# Leaderless State Machine Replication: Specification, Properties, Limits

Tuanir França Rezende [1]     Pierre Sutra [1]

[1] Telecom SudParis

14 Oct. 2020, DISC'20

# Context - Principles

- Modern internet services often maintain more than one **copy** of an object.
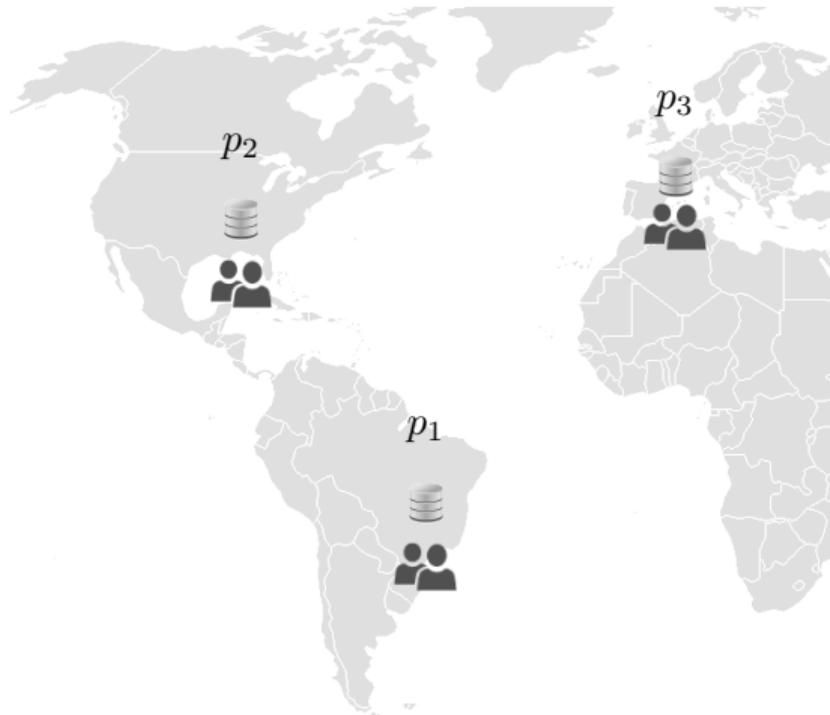
# Context - Principles



- Each copy is located on a separate server, or **replica**.

# Context - Principles



- Clients interact with these replicas by sending **commands** and **share logically** the **objects**.
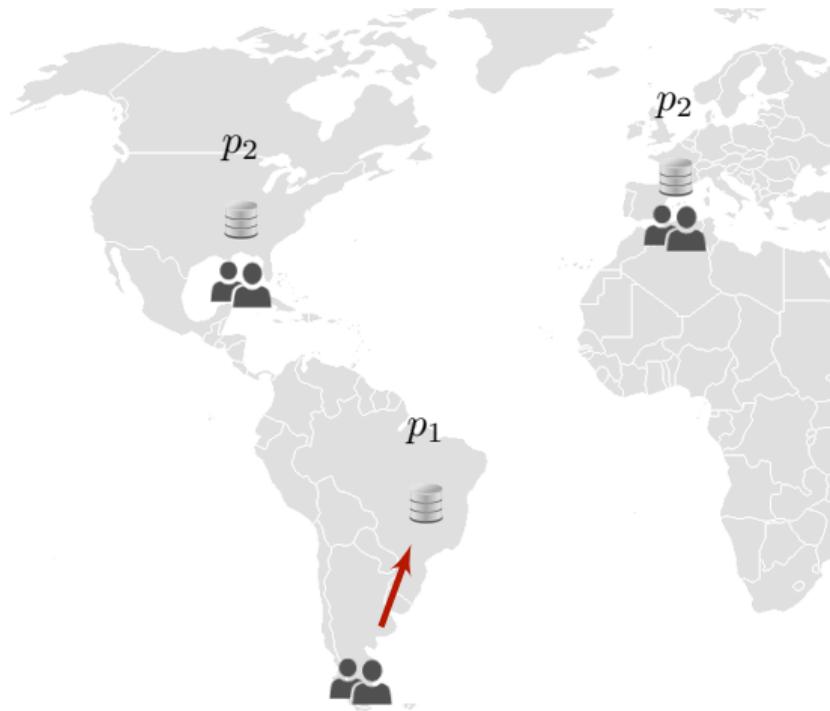
# Context - Why to replicate data?

- **Performance:** Improve latency and/or throughput.

# Context - Why to replicate data?

- **Performance:** Improve latency and/or throughput.
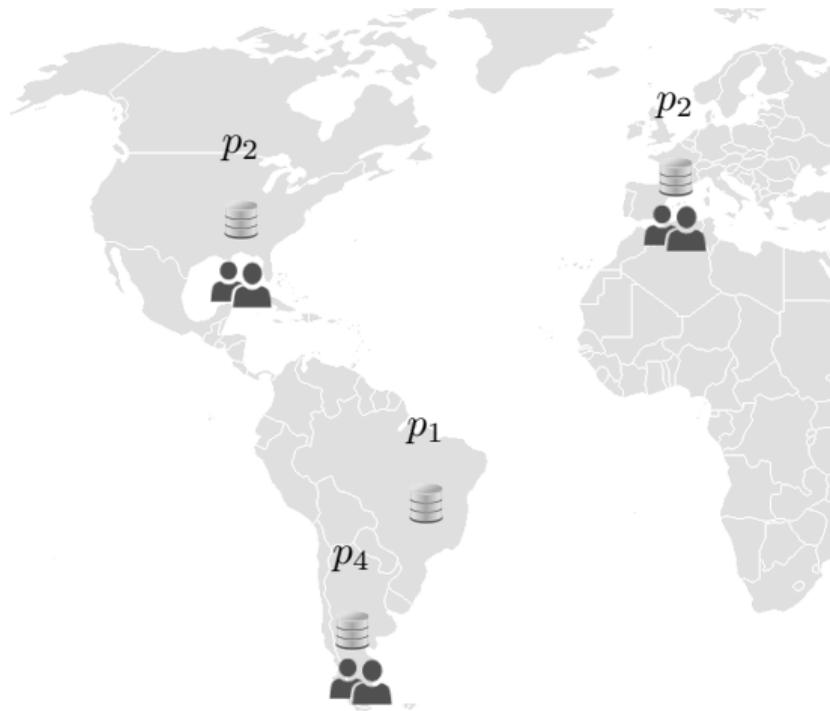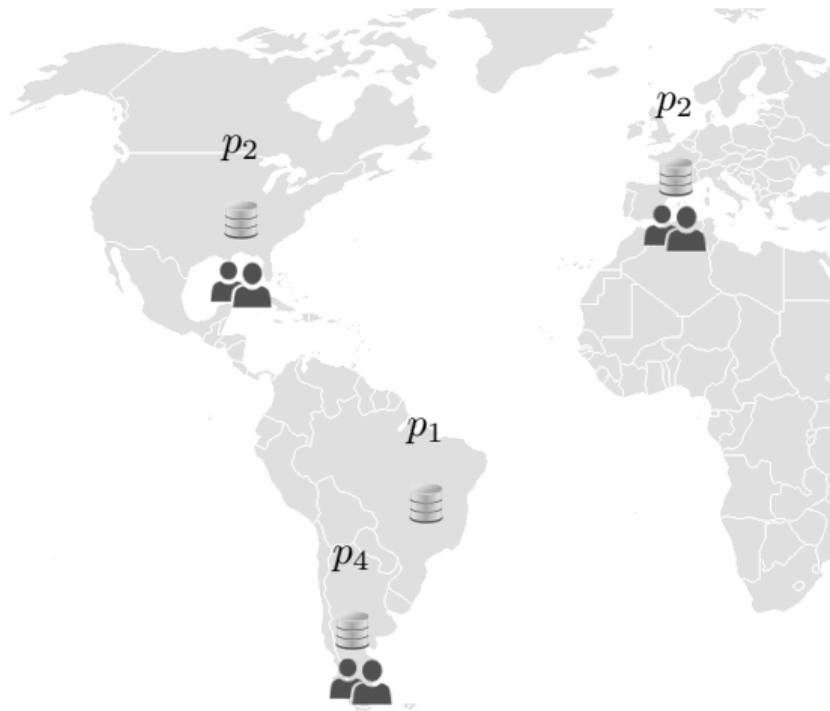
# Context - Why to replicate data?



- **Performance:** Improve latency and/or throughput.

# Context - Why to replicate data?

- **Performance:** Improve latency and/or throughput.



$p_2$

$p_2$

$p_1$

$p_4$

# Context - Why to replicate data?

- **Reliability:** Mask replica and network failures.
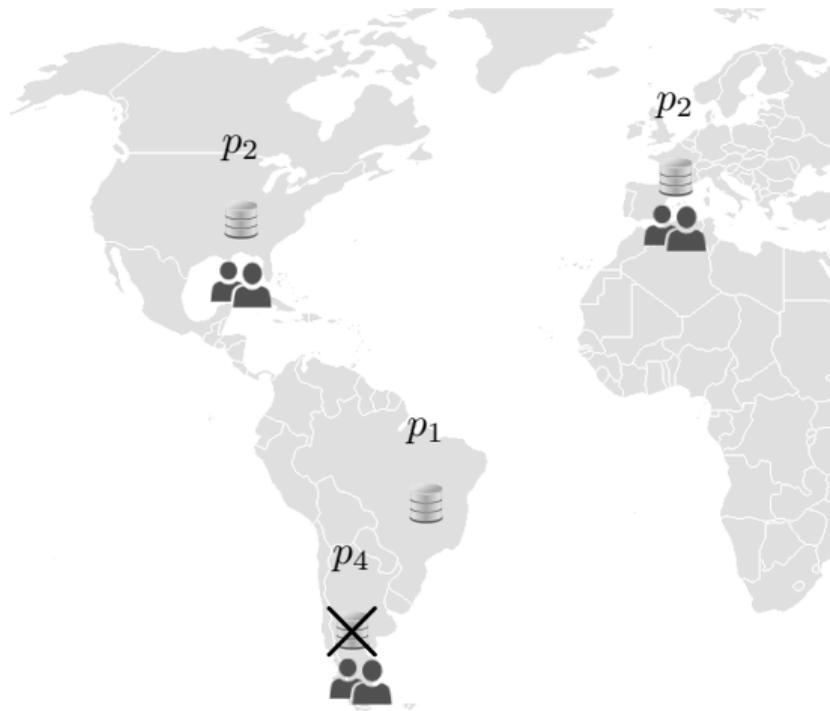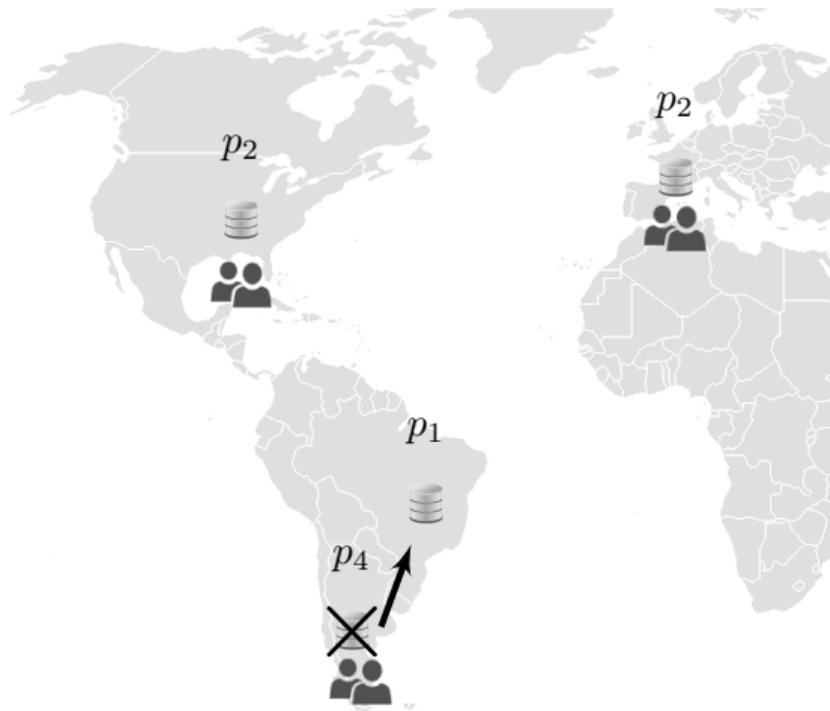
$p_2$

$p_2$

$p_1$

$p_4$

# Context - Why to replicate data?

- **Reliability:** Mask replica and network failures.

# Context - Why to replicate data?

- **Reliability:** Mask replica and network failures.
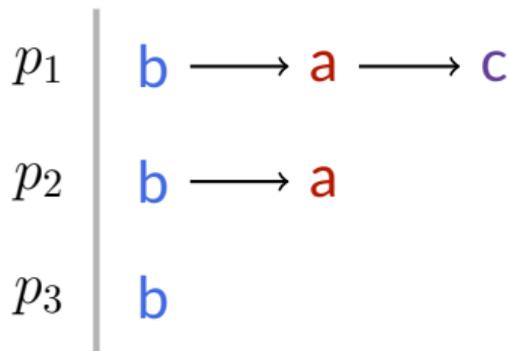
# Context - How to replicate data?
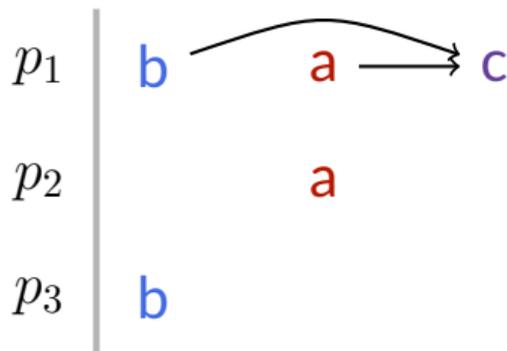
## Through **State Machine Replication (SMR)**:

- Classic SMR (Paxos). Execute commands in the same order.

$p_1$ | b $\longrightarrow$ a $\longrightarrow$ c

$p_2$ | b $\longrightarrow$ a

$p_3$ | b

# Context - How to replicate data?

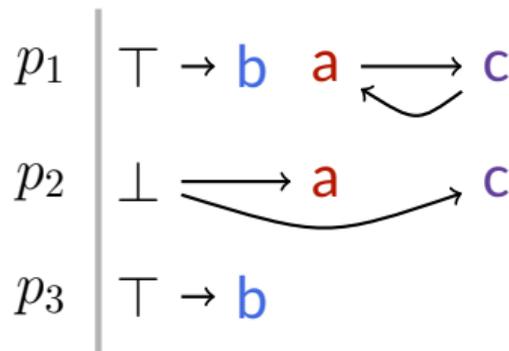## Through **State Machine Replication (SMR)**:

- Classic SMR (Paxos). Execute commands in the same order.

- Generic SMR. Execute *non-commuting* commands in the same order.

# Context - How to replicate data?

## Through **State Machine Replication (SMR)**:

- Classic SMR (Paxos). Execute commands in the same order.

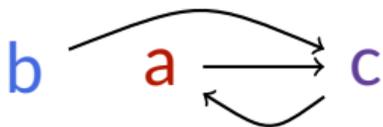- Generic SMR. Execute *non-commuting* commands in the same order.
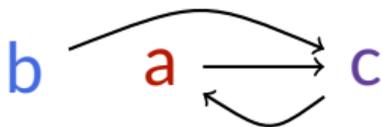
- Leaderless SMR?



$p_1 \quad \top \to \mathsf{b} \quad \mathsf{a} \rightleftarrows \mathsf{c}$

$p_2 \quad \bot \longrightarrow \mathsf{a} \longrightarrow \mathsf{c}$

$p_3 \quad \top \to \mathsf{b}$
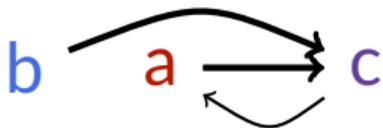
**Leaderless SMR - Dependency Graph**



Relies on the notion of **dependency graph** instead of a partially ordered log as found in Generic SMR.

# Leaderless SMR - Dependency Graph



A **dependency graph** is a directed graph that records the constraints defining how commands are executed.

# Leaderless SMR - Dependency Graph



For some command c, the **incoming neighbors** of c in the dependency graph are its **dependencies**.
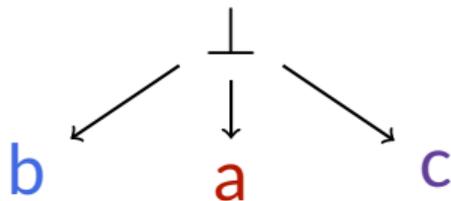
# Leaderless SMR - Understanding the abstraction

In Leaderless SMR, a process holds two mapping: $deps$ and $phase$.

## Leaderless SMR - Understanding the abstraction

The mapping $deps$ is a dependency graph storing a relation from $\mathcal{C}$ to $2^{\mathcal{C}} \cup \{\bot, \top\}$. [1]
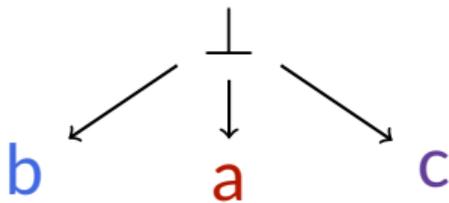
---

[1] $\mathcal{C}$ is the set of commands.

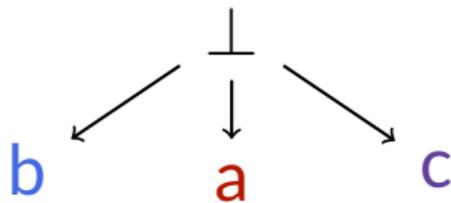# Leaderless SMR - Understanding the abstraction



Initially, for every command c, $deps(\mathsf{c})$ is set to $\perp$. This corresponds to the **pending** phase.

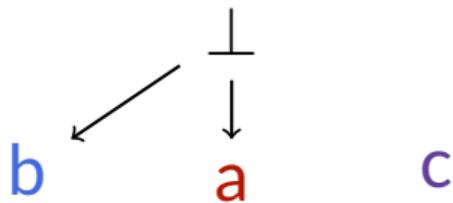# Leaderless SMR - Understanding the abstraction



When a process decides a command c, it changes the mapping $deps(\mathsf{c})$ to a non-$\perp$ value. Operation $commit(\mathsf{c}, D)$ assigns $D$ taken in $2^{\mathcal{C}}$ to $deps(\mathsf{c})$.

# Leaderless SMR - Understanding the abstraction



$$commit(\textcolor{purple}{c}, \{\textcolor{red}{a}, \textcolor{blue}{b}\})$$

# Leaderless SMR - Understanding the abstraction



$$commit(\textsf{c}, \{\textsf{a}, \textsf{b}\})$$

# Leaderless SMR - Understanding the abstraction



$$commit(\mathsf{c}, \{\mathsf{a}, \mathsf{b}\})$$

# Leaderless SMR - Understanding the abstraction



$$commit(\mathsf{c}, \{\mathsf{a}, \mathsf{b}\})$$

# Leaderless SMR - Understanding the abstraction



$$commit(\textcolor{red}{a}, \{\textcolor{purple}{c}\})$$

# Leaderless SMR - Understanding the abstraction



$$commit(\textcolor{red}{a}, \{\textcolor{blue}{c}\})$$

# Leaderless SMR - Understanding the abstraction



$$commit(\textcolor{red}{a}, \{\textcolor{purple}{c}\})$$

# Leaderless SMR - Understanding the abstraction



Let $deps^*($c$)$ be the transitive closure of the $deps$ relation starting from a command $\{$c$\}$.

# Leaderless SMR - Understanding the abstraction



A command c is stable once it is committed and no command in $deps^*(\text{c})$ is pending.

# Leaderless SMR - Understanding the abstraction



Here, only c is stable. All others are still pending, since their $deps^*$ include d (which is pending since it has $\perp$ as a dependency).

# Leaderless SMR - Understanding the abstraction



In this example, all commands are stable.
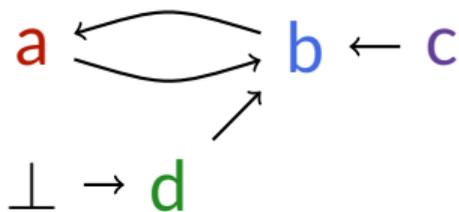
# Leaderless SMR - Understanding the abstraction

$$\perp \longrightarrow b \quad a \rightrightarrows c$$

A command c gets aborted when $deps(\mathsf{c})$ is set to $\top$.

# Leaderless SMR - Understanding the abstraction



$$\top \longrightarrow b \quad a \rightrightarrows c$$

A command c gets aborted when $deps(\textsf{c})$ is set to $\top$.

# Leaderless SMR - Understanding the abstraction



In that case, the command is removed from any $deps(\mathsf{d})$ and it will not appear later on.

# Leaderless SMR - Understanding the abstraction

$$\top \longrightarrow b \qquad a \underset{\longleftarrow}{\longrightarrow} c$$

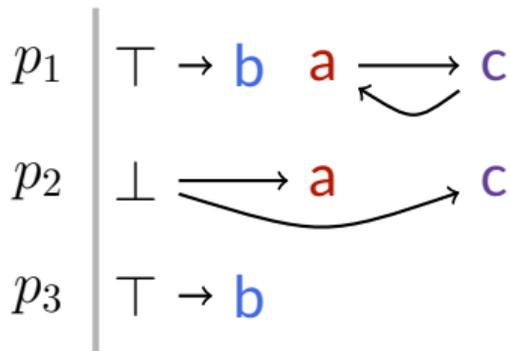In that case, the command is removed from any $deps(\mathsf{d})$ and it will not appear later on.

**Leaderless SMR - Understanding the abstraction**



$$p_1 \quad \top \to \text{b} \quad \text{a} \longrightarrow \text{c}$$
$$p_2 \quad \bot \longrightarrow \text{a} \longrightarrow \text{c}$$
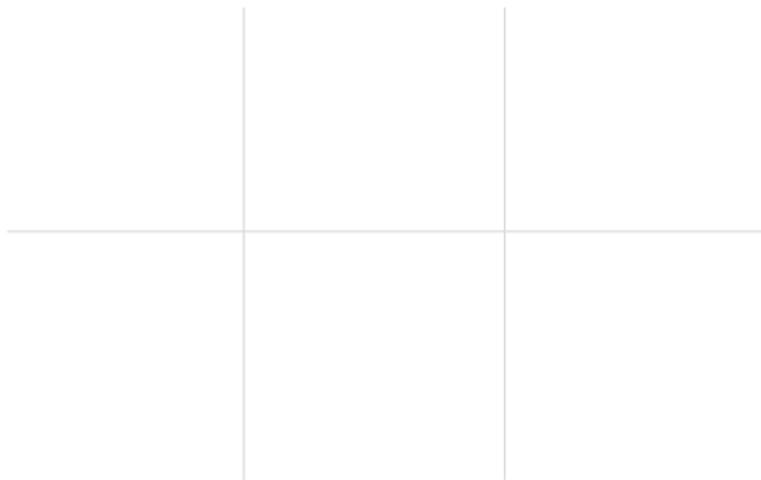$$p_3 \quad \top \to \text{b}$$

Guarantees:

**Stability**: For each command c, there exists $D$ such that if c is stable then $deps(\text{c}) = D$.

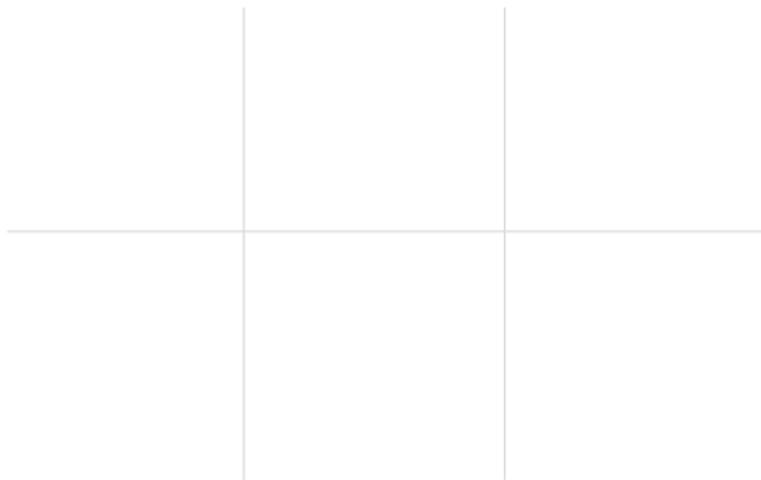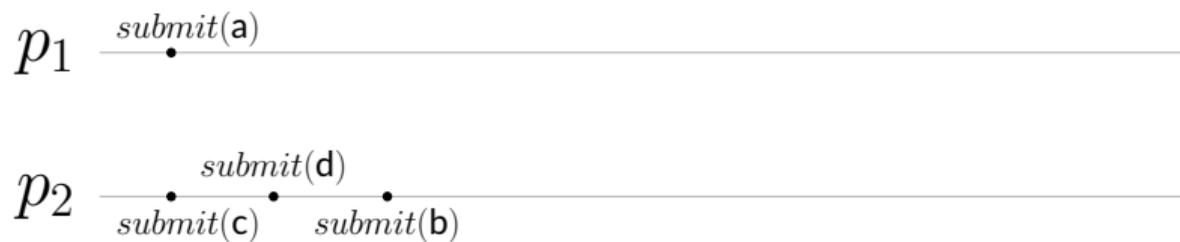**Consistency**: If a and b are both committed and conflicting, then $\text{a} \in deps(\text{b})$ or $\text{b} \in deps(\text{a})$.

# Leaderless SMR - Understanding the abstraction

$p_1$ ————————————————————————

$p_2$ ————————————————————————

# Leaderless SMR - Understanding the abstraction

# Leaderless SMR - Understanding the abstraction

$p_1$ $\quad$ $submit(\mathsf{a})$ $\qquad$ $\mathbf{g_1}$

$\qquad\qquad\qquad commit(\mathsf{a}, \{\mathsf{b}\})$

$p_2$ $\quad$ $submit(\mathsf{d})$

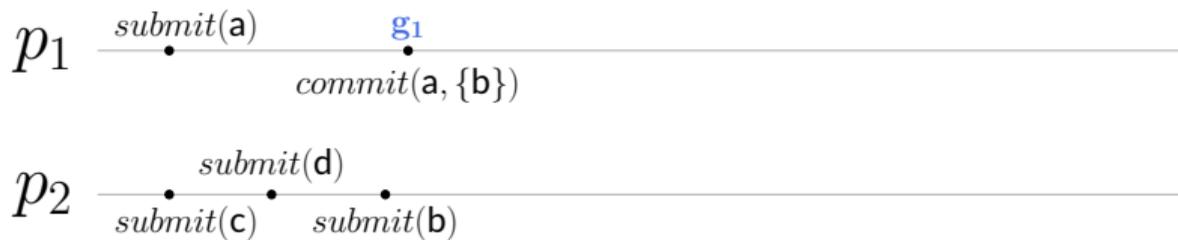$\qquad submit(\mathsf{c}) \quad submit(\mathsf{b})$
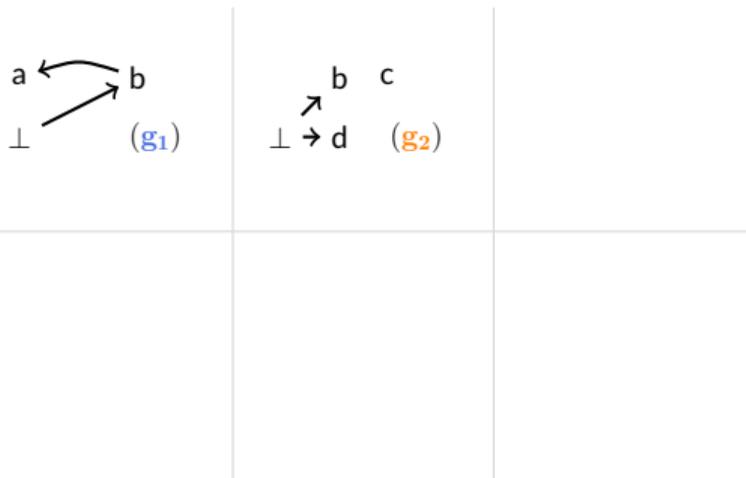
a $\leftarrow$ b

$\perp$ $\quad$ $(\mathbf{g_1})$

# Leaderless SMR - Understanding the abstraction

# Leaderless SMR - Understanding the abstraction

$$p_1 \quad \underset{\bullet}{submit(\mathsf{a})} \qquad \underset{\bullet}{\overset{\mathbf{g_1}}{\bullet}} \quad \underset{\bullet}{commit(\mathsf{c}, \{\})}$$

$$commit(\mathsf{a}, \{\mathsf{b}\})$$

$$p_2 \quad \underset{submit(\mathsf{c})}{\overset{submit(\mathsf{d})}{\bullet}} \quad \underset{submit(\mathsf{b})}{\overset{commit(\mathsf{c}, \{\})}{\bullet}} \quad \underset{\mathbf{g_2} \quad commit(\mathsf{b}, \{\mathsf{c}, \mathsf{d}, \mathsf{a}\})}{\overset{\mathbf{g_3}}{\bullet}}$$



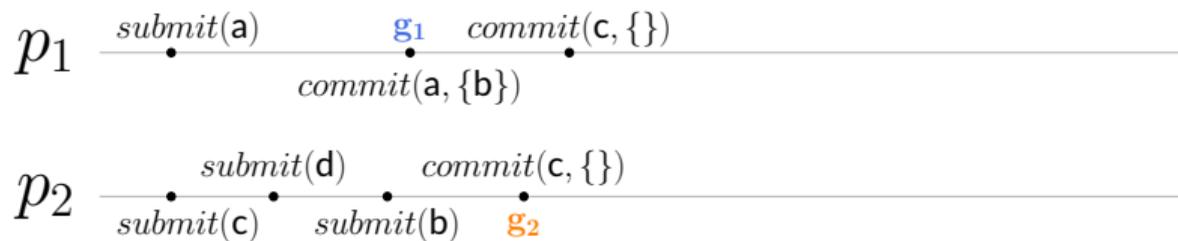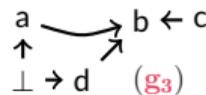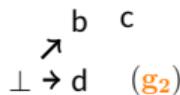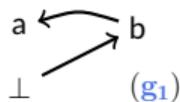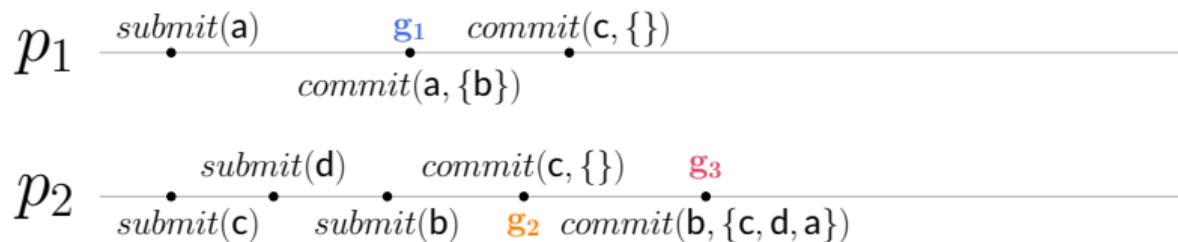| a ← b | b  c | a → b ← c |
|-------|------|-----------|
| ⊥  (g₁) | ⊥ → d  (g₂) | ⊥ → d  (g₃) |

# Leaderless SMR - Understanding the abstraction



$p_1$ — $submit(\mathsf{a})$  $g_1$  $commit(\mathsf{c}, \{\})$  $g_4$
$commit(\mathsf{a}, \{\mathsf{b}\})$  $commit(\mathsf{b}, \{\mathsf{c}, \mathsf{d}, \mathsf{a}\})$

$p_2$ — $submit(\mathsf{d})$  $commit(\mathsf{c}, \{\})$  $g_3$  $commit(\mathsf{a}, \{\mathsf{b}\})$
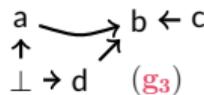$submit(\mathsf{c})$  $submit(\mathsf{b})$  $g_2$  $commit(\mathsf{b}, \{\mathsf{c}, \mathsf{d}, \mathsf{a}\})$  $g_4$

# Leaderless SMR - Understanding the abstraction

# Leaderless SMR - Understanding the abstraction

$$a \rightleftharpoons b \leftarrow c$$
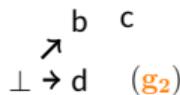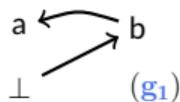$$d \quad (g_5)$$

A command gets executed once it is stable.

**Leaderless SMR - Understanding the abstraction**



All commands are stable in $g_5$, therefore they can be executed.

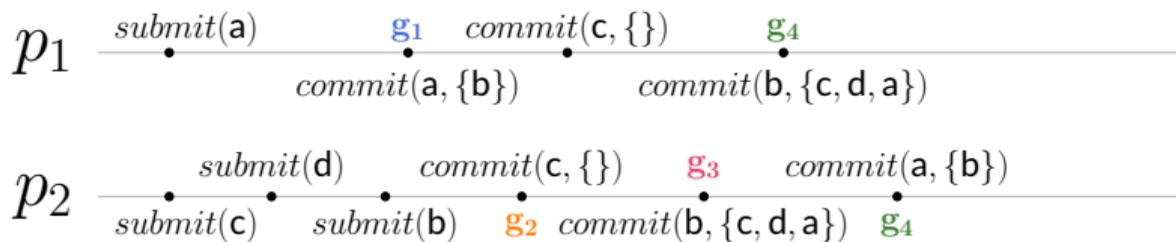# Leaderless SMR - Understanding the abstraction



Cycles are broken deterministically.

# Leaderless SMR - Understanding the abstraction



$$commit(\mathsf{c}, \{\}); \; commit(\mathsf{d}, \{\}); \; commit(\mathsf{a}, \{\mathsf{b}\}); \; commit(\mathsf{b}, \{\mathsf{a}\});$$

# Leaderless SMR - Understanding the abstraction



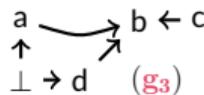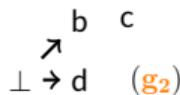$execute(\textsf{c}); execute(\textsf{d}); execute(\textsf{a}); execute(\textsf{b})$ if $\textsf{a} < \textsf{b}$

# Leaderless SMR - Understanding the abstraction



$execute(\mathsf{c}); execute(\mathsf{d}); execute(\mathsf{b}); execute(\mathsf{a})$ if $\mathsf{b} < \mathsf{a}$

# A General Framework for Leaderless SMR

- What is the essence of Leaderless SMR?
  - Leaderless SMR **does not** compute an ordering over conflicting commands.
  - Conflicting commands must simply observe one another (**Consistency** property).

- How to capture this feature?
  - Decomposition into two services: Dependency Discovery Service and Consensus Service.

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

```
1:  submit(c) :=
2:      pre:  p = coord(c) ∨ coord(c) ∈ 𝒟
3:      eff:  (D, b) ← DDS.announce(c)
4:            if b = false then D ← CONSc.propose(D)
5:            deps(c) ← D
6:            send(c, deps(c)) to Π \ {p}
7:
8:  when recv(c, D)
9:      eff:  deps(c) ← D
10:
```

- Algorithm 2 depicts an abstract protocol to decide a command.

# Abstract protocol to decide a command

■ **Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:   **pre:** $\boxed{p = coord(\mathsf{c})} \lor coord(\mathsf{c}) \in \mathcal{D}$
3:   **eff:** $(D, b) \leftarrow$ DDS.$announce(\mathsf{c})$
4:       **if** $b = false$ **then** $D \leftarrow$ CONS$_\mathsf{c}$.$propose(D)$
5:       $deps(\mathsf{c}) \leftarrow D$
6:       $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:   **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

- $coord(\mathsf{c})$ is the initial process in charge of submitting the command to the replicated state machine.

# Abstract protocol to decide a command

■ **Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         if $b = false$ then $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $\boxed{deps(\mathsf{c}) \leftarrow D}$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ to $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $\boxed{deps(\mathsf{c}) \leftarrow D}$
10:

- The algorithm computes a set of dependencies of a command abiding to properties **Stability** and **Consistency**

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:      **pre:** $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3:      **eff:** $(D, b) \leftarrow$ DDS.$announce(\mathsf{c})$
4:            **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:            $deps(\mathsf{c}) \leftarrow D$
6:            $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:      **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

- It uses a dependency discovery service (DDS).

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

- A family of consensus objects $((\mathsf{CONS_c})_{\mathsf{c} \in \mathcal{C}})$.

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \lor coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c})) \; to \; \Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

- A failure detector ($\mathcal{D}$) that returns a set of suspected processes.

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \lor coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow DDS.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow CONS_{\mathsf{c}}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

- To create a valid proposal for $(CONS_{\mathsf{c}})_{\mathsf{c} \in \mathcal{C}}$, a process $p$ relies on the DDS.

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

- This shared object offers a single operation $announce(\mathsf{c})$ that returns a pair $(D, b)$, where $D \in 2^{\mathcal{C}} \cup \{\top\}$ and $b \in \{0, 1\}$ is a flag.

# Abstract protocol to decide a command

■ **Algorithm 2** Deciding a command c – code at process $p$

1: $submit(c) :=$
2:    **pre:** $p = coord(c) \lor coord(c) \in \mathcal{D}$
3:    **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(c)$
4:        **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:        $deps(c) \leftarrow D$
6:        $send(c, deps(c))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(c, D)$
9:    **eff:** $deps(c) \leftarrow D$
10:

- When the return value is in $2^{\mathcal{C}}$, the service suggests to **commit** the command. Otherwise, the command should be **aborted**.
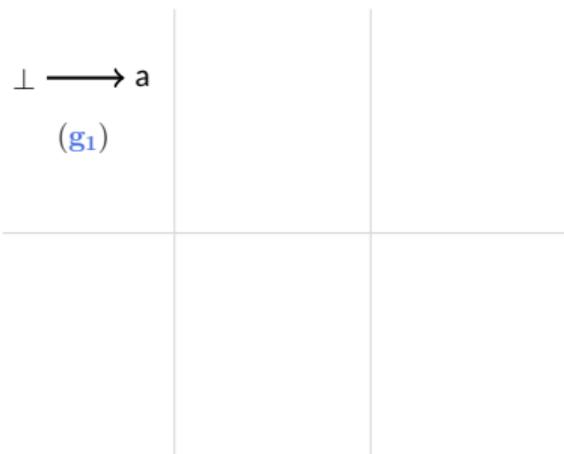
# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(c) :=$
2:     **pre:**  $p = coord(c) \lor coord(c) \in \mathcal{D}$
3:     **eff:**  $(D, b) \leftarrow \mathsf{DDS}.announce(c)$
4:          **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:          $deps(c) \leftarrow D$
6:          $send(c, deps(c))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(c, D)$
9:     **eff:**  $deps(c) \leftarrow D$
10:

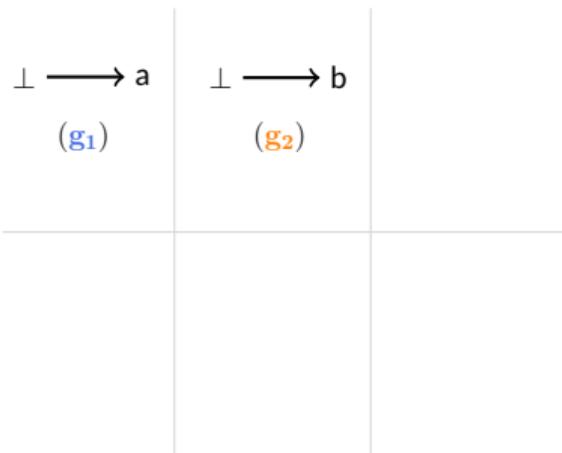- When the flag is set, the service indicates that a **spontaneous agreement occurs**.
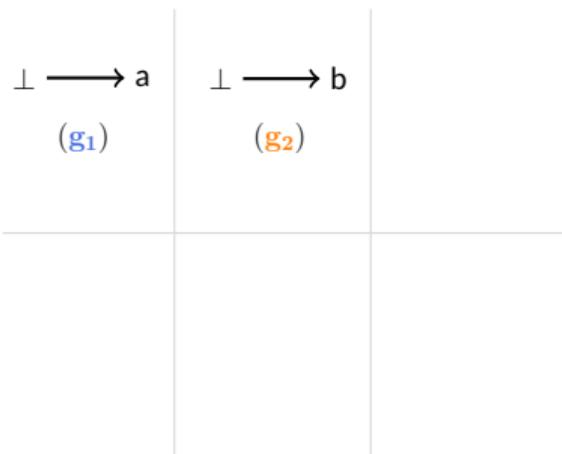
# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process $p$

```
1:  submit(c) :=
2:      pre:  p = coord(c) ∨ coord(c) ∈ 𝒟
3:      eff:  (D, b) ← DDS.announce(c)
4:            if b = false then D ← CONS_c.propose(D)
5:            deps(c) ← D
6:            send(c, deps(c)) to Π \ {p}
7:
8:  when recv(c, D)
9:      eff:  deps(c) ← D
10:
```

- In such a case, $p$ can directly commit c with the return value of the DDS service and bypass CONS$_c$; this is called a fast path.

# Abstract protocol to decide a command

$p_1$ ————————————————

$p_2$ ————————————————



■ **Algorithm 2** Deciding a command c – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

# Abstract protocol to decide a command

$p_1$ ——•—— $submit(\mathsf{a})$

$p_2$ ——•—— $submit(\mathsf{b})$

**Algorithm 2** Deciding a command $\mathsf{c}$ – code at process $p$

1: $submit(\mathsf{c})$ :
2:     **pre:** $p = coord(\mathsf{c}) \lor coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

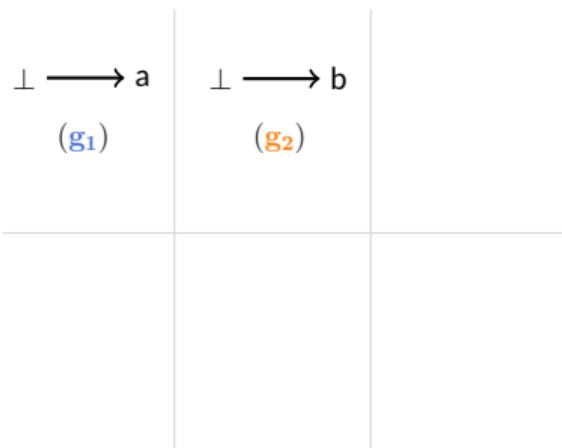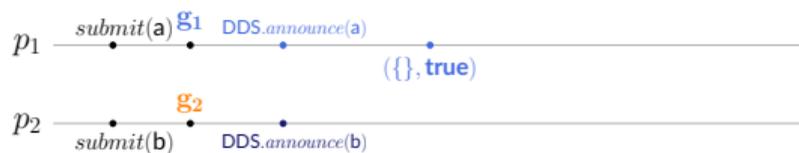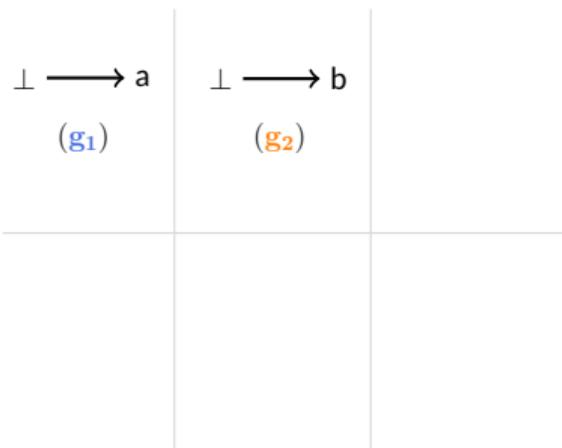# Abstract protocol to decide a command



**Algorithm 2** Deciding a command $c$ – code at process $p$

1:  $submit(c) :=$
2:      **pre:**  $p = coord(c) \lor coord(c) \in \mathcal{D}$
3:      **eff:**  $(D, b) \leftarrow \mathsf{DDS}.announce(c)$
4:              **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:              $deps(c) \leftarrow D$
6:              $send(c, deps(c))$ $to$ $\Pi \setminus \{p\}$
7:
8:  **when** $recv(c, D)$
9:      **eff:**  $deps(c) \leftarrow D$
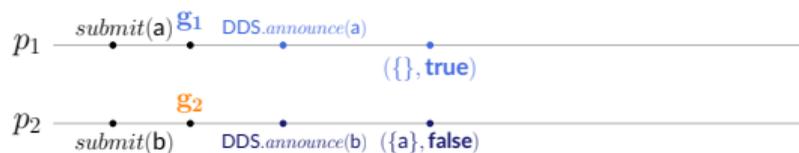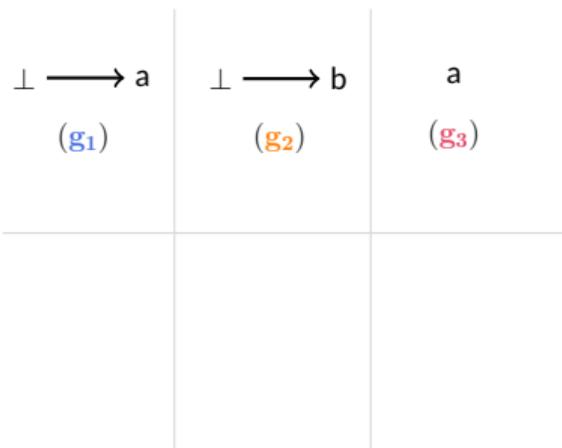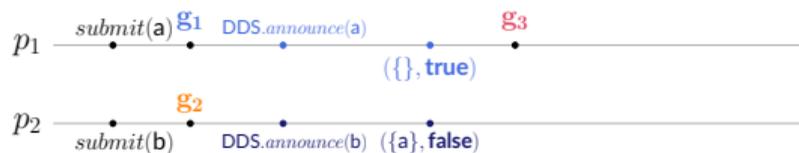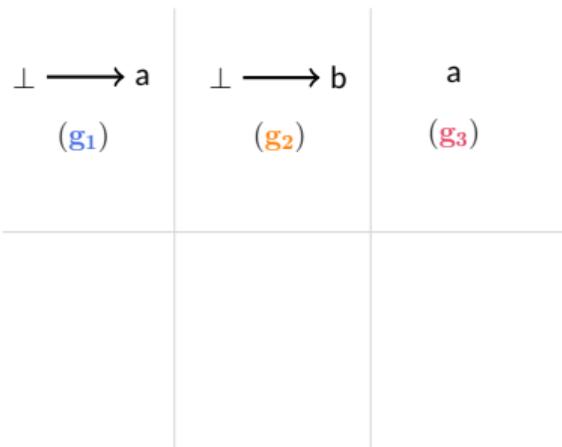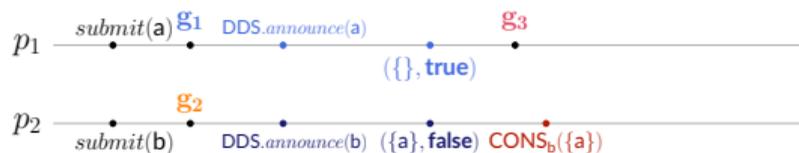10:

# Abstract protocol to decide a command



$p_1$ — $submit(\mathsf{a})$ — $\mathsf{g_1}$

$p_2$ — $\mathsf{g_2}$ — $submit(\mathsf{b})$

$\perp \longrightarrow \mathsf{a}$      $\perp \longrightarrow \mathsf{b}$

$(\mathsf{g_1})$              $(\mathsf{g_2})$

**Algorithm 2** Deciding a command $\mathsf{c}$ – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:**   $p = coord(\mathsf{c}) \lor coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:**   $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:          **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:          $deps(\mathsf{c}) \leftarrow D$
6:          $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:**   $deps(\mathsf{c}) \leftarrow D$
10:

# Abstract protocol to decide a command

# Abstract protocol to decide a command



$p_1$ $\quad$ $submit(\mathsf{a})$ $\overset{\mathbf{g_1}}{\bullet}$ $\quad$ DDS.$announce(\mathsf{a})$

$p_2$ $\quad$ $\overset{\mathbf{g_2}}{\bullet}$ $\quad$ $submit(\mathsf{b})$ $\quad$ DDS.$announce(\mathsf{b})$

$\perp \longrightarrow \mathsf{a}$ $\qquad$ $\perp \longrightarrow \mathsf{b}$

$(\mathbf{g_1})$ $\qquad\qquad$ $(\mathbf{g_2})$

---

■ **Algorithm 2** Deciding a command $\mathsf{c}$ – code at process $p$

1: $submit(\mathsf{c}) :=$
2: $\quad$ **pre:** $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3: $\quad$ **eff:** $(D, b) \leftarrow$ DDS.$announce(\mathsf{c})$
4: $\qquad$ **if** $b = false$ **then** $D \leftarrow$ CONS$_{\mathsf{c}}.propose(D)$
5: $\qquad$ $deps(\mathsf{c}) \leftarrow D$
6: $\qquad$ $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9: $\quad$ **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

# Abstract protocol to decide a command

# Abstract protocol to decide a command



$p_1$ $\quad$ $submit(\mathsf{a})$ $\mathbf{g_1}$ $\quad$ DDS.$announce(\mathsf{a})$

$(\{\}, \mathbf{true})$

$p_2$ $\quad$ $\mathbf{g_2}$

$submit(\mathsf{b})$ $\quad$ DDS.$announce(\mathsf{b})$ $(\{\mathsf{a}\}, \mathbf{false})$

$\perp \longrightarrow \mathsf{a}$ $\quad$ $\perp \longrightarrow \mathsf{b}$

$(\mathbf{g_1})$ $\quad\quad$ $(\mathbf{g_2})$

**Algorithm 2** Deciding a command $\mathsf{c}$ – code at process $p$

1:  $submit(\mathsf{c}) :=$
2:  $\quad$ **pre:** $p = coord(\mathsf{c}) \lor coord(\mathsf{c}) \in \mathcal{D}$
3:  $\quad$ **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:  $\quad\quad$ **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:  $\quad\quad$ $deps(\mathsf{c}) \leftarrow D$
6:  $\quad\quad$ $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8:  **when** $recv(\mathsf{c}, D)$
9:  $\quad$ **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

# Abstract protocol to decide a command



$p_1$    $submit(\mathsf{a})$   $\mathbf{g_1}$   DDS.*announce*(a)    $\mathbf{g_3}$

$(\{\}, \mathbf{true})$

$p_2$    $\mathbf{g_2}$

$submit(\mathsf{b})$    DDS.*announce*(b)   $(\{a\}, \mathbf{false})$

$\bot \longrightarrow \mathsf{a}$    $\bot \longrightarrow \mathsf{b}$    $\mathsf{a}$

$(\mathbf{g_1})$     $(\mathbf{g_2})$     $(\mathbf{g_3})$

---

**Algorithm 2** Deciding a command c – code at process $p$

1:   $submit(\mathsf{c}) :=$
2:     **pre:**   $p = coord(\mathsf{c}) \vee coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:**   $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:        **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:        $deps(\mathsf{c}) \leftarrow D$
6:        $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:**   $deps(\mathsf{c}) \leftarrow D$
10:

# Abstract protocol to decide a command



**Algorithm 2** Deciding a command c – code at process $p$

1: $submit(c) :=$
2:      **pre:** $p = coord(c) \lor coord(c) \in \mathcal{D}$
3:      **eff:** $(D, b) \leftarrow \text{DDS}.announce(c)$
4:          **if** $b = false$ **then** $D \leftarrow \text{CONS}_c.propose(D)$
5:          $deps(c) \leftarrow D$
6:          $send(c, deps(c))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(c, D)$
9:      **eff:** $deps(c) \leftarrow D$
10:

# Abstract protocol to decide a command



$p_1$ — $submit(\mathsf{a})$ $g_1$ DDS.$announce(\mathsf{a})$ $g_3$
$(\{\}, \mathbf{true})$

$p_2$ — $submit(\mathsf{b})$ $g_2$ DDS.$announce(\mathsf{b})$ $(\{\mathsf{a}\}, \mathbf{false})$ $\mathsf{CONS_b}(\{\mathsf{a}\})$ $g_4$

$\perp \longrightarrow \mathsf{a}$       $\perp \longrightarrow \mathsf{b}$       $\mathsf{a}$
$(g_1)$                   $(g_2)$                   $(g_3)$

$\mathsf{a} \longrightarrow \mathsf{b}$
$\uparrow$
$\perp$  $(g_4)$

**Algorithm 2** Deciding a command $\mathsf{c}$ – code at process $p$

1: $submit(\mathsf{c}) :=$
2:     **pre:** $p = coord(\mathsf{c}) \lor coord(\mathsf{c}) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(\mathsf{c})$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(\mathsf{c}) \leftarrow D$
6:         $send(\mathsf{c}, deps(\mathsf{c}))$ $to$ $\Pi \setminus \{p\}$
7:
8: **when** $recv(\mathsf{c}, D)$
9:     **eff:** $deps(\mathsf{c}) \leftarrow D$
10:

# Abstract protocol to decide a command

**Algorithm 2** Deciding a command c – code at process p

1: $submit(c) :=$
2:     **pre:** $p = coord(c) \lor coord(c) \in \mathcal{D}$
3:     **eff:** $(D, b) \leftarrow \mathsf{DDS}.announce(c)$
4:         **if** $b = false$ **then** $D \leftarrow \mathsf{CONS_c}.propose(D)$
5:         $deps(c) \leftarrow D$
6:         $send(c, deps(c))$ to $\Pi \setminus \{p\}$
7:
8: **when** $recv(c, D)$
9:     **eff:** $deps(c) \leftarrow D$
10:

# Abstract protocol to decide a command



**Algorithm 2** Deciding a command c – code at process $p$

```
1:  submit(c) :=
2:      pre:  p = coord(c) ∨ coord(c) ∈ 𝒟
3:      eff:  (D, b) ← DDS.announce(c)
4:            if b = false then D ← CONS_c.propose(D)
5:            deps(c) ← D
6:            send(c, deps(c)) to Π \ {p}
7:
8:  when recv(c, D)
9:      eff:  deps(c) ← D
10:
```

# Core Properties

Multiple implementations of Dependency Discovery Service and Consensus Service are possible.

- We have cast different leaderless SMR protocols into our framework.

- How to further **characterize** the protocols?

# Core Properties

Multiple implementations of Dependency Discovery Service and Consensus Service are possible.

- We have cast different leaderless SMR protocols into our framework.

- How to further **characterize** the protocols?

Core properties:

*(Reliability)*
*(Optimal Latency)*
*(Load Balancing)*

## Core Properties

(*Reliability*)  In every run, if there are at most $f$ failures, every submitted command gets eventually decided at every correct process.

## Core Properties

(*Reliability*)  In every run, if there are at most $f$ failures, every submitted command gets eventually decided at every correct process.

- Rationale:
  - SMR helps to mask failures and asynchrony in a distributed system.
  - Parameter $f$ captures the largest number of failures tolerated by a protocol.

**Core Properties**

(*Optimal Latency*)  During a nice run, every call to $announce(\mathsf{c})$
returns after two message delays in the absence of
conflicting commands

## Core Properties

(*Optimal Latency*)  During a nice run, every call to $announce(c)$ returns after two message delays in the absence of conflicting commands

- Rationale:
  - Leaderless SMR protocols exploit the absence of contention to boost performance.
  - Some protocols are able to execute a command after a single round-trip (optimal).

## Core Properties

(*Load Balancing*)  During a nice run, any $n - F$ replicas can be used to announce a command c.

## Core Properties

(*Load Balancing*)  During a nice run, any $n - F$ replicas can be used to announce a command c.

- Rationale:
    - The replicas that take place in the fast path vary from one protocol to another.
    - For example, Mencius use all the processes.
    - Captures scalability, as any quorum of $n - F$ processes may be used to order a command.

# Complexity result - ROLL theorem

## Theorem
*Consider an SMR protocol that satisfies the ROLL properties. Then, it is true that $2F + f - 1 \leq n$.*

- The inequation presented in the theorem expresses the trade-off between scalability and fault-tolerance.

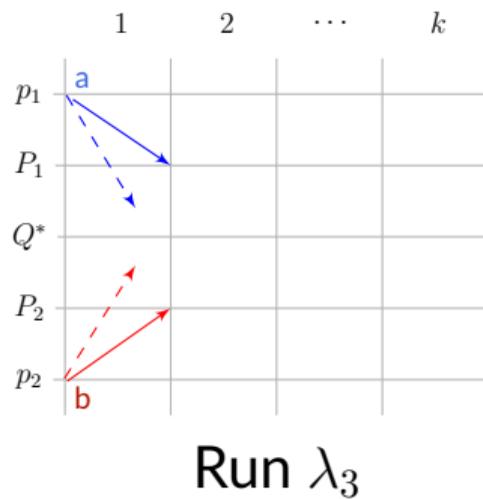- Differently than Lamport's lower bound (*Fast Learning Theorem*), ROLL captures more accurately this trade-off in leaderless protocols.

# ROLL theorem - Proof sketch

By contradiction, using a round-based reasoning. Let us assume a protocol $\mathcal{P}$ that satisfies all the ROLL properties with $2F + f - 1 > n$.



Quorums in use

# ROLL theorem - Proof sketch



Quorums in use

Run $\lambda_1$
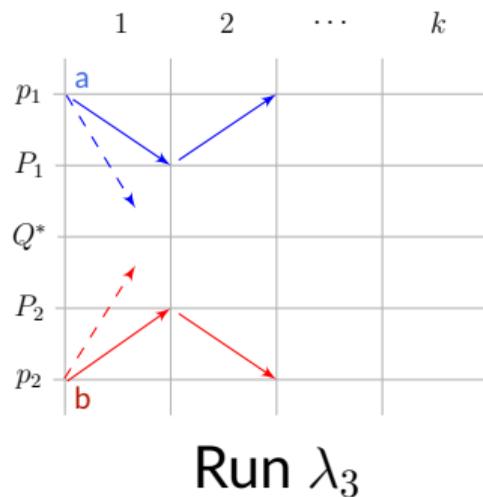
# ROLL theorem - Proof sketch



Quorums in use

Run $\lambda_1$

Run $\lambda_2$

# ROLL theorem - Proof sketch



Run $\lambda_1$       Run $\lambda_3$

# ROLL theorem - Proof sketch



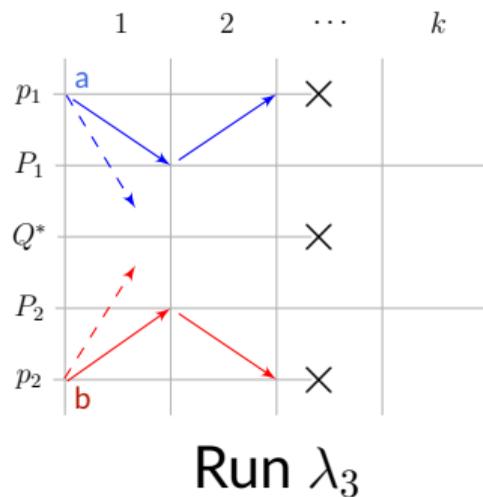Run $\lambda_1$          Run $\lambda_3$

# ROLL theorem - Proof sketch



Run $\lambda_1$                    Run $\lambda_3$
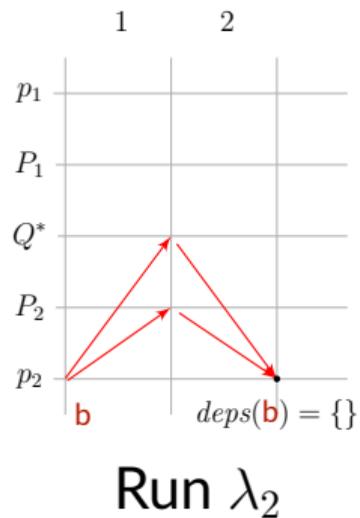
# ROLL theorem - Proof sketch



Run $\lambda_1$

Run $\lambda_3$

# ROLL theorem - Proof sketch



Run $\lambda_1$          Run $\lambda_3$

# ROLL theorem - Proof sketch



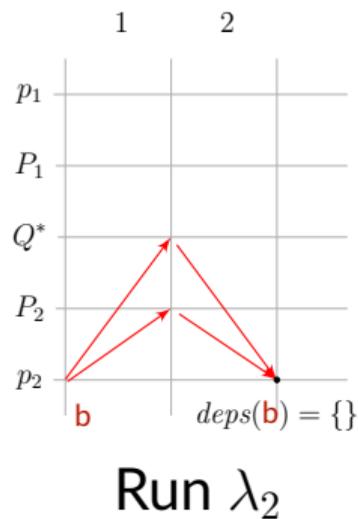Run $\lambda_1$      Run $\lambda_3$      Run $\lambda_4$

# ROLL theorem - Proof sketch



Run $\lambda_1$                 Run $\lambda_3$                 Run $\lambda_4$

# ROLL theorem - Proof sketch



Run $\lambda_1$

Run $\lambda_3$

Run $\lambda_4$

# ROLL theorem - Proof sketch



Run $\lambda_1$      Run $\lambda_3$      Run $\lambda_4$

# ROLL theorem - Proof sketch



Run $\lambda_1$      Run $\lambda_3$      Run $\lambda_4$

# ROLL theorem - Proof sketch



Run $\lambda_1$      Run $\lambda_3$      Run $\lambda_4$

# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

Run $\lambda_2$

Run $\lambda_3$

# ROLL theorem - Proof sketch



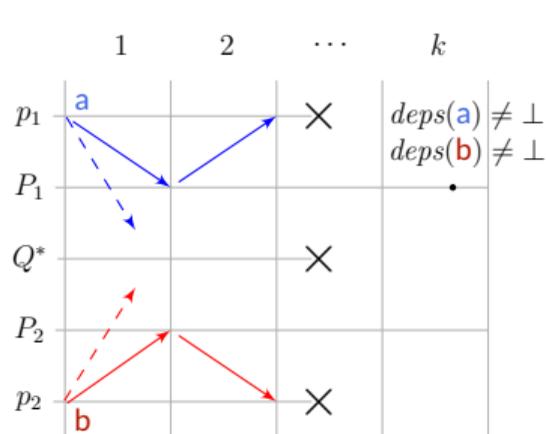Run $\lambda_2$

Run $\lambda_3$

# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

Run $\lambda_5$

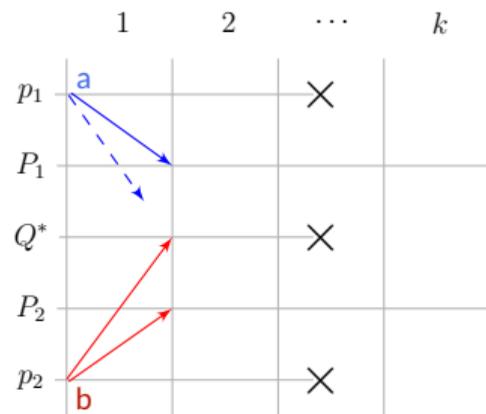# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

Run $\lambda_5$

# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

Run $\lambda_5$

# ROLL theorem - Proof sketch



Run $\lambda_2$        Run $\lambda_3$        Run $\lambda_5$

# ROLL theorem - Proof sketch



Run $\lambda_2$

Run $\lambda_3$

Run $\lambda_5$

## Complexity result - Optimality

## Theorem
*Consider an SMR protocol that satisfies the ROLL properties. Then, it is true that $2F + f - 1 \leq n$.*

- A protocol is *ROLL-optimal* when the parameters $F$ and $f$ cannot be improved according to the theorem.

## For example:

- With $n = 5$, there is a single such tuple: $(F, f) = (2, 2)$.

## Complexity result - Optimality

## Theorem
*Consider an SMR protocol that satisfies the ROLL properties. Then, it is true that $2F + f - 1 \leq n$.*

- A protocol is *ROLL-optimal* when the parameters $F$ and $f$ cannot be improved according to the theorem.

For example:

- With $n = 7$, there are two tuples possible: $(F, f) = (2, 3)$ and $(3, 2)$.

# Complexity result - Optimality

## Theorem
*Consider an SMR protocol that satisfies the ROLL properties. Then, it is true that $2F + f - 1 \leq n$.*

- A protocol is *ROLL-optimal* when the parameters $F$ and $f$ cannot be improved according to the theorem.

## For example:

- Epaxos is optimal for $n = 5$.

## Complexity result - Optimality

## Theorem
*Consider an SMR protocol that satisfies the ROLL properties. Then, it is true that $2F + f - 1 \leq n$.*

- A protocol is *ROLL-optimal* when the parameters $F$ and $f$ cannot be improved according to the theorem.

For example:

- No known optimal protocol when $n = 7$.

# Conclusion

- The decomposition into Dependency Discovery Service and Consensus Service.

# Conclusion

- The decomposition into Dependency Discovery Service and Consensus Service.

- Leaderless SMR Properties.

# Conclusion

- The decomposition into Dependency Discovery Service and Consensus Service.

- Leaderless SMR Properties.

- Complexity results:
    - ROLL theorem. Explains the trade-off between reliability and scalability.
    - Chaining effect. Explains the reason for the high latency found in real-world implementation of leaderless protocols.