# Benchmarking Model Checkers with Distributed Algorithms

Étienne Coulouma-Dupont



November 24, 2011

# Introduction

*LastVoting*

Simulator

Benchmarks

## Panda Project



Industrial Benchmark  Airbus Code
Academic Benchmark  Paxos Algorithm

## The Consensus Problem

- ▶ Context: Distributed Systems
- ▶ *Consensus* is the problem of making processes agree on a common value in spite of faults

# The Consensus Problem

An Algorithms solves Consensus if and only if it satifies the following conditions:

- *Integrity*: Any decision value is the proposed value of some process.
- *Agreement*: No two different values are decided.
- *Termination*: All process eventually decide.

Consensus : application

# Consensus Algorithm with industrial applications

Chubby (Google, 2006) [Bur06][CGR07] implements a fault tolerant data-base:

- ▶ Fault-tolerance is achieved through replication
- ▶ Consistency is achieved using Paxos Algorithm

## Paxos

- ▶ Family of algorithms built to solve Consensus
- ▶ First published by Leslie Lamport in 1998 [Lam98]
- ▶ Termination is not guaranteed without additional assumptions on liveness
- ▶ Safety is guaranteed

Several implementations in different models:

*ChandraToueg* Attach to each process in the system a *failure detector*

*LastVoting* Synchronous communication (*rounds*)

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○● | ○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

Paxos

## Related Works

This algorithm has been intensively studied

- ▶ Fuzzati [Fuz08]: proof of Paxos and *ChandraToueg* with rewriting rules
- ▶ Küfner et al.: assisted proofs in `Isabelle`
- ▶ Tsuchiya & Schiper: model-checking of *LastVoting*
  - ▶ SAT solver: up to 10 processes
  - ▶ Traditional tools: up to 4 processes

Introduction

*LastVoting*

Simulator

Benchmarks

## *LastVoting* : Hypothesis
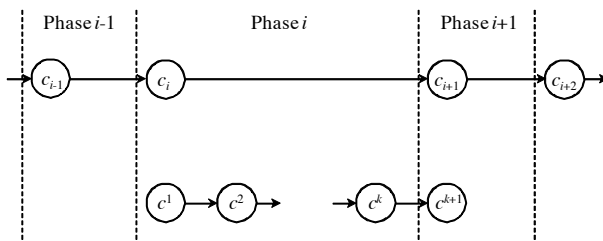
Assumptions for the environment:

- ▶ with transmission faults : in each round, the adversary chooses a set of edges in which all messages will be correctly transmitted

- ▶ pseudo-synchronous: any message which is not received in the same round as the one during which it was send is thrown out by the process which receives it

- ▶ complete network: simplifies the way algorithms are expressed

- ▶ each process has a unique identity

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○● | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○○ | |

Hypothesis

## *LastVoting* : Hypothesis

We assume that each process *p* has at any time access to the following pieces of information:

- ► The round *r* in which it is
- ► The coordinator at round *r*

Introduction
○○
○
○○○

*LastVoting*
○○
●○○○○○○○○
○○○○○○○○○○○○○○○○○○○○○○○○

Simulator
○
○
○

Benchmarks
○
○○
○○○○○○○○○○○
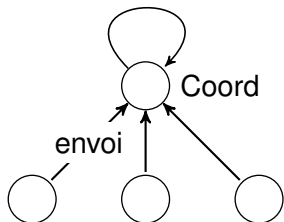
Conclusion

The Algorithm

# The Algorithm's Procedure



Figure: Transitions of configurations at the phase level (top) and at the round level (bottom) [TS11].

## The Algorithm's Procedure

In each round $r$, each process $p$:

- can send a messages according to a sending function $S_p^r$, depending on the process' state
- then can compute a new state according to a state transition function $T_p^r$, depending on the messages received

**Selection** ($r = 4\phi - 3$)



Coord

envoi

**Commit** ($r = 4\phi - 2$)



Coord

envoi

**Ack** ($r = 4\phi - 1$)



Coord

envoi

**Decision** ($r = 4\phi$)



Coord

envoi

## Notations

- $co = Coord(\phi)$
- $\mathfrak{D}(\phi) = \{p \in \Pi \mid d_p^{4\phi-3} \neq ?\}$
- for all $p \in \Pi$ and $r \in [4\phi - 3, 4\phi - 1]$, $RCV(p, r)$ is the set of processes from which $p$ receives a *non empty* message in round $r$

**Initialization :**

$x_p \in Val$, initially $v_p$

// $v_p$ is the proposed value of $p$.

$vote_p \in Val \cup \{?\}$, initially ?

// $Val$ is the set of values that may be proposed.

$commit_p$ a Boolean, initially false

$ready_p$ a Boolean, initially false

$ts_p \in \mathbb{N}$, initially 0

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
| :--- | :--- | :--- | :--- | :--- |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○●○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

The Algorithm

**Round** $r = 4\phi - 3$ :             // $\#RCV(co, 4\phi - 3) > n/2$

$S_p^r$ :
  **send** $\langle x_p, ts_p \rangle$ to $Coord(p, \phi)$

$T_p^r$ :
  **if** $p = Coord(p, \phi)$ **and** *number of* $\langle \nu, \theta \rangle$ *received* $> n/2$
  **then**
    let be the largest $\overline{\theta}$ from $\langle \nu, \theta \rangle$ received;
    $vote_p := $ one $\nu$ such that $\langle \nu, \overline{\theta} \rangle$ is received;
    $commit_p := \texttt{true};$

**Round** $r = 4\phi - 2$ **:** // $\#\{p \in \Pi \mid co \in RCV(p, 4\phi - 2)\} > n/2$

> $S_p^r$ :
> > **if** $p = Coord(p, \phi)$ **and** $commit_p$ **then**
> > > **send** $\langle vote_p \rangle$ to all processes;
>
> $T_p^r$ :
> > **if** $received\langle v \rangle$ *from* $Coord(p, \phi)$ **then**
> > > $x_p := v$;
> > > $ts_p := \phi$;

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
| :-- | :-- | :-- | :-- | :-- |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○●○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

The Algorithm

**Round** $r = 4\phi - 1$ **:**         // $\#RCV(co, 4\phi - 1) > n/2$

  $S_p^r$ :
     **if** $ts_p = \phi$ **then**
        **send** $\langle ack \rangle$ to $Coord(p, \phi)$;

  $T_p^r$ :
     **if** $p = Coord(p, \phi)$ **and** *number of* $\langle ack \rangle$ *received* $> n/2$
     **then**
        $ready_p := \mathtt{true}$;

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○● | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

The Algorithm

**Round** $r = 4\phi$ **:**          $// \ \forall p \in \Pi \setminus \mathfrak{D}(\phi) : co \in RCV(p, 4\phi)$

$S_p^r :$
    **if** $Coord(p, \phi)$ **and** $ready_p$ **then**
        **send** $\langle vote_p \rangle$ to all processes;

$T_p^r :$
    **if** $received\langle v \rangle$ *from* $Coord(p, \phi)$ **then**
        DECIDE $(v)$

    **if** $p = Coord(p, \phi)$ **then**
        $ready_p := \texttt{false};$
        $commit_p := \texttt{false};$

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ●○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○○ | |

Analysis of the Algorithm

# Terminology

Given a partial execution, we define

0-Valence if the value 0 is to be decided

1-Valence if the value 1 is to be decided

Univalence if one the above holds

Bivalence if none of the above holds

Formally, Univalence is defined for *LastVoting* as follows:

$$\exists\, v \in \mathit{Val},\ \exists\, \mathfrak{P} \subseteq \Pi :$$
$$\wedge\, \#\mathfrak{P} \geq n/2$$
$$\wedge\, \mathfrak{P} = \{p \in \Pi \,|\, pr_p = v\} \wedge (\forall q \in \Pi \setminus \mathfrak{P} : ts_q < ts_p)$$

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○●○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○○ | |

Analysis of the Algorithm

# Termination

Termination holds iff $P^{sync}(\phi)$ eventually holds:

$P^{sync}(\phi) \triangleq \exists \phi > 0, \exists co \in \Pi :$
*ronde*1 : $\quad \wedge \; \#RCV(co, 4\phi - 3) > n/2 \qquad\qquad (CPROP(\phi))$

Introduction  **LastVoting**  Simulator  Benchmarks  Conclusion
○○  ○○  ○  ○
○  ○○○○○○○○○○  ○  ○○
○○○  ○○○○○○○○○○○○○○○○○○○○○○○○  ○○○○○○○○○○○○

Analysis of the Algorithm

# Termination
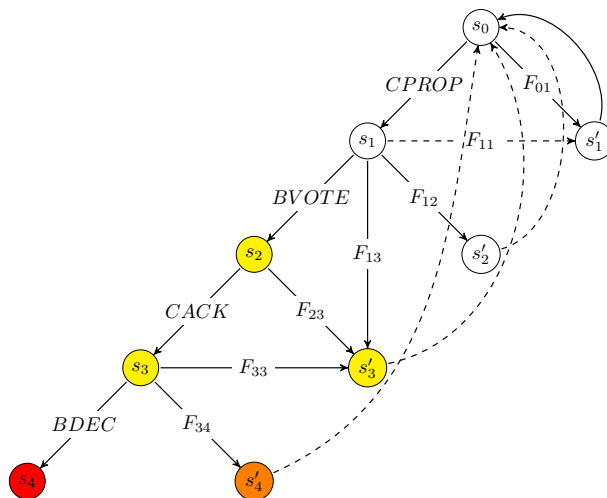
Termination holds iff $P^{sync}(\phi)$ eventually holds:

$P^{sync}(\phi) \triangleq \exists \phi > 0, \exists co \in \Pi :$

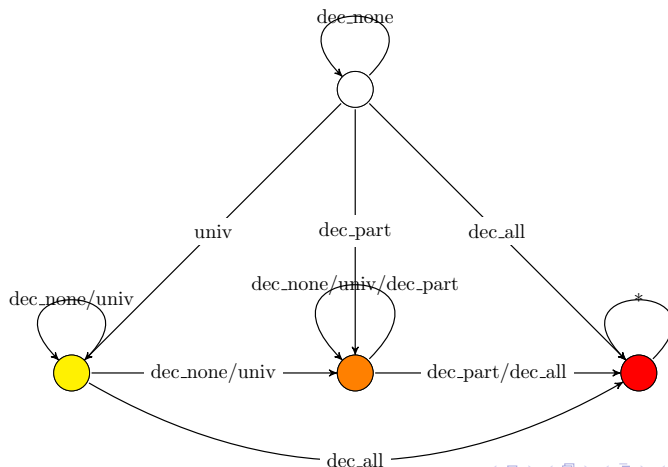*ronde*1 : $\quad \wedge \ \#RCV(co, 4\phi - 3) > n/2 \qquad (CPROP(\phi))$

*ronde*2 : $\quad \wedge \ \#\{p \in \Pi \,|\, co \in RCV(p, 4\phi - 2)\} > n/2 \ (BVOTE(\phi))$

Introduction
○○
○
○○○

*LastVoting*
○○
○○○○○○○○○
○○○○○○○○○○○○○○○○○○○○○○○

Simulator
○
○

Benchmarks
○
○○
○○○○○○○○○○○

Conclusion

Analysis of the Algorithm

# Termination

Termination holds iff $P^{sync}(\phi)$ eventually holds:

$P^{sync}(\phi) \triangleq \exists \phi > 0, \exists co \in \Pi :$
$ronde1 : \quad \wedge \ \#RCV(co, 4\phi - 3) > n/2 \qquad (CPROP(\phi))$
$ronde2 : \quad \wedge \ \#\{p \in \Pi \mid co \in RCV(p, 4\phi - 2)\} > n/2 \ (BVOTE(\phi))$
$ronde3 : \quad \wedge \ \#RCV(co, 4\phi - 1) > n/2 \qquad (CACK(\phi))$

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○●○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○○ | |

Analysis of the Algorithm

# Termination

Termination holds iff $P^{sync}(\phi)$ eventually holds:

$P^{sync}(\phi) \triangleq \exists \phi > 0, \exists co \in \Pi :$

$ronde1 : \quad \wedge \ \#RCV(co, 4\phi - 3) > n/2 \qquad (CPROP(\phi))$

$ronde2 : \quad \wedge \ \#\{p \in \Pi \mid co \in RCV(p, 4\phi - 2)\} > n/2 \ (BVOTE(\phi))$

$ronde3 : \quad \wedge \ \#RCV(co, 4\phi - 1) > n/2 \qquad (CACK(\phi))$

$ronde4 : \quad \wedge \ \forall p \in \Pi \setminus \mathfrak{D}(\phi) : co \in RCV(p, 4\phi) \quad (BDEC(\phi))$

Analysis of the Algorithm

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○●○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

Analysis of the Algorithm

# Execution-tree (graph. . . )

Analysis of the Algorithm

# Example of execution

| $\phi = 1$ | | |
|---|---|---|
| $\downarrow$ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $pr_1 = 0$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 0$ | $ts_2 = 0$ | $ts_3 = 0$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

$s_0$

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
| --- | --- | --- | --- | --- |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○●○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○○ | |

Analysis of the Algorithm

| $\phi = 1$ | | |
| --- | --- | --- |
| $\downarrow$ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 0$ | $ho_3 = 1$ |
| $pr_1 = 0$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 0$ | $ts_2 = 0$ | $ts_3 = 0$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

$s_0$

$CPROP$

$s_1$

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○●○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

Analysis of the Algorithm

| $\phi = 1$ | | |
|---|---|---|
| $\downarrow$ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 0$ | $ho_2 = 0$ | $ho_3 = 1$ |
| $pr_1 = 0$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 0$ | $ts_2 = 0$ | $ts_3 = 1$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

Analysis of the Algorithm



| $\phi = 2$ | | |
|:---:|:---:|:---:|
| | $\downarrow$ | |
| $p_1$ | $p_2$ | $p_3$ |
| $pr_1 = 0$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 0$ | $ts_2 = 0$ | $ts_3 = 1$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

Analysis of the Algorithm

| $\phi = 2$ | | |
|:---:|:---:|:---:|
| | $\downarrow$ | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 0$ | $ho_2 = 1$ | $ho_3 = 1$ |
| $pr_1 = 0$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 0$ | $ts_2 = 0$ | $ts_3 = 1$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

Analysis of the Algorithm

| $\phi = 2$ | | |
| :--: | :--: | :--: |
| | $\downarrow$ | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 0$ | $ho_3 = 0$ |
| $pr_1 = 1$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 0$ | $ts_3 = 1$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

| $\phi = 3$ | | |
|:---:|:---:|:---:|
| | | $\downarrow$ |
| $p_1$ | $p_2$ | $p_3$ |
| $pr_1 = 1$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 0$ | $ts_3 = 1$ |
| $d_1 =?$ | $d_2 =?$ | $d_3 =?$ |

| $\phi = 3$ | | |
|---|---|---|
| | | $\downarrow$ |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 1$ | $ho_3 = 0$ |
| $pr_1 = 1$ | $pr_2 = 0$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 0$ | $ts_3 = 1$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

$s_0$

$CPROP$

$s_1$

| $\phi = 3$ | | |
|---|---|---|
| | | $\downarrow$ |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 0$ | $ho_2 = 1$ | $ho_3 = 1$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 3$ | $ts_3 = 3$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

$s_0$

$CPROP$

$s_1$

$BVOTE$

$s_2$

| $\phi = 3$ | | |
|---|---|---|
| | | $\downarrow$ |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 0$ | $ho_3 = 1$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 3$ | $ts_3 = 3$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = ?$ |

$s_0$

$CPROP$

$s_1$

$BVOTE$

$s_2$

$CACK$

$s_3$

Analysis of the Algorithm



| $\phi = 3$ | | |
|---|---|---|
| | | $\downarrow$ |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 0$ | $ho_3 = 1$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 3$ | $ts_3 = 3$ |
| $d_1 =?$ | $d_2 =?$ | $d_3 = 1$ |

Analysis of the Algorithm

| $\phi = 4$ | | |
|---|---|---|
| $\downarrow$ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 2$ | $ts_2 = 3$ | $ts_3 = 3$ |
| $d_1 =?$ | $d_2 =?$ | $d_3 = 1$ |

| $\phi = 4$ | | |
|:---:|:---:|:---:|
| $\downarrow$ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 1$ | $ho_3 = 1$ |
| $pr_1 = 1$ $ts_1 = 2$ | $pr_2 = 1$ $ts_2 = 3$ | $pr_3 = 1$ $ts_3 = 3$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = 1$ |

| $\phi = 4$ | | |
|---|---|---|
| $\downarrow$ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 0$ | $ho_3 = 1$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 4$ | $ts_2 = 3$ | $ts_3 = 4$ |
| $d_1 =?$ | $d_2 =?$ | $d_3 = 1$ |



$s_0$

$CPROP$

$s_1$

$BVOTE$

$s_2$

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○●○○○○ | | ○○○○○○○○○○○○ | |

Analysis of the Algorithm

| $\phi = 4$ | | |
|---|---|---|
| ↓ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 0$ | $ho_3 = 1$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 4$ | $ts_2 = 3$ | $ts_3 = 4$ |
| $d_1 = ?$ | $d_2 = ?$ | $d_3 = 1$ |

Introduction     *LastVoting*     Simulator     Benchmarks     Conclusion
○○               ○○               ○            ○
○                ○○○○○○○○○        ○            ○○
○○○              ○○○○○○○○○○○○○○○○○○○●○○○○                      ○○○○○○○○○○○○

Analysis of the Algorithm

| $\phi = 4$ | | |
|---|---|---|
| ↓ | | |
| $p_1$ | $p_2$ | $p_3$ |
| $ho_1 = 1$ | $ho_2 = 1$ | $ho_3 = 0$ |
| $pr_1 = 1$ | $pr_2 = 1$ | $pr_3 = 1$ |
| $ts_1 = 4$ | $ts_2 = 3$ | $ts_3 = 4$ |
| $d_1 = 1$ | $d_2 = 1$ | $d_3 = 1$ |

$s_0$

$CPROP$

$s_1$

$BVOTE$

$s_2$

$CACK$

$s_3$

$BDEC$

$s_4$

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○●○○ | | ○○○○○○○○○○○ | |

Analysis of the Algorithm

# Comparison with the asynchronous version

- ▶ Each process has a infinite set of integers which is disjoint from the sets of the other processes
- ▶ This set is used to uniquely identify the consensus requests
- ▶ In the asynchronous version, the coordinator first broadcasts the request number of the new request
- ▶ The last broadcast is a flooding broadcast procedure

Introduction     *LastVoting*     Simulator     Benchmarks     Conclusion
○○               ○○              ○           ○
○                ○○○○○○○○○        ○           ○○
○○○              ○○○○○○○○○○○○○○○○○○○○○○○○○●○      ○○○○○○○○○○○○

Analysis of the Algorithm

# Computability of Paxos

- ▶ Paxos does not terminate [FLP85][Fuz08]
- ▶ Optimal condition not known for consensus algorithms

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○● | | ○○○○○○○○○○○ | |

Analysis of the Algorithm

## Fuzzati's proof of Paxos

A general pattern of rule:

$$(\text{RULE}) \ \frac{\text{Condition(s)}}{\text{Action(s) or Event(s)}}$$

Example of rule from Fuzzati's Paxos' rules:

$$(\text{CRASH}) \ \frac{S(i) = (a, r, p, b, (\top, \iota), \bot)}{\langle B, C, S \rangle \to \langle B, C, S\{i \mapsto (a, r, p, b, (\bot, \iota), \bot)\} \rangle}$$

Introduction

*LastVoting*

Simulator

Benchmarks

Motivations

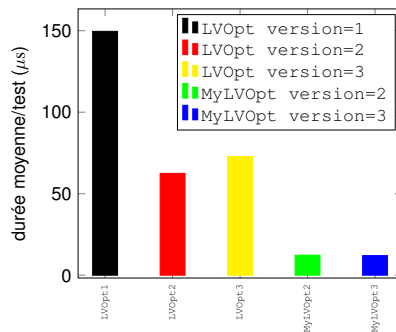## Motivations

- ► Correct Algorithms
- ► Benchmark Algorithms' performances

Benchmark

# Benchmark



**Banc d'essai pour $N = 3$**

**Banc d'essai pour $N = 4$**

Introduction

*LastVoting*

Simulator

Benchmarks

## Motivations

Benchmarking with Paxos

- ▶ Tsuchiya & Schiper:
  - ▶ SAT Solvers: up to 10 processes
  - ▶ NuSMV: up to 4 processes
  - ▶ SPIN: up to 3 processes

## Motivations

Benchmarking with Paxos

- Tsuchiya & Schiper:
    - SAT Solvers: up to 10 processes
    - NuSMV: up to 4 processes
    - SPIN: up to 3 processes
- Our work in progress:
    - POEM (Peter Niebert, Marseille):
      frontend for different tools with partial order techniques
      preprocessing
    - ALCOOL (CEA/List, Panda):
      Topological techniques

| Introduction | *LastVoting* | Simulator | **Benchmarks** | Conclusion |
|---|---|---|---|---|
| OO | OO | O | ● | |
| O | OOOOOOOOO | O | OO | |
| OOO | OOOOOOOOOOOOOOOOOOOOOOOO | | OOOOOOOOOOO | |

Various parameters

## Toward Benchmarking of Model-Checkers

- ▶ Number *n* of processes of the system

| Introduction | *LastVoting* | Simulator | **Benchmarks** | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ● | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○○○ | |

Various parameters

# Toward Benchmarking of Model-Checkers

- ▶ Number *n* of processes of the system
- ▶ Number of losses for the system:
  - ▶ for one round
  - ▶ for one phase
  - ▶ for *k* phases
  - ▶ for one execution
- ▶ Number of losses per process:
  - ▶ for one phase
  - ▶ for *k* phases
  - ▶ for one execution

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
|---|---|---|---|---|
| oo | oo | o | o | |
| o | ooooooooo | o | ●o | |
| ooo | ooooooooooooooooooooooooo | | ooooooooooo | |

POEM

## POEM

Checking is done with IF requests (first order logic):
(d[0]<>NOT_DEF) and (d[1]<>NOT_DEF) and (d[0]<>d[1]
Current problem: false positives (Bug in the frontend?)

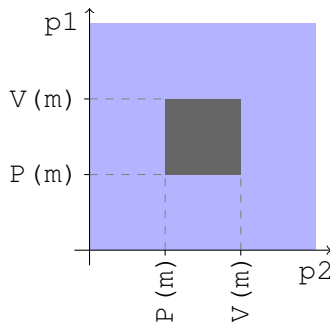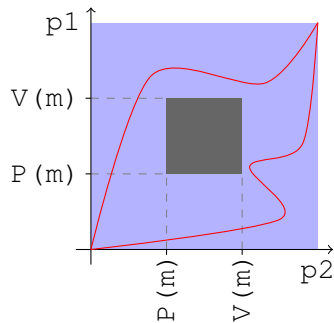Introduction          LastVoting          Simulator          **Benchmarks**          Conclusion
○○                    ○○                  ○                   ○
○                     ○○○○○○○○○           ○                   ○●
○○○                   ○○○○○○○○○○○○○○○○○○○○○○○○○○○                ○○○○○○○○○○○○○

POEM

# First Tests

| Introduction | *LastVoting* | Simulator | **Benchmarks** | Conclusion |
|---|---|---|---|---|
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○ | | ●○○○○○○○○○○○ | |

ALCOOL

## ALCOOL

- ▶ ALCOOL can detect deadlocks and forbidden zones
- ▶ We transform the properties which are to be verified into deadlock properties
- ▶ Example:
  $B_0$ and $B_1$ two n-ary synchronization barriers
  Agreement $\longrightarrow$ process $p$ waits on $B_{d_p}$
- ▶ Current problem: false positive (feature)

## Example 2

```
#mutex m
processes:
p1 = P(m).V(m)
p2 = P(m).V(m)
init:
p1 p2
```
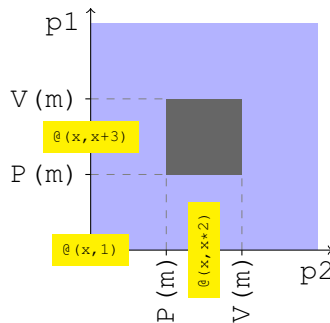
Introduction       *LastVoting*       Simulator       **Benchmarks**       Conclusion
○○                 ○○                  ○                ○
○                  ○○○○○○○○○           ○                ○○
○○○                ○○○○○○○○○○○○○○○○○○○○○○○○○○   ○        ○○●○○○○○○○○○○○

ALCOOL

| Introduction | *LastVoting* | Simulator | **Benchmarks** | Conclusion |
| :-- | :-- | :-- | :-- | :-- |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○●○○○○○○○ | |

ALCOOL

Basically only 2 executions:

- `p1.p2`
- `p2.p1`

- `p1.p2 :   x=8`
- `p2.p1 :   x=5`

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
| oo | oo | o | o | |
| o | ooooooooo | o | oo | |
| ooo | oooooooooooooooooooooooo | | ooooo●oooooooo | |

ALCOOL

## Example 3

```
#mutex a b
processes:
p1 = P(a).P(b).V(b).V(a)
p2 = P(b).P(a).V(a).V(b)
init:
p1 p2
```
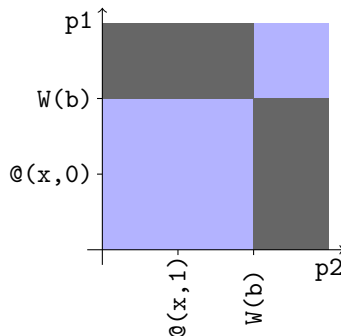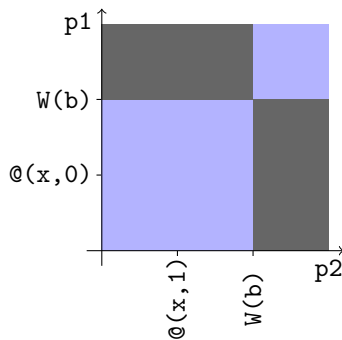
## Example 4

```
#synchronization 2 b
processes:
p1 = @(x,1).
(W(b) + [x==1] + void)
p2 = @(x,0).W(b)
init:
p1 p2
```
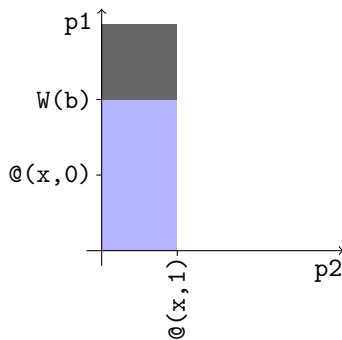
Introduction · · · · · · · · LastVoting · · Simulator · · Benchmarks · · Conclusion
○○ · · · · · · · · · · · · · · · · ○○ · · · · · · · · · ○ · · · · · · · · · ○
○ · · · · · · · · · · · · · · · · · ○○○○○○○○○ · · · ○ · · · · · · · · · ○○
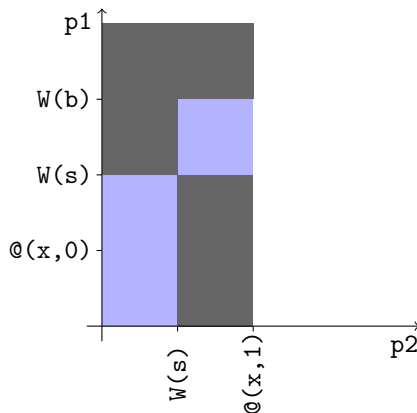○○○ · · · · · · · · · · · · · · · ○○○○○○○○○○○○○○○○○○○○○○○○ · · · ○○○○○○○●○○○○○

ALCOOL

## Example 4

```
#synchronization 2 b
processes:
p1 = @(x,1).W(b)
p2 = @(x,0).W(b)
init:
p1 p2
```

| Introduction | *LastVoting* | Simulator | Benchmarks | Conclusion |
| --- | --- | --- | --- | --- |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○●○○○○ | |

ALCOOL
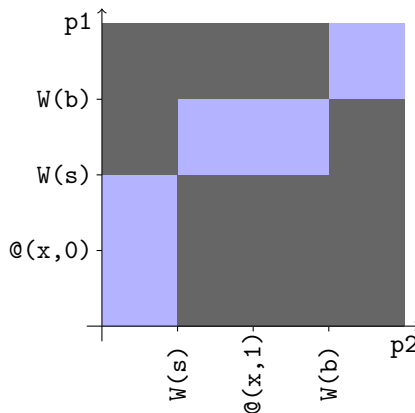
## Example 5

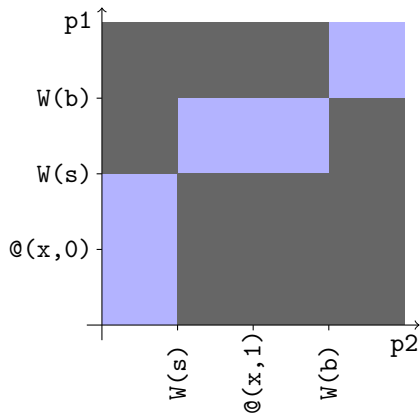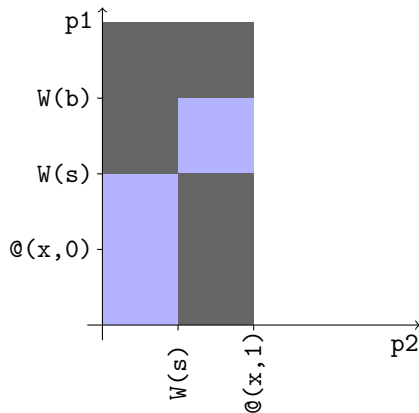```
#synchronization 2 b s
processes:
p1 = void.W(s).@(x,1).
(W(b) + [x==1] + void)
p2 = @(x,0).W(s).W(b)
init:
p1 p2
```

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
| --- | --- | --- | --- | --- |
| ○○ | ○○ | ○ | ○ | |
| ○ | ○○○○○○○○○ | ○ | ○○ | |
| ○○○ | ○○○○○○○○○○○○○○○○○○○○○○ | | ○○○○○○○○○●○○ | |

ALCOOL

## Example 5

```
#synchronization 2 b s
processes:
p1 = W(s).@(x,1).W(b)
p2 = @(x,0).W(s).W(b)
init:
p1 p2
```

## Benchmarks: ALCOOL

- ▶ Analysis of numerical variables
- ▶ Analysis of potential deadlocks
- ▶ ALCOOL +POEM?

| Introduction | LastVoting | Simulator | Benchmarks | Conclusion |
| oo | oo | o | o | |
| o | ooooooooo | o | oo | |
| ooo | ooooooooooooooooooooooo | | ooooooooooo | |

📄 BURROWS, MIKE: *The Chubby lock service for loosely-coupled distributed systems*.
In *Proceedings of the 7th symposium on Operating systems design and implementation*, OSDI '06, pages 335–350, Berkeley, CA, USA, 2006. USENIX Association.

📄 CHANDRA, TUSHAR D., ROBERT GRIESEMER and JOSHUA REDSTONE: *Paxos made live: an engineering perspective*.
In *Proceedings of the twenty-sixth annual ACM symposium on Principles of distributed computing*, PODC '07, pages 398–407, New York, NY, USA, 2007. ACM.

📄 FISCHER, MICHAEL J., NANCY A. LYNCH and MICHAEL S. PATERSON: *Impossibility of distributed consensus with one faulty process*.
J. ACM, 32:374–382, April 1985.

📄 FUZZATI, RACHELE: *A formal approach to fault tolerant distributed consensus*.
PhD thesis, École Polytechnique Fédérale de Lausanne, Lausanne, 2008.

📄 LAMPORT, LESLIE: *The part-time parliament*.
ACM Transactions on Computer Systems, 16:133–169, 1998.

📄 TSUCHIYA, TATSUHIRO and ANDRÉ SCHIPER: *Verification of consensus algorithms using satisfiability solving*.
Distributed Computing, 23:341–358, 2011.