

# Identifying the Right Replication Level to Detect and Correct Silent Errors at Scale

Anne Benoit<sup>1</sup>   Aurélien Cavelan<sup>1</sup>   Franck Cappello<sup>2</sup>  
Padma Raghavan<sup>4</sup>   Yves Robert<sup>1,3</sup>   Hongyang Sun<sup>4</sup>

<sup>1</sup>ENS Lyon & INRIA, France

<sup>2</sup>Argonne National Laboratory, USA

<sup>3</sup>University of Tennessee Knoxville, USA

<sup>4</sup>Vanderbilt University, USA

[hongyang.sun@vanderbilt.edu](mailto:hongyang.sun@vanderbilt.edu)

Fault Tolerance for HPC at eXtreme Scale (FTXS) Workshop  
June 26, 2017@Washington, D.C., USA

# An Inconvenient Truth

Top ranked supercomputers in the US (June 2017)

Rank	Name	Laboratory	Technology	Cores	PFlops/s	MTBF
4	Titan	ORNL	Cray XK7	560,640	17.59	≈ 1 day
5	Sequoia	LLNL	BG/Q	1,572,864	17.17	≈ 1 day
6	Cori	LBNL	Cray XC40	622,336	14.01	≈ 1 day
9	Mira	ANL	BG/Q	786,432	8.59	≈ 1 day

**Fail-stop errors:** Node failure, resource crashes

**Silent errors or silent data corruptions (SDCs):** Double bit flips, soft faults

# An Inconvenient Truth

Top ranked supercomputers in the US (June 2017)

Rank	Name	Laboratory	Technology	Cores	PFlops/s	MTBF
4	Titan	ORNL	Cray XK7	560,640	17.59	≈ 1 day
5	Sequoia	LLNL	BG/Q	1,572,864	17.17	≈ 1 day
6	Cori	LBNL	Cray XC40	622,336	14.01	≈ 1 day
9	Mira	ANL	BG/Q	786,432	8.59	≈ 1 day

**Fail-stop errors:** Node failure, resource crashes

**Silent errors or silent data corruptions (SDCs):** Double bit flips, soft faults

**Exascale computing** (1000 PFlops/s):

- ▶ Larger core count: millions or even billions of cores
- ▶ Shorter Mean Time Between Failures (MTBF)  $\mu$

# An Inconvenient Truth

Top ranked supercomputers in the US (June 2017)

Rank	Name	Laboratory	Technology	Cores	PFlops/s	MTBF
4	Titan	ORNL	Cray XK7	560,640	17.59	≈ 1 day
5	Sequoia	LLNL	BG/Q	1,572,864	17.17	≈ 1 day
6	Cori	LBNL	Cray XC40	622,336	14.01	≈ 1 day
9	Mira	ANL	BG/Q	786,432	8.59	≈ 1 day

**Fail-stop errors:** Node failure, resource crashes

**Silent errors or silent data corruptions (SDCs):** Double bit flips, soft faults

**Exascale computing** (1000 PFlops/s):

- ▶ Larger core count: millions or even billions of cores
- ▶ Shorter Mean Time Between Failures (MTBF)  $\mu$

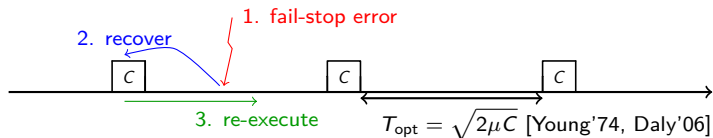
**Coping with faults:**

- ▶ Build more reliable hardware!
- ▶ Make applications more fault tolerant!
- ▶ Design better resilience techniques/algorithms!

# Resilience Techniques for HPC

Fail-stop errors (instantaneous error detection)

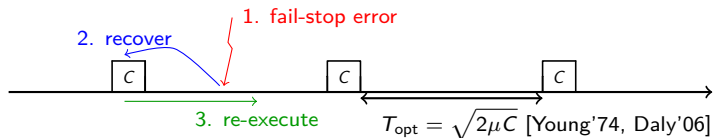
Standard approach: periodic checkpointing, rollback and recovery



# Resilience Techniques for HPC

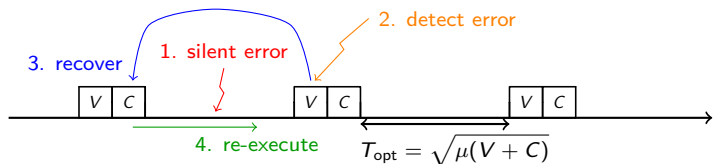
Fail-stop errors (instantaneous error detection)

Standard approach: periodic checkpointing, rollback and recovery



Silent errors (arbitrary detection latency)

Promising approach: checkpointing + verification (error detection)



# Approaches for Detecting Silent Errors

## Application-specific approaches

- ▶ Algorithm-based fault tolerance (ABFT): checksums in dense matrices, limited to one error detection and/or correction in practice [*Huang and Abraham 1984*]
- ▶ Partial differential equations (PDE): use lower-order scheme as verification mechanism [*Benson, Schmit and Schreiber 2014*]
- ▶ Generalized minimal residual method (GMRES): inner-outer iterations [*Hoemmen and Heroux 2011*]
- ▶ Preconditioned conjugate gradients (PCG): orthogonalization check iteratively, re-orthogonalization if error detected [*Sao and Vuduc 2013, Chen 2013*]

## Data-analytics/machine learning approaches

- ▶ Dynamic monitoring of datasets based on physical laws (e.g., temperature/speed limit) and space or temporal proximity [*Bautista-Gomez and Cappello 2014*]
- ▶ Time-series prediction, spatial multivariate interpolation [*Di et al. 2014*]
- ▶ Offline training, online detection based on SDC signature for convergent iterative applications [*Liu and Agrawal 2016*]
- ▶ Spatial regression based on support vector machines [*Subasi et al. 2016*]

## General-purpose approaches

- ▶ Process replication [*Fiala et al. 2012*]
- ▶ Group replication [*Casanova et al. 2014*]
- ▶ Triple modular redundancy (TMR) and voting [*Lyons and Vanderkulk 1962*]

## Focus:

Analytical model for applying replication/redundancy (general purpose approaches) in combination with checkpointing to detect and correct silent errors for HPC!

## Question:

How to *optimally* execute a parallel job obeying Amdahl's law on an error-prone platform?

What is the optimal *error-aware speedup*?



# When Amdahl Meets Young/Daly

*Error-free speedup* with  $P$  processors and  $\alpha$  sequential fraction:

$$\text{Amdahl's Law: } S(P) = \frac{1}{\alpha + \frac{1-\alpha}{P}}$$

- ▶ Bounded above by  $1/\alpha$
- ▶ Strictly increasing function of  $P$

# When Amdahl Meets Young/Daly

*Error-free speedup* with  $P$  processors and  $\alpha$  sequential fraction:

$$\text{Amdahl's Law: } S(P) = \frac{1}{\alpha + \frac{1-\alpha}{P}}$$

- ▶ Bounded above by  $1/\alpha$
- ▶ Strictly increasing function of  $P$

Allocating more processors on an error-prone platform?

- ▶ Higher error-free speedup 😊
- ▶ More errors/faults ☹️
  - ▶ More frequent checkpointing ☹️
    - ▶ More resilience overhead ☹️

# When Amdahl Meets Young/Daly

*Error-free speedup* with  $P$  processors and  $\alpha$  sequential fraction:

$$\text{Amdahl's Law: } S(P) = \frac{1}{\alpha + \frac{1-\alpha}{P}}$$

- ▶ Bounded above by  $1/\alpha$
- ▶ Strictly increasing function of  $P$

Allocating more processors on an error-prone platform?

- ▶ Higher error-free speedup 😊
- ▶ More errors/faults ☹️
  - ▶ More frequent checkpointing ☹️
    - ▶ More resilience overhead ☹️

Optimal processor allocation and checkpointing interval?

# How Is Replication Used?

On a  $Q$ -processor platform, application is replicated  $n$  times:

- ▶ **Duplication:** each replica has  $P = Q/2$  processors
- ▶ **Triplication:** each replica has  $P = Q/3$  processors
- ▶ **General case:** each replica has  $P = Q/n$  processors

# How Is Replication Used?

On a  $Q$ -processor platform, application is replicated  $n$  times:

- ▶ **Duplication:** each replica has  $P = Q/2$  processors
- ▶ **Triplication:** each replica has  $P = Q/3$  processors
- ▶ **General case:** each replica has  $P = Q/n$  processors

Having more replicas on an error-prone platform?

- ▶ Lower error-free speedup 😞
- ▶ More resilient 😊
  - ▶ Smaller checkpointing frequency 😊
    - ▶ Less resilience overhead 😊

# How Is Replication Used?

On a  $Q$ -processor platform, application is replicated  $n$  times:

- ▶ **Duplication:** each replica has  $P = Q/2$  processors
- ▶ **Triplcation:** each replica has  $P = Q/3$  processors
- ▶ **General case:** each replica has  $P = Q/n$  processors

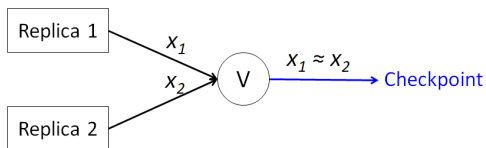
Having more replicas on an error-prone platform?

- ▶ Lower error-free speedup 😞
- ▶ More resilient 😊
  - ▶ Smaller checkpointing frequency 😊
    - ▶ Less resilience overhead 😊

Optimal replication level, processor allocation per replica  
and checkpointing interval?

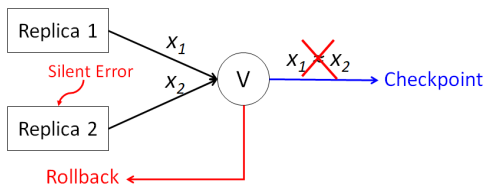
# Why Is Replication Useful?

► **Error detection (duplication):**



# Why Is Replication Useful?

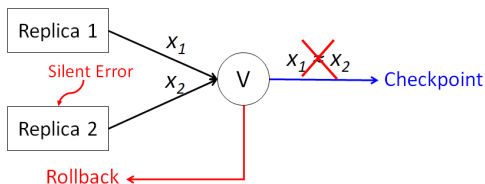
► **Error detection (duplication):**



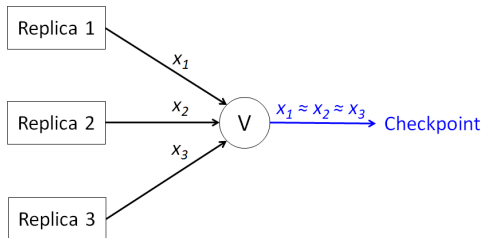


# Why Is Replication Useful?

## ▶ Error detection (duplication):

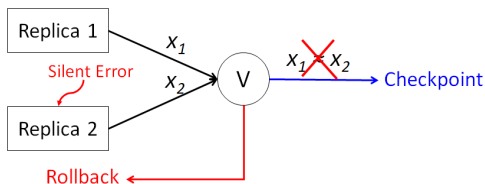


## ▶ Error correction (triplication):

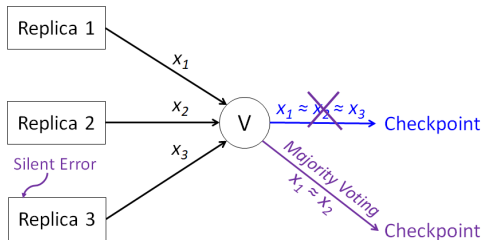


# Why Is Replication Useful?

## ▶ Error detection (duplication):

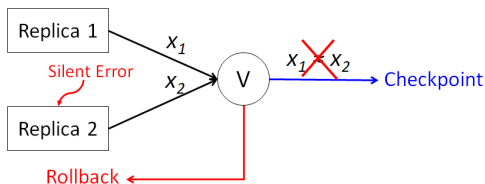


## ▶ Error correction (triplication):

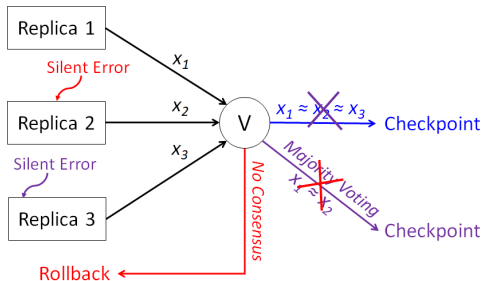


# Why Is Replication Useful?

## ▶ Error detection (duplication):

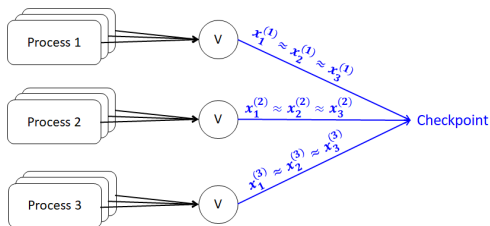


## ▶ Error correction (triplication):

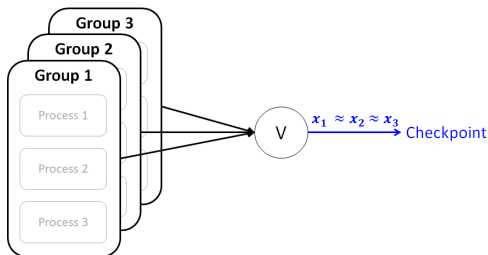


# Two Replication Modes

## ▶ Process Replication:

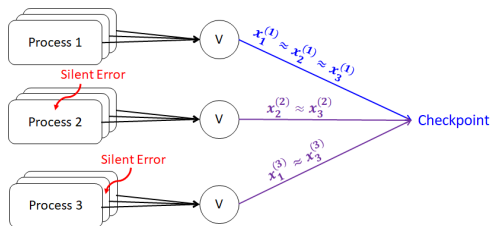


## ▶ Group Replication:

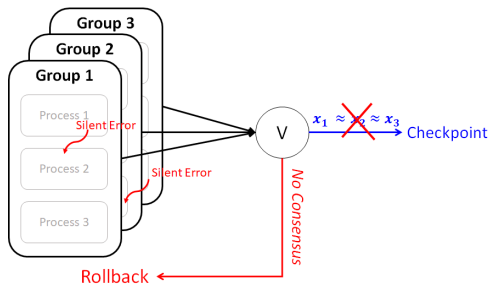


# Two Replication Modes

## ▶ Process Replication:



## ▶ Group Replication:



# Probability of Failure

Independent process error distribution

- ▶ Exponential  $Exp(\lambda)$ ,  $\lambda = 1/\mu$  (Memoryless)
- ▶ *Error probability of one process during  $T$  time of computation:*

$$\mathbb{P}(T) = 1 - e^{-\lambda T}$$

# Probability of Failure

Independent process error distribution

- ▶ Exponential  $Exp(\lambda)$ ,  $\lambda = 1/\mu$  (Memoryless)
- ▶ Error probability of one process during  $T$  time of computation:

$$\mathbb{P}(T) = 1 - e^{-\lambda T}$$

**Process Triplication:**

- ▶ Failure probability of any triplicated process:

$$\begin{aligned}\mathbb{P}_3^{\text{prc}}(T, 1) &= \binom{3}{2} (1 - \mathbb{P}(T)) \mathbb{P}(T)^2 + \mathbb{P}(T)^3 \\ &= 3e^{-\lambda T} (1 - e^{-\lambda T})^2 + (1 - e^{-\lambda T})^3 \\ &= 1 - 3e^{-2\lambda T} + 2e^{-3\lambda T}\end{aligned}$$

# Probability of Failure

Independent process error distribution

- ▶ Exponential  $Exp(\lambda)$ ,  $\lambda = 1/\mu$  (Memoryless)
- ▶ Error probability of one process during  $T$  time of computation:

$$\mathbb{P}(T) = 1 - e^{-\lambda T}$$

**Process Triplication:**

- ▶ Failure probability of any triplicated process:

$$\begin{aligned}\mathbb{P}_3^{\text{prc}}(T, 1) &= \binom{3}{2} (1 - \mathbb{P}(T)) \mathbb{P}(T)^2 + \mathbb{P}(T)^3 \\ &= 3e^{-\lambda T} (1 - e^{-\lambda T})^2 + (1 - e^{-\lambda T})^3 \\ &= 1 - 3e^{-2\lambda T} + 2e^{-3\lambda T}\end{aligned}$$

- ▶ Failure probability of  $P$ -process application:

$$\begin{aligned}\mathbb{P}_3^{\text{prc}}(T, P) &= 1 - \mathbb{P}(\text{"No process fails"}) \\ &= 1 - (1 - \mathbb{P}_3^{\text{prc}}(T, 1))^P \\ &= 1 - (3e^{-2\lambda T} - 2e^{-3\lambda T})^P\end{aligned}$$



# Probability of Failure

## Group Triplication:

- ▶ Failure probability of any  $P$ -process group:

$$\begin{aligned}\mathbb{P}_1^{\text{grp}}(T, P) &= 1 - \mathbb{P}(\text{"No process in group fails"}) \\ &= 1 - (1 - \mathbb{P}(T))^P \\ &= 1 - e^{-\lambda PT}\end{aligned}$$

# Probability of Failure

## Group Triplication:

- ▶ Failure probability of any  $P$ -process group:

$$\begin{aligned}\mathbb{P}_1^{\text{grp}}(T, P) &= 1 - \mathbb{P}(\text{"No process in group fails"}) \\ &= 1 - (1 - \mathbb{P}(T))^P \\ &= 1 - e^{-\lambda PT}\end{aligned}$$

- ▶ Failure probability of three-group application:

$$\begin{aligned}\mathbb{P}_3^{\text{grp}}(T, P) &= \binom{3}{2} (1 - \mathbb{P}_1^{\text{grp}}(T, 1)) \mathbb{P}_1^{\text{grp}}(T, 1)^2 + \mathbb{P}_1^{\text{grp}}(T, 1)^3 \\ &= 3e^{-\lambda PT} (1 - e^{-\lambda PT})^2 + (1 - e^{-\lambda PT})^3 \\ &= 1 - 3e^{-2\lambda PT} + 2e^{-3\lambda PT}\end{aligned}$$

# Probability of Failure

## Group Triplication:

- ▶ Failure probability of any  $P$ -process group:

$$\begin{aligned}\mathbb{P}_1^{\text{grp}}(T, P) &= 1 - \mathbb{P}(\text{"No process in group fails"}) \\ &= 1 - (1 - \mathbb{P}(T))^P \\ &= 1 - e^{-\lambda PT}\end{aligned}$$

- ▶ Failure probability of three-group application:

$$\begin{aligned}\mathbb{P}_3^{\text{grp}}(T, P) &= \binom{3}{2} (1 - \mathbb{P}_1^{\text{grp}}(T, 1)) \mathbb{P}_1^{\text{grp}}(T, 1)^2 + \mathbb{P}_1^{\text{grp}}(T, 1)^3 \\ &= 3e^{-\lambda PT} (1 - e^{-\lambda PT})^2 + (1 - e^{-\lambda PT})^3 \\ &= 1 - 3e^{-2\lambda PT} + 2e^{-3\lambda PT} \\ &> 1 - (3e^{-2\lambda T} - 2e^{-3\lambda T})^P = \mathbb{P}_3^{\text{prc}}(T, P)\end{aligned}$$

# Probability of Failure

## Group Triplication:

- ▶ Failure probability of any  $P$ -process group:

$$\begin{aligned}\mathbb{P}_1^{\text{grp}}(T, P) &= 1 - \mathbb{P}(\text{"No process in group fails"}) \\ &= 1 - (1 - \mathbb{P}(T))^P \\ &= 1 - e^{-\lambda PT}\end{aligned}$$

- ▶ Failure probability of three-group application:

$$\begin{aligned}\mathbb{P}_3^{\text{grp}}(T, P) &= \binom{3}{2} (1 - \mathbb{P}_1^{\text{grp}}(T, 1)) \mathbb{P}_1^{\text{grp}}(T, 1)^2 + \mathbb{P}_1^{\text{grp}}(T, 1)^3 \\ &= 3e^{-\lambda PT} (1 - e^{-\lambda PT})^2 + (1 - e^{-\lambda PT})^3 \\ &= 1 - 3e^{-2\lambda PT} + 2e^{-3\lambda PT} \\ &> 1 - (3e^{-2\lambda T} - 2e^{-3\lambda T})^P = \mathbb{P}_3^{\text{prc}}(T, P)\end{aligned}$$

**What about duplication?** (any error kills both cases)

$$\mathbb{P}_2^{\text{prc}}(T, P) = \mathbb{P}_2^{\text{grp}}(T, P) = 1 - e^{-2\lambda TP}$$

# Two Observations

## Observation 1 (Implementation)

- ▶ **Process replication** is more resilient than group replication (assuming same overhead)
- ▶ **Group replication** is easier to implement by treating an application as a blackbox

# Two Observations

## Observation 1 (Implementation)

- ▶ **Process replication** is more resilient than group replication (assuming same overhead)
- ▶ **Group replication** is easier to implement by treating an application as a blackbox

## Observation 2 (Analysis)

Following two scenarios are equivalent w.r.t. failure probability:

- ▶ **Group replication** with  $n$  replicas, where each replica has  $P$  processes and each process has error rate  $\lambda$
- ▶ **Process replication** with one process, which has error rate  $\lambda P$  and which is replicated  $n$  times

Benefit of analysis:  $\text{Group}(n, P, \lambda) \rightarrow \text{Process}(n, 1, \lambda P)$

# Analysis Steps

Maximize error-aware speedup

$$\mathbb{S}_n(T, P) = \frac{S(P)}{\mathbb{E}_n(T, P)/T}$$

1. Derive failure probability  $\mathbb{P}_n^{\text{prc}}(T, P)$  or  $\mathbb{P}_n^{\text{grp}}(T, P)$  — **exact**
2. Compute expected execution time  $\mathbb{E}_n(T, P)$  — **exact**
3. Compute **first-order approx.** of error-aware speedup  $\mathbb{S}_n(T, P)$
4. Derive optimal  $T_{\text{opt}}, P_{\text{opt}}$  and get  $\mathbb{S}_n(T_{\text{opt}}, P_{\text{opt}})$
5. Choose right replication level  $n$

# Analytical Results

## Duplication:

On a platform with  $Q$  processors and checkpointing cost  $C$ , the optimal resilience parameters for *process/group duplication* are:

$$P_{\text{opt}} = \min \left\{ \frac{Q}{2}, \left( \frac{1}{2} \left( \frac{1-\alpha}{\alpha} \right)^2 \frac{1}{C\lambda} \right)^{\frac{1}{3}} \right\}$$

$$T_{\text{opt}} = \left( \frac{C}{2\lambda P_{\text{opt}}} \right)^{\frac{1}{2}}$$

$$S_{\text{opt}} = \frac{S(P_{\text{opt}})}{1 + 2(2\lambda C P_{\text{opt}})^{\frac{1}{2}}}$$



# Analytical Results

## Duplication:

On a platform with  $Q$  processors and checkpointing cost  $C$ , the optimal resilience parameters for *process/group duplication* are:

$$P_{\text{opt}} = \min \left\{ \frac{Q}{2}, \left( \frac{1}{2} \left( \frac{1-\alpha}{\alpha} \right)^2 \frac{1}{C\lambda} \right)^{\frac{1}{3}} \right\}$$

$$T_{\text{opt}} = \left( \frac{C}{2\lambda P_{\text{opt}}} \right)^{\frac{1}{2}}$$

$$S_{\text{opt}} = \frac{S(P_{\text{opt}})}{1 + 2(2\lambda C P_{\text{opt}})^{\frac{1}{2}}}$$

## Triplexation & $(n, k)$ -replication ( $k$ -out-of- $n$ replica consensus):

similar results but different for process and group, less practical for  $n > 3$

# Analytical Results

## Duplication:

On a platform with  $Q$  processors and checkpointing cost  $C$ , the optimal resilience parameters for *process/group duplication* are:

$$P_{\text{opt}} = \min \left\{ \frac{Q}{2}, \left( \frac{1}{2} \left( \frac{1-\alpha}{\alpha} \right)^2 \frac{1}{C\lambda} \right)^{\frac{1}{3}} \right\}$$

$$T_{\text{opt}} = \left( \frac{C}{2\lambda P_{\text{opt}}} \right)^{\frac{1}{2}}$$

$$S_{\text{opt}} = \frac{S(P_{\text{opt}})}{1 + 2(2\lambda C P_{\text{opt}})^{\frac{1}{2}}}$$

## Triplication & $(n, k)$ -replication ( $k$ -out-of- $n$ replica consensus):

similar results but different for process and group, less practical for  $n > 3$

- ▶ For  $\alpha > 0$ , not necessarily use up all available  $Q$  processors
- ▶ Checkpointing interval  $T_{\text{opt}}$  nicely extends Young/Daly's result
- ▶ Error-aware speedup  $S_{\text{opt}}$  minimally affected for small  $\lambda$

# Results Comparison

For fully parallel jobs, i.e.,  $\alpha = 0$  (similar for  $\alpha > 0$ )

► Duplication v.s. Process triplication

$$P_{\text{opt}} = \frac{Q}{2}$$

$$P_{\text{opt}} = \frac{Q}{3}$$

(Processors ↓)

$$T_{\text{opt}} = \sqrt{\frac{C}{\lambda Q}}$$

$$T_{\text{opt}} = \sqrt[3]{\frac{C}{2\lambda^2 Q}}$$

(Chkpt interval ↑)

$$S_{\text{opt}} = \frac{Q/2}{1 + 2\sqrt{\lambda C Q}}$$

$$S_{\text{opt}} = \frac{Q/3}{1 + 3\sqrt[3]{\left(\frac{\lambda C}{2}\right)^2 Q}}$$

(Exp. speedup??)

# Results Comparison

For fully parallel jobs, i.e.,  $\alpha = 0$  (similar for  $\alpha > 0$ )

► Duplication v.s. Process triplication

$$P_{\text{opt}} = \frac{Q}{2}$$

$$P_{\text{opt}} = \frac{Q}{3}$$

(Processors ↓)

$$T_{\text{opt}} = \sqrt{\frac{C}{\lambda Q}}$$

$$T_{\text{opt}} = \sqrt[3]{\frac{C}{2\lambda^2 Q}}$$

(Chkpt interval ↑)

$$S_{\text{opt}} = \frac{Q/2}{1 + 2\sqrt{\lambda C Q}}$$

$$S_{\text{opt}} = \frac{Q/3}{1 + 3\sqrt[3]{\left(\frac{\lambda C}{2}\right)^2 Q}}$$

(Exp. speedup??)

► Process triplication v.s. Group triplication

$$P_{\text{opt}} = \frac{Q}{3}$$

$$P_{\text{opt}} = \frac{Q}{3}$$

(Processors =)

$$T_{\text{opt}} = \sqrt[3]{\frac{C}{2\lambda^2 Q}}$$

$$T_{\text{opt}} = \sqrt[3]{\frac{3C}{2(\lambda Q)^2}}$$

(Chkpt interval ↓)

$$S_{\text{opt}} = \frac{Q/3}{1 + 3\sqrt[3]{\left(\frac{\lambda C}{2}\right)^2 Q}}$$

$$S_{\text{opt}} = \frac{Q/3}{1 + 3\sqrt[3]{\frac{1}{3}\left(\frac{\lambda C Q}{2}\right)^2}}$$

(Exp. speedup ↓)

# Results Comparison

For fully parallel jobs, i.e.,  $\alpha = 0$  (similar for  $\alpha > 0$ )

- ▶ Duplication v.s. Process triplication

$$P_{\text{opt}} = \frac{Q}{2}$$

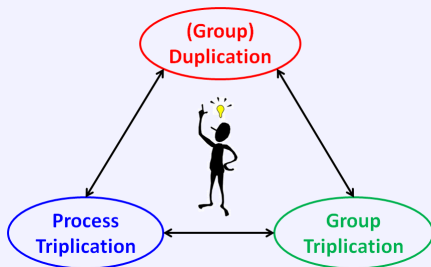
$$P_{\text{opt}} = \frac{Q}{2}$$

(Processors ↓)

## Choosing Right Mode & Level of Replication

Based on analytical model and whether process replication is supported

- ▶ P



$$1 + 3\sqrt[3]{\left(\frac{\lambda C}{2}\right)^2} Q$$

$$1 + 3\sqrt[3]{\frac{1}{3} \left(\frac{\lambda C Q}{2}\right)^2}$$

# Limitation of First-Order Approximation

## Observation 3 (First-Order)

Suppose  $P = \Theta(\lambda^{-x})$  and  $T = \Theta(\lambda^{-y})$ . Then, for first-order approximation to accurately estimate error probabilities (e.g.,  $1 - e^{-\lambda PT} \approx \lambda PT$ ), we need:

$$x + y < 1$$

or  $P \cdot T = o(\mu)$

e.g.,  $\mu = 10$  years  $\Rightarrow P \cdot T < 3 \cdot 10^8$  processor-seconds  
Generally accurate for platform MTBF  $\mu_P = \Theta(\text{days})$  or  
 $\mu_P = \Theta(\text{hours})$  depending on checkpointing cost  $C$

# Limitation of First-Order Approximation

## Observation 3 (First-Order)

Suppose  $P = \Theta(\lambda^{-x})$  and  $T = \Theta(\lambda^{-y})$ . Then, for first-order approximation to accurately estimate error probabilities (e.g.,  $1 - e^{-\lambda PT} \approx \lambda PT$ ), we need:

$$x + y < 1$$

or  $P \cdot T = o(\mu)$

e.g.,  $\mu = 10$  years  $\Rightarrow P \cdot T < 3 \cdot 10^8$  processor-seconds  
Generally accurate for platform MTBF  $\mu_P = \Theta(\text{days})$  or  
 $\mu_P = \Theta(\text{hours})$  depending on checkpointing cost  $C$

## What about larger systems?

One solution: **multi-level checkpointing**  $\Rightarrow$  error separation

[3] A. Benoit, A. Cavelan, V. Le Fèvre, Y. Robert and H. Sun. Towards Optimal Multi-Level Checkpointing. IEEE Transactions on Computers, 2017.

Consider an platform with  $Q = 10^6$ , and study

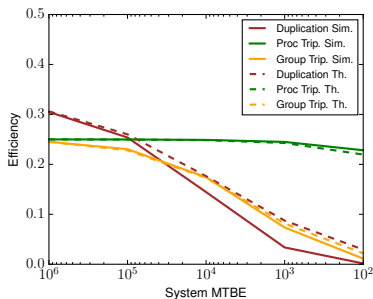
$$Efficiency = \frac{S_{opt}}{Q}$$

- ▶ Impact of MTBE and checkpointing cost  $C$
- ▶ Impact of sequential fraction  $\alpha$
- ▶ Impact of number of processes  $P$

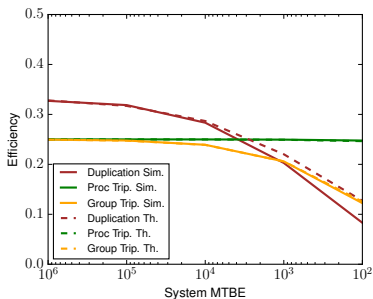


# Impact of MTBE and Checkpointing Cost

$$\alpha = 10^{-6}$$



(a)  $C = 1800s$

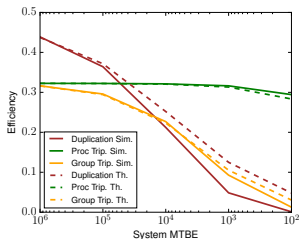


(b)  $C = 60s$

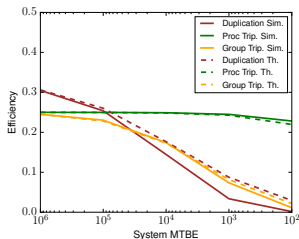
- ▶ First-order accurate except for duplication (where  $P$  is larger) and with small MTBE
- ▶ Duplication can be sufficient for large MTBE, especially for small checkpointing cost

# Impact of Sequential Fraction

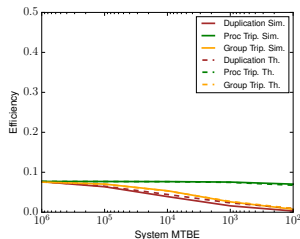
$C = 1800s$



(c)  $\alpha = 10^{-7}$



(d)  $\alpha = 10^{-6}$

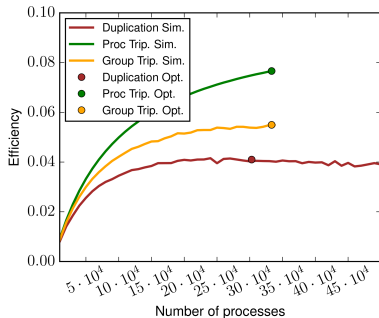


(e)  $\alpha = 10^{-5}$

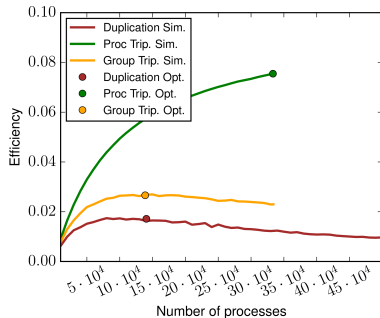
- ▶ Increased  $\alpha$  reduces efficiency
- ▶ Increased  $\alpha$  increases minimum MTBE for which duplication is sufficient

# Impact of Number of Processes

$$\alpha = 10^{-5}, C = 1800s$$



(f)  $MTBE = 10^4$



(g)  $MTBE = 10^3$

- ▶ Efficiency/error-aware speedup no longer strictly increasing function of  $P$
- ▶ First-order  $P_{\text{opt}}$  close to actual optimum

## What to Remember

- ▶ “Replication + checkpointing” as a general-purpose fault-tolerance protocol for coping with silent errors in HPC
- ▶ Process replication is more resilient than group replication, but group replication is easier to implement
- ▶ Analytical solution for  $P_{opt}$ ,  $T_{opt}$ , and  $S_{opt}$  and for choosing right replication mode and level

## Future Work

- ▶ Analyzing partial replication paradigm: different replication modes and levels for tasks with different criticality
- ▶ Dealing with co-existence of fail-stop errors and silent errors
- ▶ Experimenting with real applications/platforms