# Leaderless State-Machine Replication: An Overview
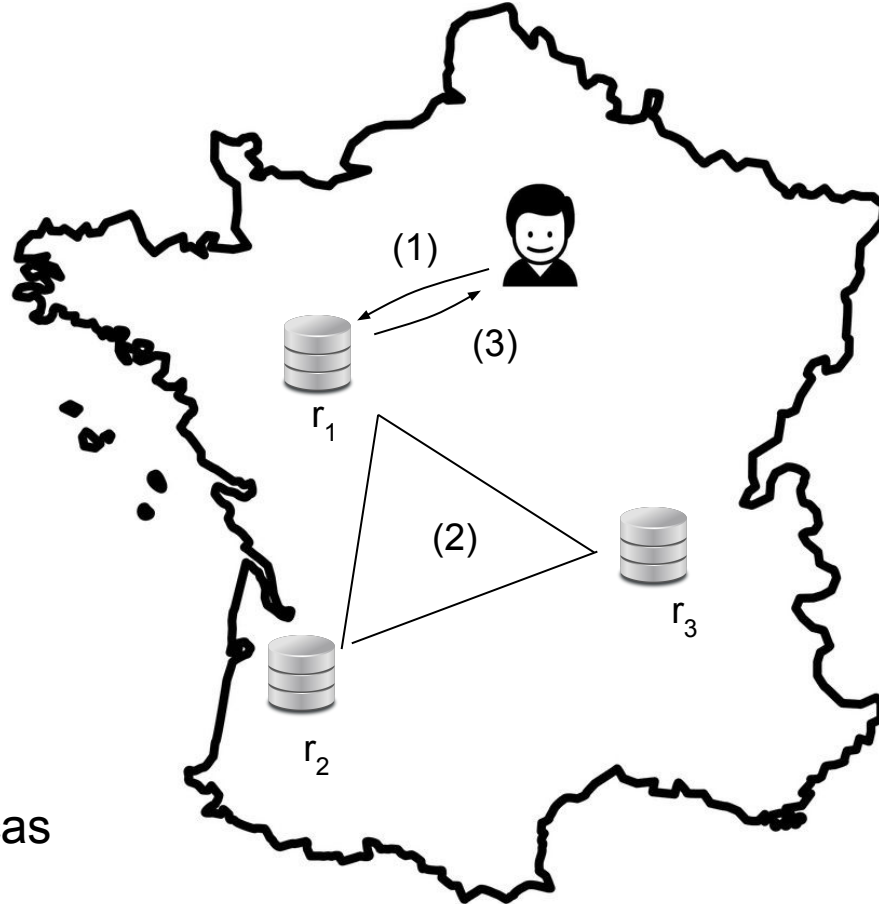
Pierre Sutra

Télécom SudParis
Institut Polytechnique de Paris

*Per3S Workshop, 13.06.2022*
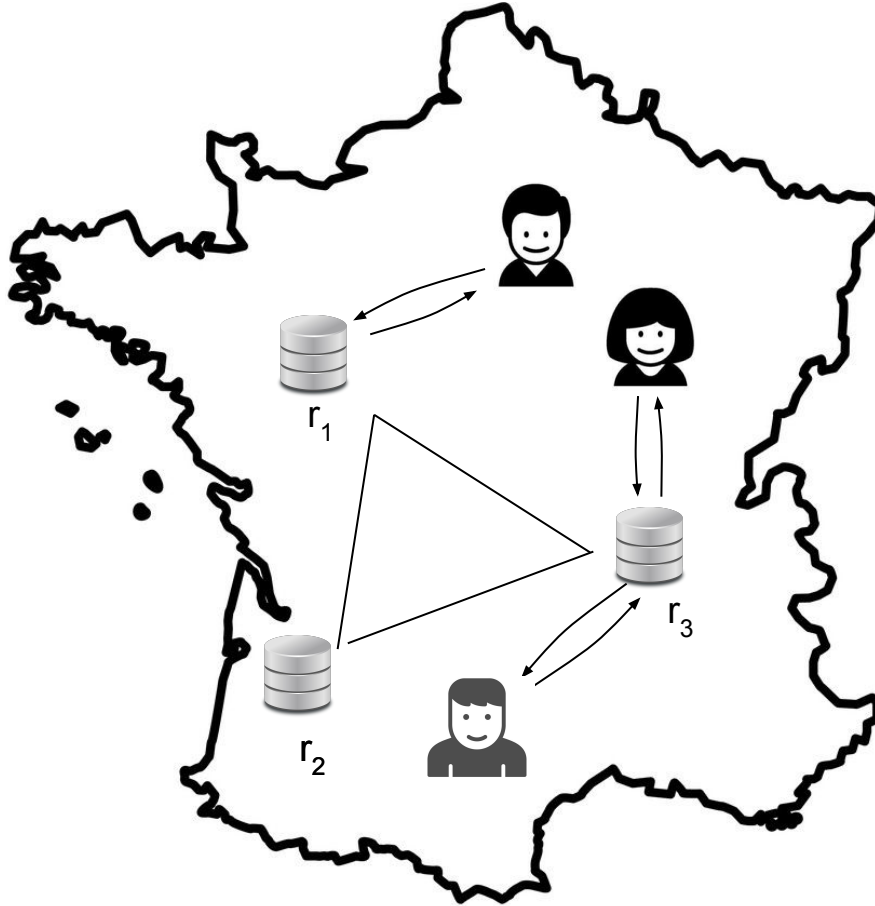
(1) command
(2) some protocol
(3) response

$r_1, r_2, ..$ = data replicas
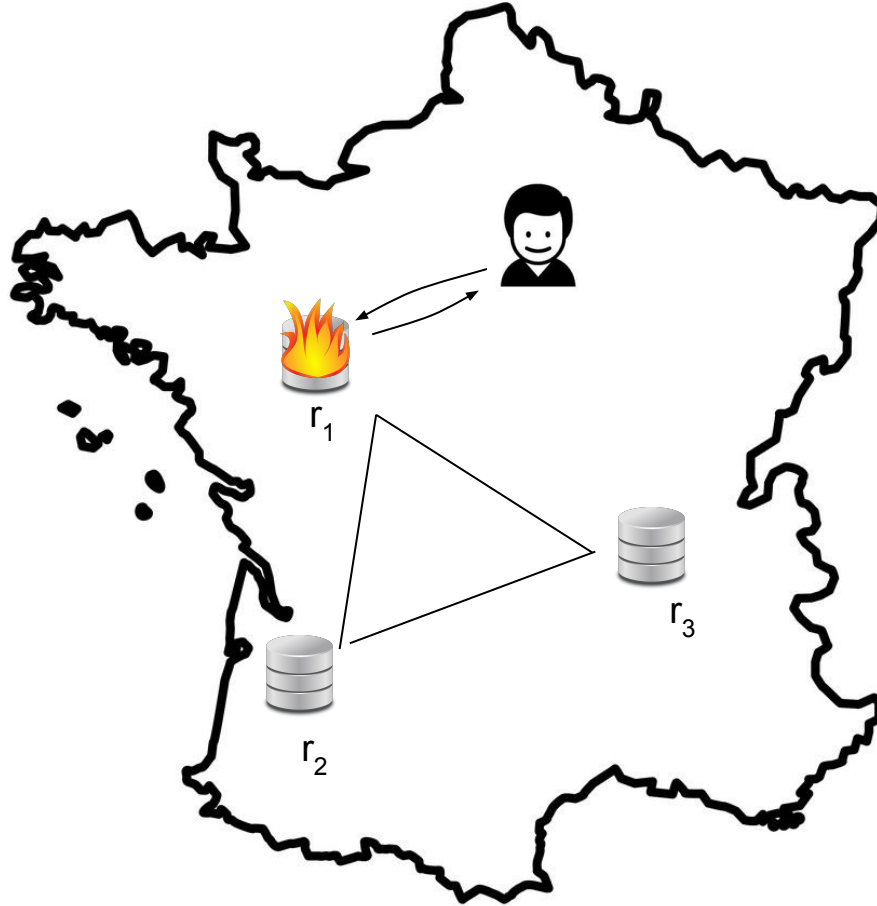
**Problematic**:
*transparent efficient*
geo-replication

3

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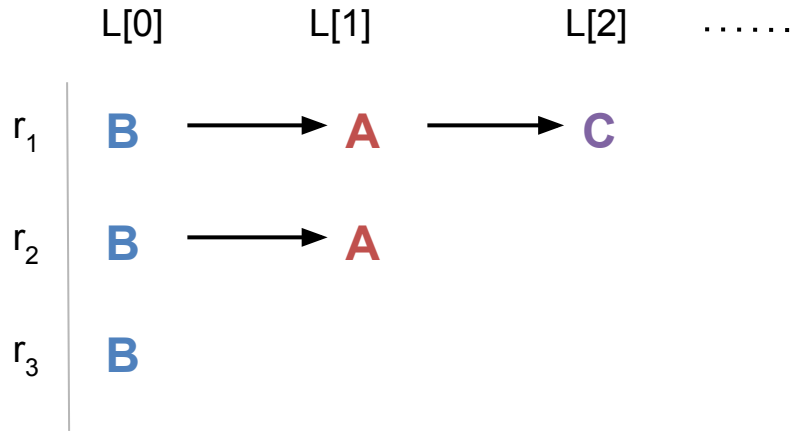geo-replication

$r_1$

$r_2$

$r_3$

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

5

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

Each replica holds a log L
For each i, *agree* on command L[i]
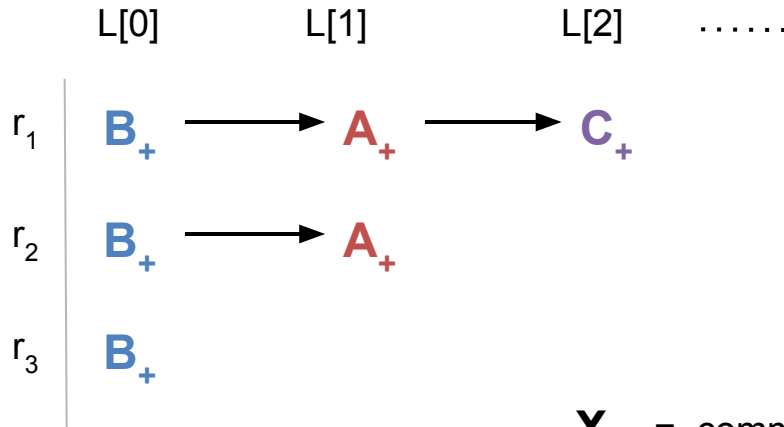
L[0]   L[1]   L[2]  ······

$r_1$   **B** ⟶ **A** ⟶ **C**

$r_2$   **B** ⟶ **A**

$r_3$   **B**

# Classic SMR

Each replica holds a log L
For each i, *agree* on command L[i]
Execute commands in log order



$X_+$ = command is executed

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

$r_1$  |  **B**+      **A**+      **C**+

$r_2$  |           **A**+

$r_3$  |  **B**+

**A** = x ← 42

**B** = y ← 7

**C** = z ← x + y

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



dep(**A**) = {**C**}

dep(**C**) = {**B**, **A**}

dep(**B**) = ∅

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



- operation **X** executed once dep(**X**) transitively closed
- cycles are broken deterministically

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



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

# Leaderless SMR

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



## Properties

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

12

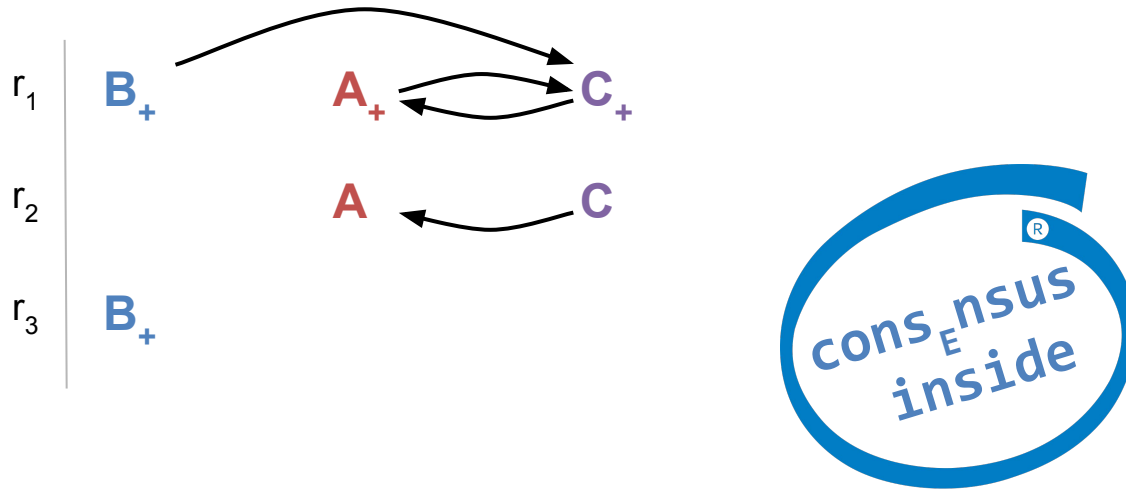# Egalitarian Paxos [*SOSP'13*]

EPaxos uses *2f+1* replicas.

When a client executes command **X**
- pick a replica
- this replica is the _coordinator_ for **X**, coord(**X**)
- coord(**X**) runs consensus over dep(**X**)

To do consensus on dep(**X**)
- try to agree spontaneously by contacting a <u>fast</u> quorum *(f+f/2 replicas)*
- if this fails, ask a <u>slow</u> quorum *(f+1 replicas)*

# EPaxos



*(n=5, f=2)*

14

A

$r_1$

A

$\varnothing \rightarrow$ A

$r_2$

$r_3$

$r_4$

$r_5$

*(n=5, f=2)*

15

# EPaxos



fast path

$r_1$

A

A

∅ → A

$r_2$

$r_3$

$r_4$

$r_5$

*(n=5, f=2)*

16

# EPaxos



A

dep(A)=∅

r₁

r₂

r₃

r₄

r₅

A

∅ → A

*(n=5, f=2)*

# EPaxos

A

dep(A)=∅

r₁

A

∅ →A

r₂

r₃

A →B

r₄

B

∅ →B

r₅

B

*(n=5, f=2)*

# EPaxos



$dep(A) = \varnothing$

$A$

$B$

$dep(B) = \{A\}$

*(n=5, f=2)*

19

**A**

dep(**A**)=∅

$r_1$

A

∅ →A

$r_2$

$r_3$

A →B

$r_4$

B

∅ →B

A →B

A →B

$r_5$

**B**

*disagreement!*

*slow path*

dep(**B**)={**A**}

*(n=5, f=2)*

# EPaxos - AWS experiments



VA, CA, .. = datacenters
x% = ratio of commuting commands
(≤ 2% is typical)

**Takeaways:**
- leaderless SMR is <u>faster</u>

**Takeaways:**
- leaderless SMR is faster and <u>more fair</u>

# EPaxos - AWS experiments



**Takeaways:**

- leaderless SMR is faster and more fair
- but commands *should commute*

A

$r_1$

A

∅ → A

$r_2$

$r_3$

A → B

$r_4$

B

∅ → B

A → B

A → B

$r_5$

B

*disagreement!*

*slow path*

dep(B)={A}

*(n=5, f=2)*

25

*goal*: avoid disagreement

Consider a bag of items E, the *k-threshold union* of E, written $\cup_k$ E, are the items reported at least k+1 times in the sets of E
formally,

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

Example:

let E = $\{E_1, E_2, E_3\}$ with $E_1 = \{A, B, C\}$, $E_2 = \{A, C\}$ and $E_3 = \{A\}$
then
- $\cup_1 E = \{A, C\}$,
- $\cup_2 E = \{A\}$,

26

# Atlas

_goal_: to avoid disagreement

    _EPaxos_ fast path condition:

        let Q be a fast path quorum ($f+f/2$ replicas)

        given $q \in Q$, let $dep_q$ be the dep. reported by q

        then

            fast-path **iff** $\forall\ q,p \in Q.\ dep_q = dep_p$

*goal:* to avoid disagreement

*Atlas* fast path condition:

given $q \in Q$, let $dep_q$ be the dep. reported by q
then

fast-path **iff** $\biguplus_f Q = \cup_q dep_q$
*(i.e., every dep. is reported at least f+1 times)*

*why this works?*
- if a failure occurs, the dep. reported by
<u>any</u> majority quorum in Q is exactly $\biguplus_f Q$

28

# Atlas



dep(A)=∅

dep(B)={A}

*the coordinator counts as the union of the reported deps.*

*(n=5, f=1)*

# Atlas - asynchrony in practice



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

**Takeaways:**

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

# Atlas - GCP experiments



*x% = how far from optimal*

**Takeaways:**
- Atlas better than EPaxos for large-scale deployment (n ≥ 5)

**Takeaways:**

- Tail latency in leaderless SMR protocols is a problem

*(n=5, f=1)*

*goal:* tame tail latency

   *key idea: 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 count(t) ≥ f+1

# Tempo



$r_1$

$r_2$

$r_3$

$r_4$

$r_5$

A

ts(A)=1

A

1

2

B

1

B

ts(B)=2

*(n=5, f=1)*    35

# Tempo - background stability mechanism

A command is <u>stable</u> once
- its timestamp, say t, is agreed;
- every command with a timestamp lower (or equal) to t is stable; and
- a quorum reports promises higher (or equal) to t.

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

Here, **A**;**B**
as ts(**A**) = ts(**B**) and **A** < **B**

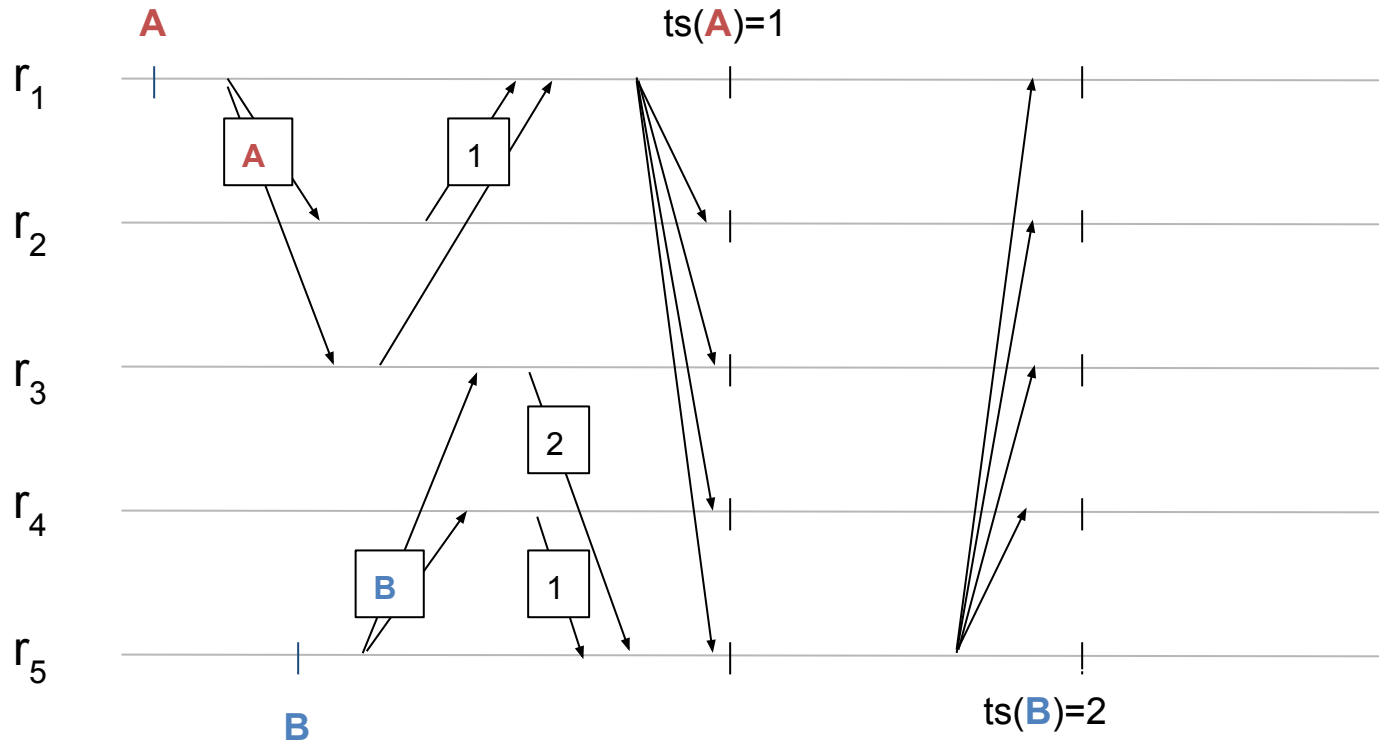| promises | | | | | |
|---|---|---|---|---|---|
| ⋮ | | | | | |
| 3 | | | | | |
| 2 | | **C** | <u>**A**</u> | <u>**B**</u> | |
| 1 | <u>**A**</u> | <u>**A**</u> | <u>**B**</u> | **C** | <u>**B**</u> |
| | $r_1$ | $r_2$ | $r_3$ | $r_4$ | $r_5$ |

replicas

<u>**X**</u>  = command
      is stable

36

# Tempo - background stability mechanism



|   | $r_1$ | $r_2$ | $r_3$ | $r_4$ | $r_5$ |
|---|-------|-------|-------|-------|-------|
| ⋮ |       |       |       |       |       |
| 3 | **C** |       | **D** |       | **D** |
| 2 | **D** | **C** | **A** | **B** | **E** |
| 1 | **A** | **A** | **B** | **C** | **B** |

**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 dep(**X**)
- better than Paxos/Raft

Future directions

- higher scalability
- byzantine failures

# References

**[Paxos]** The Part-Time Parliament, *L. Lamport, ACM TOCS, 1998.*

**[BGcast]** Generic Broadcast, *F. Pedone, A. Schiper, DISC, 1999.*

**[GPaxos]** Generalized Consensus and Paxos, *L. Lamport, Tech Report, 2005.*

**[DISC'05]** Generic Optimistic Broadcast, *F. Pedone, P. Zieliński, DISC, 2005.*

**[SOSP'13]** There is more consensus in Egalitarian parliaments, *I. Moraru, D. G. Andersen, M. Kaminsky, SOSP, 2013.*

**[Raft]** In search of an understandable consensus algorithm, *D. Ongaro, J. Ousterhout, USENIX ATC, 2014.*

**[DISC'20]** Leaderless State-Machine Replication: Specification, Properties, Limits, *T. F. Rezende, P. Sutra, DISC, 2020.*

**[Eurosys'20]** State-Machine Replication for Planet-Scale Systems, *V. Enes, M. Perrin, C. Baquero, A. Gotsman, P. Sutra, Eurosys, 2020.*

**[NSDI'21]** EPaxos Revisited, *S. Tollman, J. Ousterhout, NSDI, 2021.*

**[Eurosys'21]** Efficient Replication via Timestamp Stability, *V. Enes, C. Baquero, A. Gotsman, P. Sutra, Eurosys, 2021.*