

Multi-criteria Checkpointing Strategies: Response-time versus Resource Utilization

Aurélien BOUTEILLER¹,
Franck CAPPELLO²,
Jack DONGARRA¹, Amina GUERMOUCHE³,
Thomas HÉRAULT¹, Yves ROBERT^{1,4},

1. University of Tennessee Knoxville, USA
2. INRIA & University of Illinois at Urbana Champaign, USA
3. Université de Versailles Saint Quentin, France
4. Ecole Normale Supérieure de Lyon & INRIA, France

June 13th, 2013

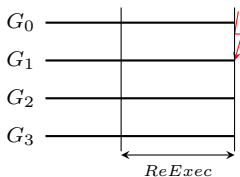
Introduction

- **Very very** large number of processing elements (e.g., 2^{20})
⇒ Probability of failures dramatically increases
- Large application to be executed on whole platform
⇒ Failure(s) will most likely occur before completion!
- Resilience provided through checkpointing
 - ① Coordinated checkpointing protocols
 - ☹ I/O overhead
 - ② Uncoordinated checkpointing protocols with message logging
 - Hierarchical checkpointing protocols

Introduction

Hierarchical protocols:

- 😊 A subset of processes roll back
 - Overlap recovery and normal execution
 - Tightly coupled applications:
 - ☹ Non-rolled back processes have to wait
 - Execute another application during recovery
 - The failed group recovers
 - The non failed groups load and execute another application
 - 😊 Improve platform efficiency



Introduction

① Application-oriented scenario:

- Non failed processes wait for the recovering ones
- The application is executed on $G + 1$ groups

② Platform-oriented scenario:

- The application is executed on G groups
- A spare group is used for recovery
- The G groups are used to execute another application while the spare group is recovering

Framework

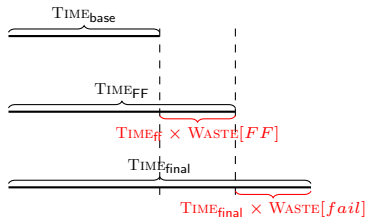
- Periodic checkpointing policies (of period T)
- Independent and identically distributed failures
- Platform failure inter-arrival time: μ
- Tightly-coupled application:
 progress \Leftrightarrow all processors available
- First-order approximation: at most one failure within a period

Waste: fraction of time not spent for useful computations

- ① Application waste: fraction of time the processes do not execute the application
- ② Platform waste: fraction of time the resources are not used to perform useful work

Waste

- $\text{TIME}_{\text{base}}$: application base time
- TIME_{FF} : with periodic checkpoints but failure-free
- $\text{TIME}_{\text{final}}$: expectation of time with failures



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

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

$$\text{WASTE} = \frac{\text{TIME}_{\text{final}} - \text{TIME}_{\text{base}}}{\text{TIME}_{\text{final}}}$$

$$\text{WASTE} = 1 - (1 - \text{WASTE}[FF])(1 - \text{WASTE}[fail])$$

Outline

① Protocol overhead

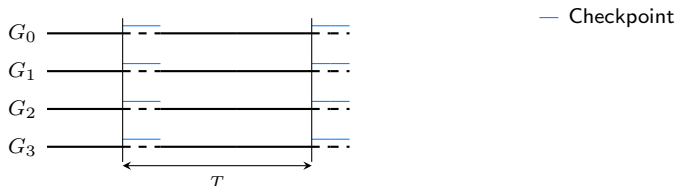
② Application scenario

③ Platform scenario

④ Instantiating the model

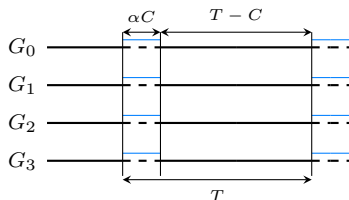
⑤ Simulations

Checkpoint overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)

Checkpoint overhead

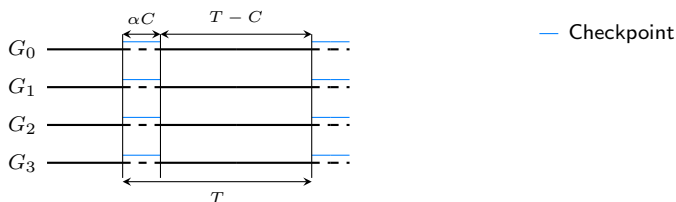


— Checkpoint

- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :

$$Work = \alpha C + T - C$$

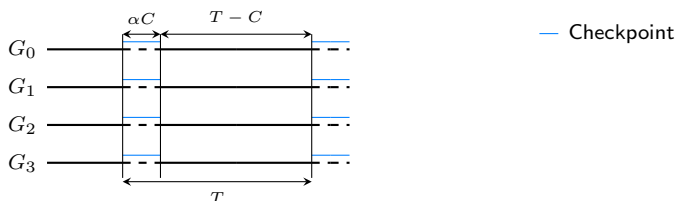
Message logging overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :

$$Work = \alpha C + T - C$$

Message logging overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :

$$Work = \alpha C + T - C$$
- Message logging slows down the execution with a factor λ ($0 < \lambda < 1$):
 The amount of computation executed in a Period T : $\lambda Work$

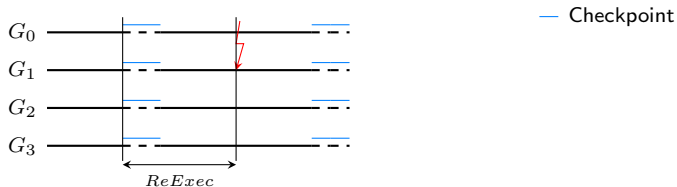
Message logging overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :

$$\lambda Work = \lambda(\alpha C + T - C)$$

Message logging overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :

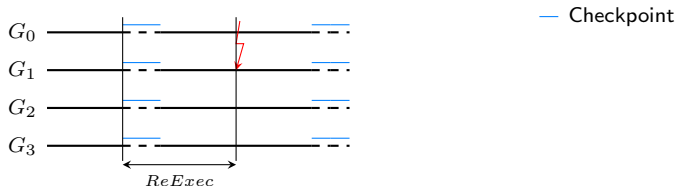
$$\lambda Work = \lambda(\alpha C + T - C)$$

Message logging overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :
 $\lambda Work = \lambda(\alpha C + T - C)$
- The amount of work to re-execute is $ReExec$

Message logging overhead



- During a checkpoint C , αC work is done ($0 \leq \alpha \leq 1$)
- The amount of computation executed in a Period T :
 $\lambda Work = \lambda(\alpha C + T - C)$
- The amount of work to re-execute is $ReExec$
- Message logging speeds up the re-execution with a factor ρ
 $(1 < \rho < 2)$: It takes $\frac{ReExec}{\rho}$ to execute $ReExec$

Outline

① Protocol overhead

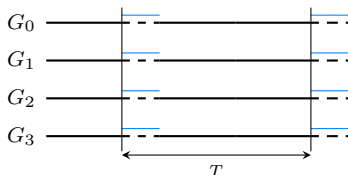
② Application scenario

③ Platform scenario

④ Instantiating the model

⑤ Simulations

Waste in the absence of failures



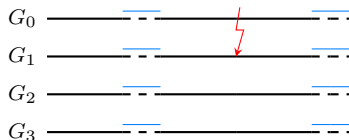
- T work should be done
- $\lambda Work = \lambda(T - C + \alpha C)$ is done

$$WASTE[ff] = \frac{T - \lambda Work}{T}$$

Waste in case of failures



Waste in case of failures



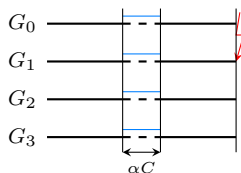
- ① Failure during work:

Waste in case of failures



- ① Failure during work:

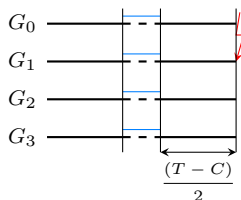
Waste in case of failures



- ① Failure during work:

$$ReExec1 = \alpha C$$

Waste in case of failures

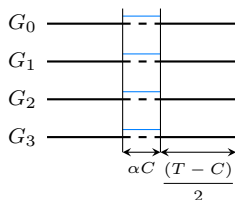


① Failure during work:

$$ReExec1 = \alpha C + \frac{T - C}{2}$$

Probability: $\frac{T - C}{T}$

Waste in case of failures



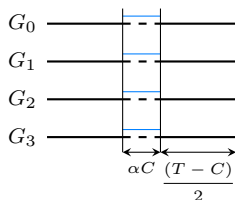
- ① Failure during work:

$$ReExec1 = \alpha C + \frac{T - C}{2}$$

Probability: $\frac{T - C}{T}$

- ② Failure during checkpoint:

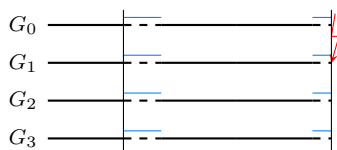
Waste in case of failures



- ① Failure during work:

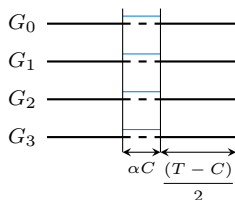
$$ReExec1 = \alpha C + \frac{T-C}{2}$$

Probability: $\frac{T-C}{T}$



- ② Failure during checkpoint:

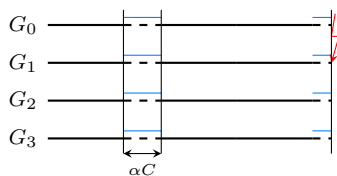
Waste in case of failures



- ① Failure during work:

$$ReExec1 = \alpha C + \frac{T - C}{2}$$

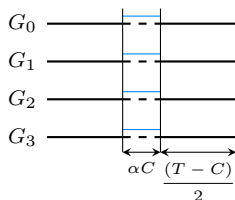
Probability: $\frac{T - C}{T}$



- ② Failure during checkpoint:

$$ReExec2 = \alpha C$$

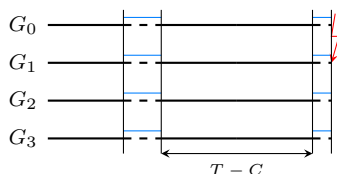
Waste in case of failures



① Failure during work:

$$ReExec1 = \alpha C + \frac{T - C}{2}$$

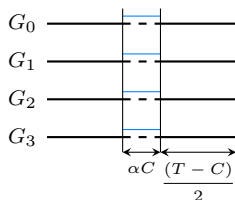
Probability: $\frac{T - C}{T}$



② Failure during checkpoint:

$$ReExec2 = \alpha C + T - C$$

Waste in case of failures



① Failure during work:

$$ReExec1 = \alpha C + \frac{T-C}{2}$$

$$\text{Probability: } \frac{T-C}{T}$$

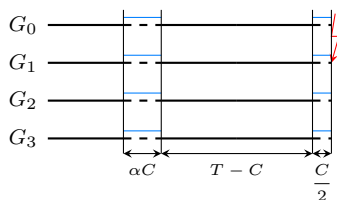
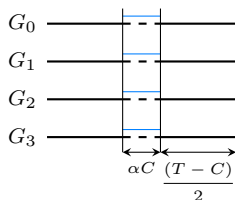


② Failure during checkpoint:

$$ReExec2 = \alpha C + T - C + \frac{C}{2}$$

$$\text{Probability: } \frac{C}{T}$$

Waste in case of failures



① Failure during work:

$$ReExec1 = \alpha C + \frac{T-C}{2}$$

$$\text{Probability: } \frac{T-C}{T}$$

② Failure during checkpoint:

$$ReExec2 = \alpha C + T - C + \frac{C}{2}$$

$$\text{Probability: } \frac{C}{T}$$

$$\text{WASTE}[fail] = \frac{1}{\mu} \left[D + R + \frac{T-C}{T} \times \frac{ReExec1}{\rho} + \frac{T}{C} \times \frac{ReExec2}{\rho} \right]$$

Outline

① Protocol overhead

② Application scenario

③ Platform scenario

④ Instantiating the model

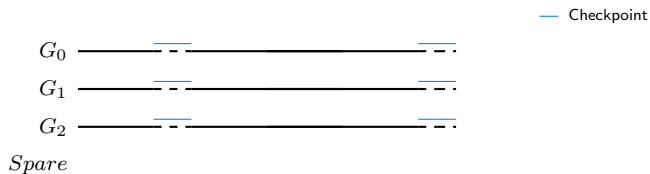
⑤ Simulations

Framework

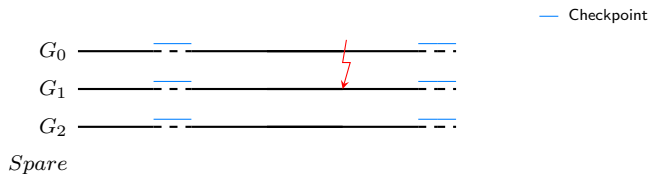
- $G + 1$ available groups
- G groups used to run the application
- A spare group used for the rollback
 - Longer time needed to compute the same amount of work
 - Longer checkpoint duration (by a factor of $\frac{G+1}{G}$)

$$\text{WASTE} = \frac{1}{G+1} + \frac{G}{G+1} (\text{WASTE}[ff] + \text{WASTE}[fail] - \text{WASTE}[ff]\text{WASTE}[fail])$$

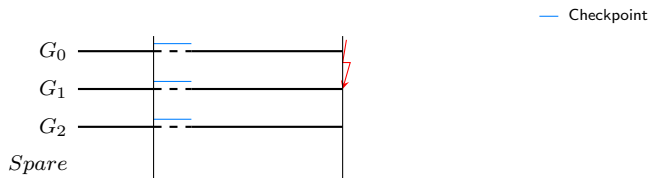
Waste in case of failures



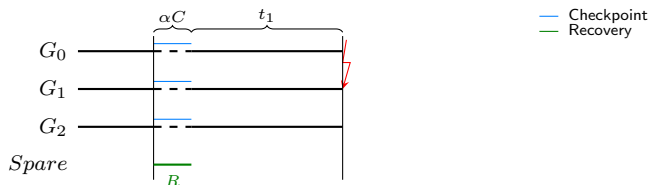
Waste in case of failures



Waste in case of failures

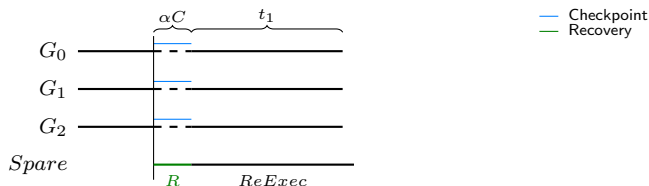


Waste in case of failures



- The spare group restarts from the checkpoint

Waste in case of failures



- The spare group restarts from the checkpoint
- $ReExec = \alpha C + t_1$

Waste in case of failures



- The non failed groups checkpoint

Waste in case of failures

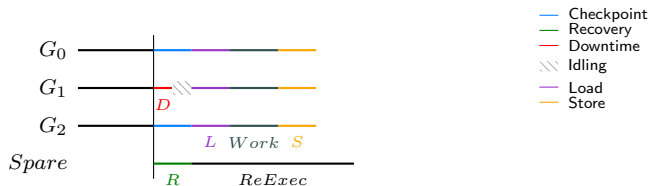


- They load the other application from its checkpoint

Waste in case of failures

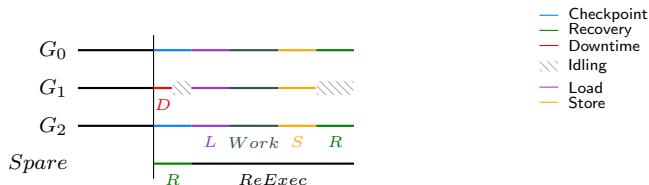


Waste in case of failures



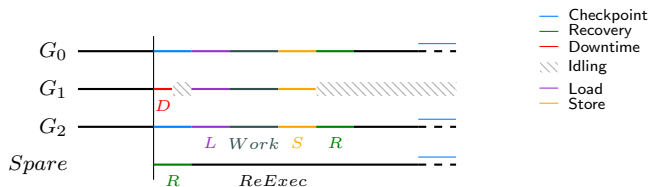
- They store the work they did

Waste in case of failures

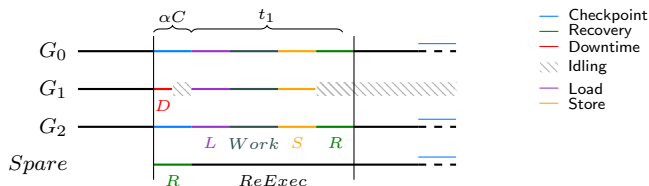


- They reload the other application
- The faulty group becomes the spare group

Waste in case of failures

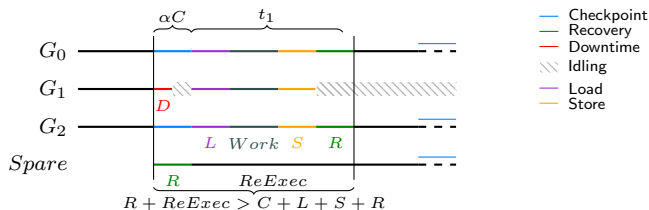


Platform scenario



$$ReExec = \frac{\alpha C + t_1}{\rho}$$

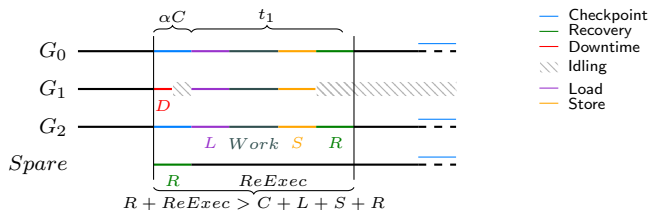
Platform scenario



$$ReExec = \frac{\alpha C + t_1}{\rho}$$

$$R + ReExec > C + L + S + R$$

Platform scenario

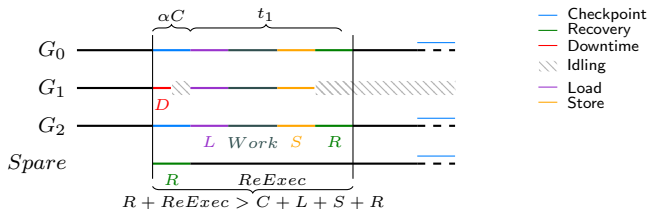


$$ReExec = \frac{\alpha C + t_1}{\rho}$$

$$R + ReExec > C + L + S + R$$

$$ReExec > C + L + S$$

Platform scenario



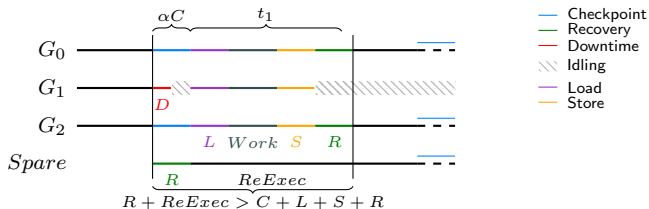
$$ReExec = \frac{\alpha C + t_1}{\rho}$$

$$R + ReExec > C + L + S + R$$

$$ReExec > C + L + S$$

$$t_1 > \rho(C + L + S) - \alpha C = Z$$

Platform scenario



$$ReExec = \frac{\alpha C + t_1}{\rho}$$

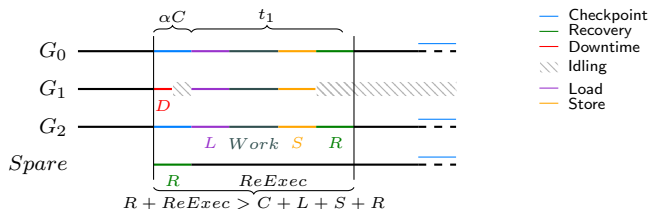
$$R + ReExec > C + L + S + R$$

$$ReExec > C + L + S$$

$$t_1 > \rho(C + L + S) - \alpha C = Z$$

$$t_1 > Z$$

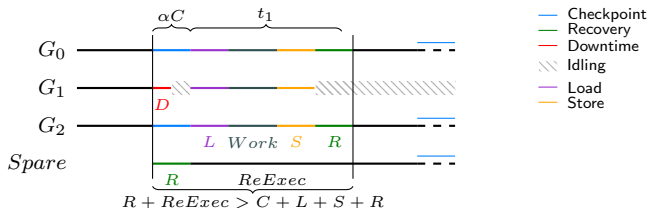
Platform scenario



- $t_1 > Z$ with a probability $\frac{T - Z}{T}$

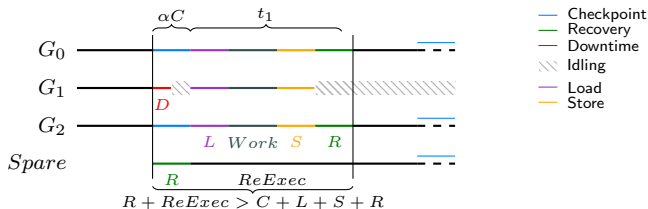
The G groups loose $L + C + S + R = X$

Platform scenario



- $t_1 > Z$ with a probability $\frac{T - Z}{T}$
The G groups loose $L + C + S + R = X$
- $t_1 \leq Z$ with a probability $\frac{Z}{T}$:
 - The G groups loose $R + ReExec$
 - The expectation of t_1 is $\frac{Z}{2}$

Platform scenario



- $t_1 > Z$ with a probability $\frac{T - Z}{T}$

The G groups loose $L + C + S + R = X$

- $t_1 \leq Z$ with a probability $\frac{Z}{T}$:

- The G groups loose $R + ReExec$
- The expectation of t_1 is $\frac{Z}{2}$

$$\text{WASTE}[fail] = \frac{1}{\mu} \left[\frac{T - Z}{T} \times X + \frac{Z}{T} \times \left(R + \frac{X - R}{2} + \frac{\alpha C}{2\rho} \right) \right]$$

Outline

- ① Protocol overhead
- ② Application scenario
- ③ Platform scenario
- ④ Instantiating the model
- ⑤ Simulations

Model instantiation

- ① Applications
 - 2D-stencil
 - Matrix product
- ② Platforms
 - K-Computer
 - ExaScaleFat

Impact of message logging on checkpoint size

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

$$C = C_0(1 + \beta \text{WORK}) \Leftrightarrow \beta = \frac{C - C_0}{C_0 \text{WORK}}$$

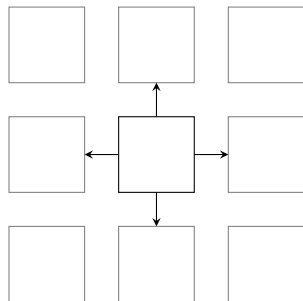
$$\text{WORK} = \lambda(T - (1 - \alpha)GC)$$

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

Computing β for Stencil-2D

$$C = C_0 + \textit{Logged_Msg} = C_0(1 + \beta \textit{WORK})$$

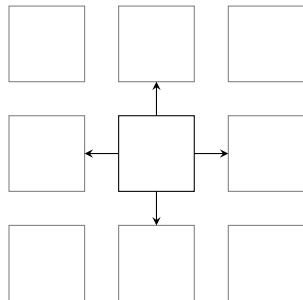
- $C_0 = \frac{\textit{Mem}}{G}$
- Real matrix $n \times n$
- $\textit{Mem} = 8n^2$



Computing β for Stencil-2D

$$C = C_0 + \text{Logged_Msg} = C_0(1 + \beta \text{WORK})$$

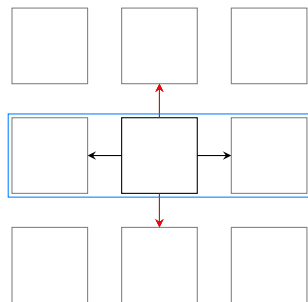
- $C_0 = \frac{\text{Mem}}{G}$
- s_p : speed of the process
- Block update: 9 floating points operations
- Each process holds a block of size b
- $\text{Work} = \frac{9b^2}{s_p}$



Computing β for Stencil-2D

$$C = C_0 + \text{Logged_Msg} = C_0(1 + \beta \text{WORK})$$

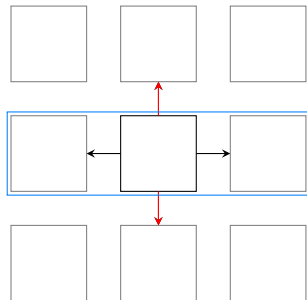
- $C_0 = \frac{\text{Mem}}{G}$
- $\text{Work} = \frac{9b^2}{s_p}$
- 1 group = 1 line
- Each process sends a block to its 4 neighbors
 - 2 out of the 4 messages are logged



Computing β for Stencil-2D

$$C = C_0 + \text{Logged_Msg} = C_0(1 + \beta \text{WORK})$$

- $C_0 = \frac{\text{Mem}}{G}$
- $\text{Work} = \frac{9b^2}{s_p}$
 - 2 out of the 4 messages are logged
 - $\beta = \frac{2s_p}{9b^3}$



Outline

① Protocol overhead

② Application scenario

③ Platform scenario

④ Instantiating the model

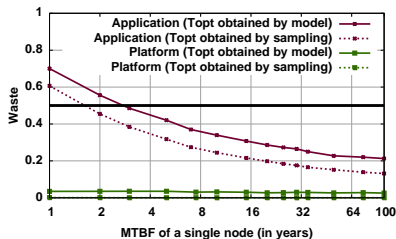
⑤ Simulations

Simulation parameters

- Failure distribution: Weibull, $k = 0.7$
- Failure free execution on each process: 4 days
- Time-out: 1 year
- No assumption on failures
- $\alpha = 0.3$, $\rho = 1.5$, $\lambda = 0.98$
- Each point: average over 20 randomly generated instances
- Computed period and best period:
 - Generate 480 periods in the neighborhood of the period from the model
 - Numerically evaluate the best one through simulations

Waste comparison

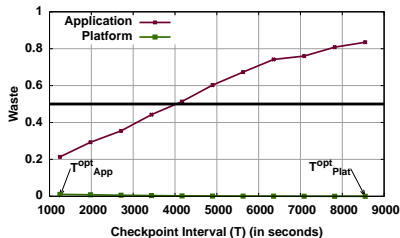
- Solid line: Computed period
- Dotted line: Best Period



Matrix-Product waste on K-Computer

Checkpointing period impact

- T_{App}^{opt} : Application-scenario optimal period
- T_{Plat}^{opt} : Platform-scenario optimal period



2D-Stencil waste on ExaScale Fat for an MTBF of 20 years

Outline

- ① Protocol overhead
- ② Application scenario
- ③ Platform scenario
- ④ Instantiating the model
- ⑤ Simulations

Conclusion and future work

- Overlap the idle time created by recovery periods
- Analytical model
 - ① Application scenario:
 - No overlap
 - $G + 1$ groups
 - ② Platform scenario:
 - Overlap with another application
 - G groups + 1 spare group
- Improve platform efficiency
 - A spare group
 - Application optimal period
- Future work:
 - Energy efficiency