## Models for fault-tolerance at very large scale

### Yves Robert

### ENS Lyon & Institut Universitaire de France University of Tennessee Knoxville

yves.robert@ens-lyon.fr
http://graal.ens-lyon.fr/~yrobert/espas2013.pdf

### **ESPAS 2013**

(B)

### Thanks

### INRIA & ENS Lyon

- Frédéric Vivien
- PhD students (Guillaume Aupy, Dounia Zaidouni)

### UT Knoxville

- George Bosilca
- Aurélien Bouteiller
- Jack Dongarra
- Thomas Hérault

### Others

- Franck Cappello, UIUC-Inria joint lab
- Henri Casanova, Univ. Hawai'i
- Amina Guermouche, Univ. Versailles Paris



- Large-scale computing platforms
- Faults and failures



Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



3

### Introduction

- Large-scale computing platforms
- Faults and failures

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



<ロ> (日) (日) (日) (日) (日)

Introduction Large-scale computing platforms Eaults and failures

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



3

Probabilistic Models

Conclusion

### Exascale platforms (courtesy Jack Dongarra)

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

Systems	2011 K computer	2019	Difference 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)		

### Exascale platforms

### • Hierarchical

- $\bullet~10^5~{\rm or}~10^6~{\rm nodes}$
- Each node equipped with  $10^4$  or  $10^3$  cores

### • Failure-prone

MTBF – one node	1 year	10 years	120 years
MTBF – platform	30sec	5mn	1h
of 10 <sup>6</sup> nodes			

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

Introduction

### Exascale platforms



### Scenario for 2015

- Phase-Change memory, read 100GB/sec, write 10GB/sec
- Checkpoint size 128GB
- C: checkpoint save time: C = 12sec
- R: checkpoint recovery time: R = 1.2sec
- D: down/reboot time: D = 15sec
- Steady-state job utilization of whole platform
- p: total number of (multicore) nodes:  $p = 2^8$  to  $p = 2^{20}$
- MTBF  $\mu = 1$  week, 1 month, 1|10|100|1000 years (per node)

Introduction

Probabilistic Models

## Platform throughput with optimal checkpointing period

	р	Throughput			р	Throughput	]		р	Throughput
	2 <sup>8</sup>	91.56%	ſ	th	2 <sup>8</sup>	96.04%		2	2 <sup>8</sup>	98.89%
vee	2 <sup>11</sup>	73.75%		nor	2 <sup>11</sup>	88.23%		yea	2 <sup>11</sup>	96.80%
	2 <sup>14</sup>	20.07%		2	2 <sup>14</sup>	62.28%		Ë.	2 <sup>14</sup>	90.59%
	2 <sup>17</sup>	2.51%		Ī	2 <sup>17</sup>	10.66%		= 7	2 <sup>17</sup>	70.46%
1	2 <sup>20</sup>	0.31%		ή	2 <sup>20</sup>	1.33%		1	2 <sup>20</sup>	15.96%

	р	Throughput			р	Throughput		р	Throughput
rs	2 <sup>8</sup>	99.65%	[	ars	2 <sup>8</sup>	99.89%	ars	2 <sup>8</sup>	99.97%
yea	2 <sup>11</sup>	99.00%		ye	2 <sup>11</sup>	99.69%	ye	2 <sup>11</sup>	99.90%
10	2 <sup>14</sup>	97.15%		00	2 <sup>14</sup>	99.11%	8	2 <sup>14</sup>	99.72%
Π.	2 <sup>17</sup>	91.63%		1	2 <sup>17</sup>	97.45%	= 10	2 <sup>17</sup>	99.20%
ή	2 <sup>20</sup>	74.01%		μ	2 <sup>20</sup>	92.56%	= 11	2 <sup>20</sup>	97.73%

. . . . . . .



#### Introduction

Large-scale computing platforms

#### Faults and failures



Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



<ロ> (日) (日) (日) (日) (日)

# Error sources (courtesy Franck Cappello)

# Sources of failures

Probabilistic Models

Analysis of error and failure logs

Introduction

 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."



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

Fault-tolerance for HPC

### Faults and failures

### A few definitions

- 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

### Failure distributionss

- Exponential (memoryless)
- Weibull (account for infant mortality)

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?

ullet Well, it depends whether with or without rejuvenation  $\textcircled{\bullet}$ 

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?

ullet Well, it depends whether with or without rejuvenation ildot

Introduction

- Rebooting only faulty processor
- Platform failure distribution
  - $\Rightarrow$  superposition of *p* IID processor distributions
- Simple formula for arbitrary distributions:

$$\mu_{p} = \frac{\mu}{p}$$

with  ${\it p}$  processors of MTBF  $\mu$ 

- Rejuvenation does not matter for Exponential
- Rejuvenation harmful for Weibull with k < 1

### Values from the literature

- 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)

- Large-scale computing platforms
- Eaults and failures

### 2

#### Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



<ロ> (日) (日) (日) (日) (日)

### Introduc

- Large-scale computing platforms
- Faults and failures



#### Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



#### Conclusion

<ロ> (日) (日) (日) (日) (日)

### Checkpointing cost



# **Blocking model:** while a checkpoint is taken, no computation can be performed

A B F A B F

< 67 ▶

### Framework

- Periodic checkpointing policy of period T
- Independent and identically distributed failures
- Applies to a single processor with MTBF  $\mu = \mu_{\textit{ind}}$
- Applies to a platform with p processors with MTBF  $\mu = \frac{\mu_{ind}}{p}$ 
  - coordinated checkpointing
  - tightly-coupled application
  - progress  $\Leftrightarrow$  all processors available

### Waste: fraction of time not spent for useful computations

Probabilistic Models

### Waste in fault-free execution



- $T_{IME_{\text{base}}}$ : application base time
- $TIME_{FF}$ : with periodic checkpoints but failure-free

$$TIME_{FF} = TIME_{base} + #checkpoints \times C$$

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

$$\text{TIME}_{\mathsf{FF}} = \text{TIME}_{\mathsf{base}} \frac{\mathcal{T}}{\mathcal{T} - \mathcal{C}} \text{ and } \text{WASTE}[\mathcal{FF}] = \frac{\text{TIME}_{\mathsf{FF}} - \text{TIME}_{\mathsf{base}}}{\text{TIME}_{\mathsf{FF}}}$$

WASTE[*FF*] =  $\frac{C}{T}$ 

- E > - E >

### Waste due to failures

- $\bullet~{\rm TIME}_{\text{base}}:$  application base time
- $\bullet\ T{\rm IME}_{FF}{\rm :}$  with periodic checkpoints but failure-free
- $\bullet \ T{\rm IME}_{{\rm final}}:$  expectation of time with failures

 $\text{TIME}_{final} = \text{TIME}_{\mathsf{FF}} + \textit{N}_{faults} \times \textit{T}_{\mathsf{lost}}$ 

 $N_{faults}$  number of failures during execution  $T_{lost}$ : average time lost par failures

$$N_{faults} = rac{\mathrm{TIME}_{\mathrm{final}}}{\mu}$$
 $T_{\mathrm{lost}}$ ?

### Waste due to failures

- $\bullet~{\rm TIME}_{\text{base}}:$  application base time
- $\bullet~T{\scriptstyle\rm IME}_{FF}{:}$  with periodic checkpoints but failure-free
- $\bullet \ T{\rm IME}_{{\rm final}}:$  expectation of time with failures

$$\mathrm{TIME}_{\mathsf{final}} = \mathrm{TIME}_{\mathsf{FF}} + \mathit{N}_{\mathsf{faults}} \times \mathit{T}_{\mathsf{lost}}$$

 $N_{faults}$  number of failures during execution  $T_{lost}$ : average time lost par failures

$$N_{faults} = rac{\mathrm{TIME_{final}}}{\mu}$$

Introduction

Probabilistic Models

### Computing $T_{\text{lost}}$



 $\Rightarrow$  Instants when periods begin and failures strike are independent  $\Rightarrow$  Valid for all distribution laws, regardless of their particular shape

→ ∃ →

### Waste due to failures

$$\begin{split} \text{TIME}_{\text{final}} &= \text{TIME}_{\text{FF}} + N_{\text{faults}} \times T_{\text{lost}} \\ \text{TIME}_{\text{final}} &= \text{TIME}_{\text{FF}} + \frac{\text{TIME}_{\text{final}}}{\mu} \times \left(D + R + \frac{T}{2}\right) \\ \text{WASTE}[fail] &= \frac{\text{TIME}_{\text{final}} - \text{TIME}_{\text{FF}}}{\text{TIME}_{\text{final}}} \\ \text{WASTE}[fail] &= \frac{1}{\mu} \left(D + R + \frac{T}{2}\right) \end{split}$$

イロト イヨト イヨト イヨト

3

### Total waste



$$WASTE = \frac{TIME_{final} - TIME_{base}}{TIME_{final}}$$

 $(1 - \text{WASTE}[fail])(1 - \text{WASTE}[FF])\text{TIME}_{final} = Time[base]$ 1 - WASTE = (1 - WASTE[FF])(1 - WASTE[fail])

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

3

< 回 ト < 三 ト < 三 ト

### Waste minimization

$$\begin{aligned} \text{WASTE} &= \frac{C}{T} + \left(1 - \frac{C}{T}\right) \frac{1}{\mu} \left(D + R + \frac{T}{2}\right) \\ \text{WASTE} &= \frac{u}{T} + v + wT \\ u &= C \left(1 - \frac{D + R}{\mu}\right) \qquad v = \frac{D + R - C/2}{\mu} \qquad w = \frac{1}{2\mu} \end{aligned}$$

WASTE minimized for  $T = \sqrt{\frac{u}{w}}$ 

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

イロト イ理ト イヨト イヨトー

### Comparison with Young/Daly



$$(1 - \text{WASTE}[fail])$$
TIME<sub>final</sub> = TIME<sub>FF</sub>  
 $\Rightarrow T = \sqrt{2(\mu - (D + R))C}$ 

**Daly**: TIME<sub>final</sub> = 
$$(1 + \text{WASTE}[fail])$$
TIME<sub>FF</sub>  
 $\Rightarrow T = \sqrt{2(\mu + (D + R))C} + C$ 

### **Young**: TIME<sub>final</sub> = (1 + WASTE[fail])TIME<sub>FF</sub> and D = R = 0 $\Rightarrow T = \sqrt{2\mu C} + C$

• • = • • = •

Technicalities

Introduction

•  $\mathbb{E}(N_{faults}) = \frac{\text{Time_{final}}}{\mu}$  and  $\mathbb{E}(T_{\text{lost}}) = D + R + \frac{T}{2}$ but expectation of product is not product of expectations (not independent RVs here)

Probabilistic Models

- Enforce  $C \leq T$  to get  $WASTE[FF] \leq 1$
- Enforce  $D + R \le \mu$  and bound T to get  $\text{WASTE}[fail] \le 1$ but  $\mu = \frac{\mu_{ind}}{p}$  too small for large p, regardless of  $\mu_{ind}$

### Several failures within same period?

- WASTE[fail] accurate only when two or more faults do not take place within same period
- Cap period:  $\mathcal{T} \leq \gamma \mu$ , where  $\gamma$  is some tuning parameter
  - Poisson process of parameter  $\theta = \frac{T}{\mu}$
  - Probability of having  $k \ge 0$  failures :  $P(X = k) = \frac{\theta^k}{k!} e^{-\theta}$

• Probability of having two or more failures:  $\pi = P(X \ge 2) = 1 - (P(X = 0) + P(X = 1)) = 1 - (1 + \theta)e^{-\theta}$ 

• 
$$\gamma = 0.27 \Rightarrow \pi \le 0.03$$

 $\Rightarrow$  overlapping faults for only 3% of checkpointing segments

Conclusion

# Validity of the approach (3/3)

Introduction

• Enforce 
$$T \leq \gamma \mu$$
,  $C \leq \gamma \mu$ , and  $D + R \leq \gamma \mu$ 

• Optimal period  $\sqrt{2(\mu - (D + R))C}$  may not belong to admissible interval  $[C, \gamma \mu]$ 

Probabilistic Models

• Waste is then minimized for one of the bounds of this admissible interval (by convexity)



- Capping periods, and enforcing a lower bound on MTBF
   ⇒ mandatory for mathematical rigor ☺
  - $\Rightarrow$  optimal strategy only known for Exp. distributions  $\ensuremath{\mathfrak{S}}$

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

• Approach surprisingly robust 😳

### Introdu

- Large-scale computing platforms
- Faults and failures



#### Probabilistic models and execution scenarios

Young/Daly's approximation

#### Failure Prediction

- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



#### Conclusion

<ロ> (日) (日) (日) (日) (日)

> = 31/ 71

### Framework

### Predictor

- Exact prediction dates (at least C seconds in advance)
- Recall r: fraction of faults that are predicted
- Precision *p*: fraction of fault predictions that are correct

### Events

- true positive: predicted faults
- *false positive*: fault predictions that did not materialize as actual faults
- false negative: unpredicted faults

$$r = rac{True_P}{True_P + False_N}$$
 and  $p = rac{True_P}{True_P + False_P}$ 

### Fault rates

•  $\mu$ : mean time between failures (MTBF)

1

- $\mu_P$  mean time between predicted events (both true positive and false positive)
- $\mu_{NP}$  mean time between unpredicted faults (false negative).
- $\mu_e$ : mean time between events (including all three event types)

$$\frac{(1-r)}{\mu} = \frac{1}{\mu_{NP}}$$
$$\frac{r}{\mu} = \frac{p}{\mu_{P}}$$
$$\frac{1}{\mu_{e}} = \frac{1}{\mu_{P}} + \frac{1}{\mu_{NP}}$$

33/71

< 3 > < 3 >


- While no fault prediction is available:
  - ullet checkpoints taken periodically with period  ${\mathcal T}$
- When a fault is predicted at time t:
  - take a checkpoint ALAP (completion right at time t)
  - after the checkpoint, complete the execution of the period

# Computing the waste

- **•** Fault-free execution: WASTE $[FF] = \frac{C}{T}$
- **②** Unpredicted faults:  $\frac{1}{\mu_{NP}} \left[ D + R + \frac{T}{2} \right]$
- S Predictions:  $\frac{1}{\mu_P} \left[ p(C+D+R) + (1-p)C \right]$

WASTE[fail] = 
$$\frac{1}{\mu} \left[ (1-r)\frac{T}{2} + D + R + \frac{r}{p}C \right]$$

$$T_{opt} \approx \sqrt{rac{2\mu C}{1-r}}$$

伺下 イヨト イヨト

# Outline



Large-scale computing platforms

Faults and failures



#### Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction

#### Checkpointing protocols

- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



#### Conclusion

<ロ> (日) (日) (日) (日) (日)

Probabilistic Models

Conclusion

# Background: coordinated checkpointing protocols

- Coordinated checkpoints over all processes
- Global restart after a failure



- ☺ No risk of cascading rollbacks
- $\bigcirc$  No need to log messages
- ☺ All processors need to roll back

# Background: message logging protocols

- Message content logging (sender memory)
- Restart of the failed process



- ☺ No cascading rollbacks
- $\bigcirc$  Number of processes to roll back
- Memory occupation
- 🙂 Overhead

# Background: hierarchical protocols

- Clusters of processes
- Coordinated checkpointing protocol within clusters
- Message logging protocols between clusters
- Only processors from failed group need to roll back



- Over the second seco
  - Slowdowns failure-free execution
  - Increases checkpoint size/time
- ☺ Faster re-execution with logged messages

# Which checkpointing protocol to use?

#### Coordinated checkpointing

- © No risk of cascading rollbacks
- ③ No need to log messages
- ☺ All processors need to roll back
- 🙂 Rumor: May not scale to very large platforms

#### Hierarchical checkpointing

- Need to log inter-groups messages
  - Slowdowns failure-free execution
  - Increases checkpoint size/time
- $\ensuremath{\textcircled{\odot}}$  Only processors from failed group need to roll back
- ☺ Faster re-execution with logged messages
- ☺ Rumor: Should scale to very large platforms

# Outline

#### Introdu

- Large-scale computing platforms
- Faults and failures



#### Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



#### Conclusion

<ロ> (日) (日) (日) (日) (日)

# Checkpointing cost



#### Blocking model: checkpointing blocks all computations

#### Checkpointing cost



**Non-blocking model:** checkpointing has no impact on computations (e.g., first copy state to RAM, then copy RAM to disk)

### Checkpointing cost



Processing the first chunk

**General model:** checkpointing slows computations down: during a checkpoint of duration C, the same amount of computation is done as during a time  $\alpha C$  without checkpointing  $(0 \le \alpha \le 1)$ 

### Waste in fault-free execution



Time elapsed since last checkpoint: T

Amount of computations executed: WORK =  $(T - C) + \alpha C$ WASTE[FF] =  $\frac{T - WORK}{T}$ 

A B F A B F

# Waste due to failures



Failure can happen

- Ouring computation phase
- 2 During checkpointing phase

Time

#### Waste due to failures





Image: A matrix

Time

### Waste due to failures





< □ > < ---->

A B F A B F

#### Waste due to failures



Coordinated checkpointing protocol: when one processor is victim of a failure, all processors lose their work and must roll back to last checkpoint

Probabilistic Models

Conclusion

# Waste due to failures in computation phase



(日) (周) (三) (三)

#### Waste due to failures in computation phase



# Coordinated checkpointing protocol: all processors must recover from last checkpoint

43/71

3 K K 3 K

Probabilistic Models

Conclusion

#### Waste due to failures in computation phase



Redo the work destroyed by the failure, that was done in the checkpointing phase before the computation phase

But no checkpoint is taken in parallel, hence this re-execution is faster than the original computation

Probabilistic Models

Conclusion

### Waste due to failures in computation phase



#### Re-execute the computation phase

.∋...>

Probabilistic Models

Conclusion

#### Waste due to failures in computation phase



#### Finally, the checkpointing phase is executed

43/71

.∋...>

#### Total waste



WASTE[fail] = 
$$\frac{1}{\mu} \left( D + R + \alpha C + \frac{T}{2} \right)$$

Optimal period  $T_{opt} = \sqrt{2(1-\alpha)(\mu - (D+R))C}$ 

イロト イ理ト イヨト イヨトー

イロト イヨト イヨト イヨト

# Outline



- Large-scale computing platforms
- Faults and failures



#### Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



Conclusio

イロト イヨト イヨト イヨト

#### Hierarchical checkpointing



- Processors partitioned into G groups
- Each group includes q processors
- Inside each group: coordinated checkpointing in time C(q)
- Inter-group messages are logged

#### Accounting for message logging: Impact on work

- $\bigcirc$  Logging messages slows down execution:  $\Rightarrow$  WORK becomes  $\lambda$ WORK, where  $0 < \lambda < 1$ Typical value:  $\lambda \approx 0.98$
- © Re-execution after a failure is faster:  $\Rightarrow$  RE-EXEC becomes  $\frac{\text{Re-EXEC}}{\rho}$ , where  $\rho \in [1..2]$ Typical value:  $\rho \approx 1.5$

$$WASTE[FF] = \frac{T - \lambda WORK}{T}$$
$$WASTE[fail] = \frac{1}{\mu} \left( D(q) + R(q) + \frac{\text{Re-Exec}}{\rho} \right)$$

#### Accounting for message logging: Impact on checkpoint size

- Inter-groups messages logged continuously
- Checkpoint size increases with amount of work executed before a checkpoint
- $C_0(q)$ : Checkpoint size of a group without message logging

$$\mathcal{C}(q) = \mathcal{C}_0(q)(1 + eta \mathrm{WORK}) \Leftrightarrow eta = rac{\mathcal{C}(q) - \mathcal{C}_0(q)}{\mathcal{C}_0(q) \mathrm{WORK}}$$

WORK = 
$$\lambda(T - (1 - \alpha)GC(q))$$
  
 $C(q) = \frac{C_0(q)(1 + \beta\lambda T)}{1 + GC_0(q)\beta\lambda(1 - \alpha)}$ 

#### Three case studies

#### Coord-IO

Coordinated approach:  $C = C_{Mem} = \frac{Mem}{b_{io}}$ where Mem is the memory footprint of the application

#### Hierarch-IO

Several (large) groups, I/O-saturated  $\Rightarrow$  groups checkpoint sequentially

$$\mathcal{C}_0(q) = rac{\mathcal{C}_{\mathsf{Mem}}}{\mathcal{G}} = rac{\mathsf{Mem}}{\mathcal{G}\mathsf{b}_{io}}$$

#### **Hierarch-Port**

Very large number of smaller groups, *port-saturated*  $\Rightarrow$  some groups checkpoint in parallel Groups of q<sub>min</sub> processors, where q<sub>min</sub>b<sub>port</sub>  $\ge$  b<sub>io</sub>

### Three applications

- 2D-stencil
- 2 Matrix product
- 3D-Stencil
  - Plane
  - Line

★ 3 > < 3 >

Probabilistic Models

# Computing $\beta$ for 2D-Stencil

 $C(q) = C_0(q) + Logged_Msg = C_0(q)(1 + \beta WORK)$ 

Real  $n \times n$  matrix and  $p \times p$  grid  $Work = \frac{9b^2}{s_p}$ , b = n/pEach process sends a block to its 4 neighbors

#### HIERARCH-IO:

- 1 group = 1 grid row
- 2 out of the 4 messages are logged

• 
$$\beta = \frac{2s_p}{9b^3}$$

HIERARCH-PORT:

•  $\beta$  doubles



# Four platforms: basic characteristics

Name	Number of	Number of	Number of cores	Memory	I/O Network Bandwidth (bio)		I/O Bandwidth (bport)
	cores	processors p <sub>total</sub>	per processor	per processor	Read	Write	Read/Write per processor
Titan	299,008	16,688	16	32GB	300GB/s	300GB/s	20GB/s
K-Computer	705,024	88,128	8	16GB	150GB/s	96GB/s	20GB/s
Exascale-Slim	1,000,000,000	1,000,000	1,000	64GB	1TB/s	1TB/s	200GB/s
Exascale-Fat	1,000,000,000	100,000	10,000	640GB	1TB/s	1TB/s	400GB/s

Name	Scenario	G (C(q))	$\beta$ for	$\beta$ for	
			2D-Stencil	MATRIX-PRODUCT	
	Coord-IO	1 (2,048s)	/	/	
Titan	HIERARCH-IO	136 (15s)	0.0001098	0.0004280	
	HIERARCH-PORT	1,246 (1.6s)	0.0002196	0.0008561	
	Coord-IO	1 (14,688s)	/	/	
K-Computer	HIERARCH-IO	296 (50s)	0.0002858	0.001113	
	HIERARCH-PORT	17,626 (0.83s)	0.0005716	0.002227	
	Coord-IO	1 (64,000s)	/	/	
Exascale-Slim	HIERARCH-IO	1,000 (64s)	0.0002599	0.001013	
	HIERARCH-PORT	200,0000 (0.32s)	0.0005199	0.002026	
	Coord-IO	1 (64,000s)	/	/	
Exascale-Fat	HIERARCH-IO	316 (217s)	0.00008220	0.0003203	
	HIERARCH-PORT	33,3333 (1.92s)	0.00016440	0.0006407	

イロト イヨト イヨト イヨト

# Checkpoint time

Name	С		
K-Computer	14,688s		
Exascale-Slim	64,000		
Exascale-Fat	64,000		

- Large time to dump the memory
- Using 1%C
- Comparing with 0.1%C for exascale platforms
- $\alpha = 0.3$ ,  $\lambda = 0.98$  and  $\rho = 1.5$

< 3 > < 3 >

Probabilistic Models 

# Plotting formulas - Platform: Titan



Waste as a function of processor MTBF  $\mu$ 

53/71

Probabilistic Models

Conclusion

### Platform: K-Computer



Waste as a function of processor MTBF  $\mu$ 

54/71

- E > - E >

Probabilistic Models

Conclusion

### Plotting formulas - Platform: Exascale

#### WASTE = 1 for all scenarios!!!

yves.robert@ens-lyon.fr

Fault-tolerance for HPC

通 ト イヨト イヨト

> = 55/71



Probabilistic Models

# Plotting formulas - Platform: Exascale



55/71

Conclusion
Probabilistic Models

Conclusion

## Plotting formulas – Platform: Exascale with C = 1,000



yves.robert@ens-lyon.fr

Probabilistic Models

Conclusion

## Plotting formulas – Platform: Exascale with C = 100



yves.robert@ens-lyon.fr

Fault-tolerance for HPC

Probabilistic Models

Conclusion

## Simulations – Platform: Titan



Makespan (in days) as a function of processor MTBF  $\mu$ 

(日) (周) (三) (三)

3

Probabilistic Models

Conclusion

## Simulations – Platform: Exascale with C = 1,000



Makespan (in days) as a function of processor MTBF  $\mu$ , C=1,000

yves.robert@ens-lyon.fr

< 🗗 🕨

3

A B F A B F

Probabilistic Models

Conclusion

### Simulations – Platform: Exascale with C = 100



Makespan (in days) as a function of processor MTBF  $\mu$ , C = 100

yves.robert@ens-lyon.fr

・ロト ・ 日 ・ ・ ヨ ・ ・ ヨ ・

\*) Q (

### Outline

#### Introduc

- Large-scale computing platforms
- Faults and failures

#### 2

#### Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication

#### Conclusi

<ロ> (日) (日) (日) (日) (日)

#### PROCESS REPLICATION



• Each process replicated  $g \ge 2$  times  $\rightarrow$  *replica-group* 

- $n_{rg}$  = number of replica-groups ( $g \times n_{rg} = N$ )
- Study for g = 2 by Ferreira et al., SC'2011

62/71

< ∃ > <

### Number of failures to bring down application

- $MNFTI^{ah}$  Count each failure hitting any of the  $g \cdot n_{rg}$  initial processors, including those *already hit* by a failure
- *MNFTI*<sup>rp</sup> Count failures that hit *running processors*, and thus effectively kill replicas.

 $MNFTI^{ah} = 1 + MNFTI^{rp}$ 

### Number of failures to bring down application

- $MNFTI^{ah}$  Count each failure hitting any of the  $g \cdot n_{rg}$  initial processors, including those *already hit* by a failure
- *MNFTI*<sup>rp</sup> Count failures that hit *running processors*, and thus effectively kill replicas.

 $\textit{MNFTI}^{\rm ah} = 1 + \textit{MNFTI}^{\rm rp}$ 

Probabilistic Models

Conclusion

## Analogy with birthday problem



イロト イヨト イヨト イヨト

► = 64/ 71

Probabilistic Models

Conclusion

## Analogy with birthday problem



#### $n = n_{rg}$ bins, throw balls until one bin gets two balls

- E - - E -

Probabilistic Models

Conclusion

## Analogy with birthday problem



 $n = n_{rg}$  bins, throw balls until one bin gets two balls

Expected number of balls to throw: Birthday(n) =  $1 + \int_0^{+\infty} e^{-x} (1 + x/n)^{n-1} dx$ 

Probabilistic Models

Conclusion

## Analogy with birthday problem



#### But second failure may hit already struck replica 😟

64/71

B ▶ < B ▶

Probabilistic Models

Conclusion

## Analogy with birthday problem



 $n = n_{rg}$  bins, red and blue balls

 $MNFTI^{ah}$  = expected number of balls to throw until one bin gets one ball of each color

## Exponential failures, g = 2

Theorem  $MNFTI^{\mathrm{ah}} = \mathbb{E}(NFTI^{\mathrm{ah}}|0)$  where

$$\mathbb{E}(NFTI^{\mathrm{ah}}|n_f) = \begin{cases} 2 & \text{if } n_f = n_{rg}, \\ \frac{2n_{rg}-n_f}{2n_{rg}-n_f} + \frac{2n_{rg}-2n_f}{2n_{rg}-n_f} \mathbb{E}\left(NFTI^{\mathrm{ah}}|n_f+1\right) & \text{otherwise.} \end{cases}$$

 $\mathbb{E}(NFTI^{ah}|n_f)$ : expectation of number of failures to kill application, knowing that

- application is still running
- failures have already hit  $n_f$  different replica-groups

Probabilistic Models

Conclusion

## Exponential failures, g = 2 (cont'd)

#### Proof

$$\mathbb{E}\left(\textit{NFTI}^{\mathrm{ah}} \left| \textit{n}_{\textit{rg}} \right.\right) = \frac{1}{2} \times 1 + \frac{1}{2} \times \left(1 + \mathbb{E}\left(\textit{NFTI}^{\mathrm{ah}} \left| \textit{n}_{\textit{rg}} \right.\right)\right).$$

$$\mathbb{E}\left(NFTI^{\mathrm{ah}}|n_{f}\right) = \frac{2n_{rg}-2n_{f}}{2n_{rg}} \times \left(1 + \mathbb{E}\left(NFTI^{\mathrm{ah}}|n_{f}+1\right)\right) \\ + \frac{2n_{f}}{2n_{rg}} \times \left(\frac{1}{2} \times 1 + \frac{1}{2}\left(1 + \mathbb{E}\left(NFTI^{\mathrm{ah}}|n_{f}\right)\right)\right).$$

## Exponential failures, g = 2 (cont'd)

#### Proof

$$\mathbb{E}\left(\textit{NFTI}^{\mathrm{ah}} \left| \textit{n}_{\textit{rg}} \right.\right) = \frac{1}{2} \times 1 + \frac{1}{2} \times \left(1 + \mathbb{E}\left(\textit{NFTI}^{\mathrm{ah}} \left| \textit{n}_{\textit{rg}} \right.\right)\right).$$

$$\mathbb{E}\left(\mathsf{NFTI}^{\mathrm{ah}}|n_{f}\right) = \frac{2n_{rg} - 2n_{f}}{2n_{rg}} \times \left(1 + \mathbb{E}\left(\mathsf{NFTI}^{\mathrm{ah}}|n_{f} + 1\right)\right) \\ + \frac{2n_{f}}{2n_{rg}} \times \left(\frac{1}{2} \times 1 + \frac{1}{2}\left(1 + \mathbb{E}\left(\mathsf{NFTI}^{\mathrm{ah}}|n_{f}\right)\right)\right).$$

VIII - system MTBF (2016) X MFT M マロン・(アン・マネン・ミン・マスン・マン・ yves.robert@ens-lyon.fr Fault-tolerance for HPC 66/71

## Exponential failures, g = 2 (cont'd)

#### Proof

$$\mathbb{E}\left(\textit{NFTI}^{\mathrm{ah}} \left| \textit{n}_{\textit{rg}} \right.\right) = \frac{1}{2} \times 1 + \frac{1}{2} \times \left(1 + \mathbb{E}\left(\textit{NFTI}^{\mathrm{ah}} \left| \textit{n}_{\textit{rg}} \right.\right)\right).$$

$$\mathbb{E}\left(\mathsf{NFTI}^{\mathrm{ah}}|n_{f}\right) = \frac{2n_{rg} - 2n_{f}}{2n_{rg}} \times \left(1 + \mathbb{E}\left(\mathsf{NFTI}^{\mathrm{ah}}|n_{f} + 1\right)\right) \\ + \frac{2n_{f}}{2n_{rg}} \times \left(\frac{1}{2} \times 1 + \frac{1}{2}\left(1 + \mathbb{E}\left(\mathsf{NFTI}^{\mathrm{ah}}|n_{f}\right)\right)\right).$$

# $MTTI = systemMTBF(2n_{rg}) \times MNFTI^{ah}$

#### yves.robert@ens-lyon.fr

#### Fault-tolerance for HPC

イロト イポト イヨト イヨト

### Failure distribution

#### R(t) probability that application still running at time t

- All replica-groups have at least one replica running
- Exponential:  $R(t) = (1 (1 e^{-\lambda t})^g)^{n_{rg}}$
- Weibull:  $R(t) = \left(1 \left(1 e^{-\left(\frac{t}{\lambda}\right)^k}\right)^g\right)^{n_{rg}}$

• Can use dynamic programming algorithms for sequential/parallel jobs

#### MTTI

- $MTTI = \int_0^{+\infty} R(t) dt \rightarrow \text{closed-form formulas}$
- Assess the impact of replication for various scenarios 😳

A B A A B A

### Failure distribution



#### Crossover point for replication when $\mu = 125$ years

(日) (周) (三) (三)

3

### Outline

#### Introdu

- Large-scale computing platforms
- Faults and failures

#### 2

Probabilistic models and execution scenarios

- Young/Daly's approximation
- Failure Prediction
- Checkpointing protocols
- Coordinated checkpointing
- Hierarchical checkpointing
- Replication



#### Conclusion

<ロ> (日) (日) (日) (日) (日)

### Conclusion

- Multiple approaches to Fault Tolerance
- Application-specific FT will always provide more benefits
- General-purpose FT will always be needed
  - Not every computer scientist needs to learn how to write fault-tolerant applications
  - Not all parallel applications can be ported to a fault-tolerant version
- Faults are a feature of the platform. Why should it be the role of the programmers to handle them?

### Conclusion

- Software/hardware techniques to reduce checkpoint, recovery, migration times and to improve failure prediction
- Multi-criteria scheduling problem execution time/energy/reliability add replication best resource usage (performance trade-offs)
- Need combine all these approaches!

Several challenging algorithmic/scheduling problems 😳

Extended version of this talk: see SC'12 tutorial