# Resilience

Guillaume Pallez
Inria

**M2 CISD, Enseirb-Matmeca,**
Automne 2019

# Plan du cours d'aujourd'hui

In a hypothetical future, in a galaxy far far away, you decided to do a PhD.

In a hypothetical future, in a galaxy far far away, you decided to do a PhD. Even more hypothetical, you made it to the step where you are actually writing it!

In a hypothetical future, in a galaxy far far away, you decided to do a PhD. Even more hypothetical, you made it to the step where you are actually writing it!

▶ Thanks to the lack of budget in research, your laptop is very old: it crashes often.

▶ You chose to write your thesis using a WYSIWYG software which takes approx 3 minutes to save (while freezing your laptop obviously).. Too many figures.

What do you do?

In a hypothetical future, in a galaxy far far away, you decided to do a PhD. Even more hypothetical, you made it to the step where you are actually writing it!

▶ Thanks to the lack of budget in research, your laptop is very old: it crashes often.

▶ You chose to write your thesis using a WYSIWYG software which takes approx 3 minutes to save (while freezing your laptop obviously).. Too many figures.

What do you do?

Solution: Write your thesis in Latex, faster to save (but it's too late now)

In a hypothetical future, in a galaxy far far away, you decided to do a PhD. Even more hypothetical, you made it to the step where you are actually writing it!

▶ Thanks to the lack of budget in research, your laptop is very old: it crashes often.

▶ You chose to write your thesis using a WYSIWYG software which takes approx 3 minutes to save (while freezing your laptop obviously).. Too many figures.

What do you do?

▶ Solution 1: You choose to save your work every 3h.

    ▶ Mid-afternoon of Day 3, your laptop crashes. You have lost 1.5h of work.

In a hypothetical future, in a galaxy far far away, you decided to do a PhD. Even more hypothetical, you made it to the step where you are actually writing it!

▶ Thanks to the lack of budget in research, your laptop is very old: it crashes often.

▶ You chose to write your thesis using a WYSIWYG software which takes approx 3 minutes to save (while freezing your laptop obviously).. Too many figures.

What do you do?

▶ Solution 1: You choose to save your work every 3h.

    ▶ Mid-afternoon of Day 3, your laptop crashes. You have lost 1.5h of work.

▶ Solution 2: You choose to save your work every half-hour.

    ▶ No crash during the next three consecutive days.

Which solution is best?

# EXASCALE PLATFORMS (COURTESY JACK DONGARRA)

## Potential System Architecture
## with a cap of $200M and 20MW

| Systems | 2011<br>K computer | 2019 | Difference<br>Today & 2019 |
|---|---|---|---|
| **System peak** | 10.5 Pflop/s | 1 Eflop/s | O(100) |
| **Power** | 12.7 MW | ~20 MW | |
| System memory | 1.6 PB | 32 - 64 PB | O(10) |
| Node performance | 128 GF | 1,2 or 15TF | O(10) – O(100) |
| Node memory BW | 64 GB/s | 2 - 4TB/s | O(100) |
| Node concurrency | 8 | O(1k) or 10k | O(100) – O(1000) |
| Total Node Interconnect BW | 20 GB/s | 200-400GB/s | O(10) |
| System size (nodes) | 88,124 | O(100,000) or O(1M) | O(10) – O(100) |
| Total concurrency | 705,024 | O(billion) | O(1,000) |
| MTTI | days | O(1 day) | - O(10) |

## Toward Exascale Computing (My Roadmap)

*Based on proposed DOE roadmap with MTTI adjusted to scale linearly*

| Systems | 2009 | 2011 | 2015 | 2018 |
|---|---|---|---|---|
| System peak | 2 Peta | 20 Peta | 100-200 Peta | 1 Exa |
| System memory | 0.3 PB | 1.6 PB | 5 PB | 10 PB |
| Node performance | 125 GF | 200GF | 200-400 GF | 1-10TF |
| Node memory BW | 25 GB/s | 40 GB/s | 100 GB/s | 200-400 GB/s |
| Node concurrency | 12 | 32 | O(100) | O(1000) |
| Interconnect BW | 1.5 GB/s | 22 GB/s | 25 GB/s | 50 GB/s |
| System size (nodes) | 18,700 | 100,000 | 500,000 | O(million) |
| Total concurrency | 225,000 | 3,200,000 | O(50,000,000) | O(billion) |
| Storage | 15 PB | 30 PB | 150 PB | 300 PB |
| IO | 0.2 TB/s | 2 TB/s | 10 TB/s | 20 TB/s |
| MTTI | 4 days | 19 h 4 min | 3 h 52 min | 1 h 56 min |
| Power | 6 MW | ~10MW | ~10 MW | ~20 MW |

▶ Hierarchical
- $10^5$ or $10^6$ nodes
- Each node equipped with $10^4$ or $10^3$ cores

▶ Failure-prone

| MTBF – one node | 1 year | 10 years | 120 years |
|---|---|---|---|
| MTBF – platform of $10^6$ nodes | 30sec | 5mn | 1h |

More nodes $\Rightarrow$ Shorter MTBF (Mean Time Between Failures)

▶ Hierarchical
- $10^5$ or $10^6$ nodes
- Each node equipped with ~~~ or $10^3$ cores

▶ Failure-prone

| MTBF – one node | 1 year | 10 years | 120 years |
|---|---|---|---|
| MTBF – platform of $10^6$ nodes | 30sec | 5mn | 1h |

**Exascale** $\neq$ **Petascale** $\times 1000$

More no  (Failures)

## Classic approach for FT: Checkpoint-Restart

Typical "Balanced Architecture" for PetaScale Computers



Balanced System Approach

RoadRunner

Compute nodes

Total memory: 100-200 TB

40 to 200 GB/s

Parallel file system (1 to 2 PB)

Network(s)

I/O nodes

Without optimization, Checkpoint-Restart needs about 1h! (~30 minutes each)

| Systems | Perf. | Ckpt time | Source |
|---------|-------|-----------|--------|
| RoadRunner | 1PF | ~20 min. | Panasas |
| LLNL BG/L | 500 TF | >20 min. | LLNL |
| LLNL Zeus | 11TF | 26 min. | LLNL |
| YYY BG/P | 100 TF | ~30 min. | YYY |

LLNL BG/L

## Sources of failures

- Analysis of error and failure logs

- In 2005 (Ph. D. of CHARNG-DA LU) : "Software halts account for the most number of outages (59-84 percent), and take the shortest time to repair (0.6-1.5 hours). Hardware problems, albeit rarer, need 6.3-100.7 hours to solve."

- In 2007 (Garth Gibson, ICPP Keynote):



- In 2008 (Oliner and J. Stearley, DSN Conf.):

| | Raw | | Filtered | |
|---|---|---|---|---|
| Type | Count | % | Count | % |
| Hardware | 174,586,516 | 98.04 | 1,999 | 18.78 |
| Software | 144,899 | 0.08 | 6,814 | 64.01 |
| Indeterminate | 3,350,044 | 1.88 | 1,832 | 17.21 |

Relative frequency of root cause by system type.

Software errors: Applications, OS bug (kernel panic), communication libs, File system error and other.

Hardware errors, Disks, processors, memory, network

Conclusion: Both Hardware and Software failures have to be considered

10

- ▶ Many types of faults: software error, hardware malfunction, memory corruption
- ▶ Many possible behaviors: silent, transient, unrecoverable
- ▶ Restrict to faults that lead to application failures
- ▶ This includes all hardware faults, and some software ones
- ▶ Will use terms *fault* and *failure* interchangeably

▶ Many types of faults: software error, hardware malfunction, memory corruption

▶ Many possible behaviors: silent, transient, unrecoverable

▶ Restrict to faults that lead to application failures

▶ This includes all hardware faults, and some software ones

▶ Will use terms *fault* and *failure* interchangeably

▶ First question: quantify the rate or frequency at which these faults strike!

# FAILURE DISTRIBUTIONS: (1) EXPONENTIAL



$Exp(\lambda)$: Exponential distribution law of parameter $\lambda$:

▶ Probability density function (pdf): $f(t) = \lambda e^{-\lambda t} dt$ for $t \geq 0$

▶ Cumulative distribution function (cdf): $F(t) = 1 - e^{-\lambda t}$

▶ Mean: $\mu = \frac{1}{\lambda}$

Sequential Machine

$X$ random variable for $Exp(\lambda)$ failure inter-arrival times:

▶ $\mathbb{P}(X \leq t) = 1 - e^{-\lambda t} dt$ (by definition)

▶ Memoryless property: $\mathbb{P}(X \geq t + s \mid X \geq s) = \mathbb{P}(X \geq t)$
   (for all $t, s \geq 0$): at any instant, time to next failure does not depend upon time
   elapsed since last failure

▶ Mean Time Between Failures (MTBF) $\mu = \mathbb{E}(X) = \frac{1}{\lambda}$

Sequential Machine

*Weibull*($k, \lambda$): Weibull distribution law of shape parameter $k$ and scale parameter $\lambda$:

▶ Pdf: $f(t) = k\lambda(t\lambda)^{k-1}e^{-(\lambda t)^k}dt$ for $t \geq 0$

▶ Cdf: $F(t) = 1 - e^{-(\lambda t)^k}$

▶ Mean: $\mu = \frac{1}{\lambda}\Gamma(1 + \frac{1}{k})$

Sequential Machine

$X$ random variable for $Weibull(k, \lambda)$ failure inter-arrival times:

▶ If $k < 1$: failure rate decreases with time
  "infant mortality": defective items fail early

▶ If $k = 1$: $Weibull(1, \lambda) = Exp(\lambda)$ constant failure time

▶ Processor (or node): any entity subject to failures
  ⇒ approach agnostic to granularity

▶ If the MTBF is $\mu$ with one processor,
  what is its value with $p$ processors?

▶ Processor (or node): any entity subject to failures
⇒ approach agnostic to granularity

▶ If the MTBF is $\mu$ with one processor,
what is its value with $p$ processors?

▶ Well, it depends ☺

- ▶ Rebooting all $p$ processors after a failure
- ▶ Platform failure distribution
  $\Rightarrow$ minimum of $p$ IID processor distributions
- ▶ With $p$ distributions $Exp(\lambda)$:

$$\min\left(Exp(\lambda_1), Exp(\lambda_2)\right) = Exp(\lambda_1 + \lambda_2)$$

$$\mu = \frac{1}{\lambda} \Rightarrow \mu_p = \frac{\mu}{p}$$

- ▶ With $p$ distributions $Weibull(k, \lambda)$:

$$\min_{1..p}\left(Weibull(k, \lambda)\right) = Weibull(k, p^{1/k}\lambda)$$

$$\mu = \frac{1}{\lambda}\Gamma(1 + \frac{1}{k}) \Rightarrow \mu_p = \frac{\mu}{p^{1/k}}$$

▶ Rebooting only faulty processor
▶ Platform failure distribution
  ⇒ superposition of $p$ IID processor distributions
  ⇒ IID only for Exponential
▶ Define $\mu_p$ by

$$\lim_{F \to +\infty} \frac{n(F)}{F} = \frac{1}{\mu_p}$$

$n(F)$ = number of platform failures until time $F$ is exceeded

**Theorem:** $\mu_p = \dfrac{\mu}{p}$ for arbitrary distributions

If three processors have around 20 faults during a time $t$ $(\mu = \frac{t}{20})$...



...during the same time, the platform has around 60 faults $(\mu_p = \frac{t}{60})$

**Theorem:** $\mu_p = \frac{\mu}{p}$ for arbitrary distributions

**With one processor:**

- $n(F) =$ number of failures until time $F$ is exceeded

- $X_i$ iid random variables for inter-arrival times, with $\mathbb{E}(X_i) = \mu$

- $\sum_{i=1}^{n(F)-1} X_i \leq F \leq \sum_{i=1}^{n(F)} X_i$

- Wald's equation:
  $(\mathbb{E}(n(F)) - 1)\mu \leq F \leq \mathbb{E}(n(F))\,\mu$

- $\lim_{F \to +\infty} \frac{\mathbb{E}(n(F))}{F} = \frac{1}{\mu}$

**Theorem:** $\mu_p = \frac{\mu}{p}$ for arbitrary distributions

**With one processor:**

- $n(F)$ = number of failures until time $F$ is exceeded

- $X_i$ iid random variables for inter-arrival times, with $\mathbb{E}(X_i) = \mu$

- $\sum_{i=1}^{n(F)-1} X_i \leq F \leq \sum_{i=1}^{n(F)} X_i$

- Wald's equation:
  $(\mathbb{E}(n(F)) - 1)\mu \leq F \leq \mathbb{E}(n(F))\mu$

- $\lim_{F \to +\infty} \frac{\mathbb{E}(n(F))}{F} = \frac{1}{\mu}$

**With *p* processors:**

- $n(F)$ = number of platform failures until time $F$ is exceeded

- $n_q(F)$ = number of those failures that strike processor $q$

- $n_q(F) + 1$ = number of failures on processor $q$ until time $F$ is exceeded (except for processor with last-failure)

- $\lim_{F \to +\infty} \frac{n_q(F)}{F} = \frac{1}{\mu}$ as above

- $\lim_{F \to +\infty} \frac{n(F)}{F} = \frac{1}{\mu_p}$ by definition

- Hence $\mu_p = \frac{\mu}{p}$ because
  $n(F) = \sum_{q=1}^{p} n_q(F)$

- MTBF of one processor: between 1 and 125 years
- Shape parameters for Weibull: $k = 0.5$ or $k = 0.7$
- Failure trace archive from INRIA
  (`http://fta.inria.fr`)
- Computer Failure Data Repository from LANL
  (`http://institutes.lanl.gov/data/fdata`)

Parallel machine ($10^6$ nodes)

After infant mortality and before aging,
instantaneous failure rate of computer platforms is almost constant

▶ MTBF key parameter and $\mu_p = \frac{\mu}{p}$ ☺
▶ Exponential distribution OK for most purposes ☺
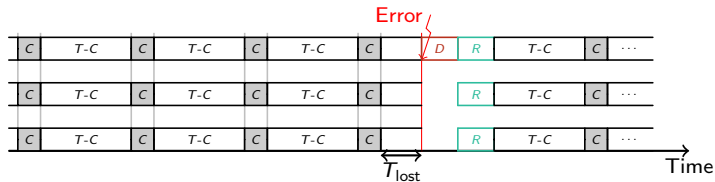▶ Assume failure independence while not (completely) true ☹

- Periodic checkpointing policy of period $T$
- Time to checkpoint $C$
- Time lost in case of a failure $T_{\text{lost}}$

- Downtime $D$
- Time for recovery $R$

- Periodic checkpointing policy of period $T$
- Time to checkpoint $C$
- Time lost in case of a failure $T_{\text{lost}}$

- Downtime $D$
- Time for recovery $R$

# PERIODIC CHECKPOINTING, DEFINITIONS



- Periodic checkpointing policy of period $T$
- Time to checkpoint $C$
- Time lost in case of a failure $T_{\text{lost}}$

- Downtime $D$
- Time for recovery $R$

- ▶ Periodic checkpointing policy of period $T$
- ▶ Time to checkpoint $C$
- ▶ Time lost in case of a failure $T_{\text{lost}}$

- ▶ Downtime $D$
- ▶ Time for recovery $R$

- Periodic checkpointing policy of period $T$
- Time to checkpoint $C$
- Time lost in case of a failure $T_{lost}$

- Downtime $D$
- Time for recovery $R$

## Strategies

1. Only one checkpoint at the end of the execution;

2. Three checkpoints during the execution, after every 10 minutes of work;

3. Five checkpoints during the execution, after every 6 minutes of work.

## Scenarios

(A) A large time between faults (in this example, no fault during the execution);

(B) A medium time between faults (only one fault at the 19th minute);

(C) A small time between faults (one fault at the 19th, 42th , and 62th minutes.).

Strategy 1

Strategy 2

Strategy 3

Time

Large MTBF: there are no or very few faults. Checkpointing is too expensive. The first strategy wins.

Strategy 1 | Strategy 2 | Strategy 3

Large MTBF: there are no or very few faults. Checkpointing is too expensive. The first strategy wins.
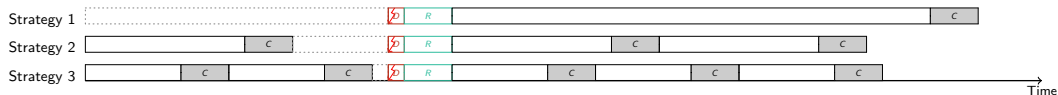


Strategy 1 | Strategy 2 | Strategy 3

Medium MTBF: there are more faults. It is good to checkpoint, but not too frequently, because of the corresponding overhead. The second strategy wins.
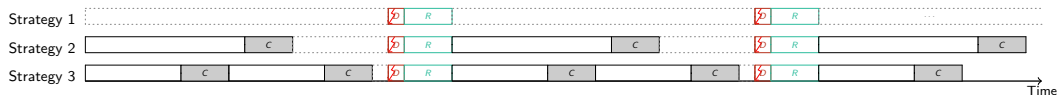
Large MTBF: there are no or very few faults. Checkpointing is too expensive. The first strategy wins.



Medium MTBF: there are more faults. It is good to checkpoint, but not too frequently, because of the corresponding overhead. The second strategy wins.



Small MTBF: there are many faults. The cost of the checkpoints is paid off because the time lost due to faults is dramatically reduced. The third strategy wins.

> **Waste:** Fraction of time not spent for useful computations. If an application needs $\text{TIME}_{\text{base}}$ volume of compute, and the final execution time is $\text{TIME}_{\text{final}}$:
>
> $$\text{WASTE} = \frac{\text{TIME}_{\text{Final}} - \text{TIME}_{\text{base}}}{\text{TIME}_{\text{Final}}}$$

Equivalent to minimizing $\text{TIME}_{\text{Final}}$: $(1 - \text{WASTE})\text{TIME}_{\text{Final}} = \text{TIME}_{\text{base}}$, but more convenient (get rid of notion of Time, and end of computation).

An execution. Black intervals correspond to work destroyed by faults, downtimes, and recoveries.

$\text{TIME}_{\mathsf{Final}}$
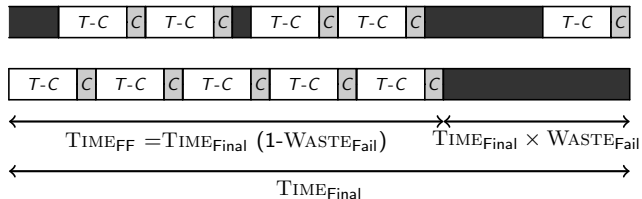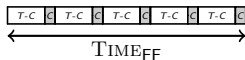
An execution. Black intervals correspond to work destroyed by faults, downtimes, and recoveries.

An execution. Black intervals correspond to work destroyed by faults, downtimes, and recoveries.

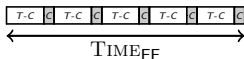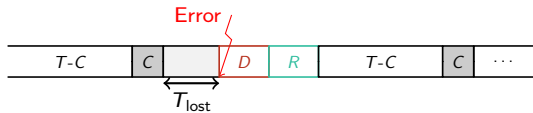| T-C | c | T-C | c | T-C | c | T-C | c | T-C | c |

$\longleftrightarrow$ TIME$_\text{FF}$

- TIME$_\text{base}$: application base time
- TIME$_\text{FF}$: with periodic checkpoints but failure-free

$$\text{TIME}_\text{FF} = \text{TIME}_\text{base} + \#\textit{checkpoints} \times C$$

| $T$-$C$ | $c$ | $T$-$C$ | $c$ | $T$-$C$ | $c$ | $T$-$C$ | $c$ | $T$-$C$ | $c$ |

$\longleftrightarrow$ $\text{TIME}_{\textsf{FF}}$

- ▶ $\text{TIME}_{\textsf{base}}$: application base time
- ▶ $\text{TIME}_{\textsf{FF}}$: with periodic checkpoints but failure-free

$$\text{TIME}_{\textsf{FF}} = \text{TIME}_{\textsf{base}} + \#checkpoints \times C$$

$$\#checkpoints = \left\lceil \frac{\text{TIME}_{\textsf{base}}}{T - C} \right\rceil \approx \frac{\text{TIME}_{\textsf{base}}}{T - C} \text{ (valid for large jobs)}$$
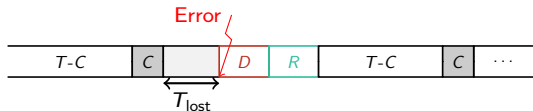
$$\text{WASTE}_{\textsf{FF}} = \frac{\text{TIME}_{\textsf{FF}} - \text{TIME}_{\textsf{base}}}{\text{TIME}_{\textsf{FF}}} = \frac{C}{T}$$

$$\mathrm{TIME_{Final}} = \mathrm{TIME_{FF}} + N_{\mathsf{faults}} \left( T_{\mathsf{lost}} + D + R \right)$$
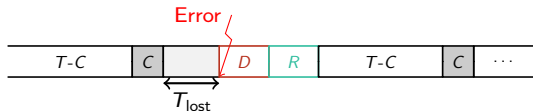
$$\text{TIME}_{\text{Final}} = \text{TIME}_{\text{FF}} + N_{\text{faults}}\left(T_{\text{lost}} + D + R\right)$$

$$\text{TIME}_{\text{Final}} = \text{TIME}_{\text{FF}} + \frac{\text{TIME}_{\text{Final}}}{\mu}\left(T/2 + D + R\right)$$

$\Rightarrow$ Instants when periods begin and failures strike are independent

$$\text{TIME}_{\textsf{Final}} = \text{TIME}_{\textsf{FF}} + N_{\textsf{faults}} \left( T_{\textsf{lost}} + D + R \right)$$

$$\text{TIME}_{\textsf{Final}} = \text{TIME}_{\textsf{FF}} + \frac{\text{TIME}_{\textsf{Final}}}{\mu} \left( T/2 + D + R \right)$$

$$\text{WASTE}_{\textsf{Fail}} = \frac{\text{TIME}_{\textsf{Final}} - \text{TIME}_{\textsf{FF}}}{\text{TIME}_{\textsf{Final}}} = \frac{1}{\mu} \left( T/2 + D + R \right)$$

- TIME$_{\mathsf{base}}$: application base time
- TIME$_{\mathsf{FF}}$: with periodic checkpoints but failure-free
- TIME$_{\mathsf{Final}}$: final time

$$\mathrm{WASTE} = \frac{\mathrm{TIME_{Final}} - \mathrm{TIME_{base}}}{\mathrm{TIME_{Final}}}$$

$$1 - \mathrm{WASTE} = (1 - \mathrm{WASTE_{FF}})(1 - \mathrm{WASTE_{Fail}})$$

$$\mathrm{WASTE} = \frac{C}{T} + \left(1 - \frac{C}{T}\right)\frac{1}{\mu}\left(D + R + \frac{T}{2}\right)$$

WASTE is minimized for

$$T = \sqrt{2\left(\mu - (D + R)\right)C}$$

▶ Capping periods, and enforcing a lower bound on MTBF
⇒ mandatory for mathematical rigor ☹

▶ Not needed for practical purposes ☺
  • actual job execution uses optimal value
  • account for multiple faults by re-executing work until success

▶ Approach surprisingly robust ☺



FINDING THE OPTIMAL CHECKPOINT INTERVAL

WASTE

CAN'T CHECKPOINT THAT OFTEN

OPTIMAL

SPEND TOO LONG CHECK-POINTING
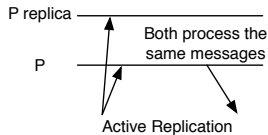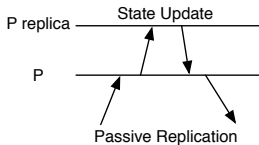
LOSE TOO MUCH COMPUTATION BECAUSE OF FAILURES

CHECKPOINT INTERVAL

- ▶ Back to your thesis
- ▶ Saving periodically is a pain: you get interrupted in important paragraphs and then lose your train of thoughts.
- ▶ You would rather save at the end of sections.
- ▶ Let's assume you know how many sections you are going to write, and what their sizes are going to be.

> ▶ Propose a model for this new problem
> ▶ Solve it!

## Replication



P replica ——————— State Update ———————

P ———————————————

Passive Replication

P replica ———————

Both process the same messages
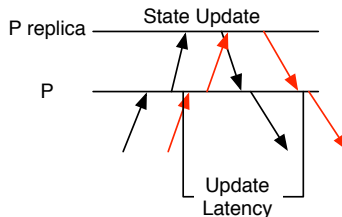
P ———————————————

Active Replication

### Idea

- Each process is replicated on a resource that has small chance to be hit by the same failure as its replica
- In case of failure, one of the replicas will continue working, while the other recovers
- Passive Replication / Active Replication

PASSIVE REPLICATION

## Replication

P replica ——— State Update

P

Update
Latency

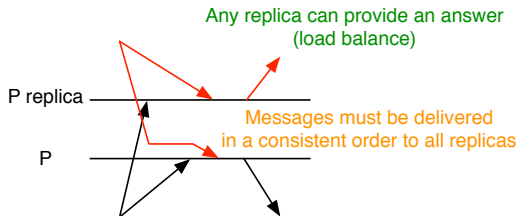### Challenges

- Passive replication: latency of state update
- Active replication: ordering of decision → internal additional communications
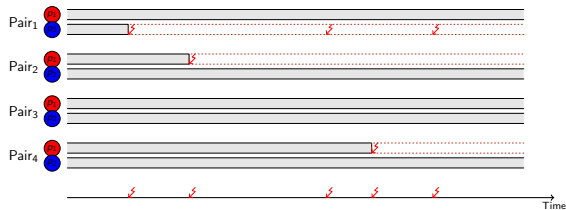
Any replica can provide an answer
(load balance)

P replica

Messages must be delivered
in a consistent order to all replicas

P

### Challenges

- Passive replication: latency of state update
- Active replication: ordering of decision $\rightarrow$ internal additional communications
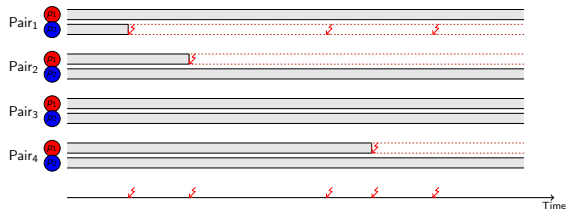
Processor pairs for replication: each blue processor is paired with a red processor and they do the same work.
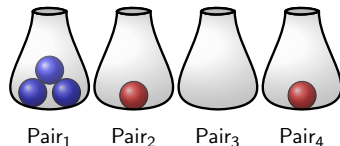
- ▶ How do you write the job model?
- ▶ How do you write the objective function?
- ▶ How do you compare to the checkpoint strategy?

Processor pairs for replication: each blue processor is paired with a red processor and they do the same work.



$Pair_1$    $Pair_2$    $Pair_3$    $Pair_4$

Modeling the state of the platform as a balls-into-bins problem. Colors of balls are important: $\neq$ birthday problem!

- ▶ How do you write the job model?
- ▶ How do you write the objective function?
- ▶ How do you compare to the checkpoint strategy?

This was not a class about resilience but a class about scheduling ☺. If you are interested, still need to read about:

▶ Technical Protocols for resilience (dealing with messages etc)

▶ Different type of checkpointing (blocking v Asynchronous, coordinated v uncoordinated, hierarchical, in memory etc)

▶ Combining replication and checkpointing

▶ ABFT

▶ ...

You can see the Tutorial by Bosilca, Bouteiller, Hérault and Robert:
http://fault-tolerance.org/2018/11/09/sc18/