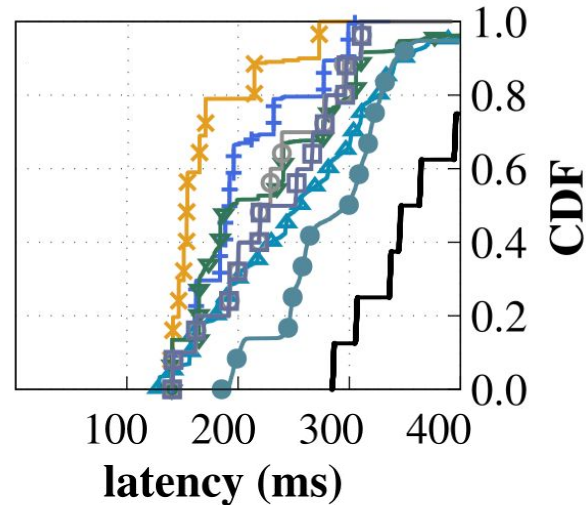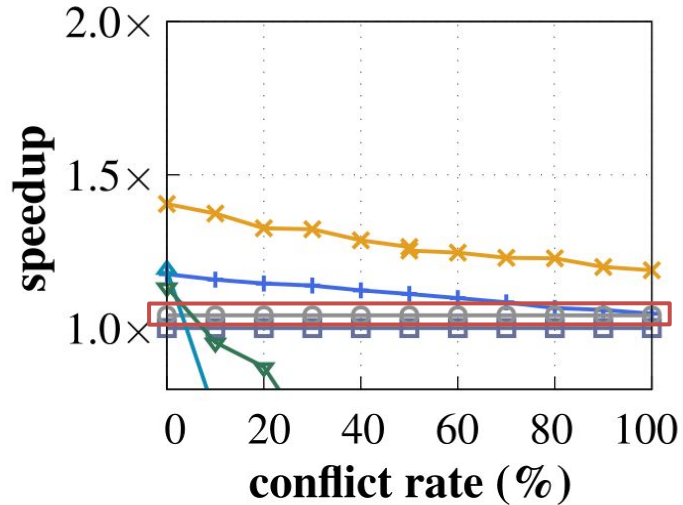*no-op* service (1KB)

- baseline

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (1KB)

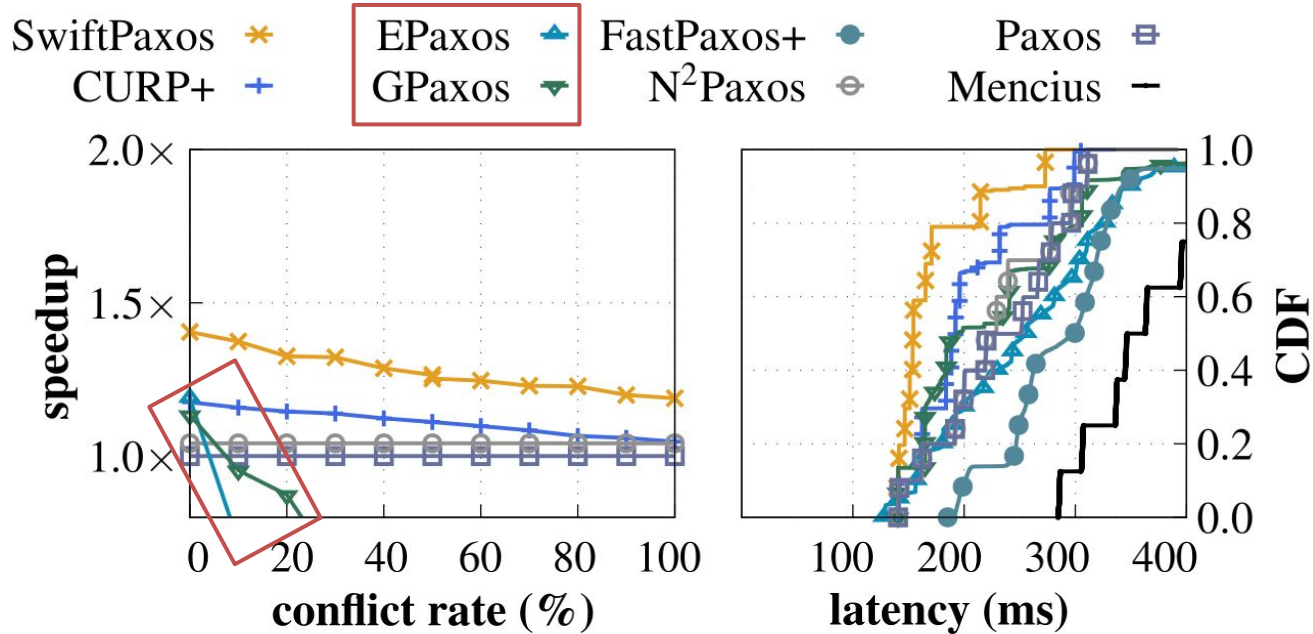- around 5% faster

# Experiments / *average latency*



SwiftPaxos ✳   EPaxos ▲   FastPaxos+ ⬤   Paxos ⊟
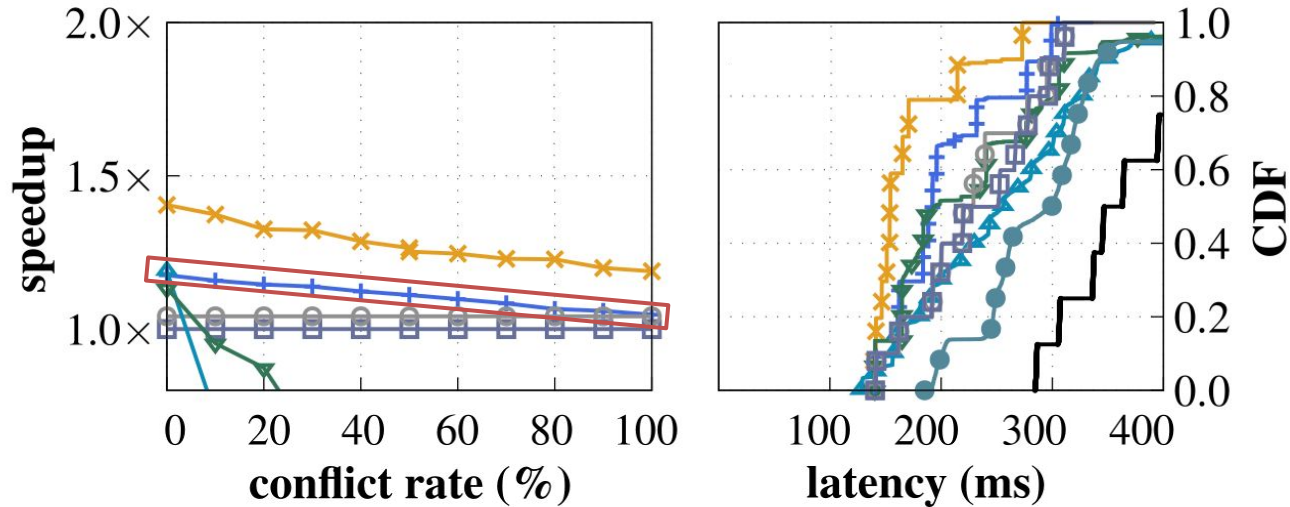CURP+ ✛   GPaxos ▼   N²Paxos ⊖   Mencius ▬

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (1KB)

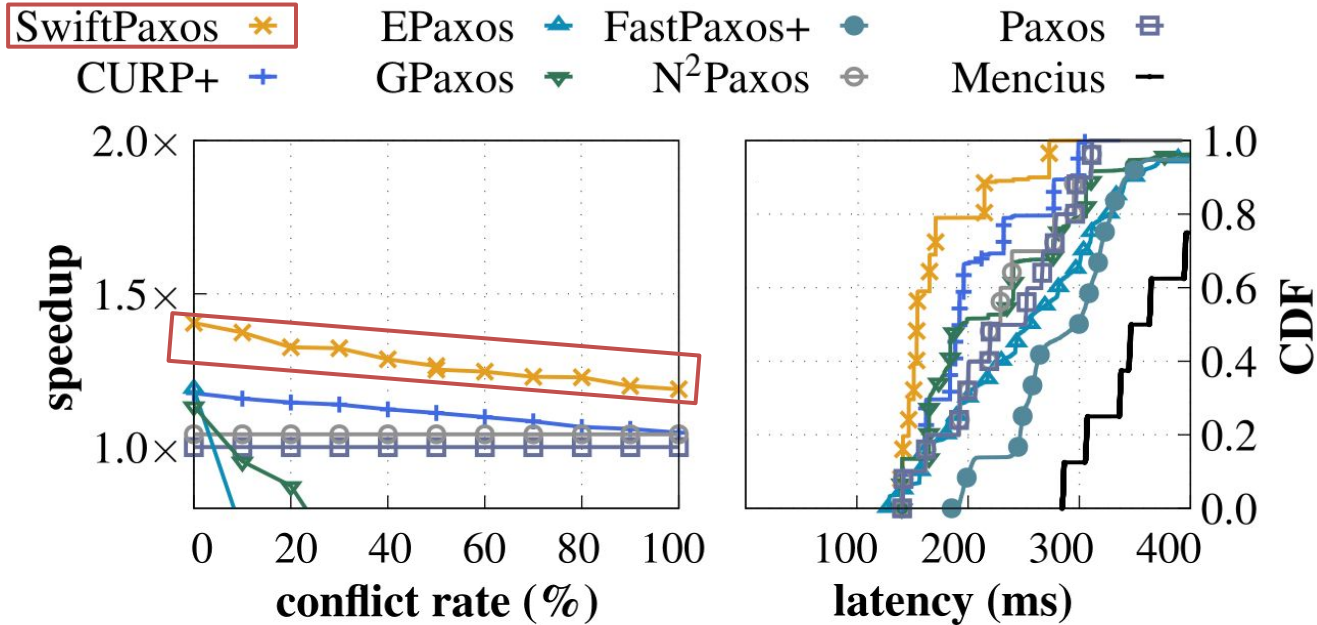- performance quickly drops w. conflicts

# Experiments / *average latency*



13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (1KB)
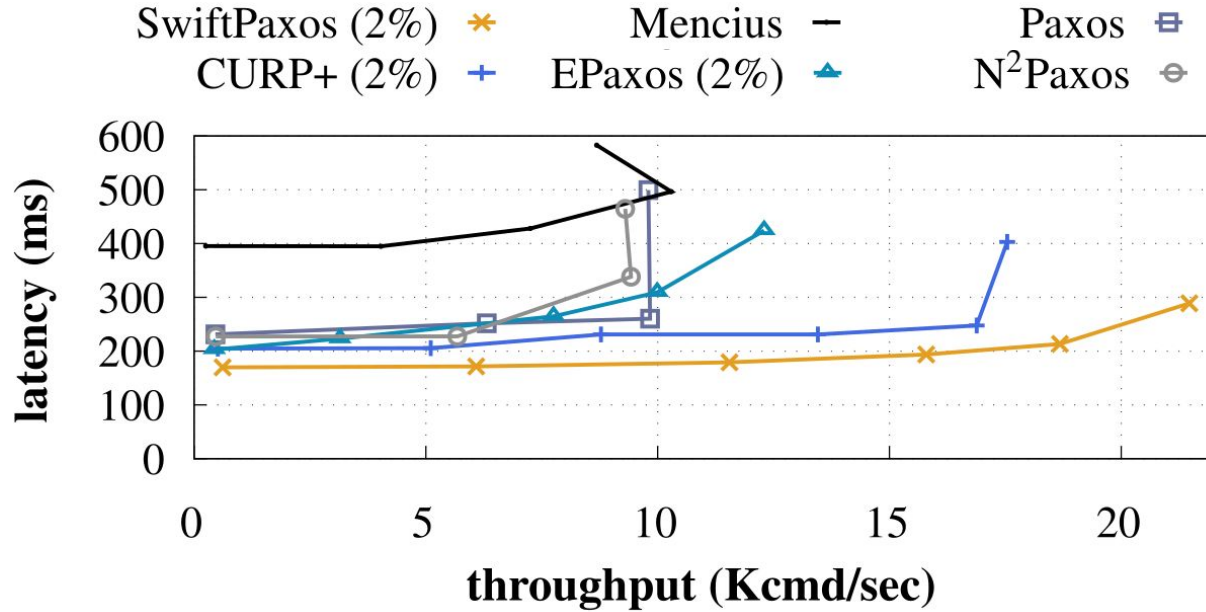
- 18% faster than Paxos on average

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (1KB)

**Takeaways:**

- fastest protocol of all (up to 1.4x)
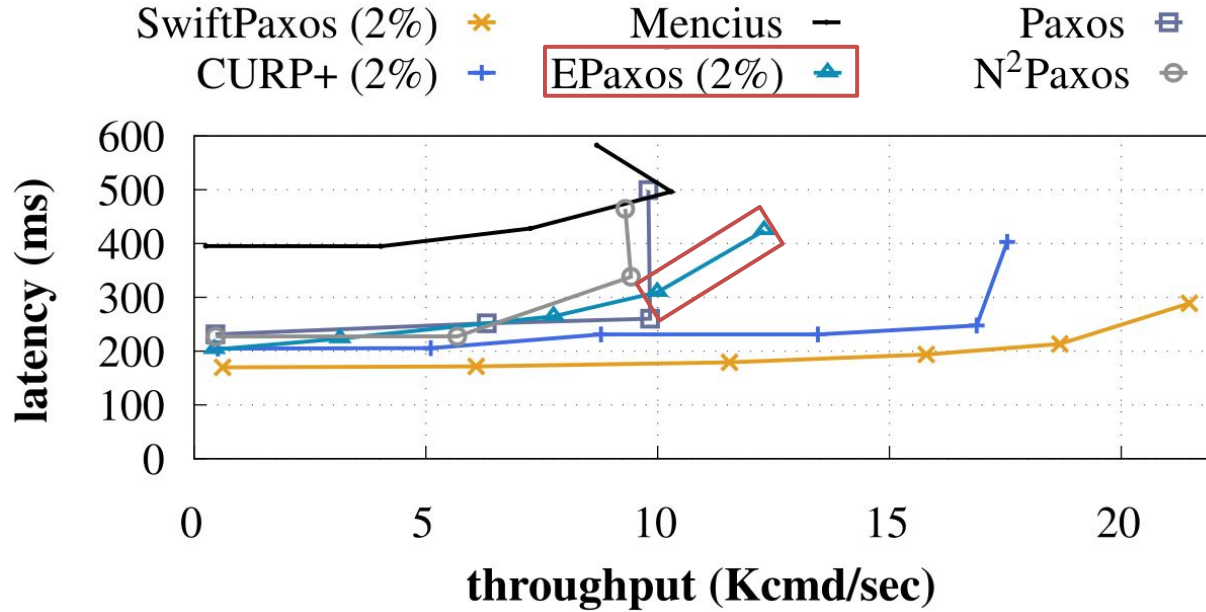- always *at least* Paxos speed

Legend: SwiftPaxos (2%), CURP+ (2%), Mencius, EPaxos (2%), Paxos, $N^2$Paxos

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (3KB)

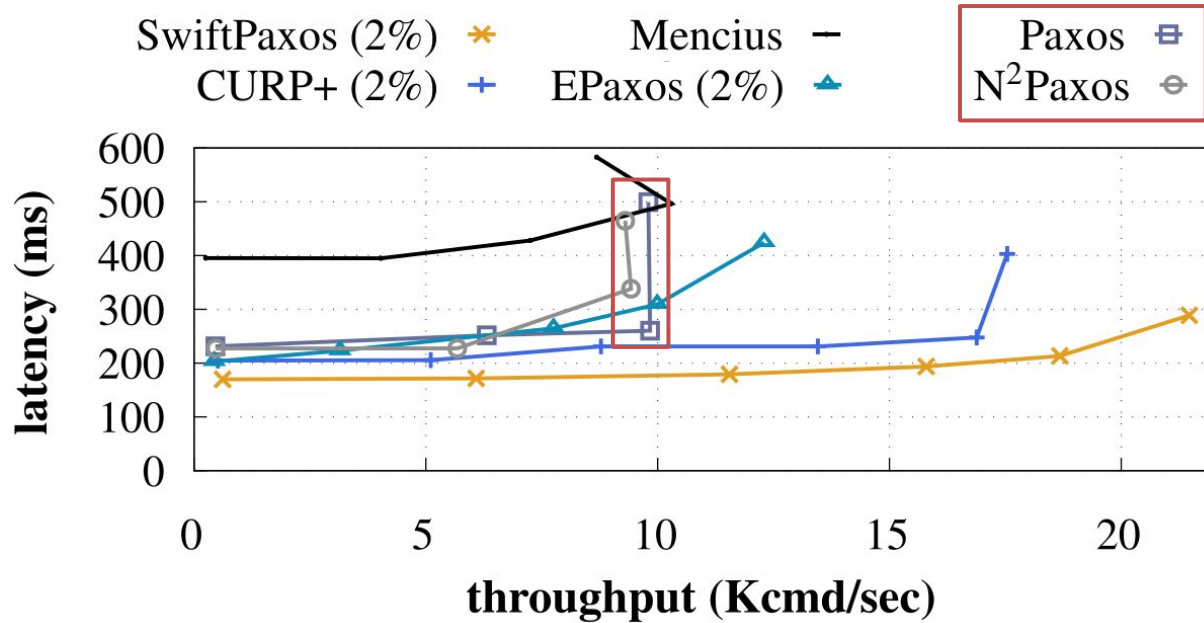progressively add new clients (up to 5,000)
at a fixed 2% conflict rate

14

SwiftPaxos (2%) ✳   Mencius ▬   Paxos ⊟
CURP+ (2%) ╋   EPaxos (2%) ▲   N²Paxos ⊖

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (3KB)

- performance saturates due to convoy effect (long chains)

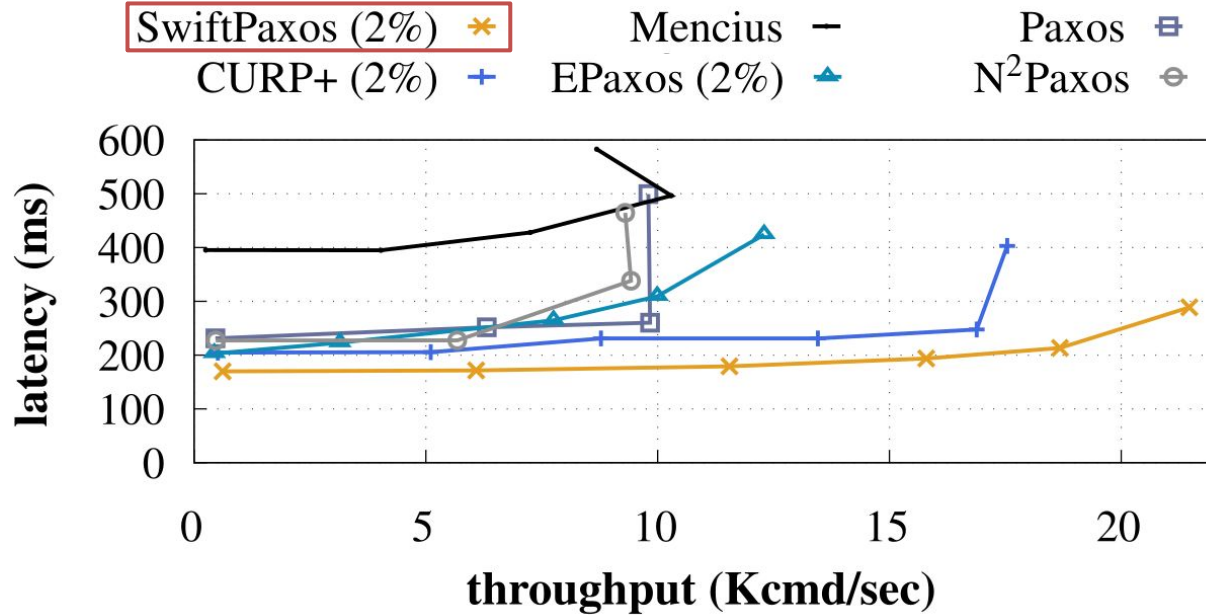SwiftPaxos (2%) ✳    Mencius —    Paxos ⊟
CURP+ (2%) ✚    EPaxos (2%) ▲    N$^2$Paxos ⊖

13 AWS EC2 sites
- 5 replica sites
- 10 client sites
*no-op* service (3KB)

- the leader is bottlenecking (because it disseminates all commands)
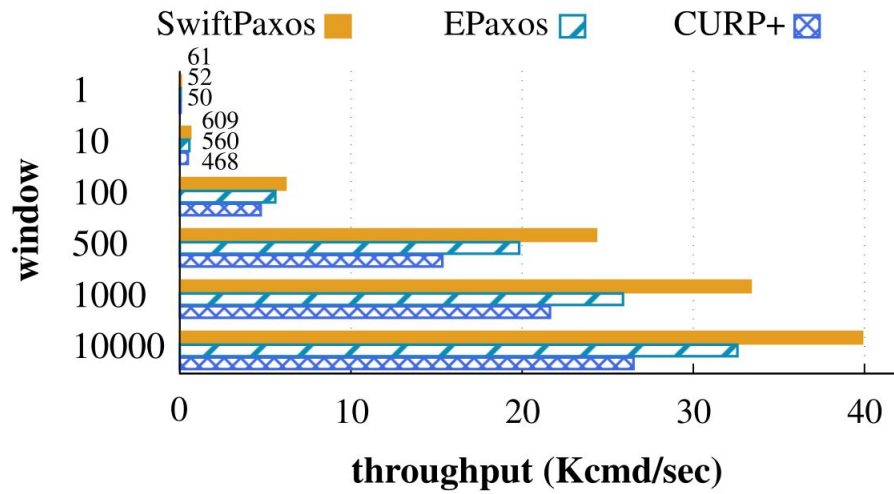
# Experiments / *scalability*



13 AWS EC2 sites
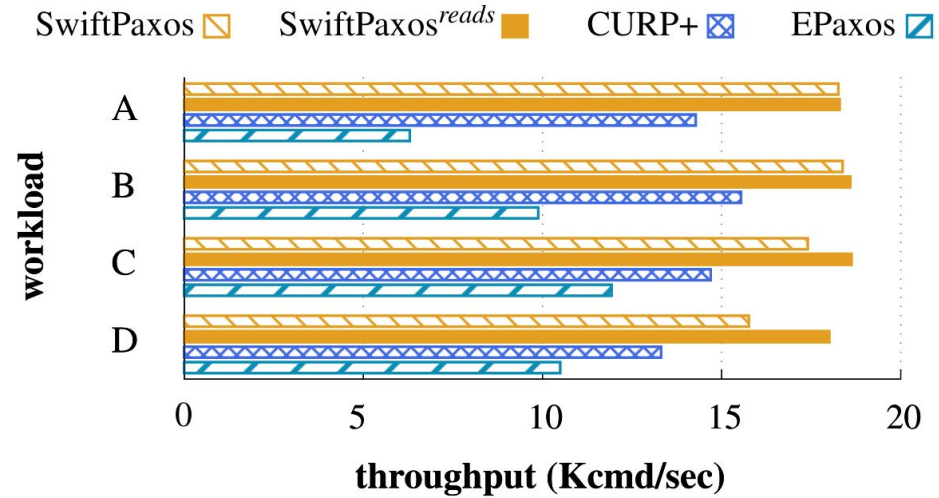- 5 replica sites
- 10 client sites
*no-op* service (3KB)

**Takeaways:**

- stable performance
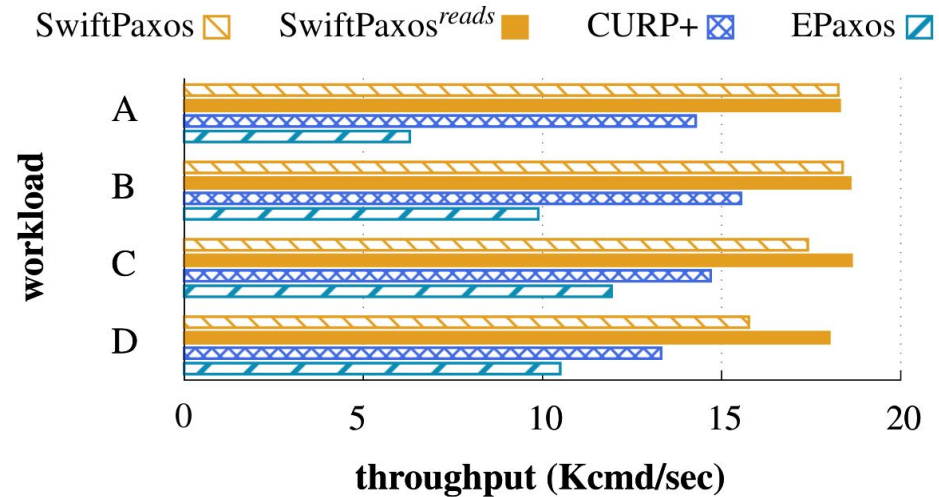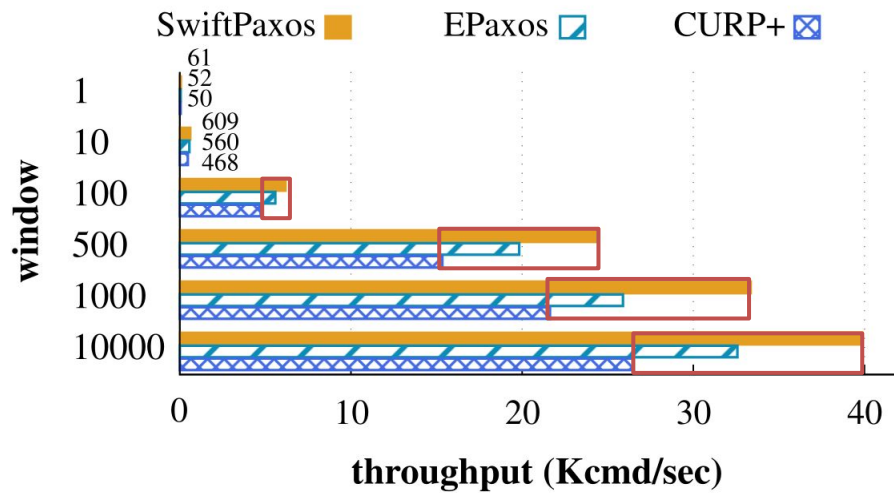- low overhead at the leader
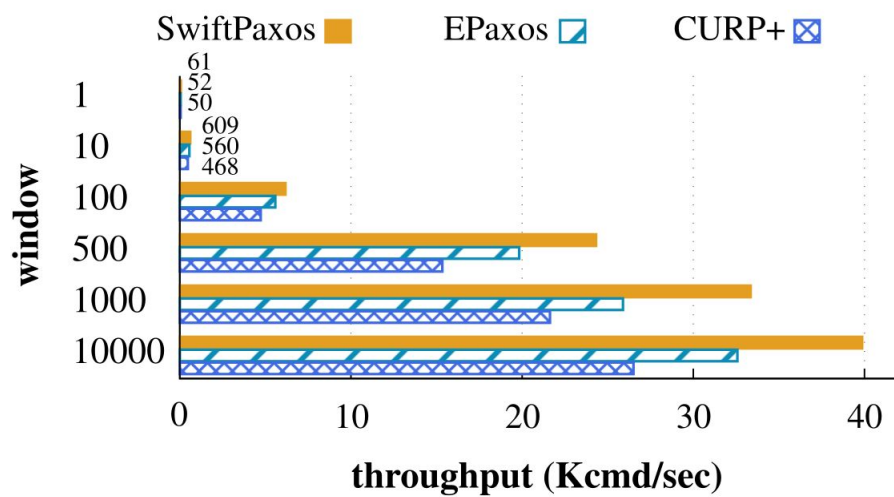
Experiments / *applications*



pipelining

YCSB

15

- up to 49% improvement over CURP [NSDI '19]

**Takeaways:**

- consistently *better performance* than competitors
- can execute *fast linearizable reads* at any replica

# Conclusion

SwiftPaxos

- A *new* (leader-driven) state-machine replication protocol
- Executes commands in optimal time:
  - 1 RTT when no contention
  - 1.5 RTT otherwise.

In practice,

- *always* faster than Paxos (16-29% better)
- up to 2.9x higher throughput than alternatives
- low metadata usage

try it!