



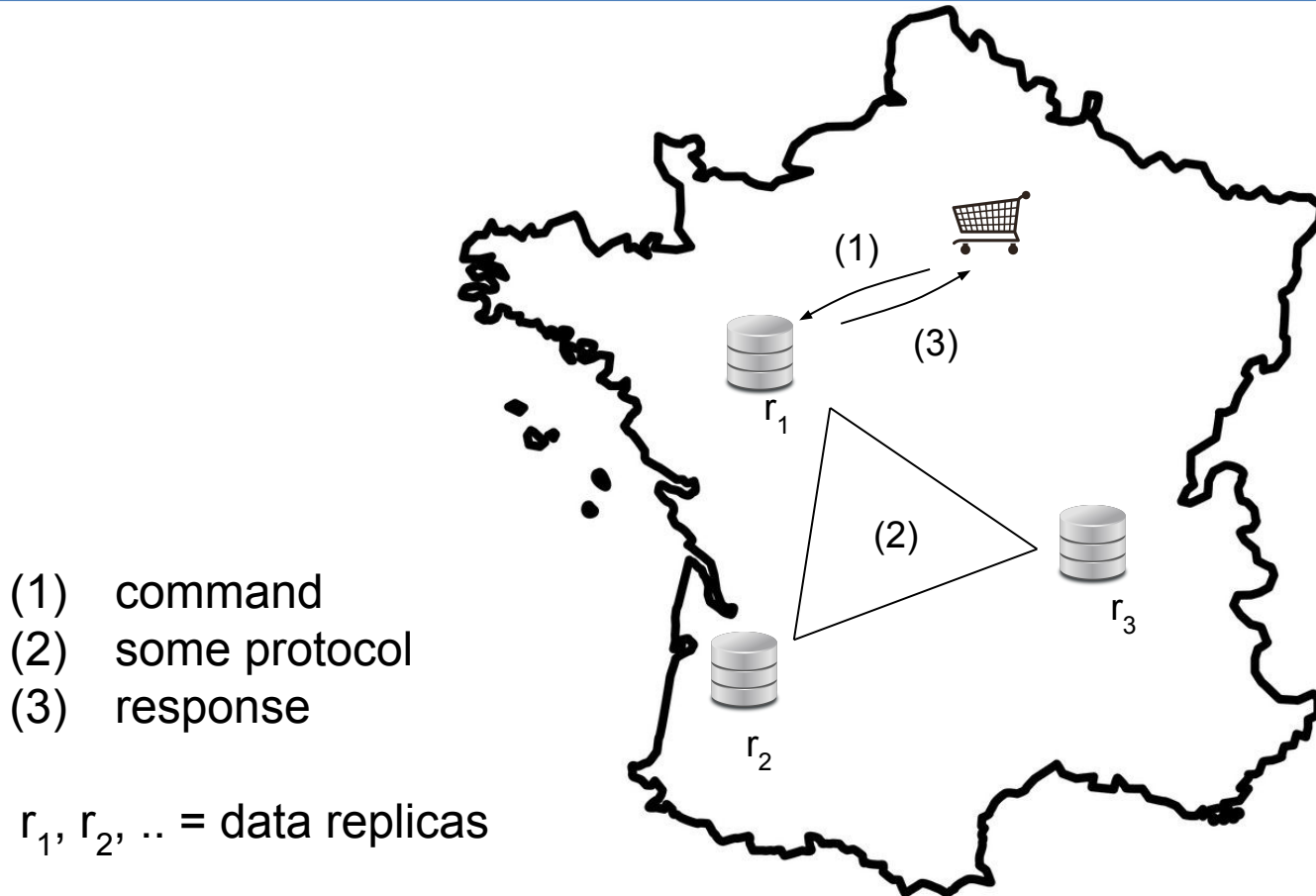
# Leaderless State-Machine Replication: An Overview

Pierre Sutra

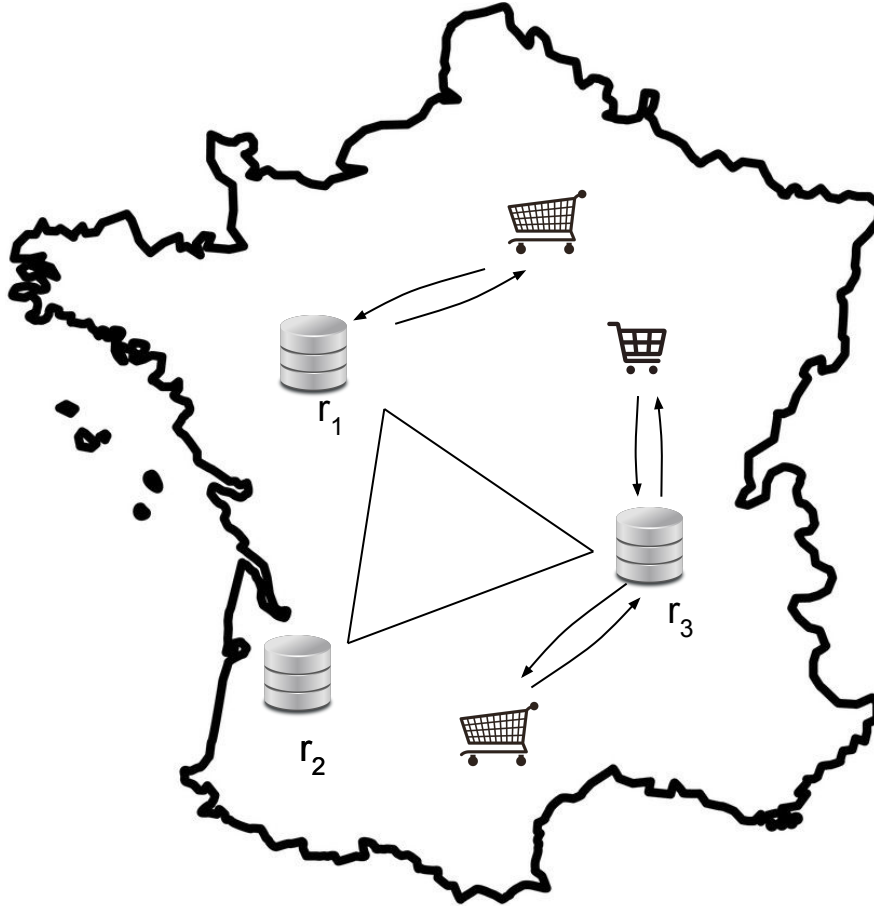
Télécom SudParis  
Institut Polytechnique de Paris

***RainbowFS Final Workshop, March 2022***

## Context: geo-replication

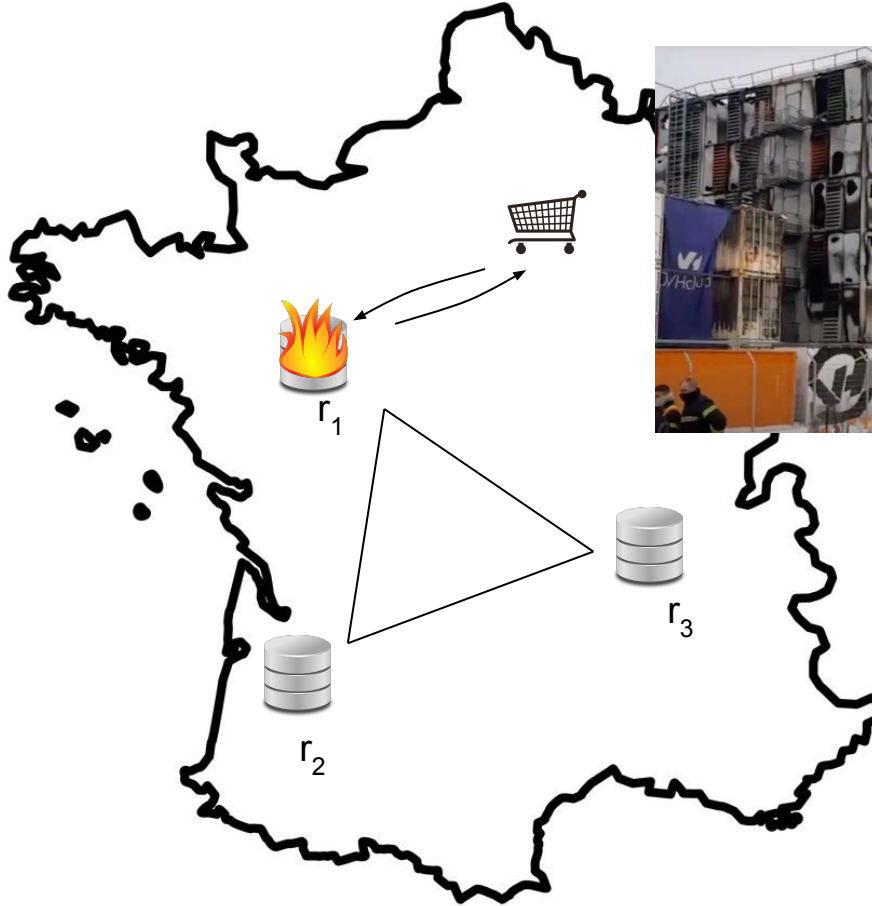


## Context: geo-replication



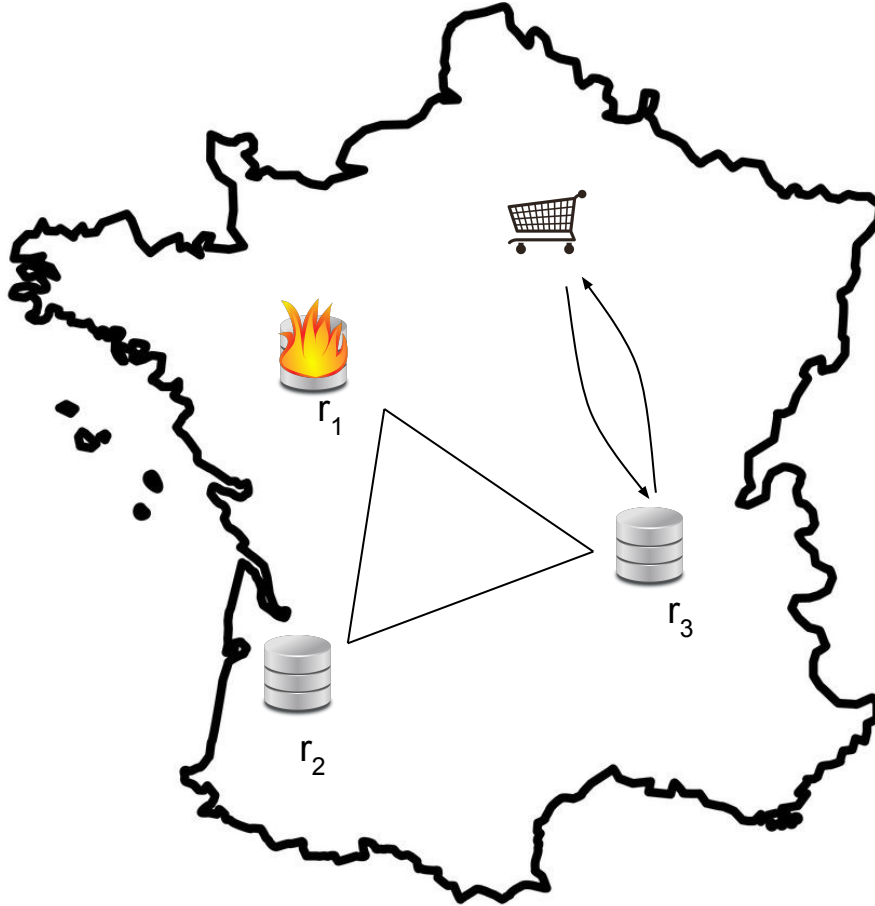
**Problematic:**  
*transparent efficient*  
geo-replication

## Context: geo-replication



**Problematic:**  
*transparent efficient*  
geo-replication

## Context: geo-replication



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

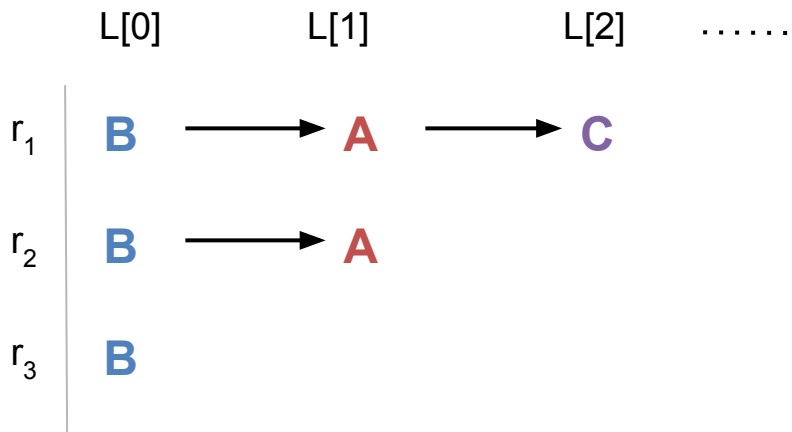
## Classic State-Machine Replication [*Paxos*, *Raft*]

Each replica holds a log  $L$

Execute commands in log order

To append a command at position  $L[i]$

- run  $i$ -th consensus



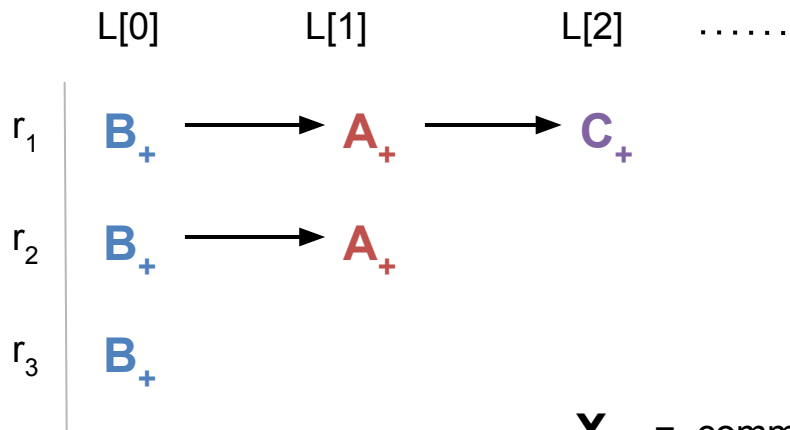
# Classic SMR

Each replica holds a log L

Execute commands in log order

To append a command at position L[i]

- run i-th consensus



X<sub>+</sub> = command  
is executed

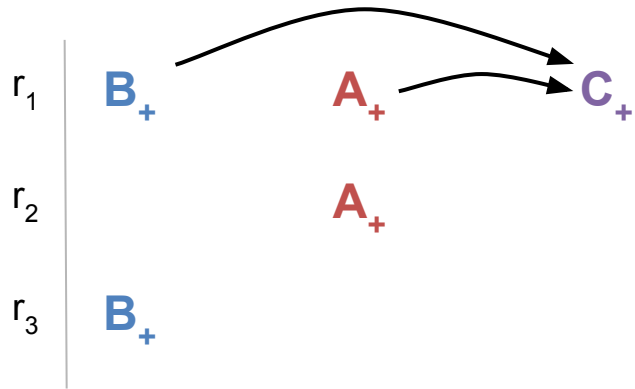


Cloud  
Spanner



## Generic SMR [*GPaxos*, *GBcast*]

Execute *non-commuting* commands in the same order in the log

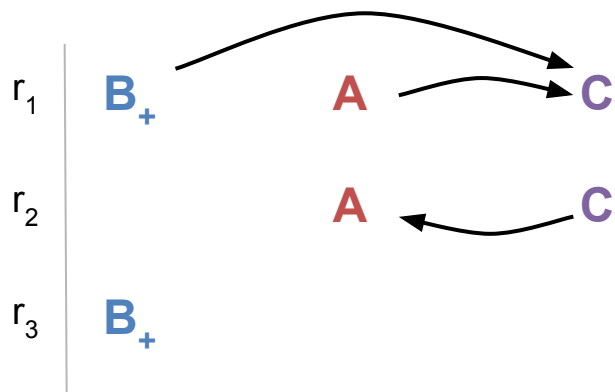


$\left\{ \begin{array}{l} \mathbf{A} = x \leftarrow 42 \\ \mathbf{B} = y \leftarrow 7 \\ \mathbf{C} = z \leftarrow x + y \end{array} \right.$



## Leaderless SMR [DISC'05, SOSP'13]

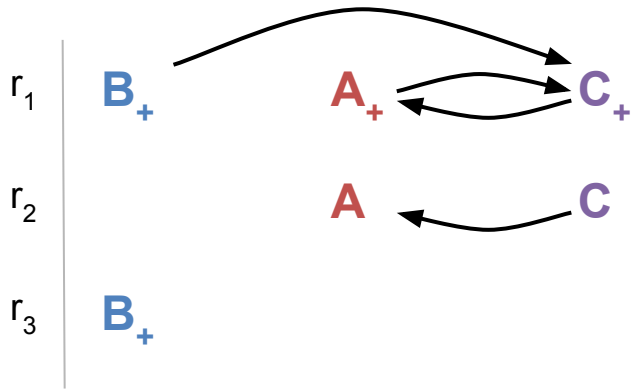
Execute *non-commuting* commands according to the same **graph**



$$\left\{ \begin{array}{l} \text{dep}(\mathbf{A}) = \{\mathbf{C}\} \\ \text{dep}(\mathbf{C}) = \{\mathbf{B}, \mathbf{A}\} \\ \text{dep}(\mathbf{B}) = \emptyset \end{array} \right.$$

## Leaderless SMR

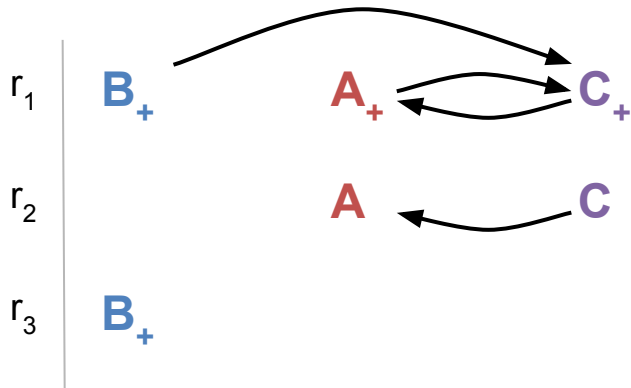
Execute *non-commuting* commands according to the same graph



- operation **X** executed once  $\text{dep}(\mathbf{X})$  transitively closed
- cycles are broken deterministically

# Leaderless SMR

Execute *non-commuting* commands according to the same graph

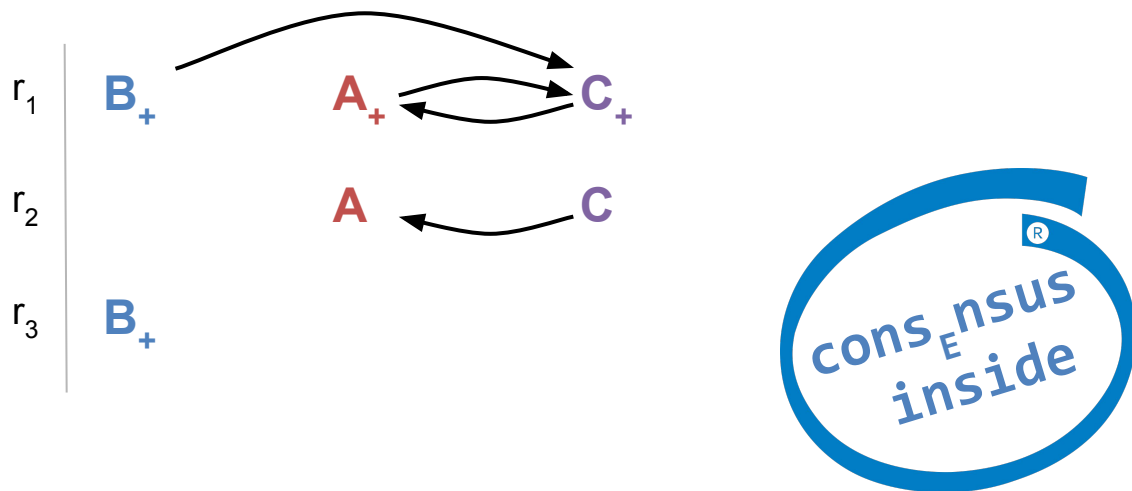


## Properties

- replicas agree *on*  $\text{dep}(\mathbf{X})$
- $(\mathbf{X}, \mathbf{Y})$  non-commuting then  $\mathbf{X} \in \text{dep}(\mathbf{Y})$  or  $\mathbf{Y} \in \text{dep}(\mathbf{X})$

# Leaderless SMR

Execute *non-commuting* commands according to the same graph



## Properties

- replicas **agree** on  $\text{dep}(\mathbf{X})$
- $(\mathbf{X}, \mathbf{Y})$  non-commuting then  $\mathbf{X} \in \text{dep}(\mathbf{Y})$  or  $\mathbf{Y} \in \text{dep}(\mathbf{X})$

## Egalitarian Paxos [SOSP'13]

---

EPaxos uses  $2f+1$  processes ( $f = \max \text{ \#failures}$ )

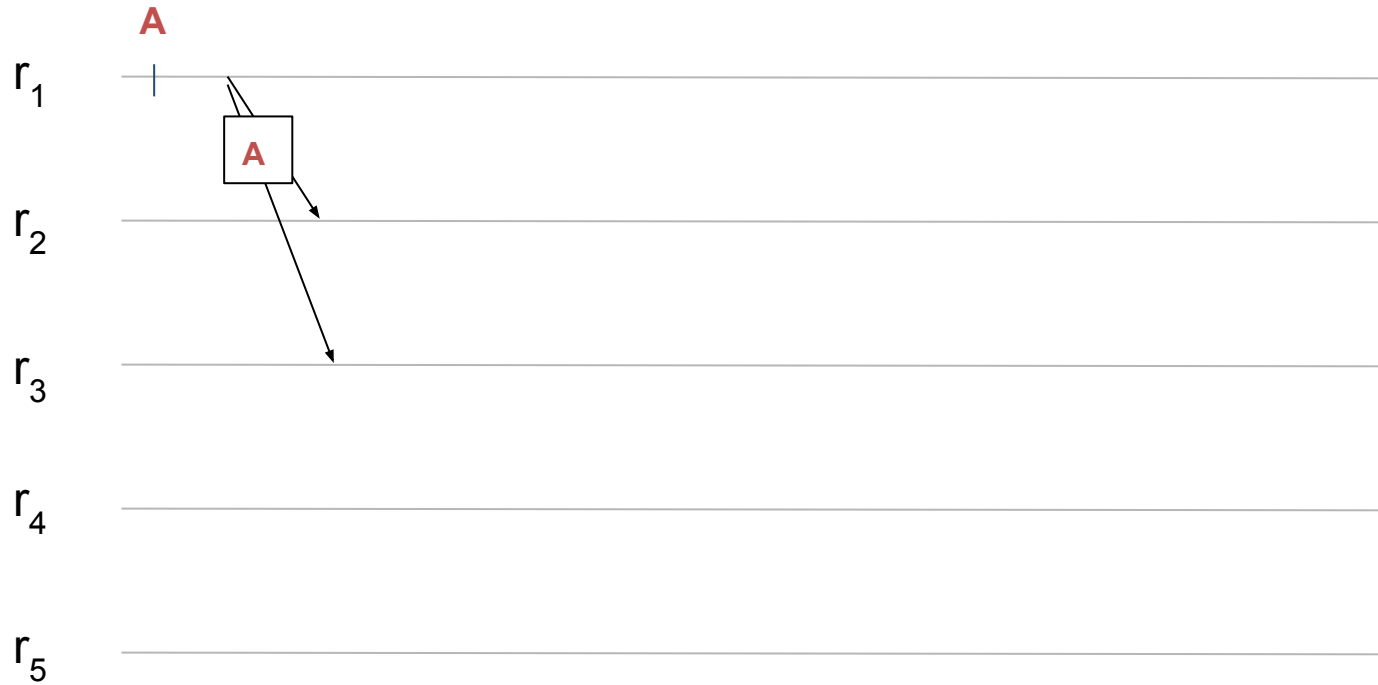
When a client executes command **X**

- pick a replica
- this replica is the *coordinator* for **X**,  $\text{coord}(\mathbf{X})$
- $\text{coord}(\mathbf{X})$  runs consensus over  $\text{dep}(\mathbf{X})$

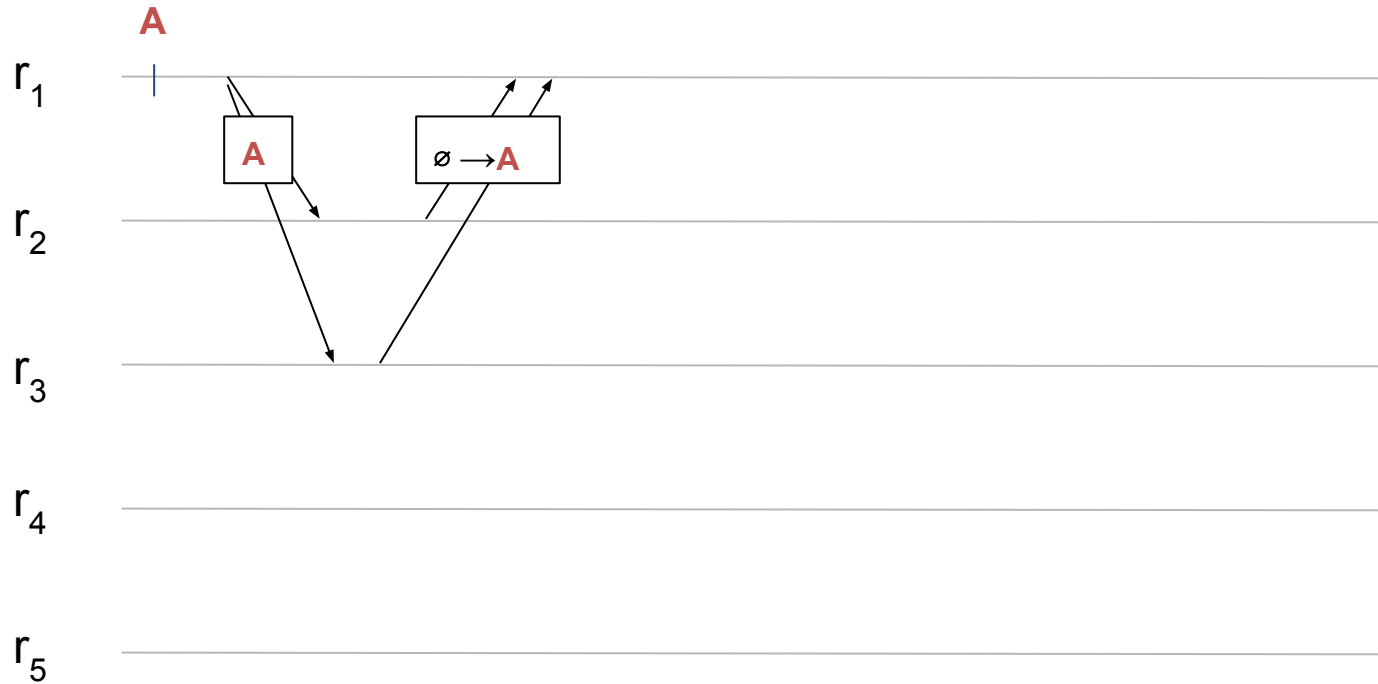
To do consensus on  $\text{dep}(\mathbf{X})$

- try to agree spontaneously by contacting a fast quorum ( $f+f/2$  replicas)
  - when contacted, a replica sent back the commands conflicting with **X** seen so far
- if this fails, ask a slow quorum ( $f+1$  replicas)
  - in this slow path, the union of the reported deps by the fast quorum is used

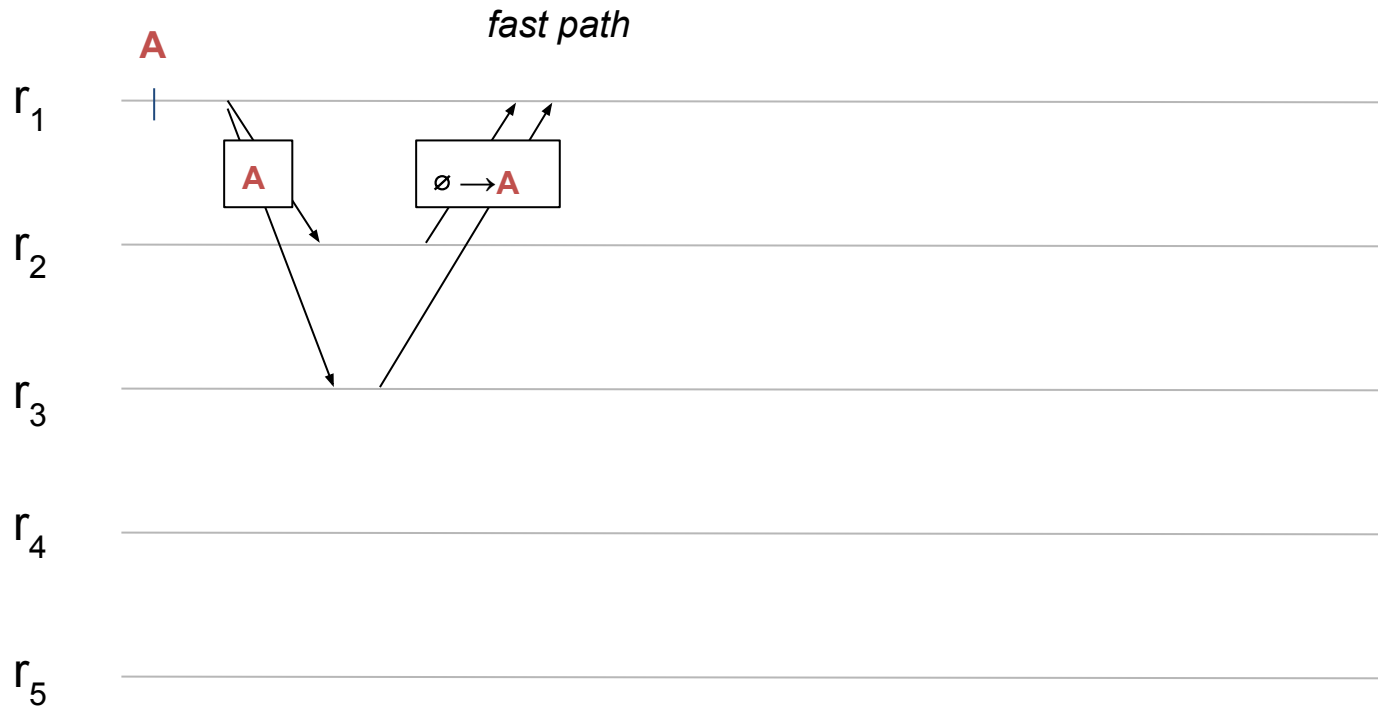
# EPaxos



# EPaxos

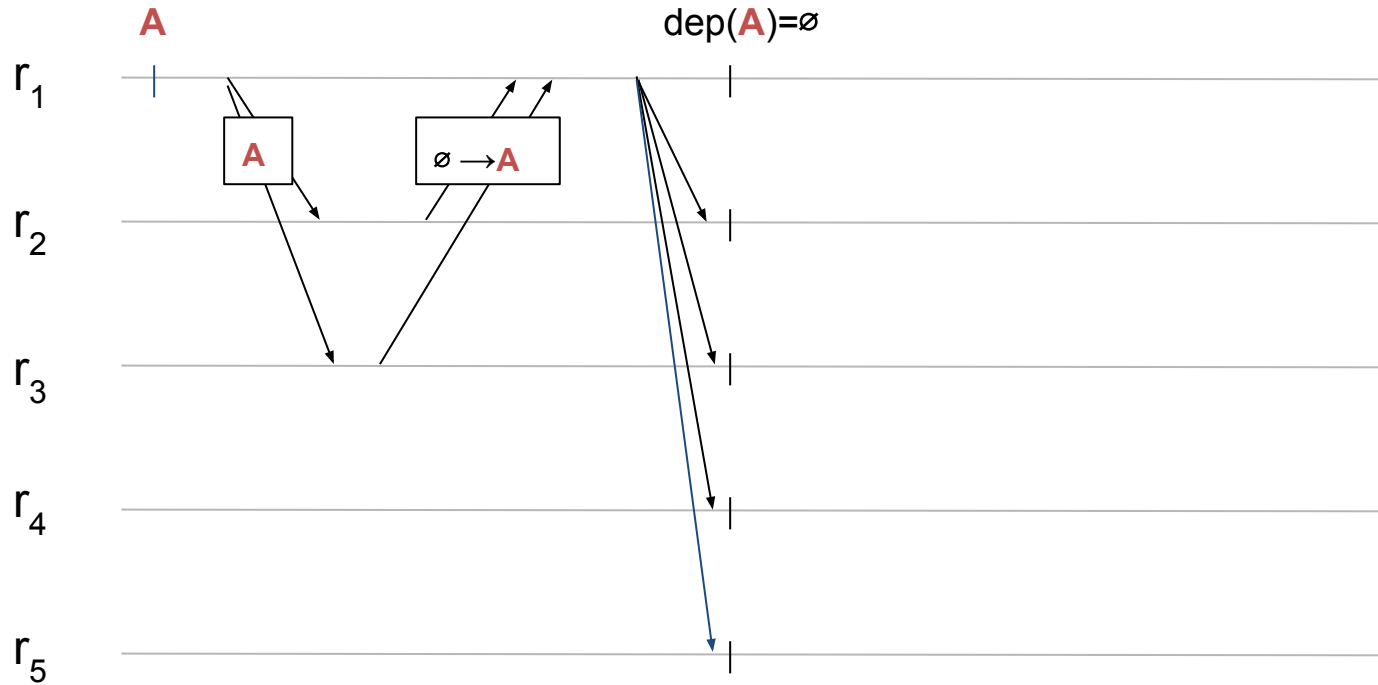


# EPaxos

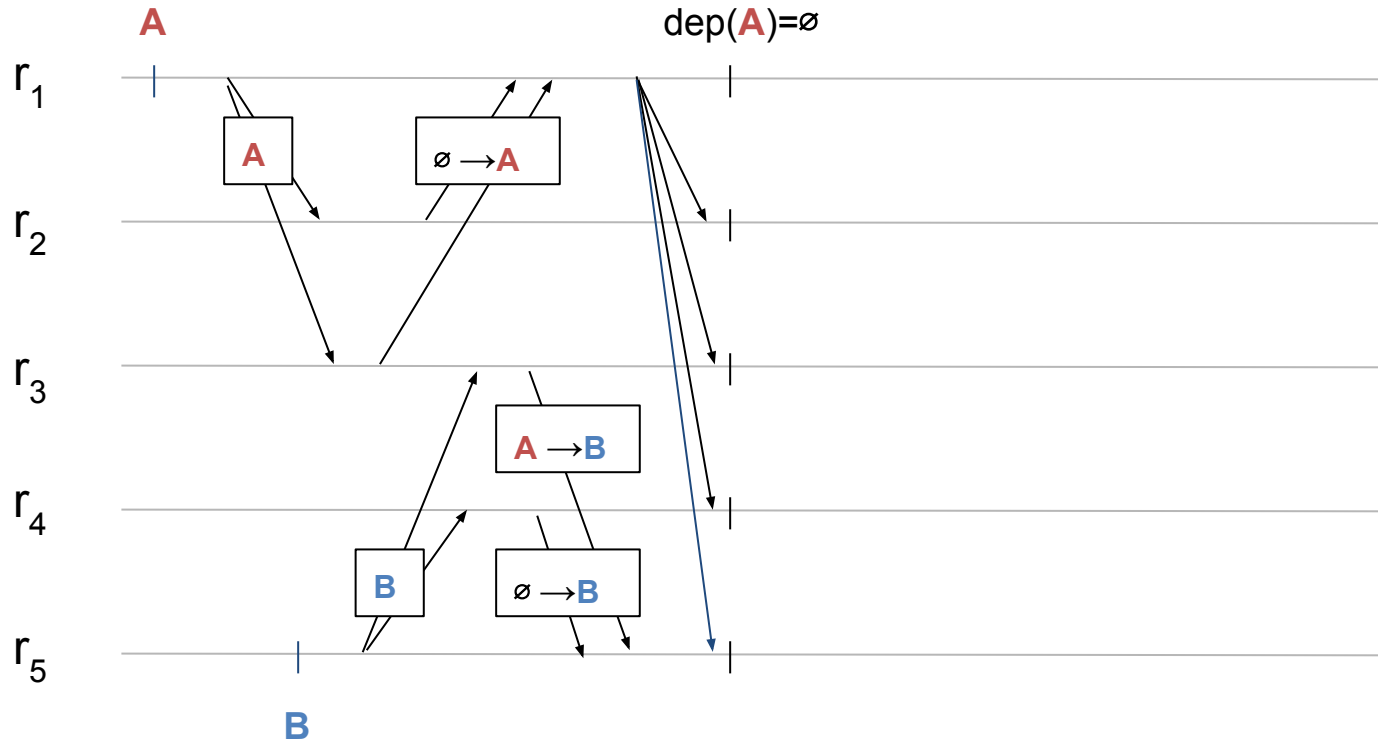




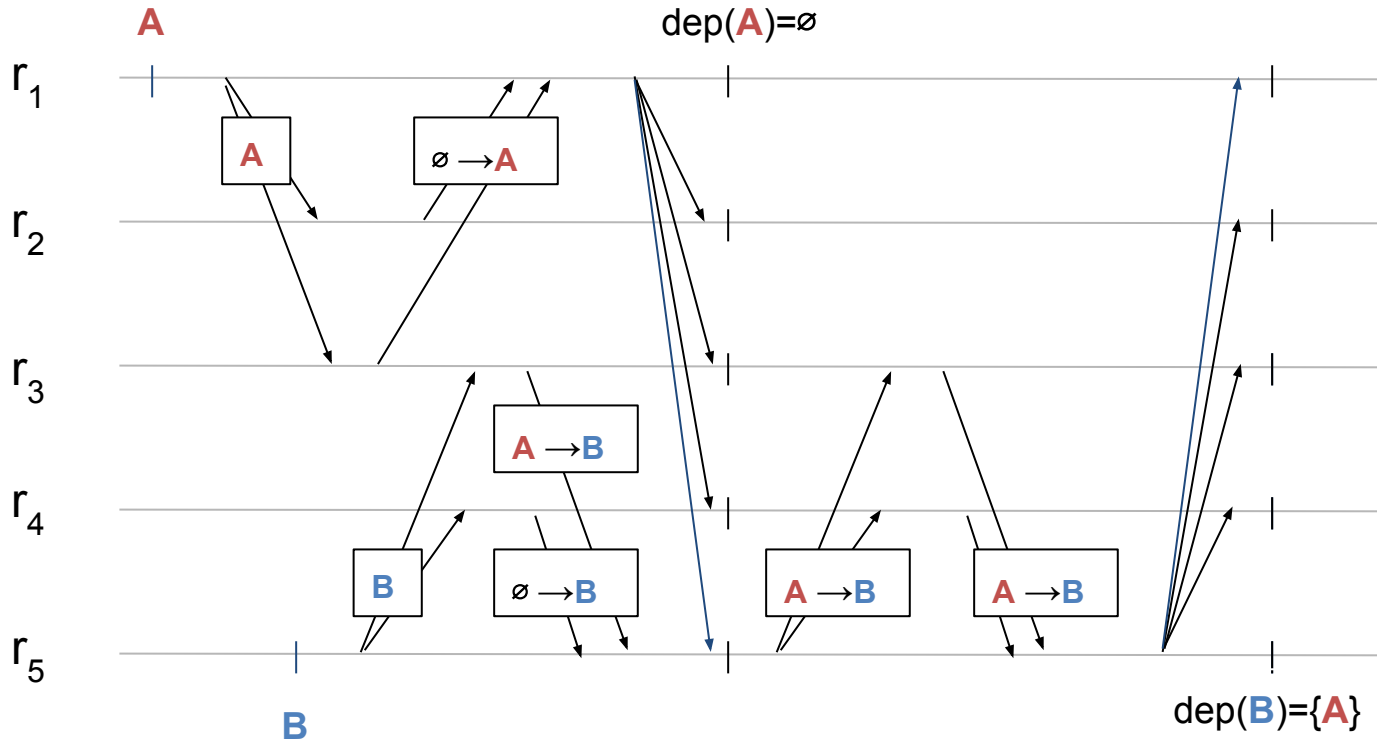
# EPaxos



# EPaxos

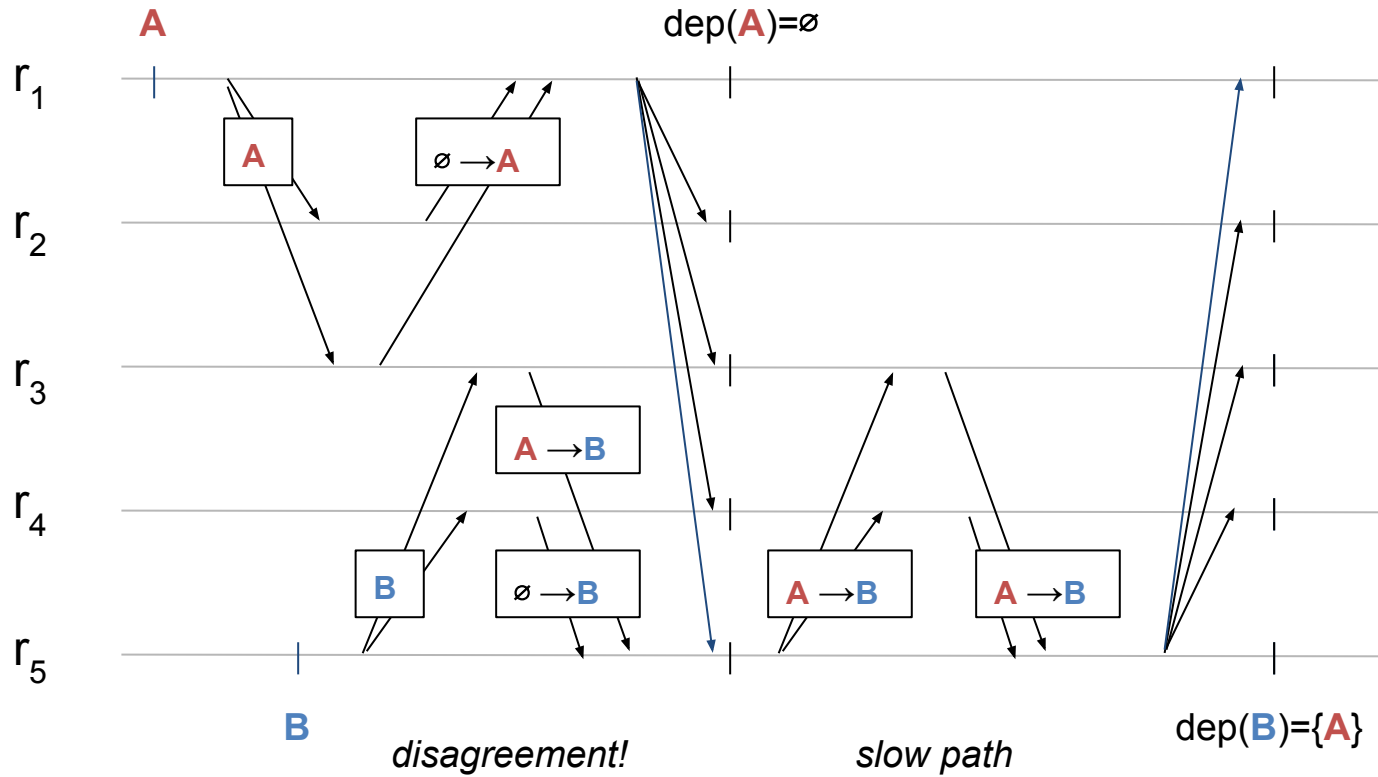


# EPaxos



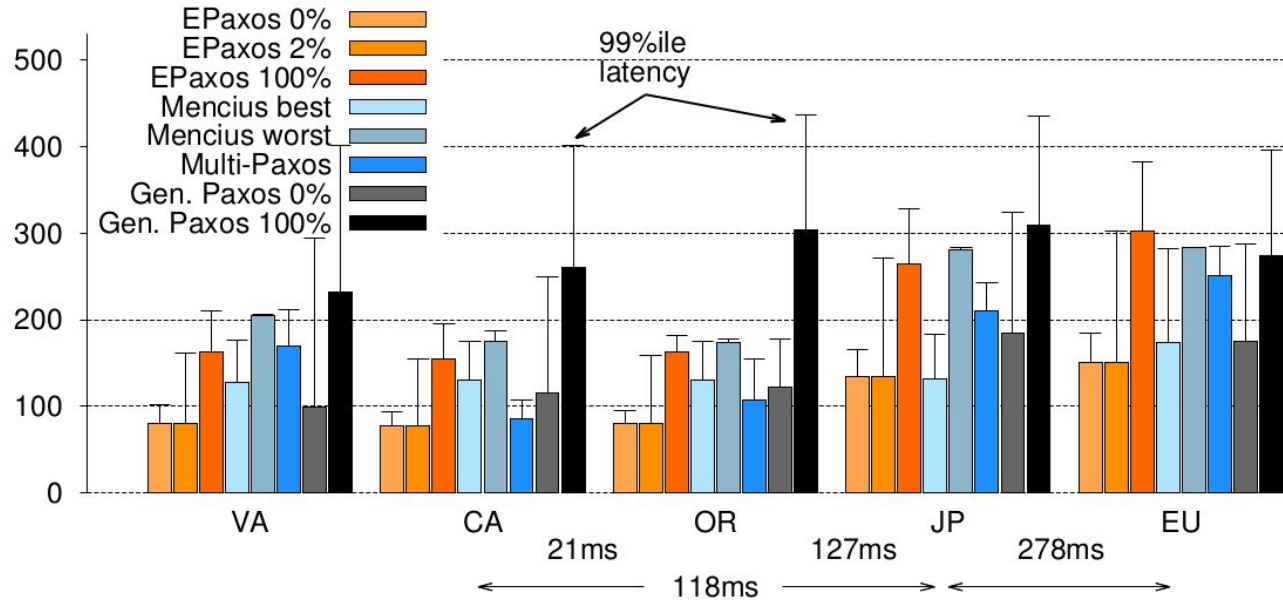
$(n=5, f=2)$

# EPaxos



( $n=5, f=2$ )

# EPaxos - AWS experiments



## Takeaways:

- leaderless SMR is faster and more fair
- but needs most commands commute (EPaxos 100% is bad)

avoid disagreement

*how?* threshold union

consider a bag of items  $E$ , the  $k$ -threshold union of  $E$ , written  $\bigcup_k E$ ,  
are the items reported at least  $k+1$  times in the sets of  $E$

formally,

$$\bigcup_k E = \{ Y : \text{count}(Y) \geq k+1 \}$$

$E = \{E_1, E_2, E_3\}$  with  $E_1 = \{\textcolor{red}{A}, \textcolor{blue}{B}, \textcolor{green}{C}\}$ ,  $E_2 = \{\textcolor{red}{A}, \textcolor{green}{C}\}$  and  $E_3 = \{\textcolor{red}{A}\}$

then

- $\bigcup_1 E = \{\textcolor{red}{A}, \textcolor{green}{C}\},$
- $\bigcup_2 E = \{\textcolor{red}{A}\},$

avoid disagreement

*how?* threshold union

EPaxos fast path condition:

given  $q \in Q$ , let  $\text{dep}_q$  be the dep. reported by  $q$

then

fast-path **iff**  $\forall q, p \in Q. \text{dep}_q = \text{dep}_p$

# Atlas

---

avoid disagreement

*how?* threshold union

Atlas fast path condition:

given  $q \in Q$ , let  $\text{dep}_q$  be the dep. reported by  $q$

then

fast-path **iff**  $\bigcup_f Q = \bigcup_q \text{dep}_q$

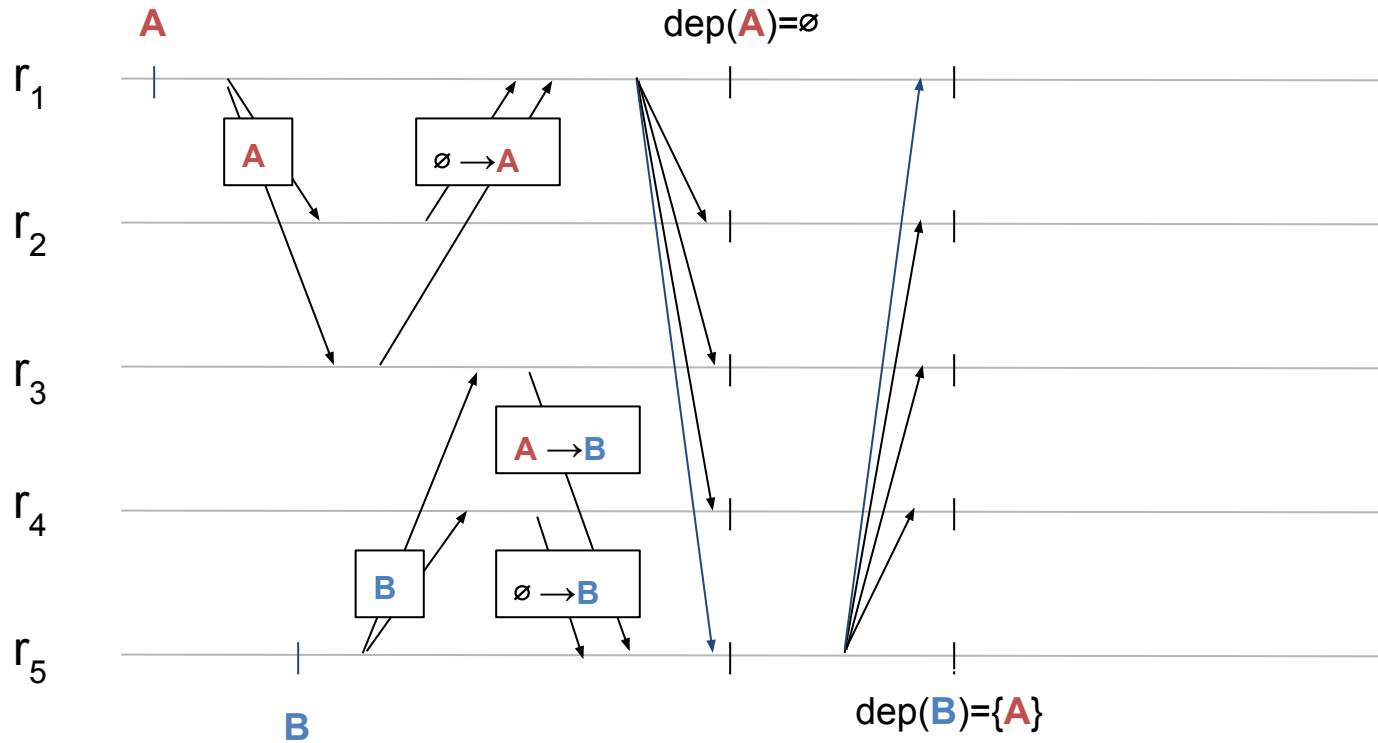
*(=every dep. is reported at least  $f+1$  times)*

why this works?

- if a failure occurs, the dep. reported by any majority quorum in  $Q$  is exactly  $\bigcup_f Q$



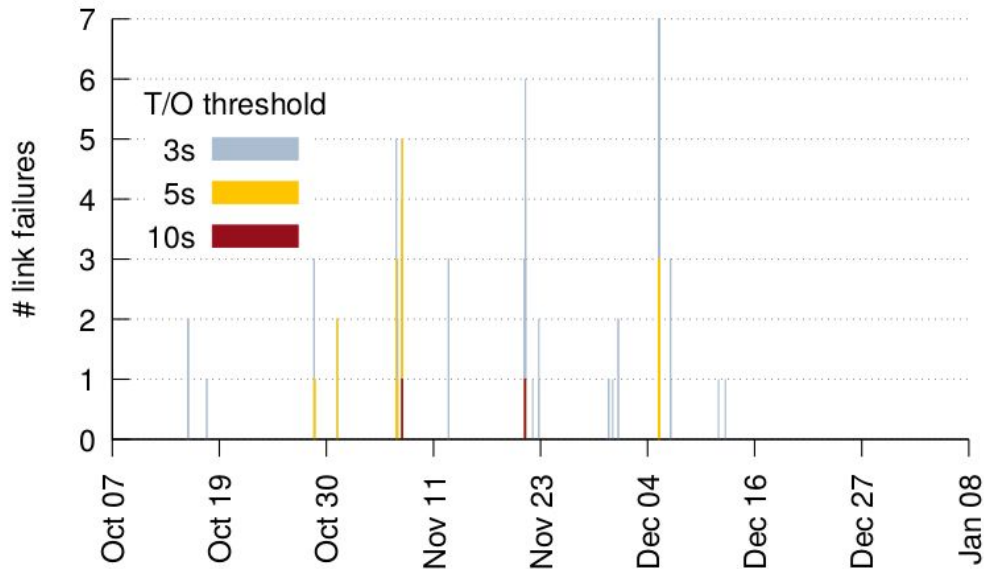
# Atlas



\* the coordinator takes the union of the reported deps.

( $n=5$ ,  $f=1$ )

## Atlas - asynchrony in practice

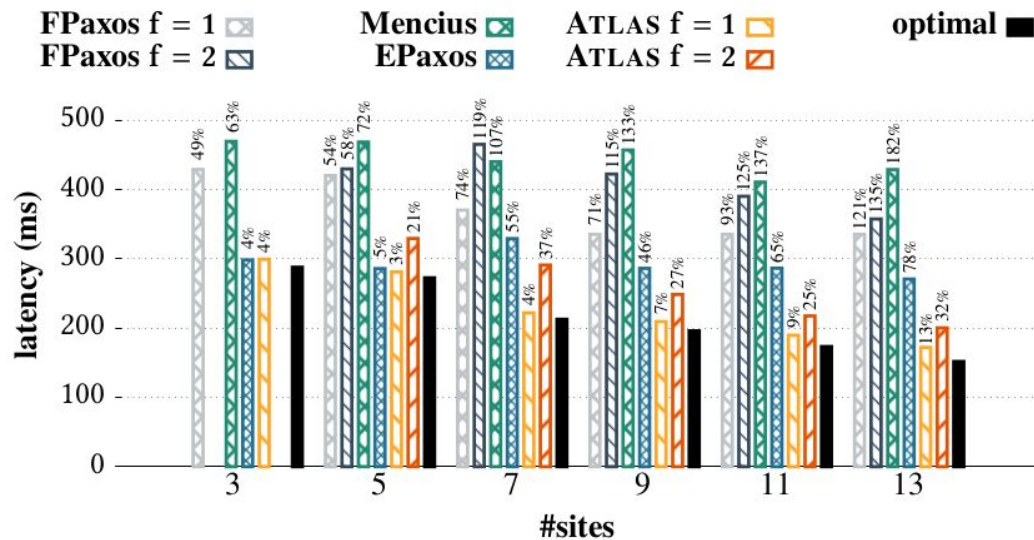


### Takeaways:

- concurrent link failures is a rare event at scale
- at most one slow site during the exp. ( $f=1$ )

*13 GCP sites  
all-to-all ping  
over 3 months*

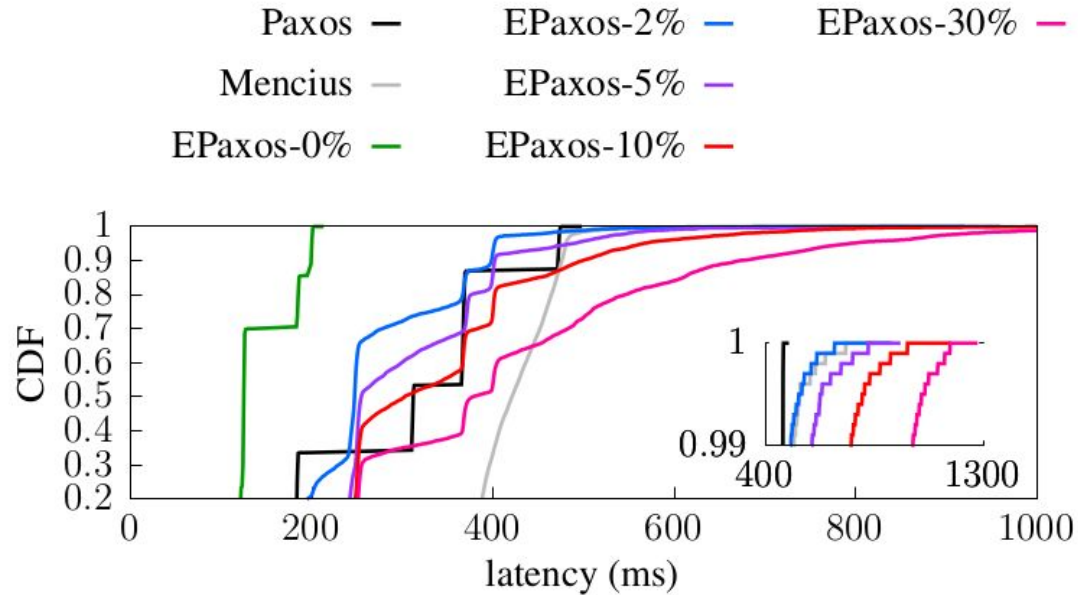
# Atlas - GCP experiments



## Takeaways:

- Atlas better than EPaxos for large-scale deployment ( $n \geq 5$ )

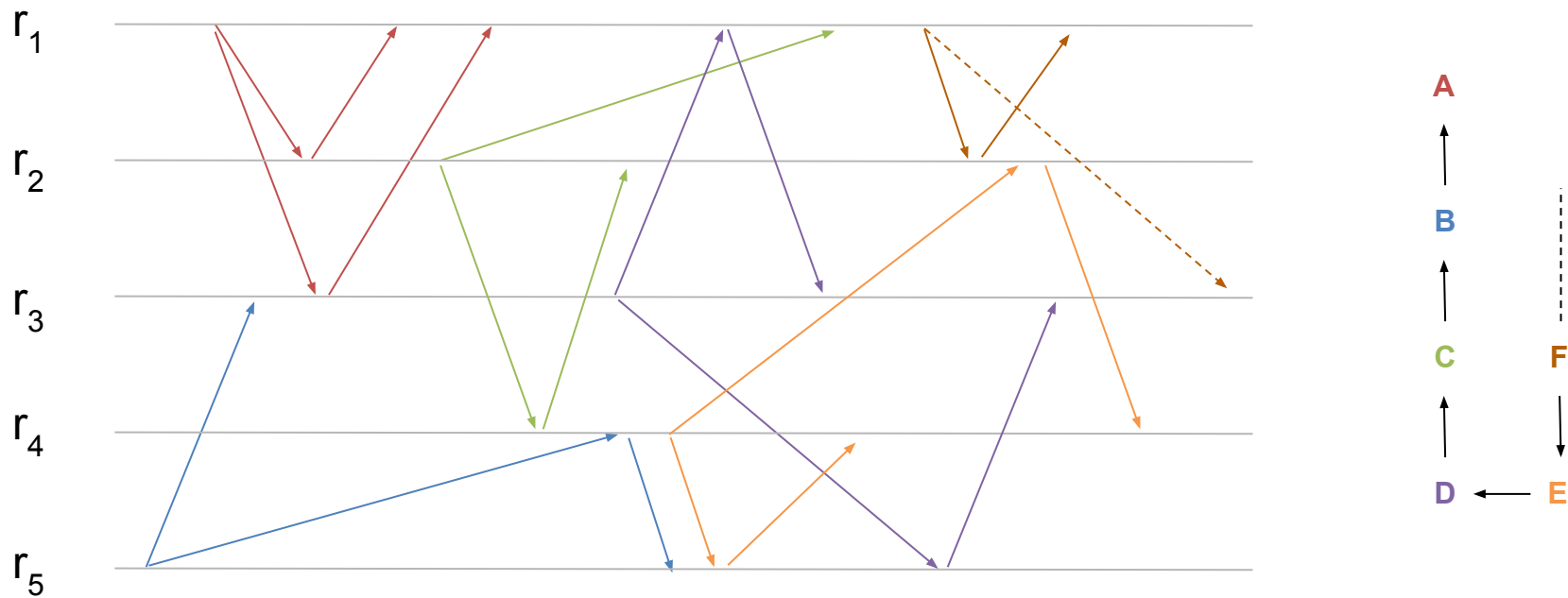
## Tail latency [*DISC'20*, *NSDI'21*]



### Takeaways:

- Tail latency in leaderless SMR protocols is a problem

# Tail latency



tame tail latency

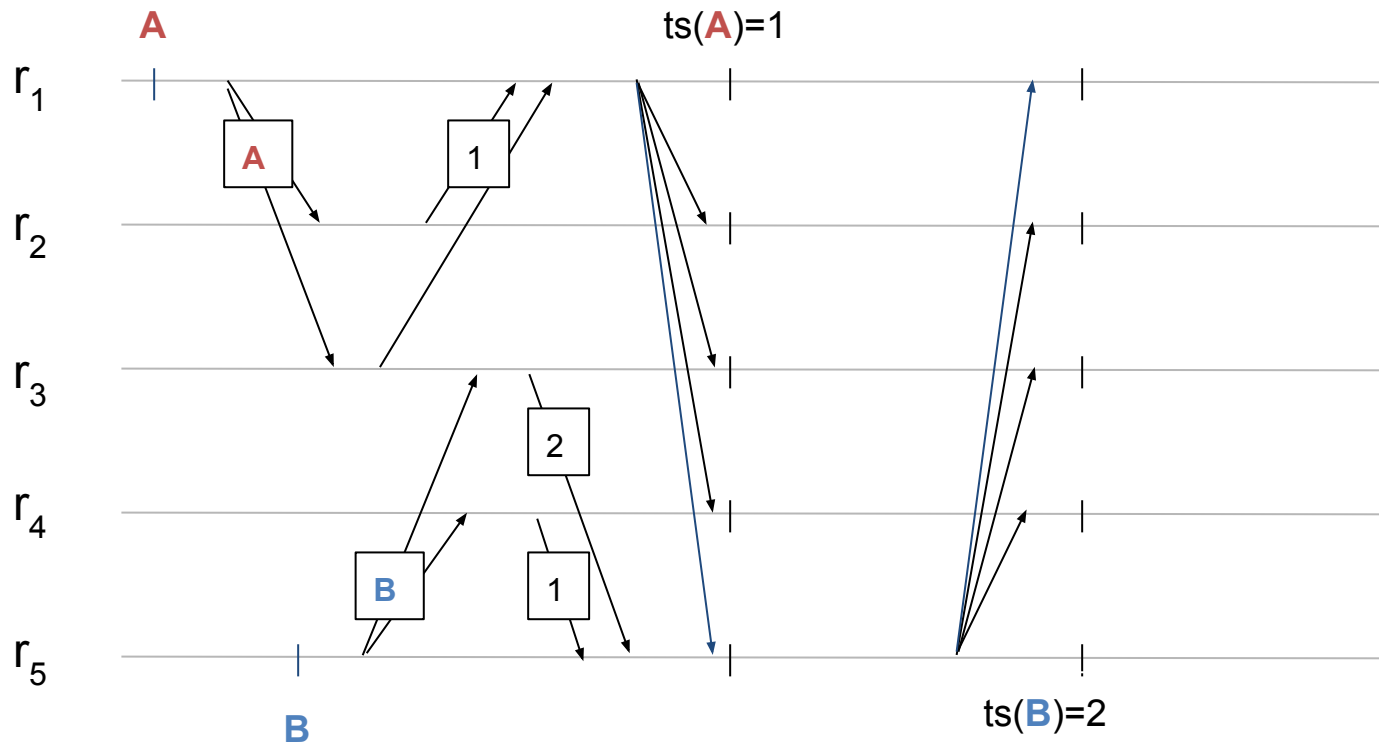
how? agree on a timestamp per command  
+ make the timestamp stable

Tempo fast path condition:

given  $q \in Q$ , let  $ts_q$  be the timestamp reported, or *promised*, by  $q$   
then

fast-path **iff** let  $t = \max\{ts_q : q \in Q\}$   
then  $\text{count}(t) \geq f+1$

# Tempo



$(n=5, f=1)$

## Tempo - background stability mechanism

A command is stable once

- its timestamp, say  $t$ , is committed,
- every command with a timestamp lower (or equal) to  $t$  is stable
- a quorum reports promises higher (or equal) to  $t$

Stable commands are executed in the order of timestamps (ties are broken arbitrarily)

Here,  $\underline{A}; \underline{B}$

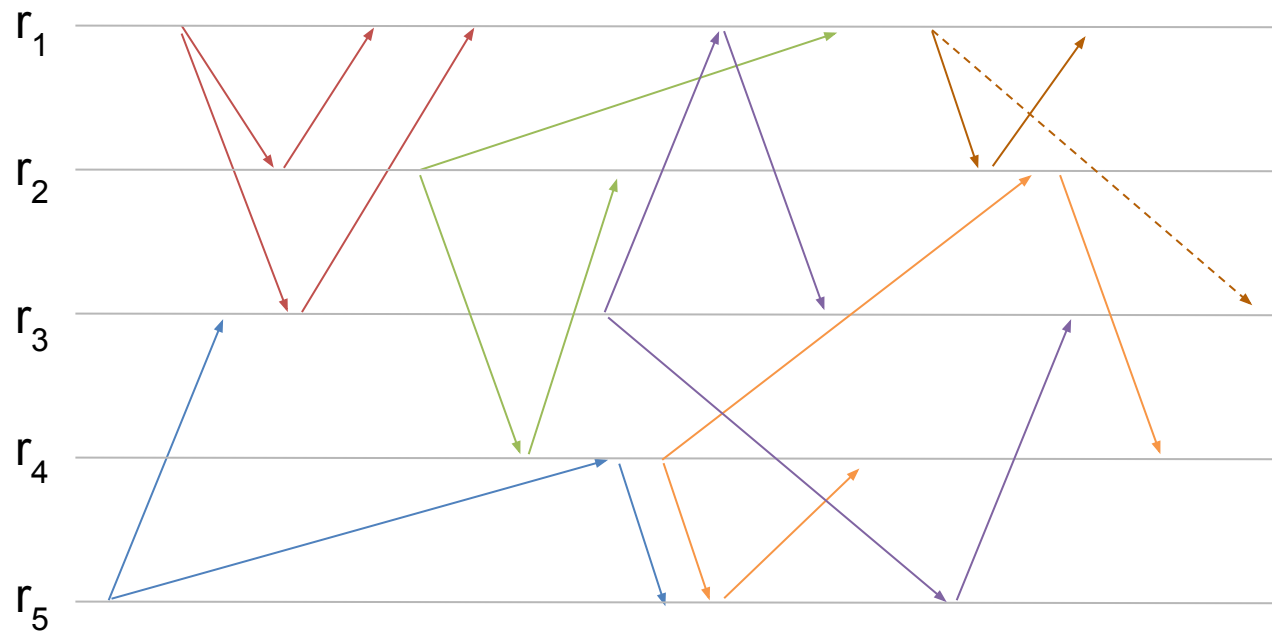
as  $ts(\underline{A}) = ts(\underline{B})$  and  $\underline{A} < \underline{B}$

promises	$\vdots$					
	3					
	2		C	<u>A</u>	<u>B</u>	
	1	<u>A</u>	<u>A</u>	<u>B</u>	C	<u>B</u>
		$r_1$	$r_2$	$r_3$	$r_4$	$r_5$
replicas						

$\underline{X}$  = command  
is stable



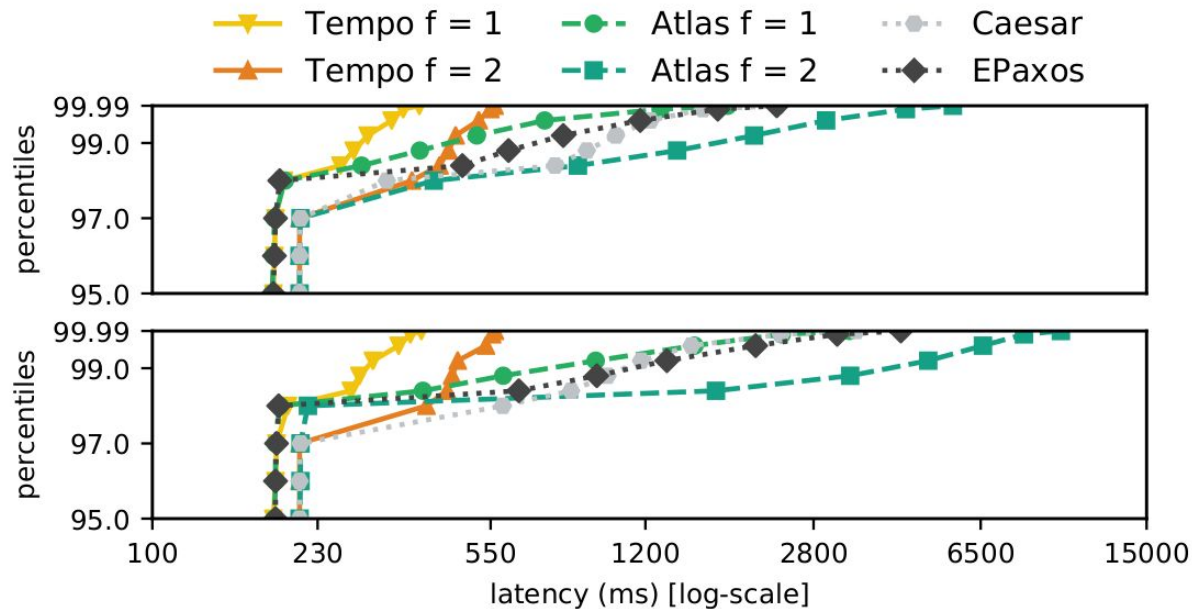
# Tempo - background stability mechanism



⋮					
3	<u>C</u>		D		D
2	D	<u>C</u>	<u>A</u>	<u>B</u>	E
1	<u>A</u>	<u>A</u>	<u>B</u>	<u>C</u>	<u>B</u>
	$r_1$	$r_2$	$r_3$	$r_4$	$r_5$

A;B;C

# Tempo



## Takeaways:

- Tempo improves tail latency in leaderless SMR

5 GCP sites  
512/256 (top/bottom)  
clients per site  
conflict rate is 2%

# Conclusion

---

## Leaderless SMR

- graph-based ordering of commands
- a coordinator per command **X**
  - runs consensus on  $\text{dep}(\mathbf{X})$
- faster and more fair than Paxos/Raft

## Future directions

- scalability
- BFT (blockchain)