

Unified Model for Assessing Checkpointing Protocols at Extreme-Scale

George BOSILCA¹, Aurélien BOUTEILLER¹,
Elisabeth BRUNET², Franck CAPPELLO³,
Jack DONGARRA¹, Amina GUERMOUCHE⁴,
Thomas HÉRAULT¹, Yves ROBERT^{1,4},
Frédéric VIVIEN⁴, and Dounia ZAIDOUNI⁴

1. University of Tennessee Knoxville, USA
2. Telecom SudParis, France
3. INRIA & University of Illinois at Urbana Champaign, USA
4. Ecole Normale Supérieure de Lyon & INRIA, France

Pittsburgh, June 28, 2012

Motivation

- **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 protocols
 - ② Hierarchical protocols

Which checkpointing protocol to use?

Coordinated checkpointing

- 😊 No risk of cascading rollbacks
- 😊 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
- 😊 Only processors from failed group need to roll back
- 😊 Faster re-execution with logged messages
- 😊 Rumor: Should scale to very large platforms

Outline

① Protocol Overhead

- Coordinated checkpointing
- Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

- Applications
- Platforms

④ Experimental results

- Plotting formulas
- Simulations

Outline

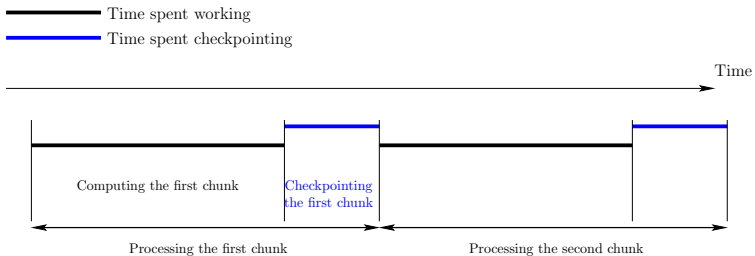
- 1 Protocol Overhead
- 2 Accounting for message logging
- 3 Instantiating the model
- 4 Experimental results

Framework

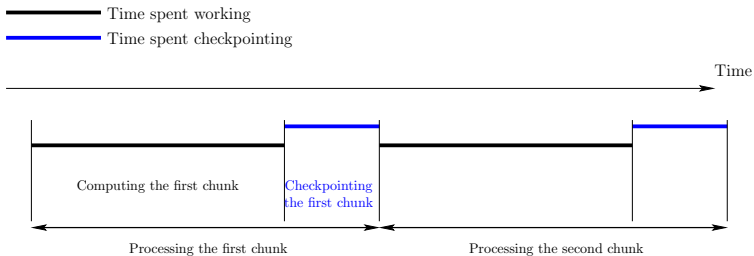
- Periodic checkpointing policies (of period T)
- Independent and identically distributed failures
- Platform failure inter-arrival time: μ
- Tightly-coupled application:
 progress \Leftrightarrow **all processors available**
- First-order approximation: at most one failure within a period

Waste: fraction of time not spent for useful computations

Checkpointing cost

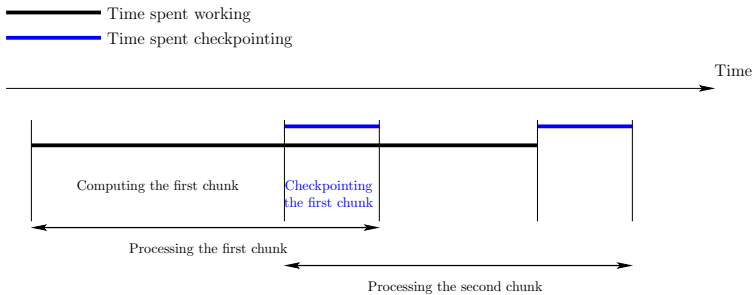


Checkpointing cost



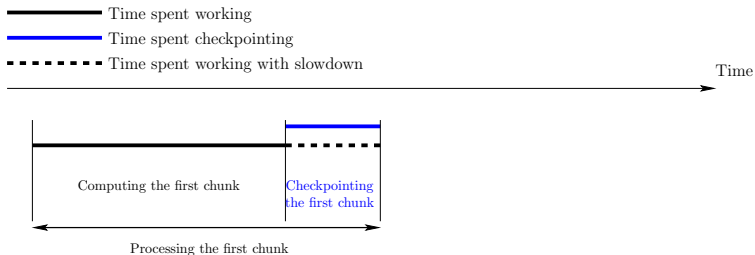
Blocking model: while a checkpoint is taken, no computation can be performed

Checkpointing cost



Non-blocking model: while a checkpoint is taken, computations are not impacted (e.g., first copy state to RAM, then copy RAM to disk)

Checkpointing cost



General model: while a checkpoint is taken, computations are slowed-down: during a checkpoint of duration C , the same amount of computation is done as during a time αC without checkpointing ($0 \leq \alpha \leq 1$).

① Protocol Overhead

- Coordinated checkpointing

- Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

- Applications

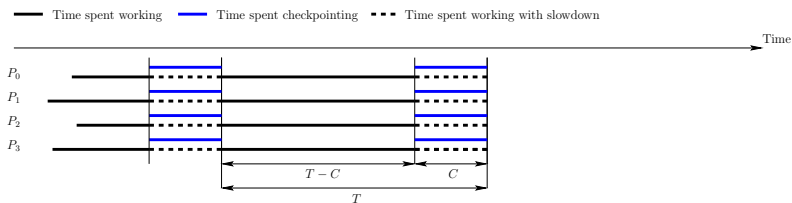
- Platforms

④ Experimental results

- Plotting formulas

- Simulations

Waste in absence of failures

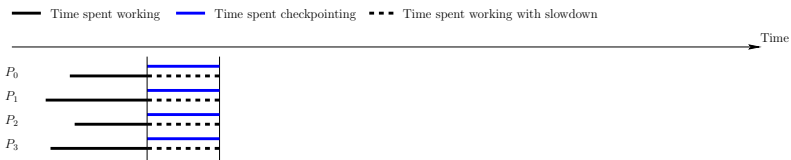


Time elapsed since last checkpoint: T

Amount of computation saved: $(T - C) + \alpha C$

$$\text{WASTE}_{\text{coord-nofailure}} = \frac{T - ((T - C) + \alpha C)}{T} = \frac{(1 - \alpha)C}{T}$$

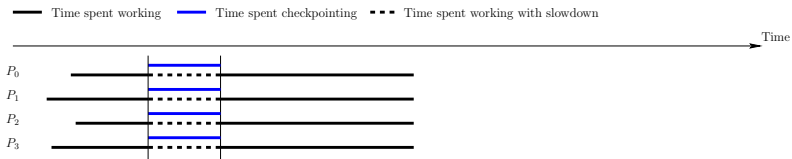
Waste due to failures



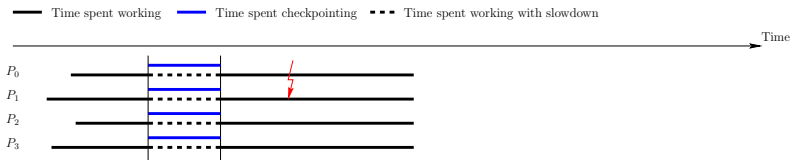
Failure can happen

- ① During computation phase
- ② During checkpointing phase

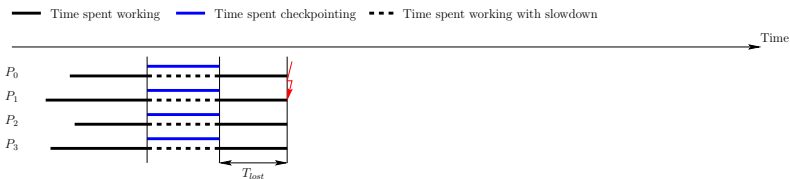
Waste due to failures



Waste due to failures

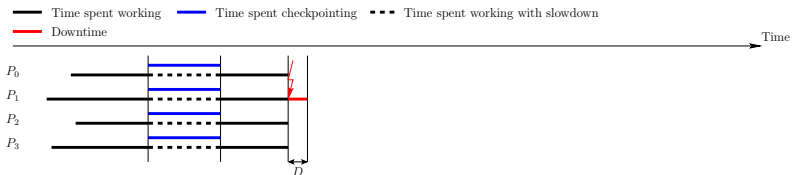


Waste due to failures

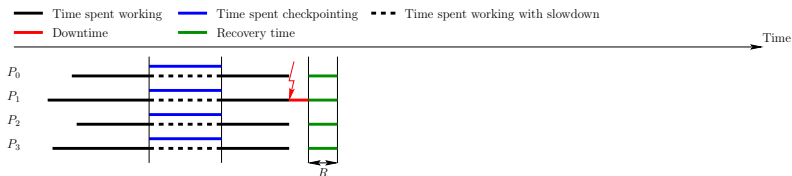


Coordinated checkpointing protocol: when one processor is victim of a failure, all processors lose their work and must roll back to last checkpoint

Waste due to failures in computation phase

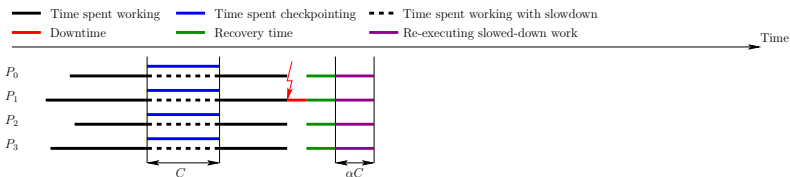


Waste due to failures in computation phase



Coordinated checkpointing protocol: All processors must recover from last checkpoint

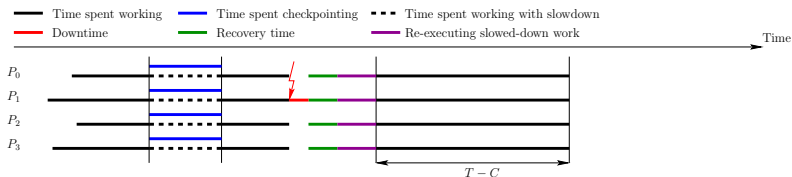
Waste due to failures in computation phase



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

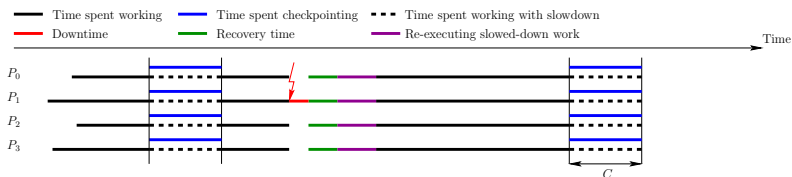
But no checkpoint is taken in parallel, hence this re-computation is faster than the original computation

Waste due to failures in computation phase



Re-execute the computation phase

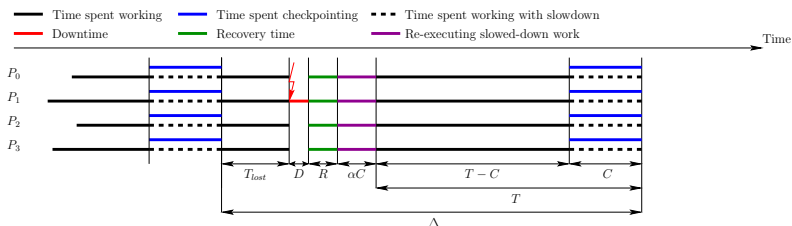
Waste due to failures in computation phase



Finally, the checkpointing phase is executed

First-order approximation: we assume that no other failure occurs during the re-execution

Waste due to failures in computation phase



$$\text{RE-EXEC: } \Delta - T = T_{lost} + \alpha C$$

$$\text{Expectation: } T_{lost} = \frac{1}{2}(T - C)$$

$$\text{RE-EXEC}_{\text{coord-fail-in-work}} = \frac{T - C}{2} + \alpha C$$

Waste due to failures

- Failure in the computation phase (probability: $\frac{T-C}{T}$)

$$\text{RE-EXEC}_{\text{coord-fail-in-work}} = \frac{T-C}{2} + \alpha C$$

- Failure in the checkpointing phase (probability: $\frac{C}{T}$)

$$\text{RE-EXEC}_{\text{coord-fail-in-checkpoint}} = T - \frac{C}{2} + \alpha C$$

$$\begin{aligned} \frac{T-C}{T} \left(\frac{T-C}{2} + \alpha C \right) + \frac{C}{T} \left(T - \frac{C}{2} + \alpha C \right) \\ = \alpha C + \frac{T}{2} \end{aligned}$$

Overall waste

$$\begin{aligned} \text{WASTE}_{\text{coord}} &= \text{WASTE}_{\text{coord-nofailure}} + \frac{1}{\mu}(D + R + \text{RE-EXEC}_{\text{coord}}) \\ &= \frac{(1 - \alpha)C}{T} + \frac{1}{\mu} \left(D + R + \alpha C + \frac{T}{2} \right) \end{aligned}$$

Minimize $\text{WASTE}_{\text{coord}}$ subject to:

- $C \leq T$ (by construction)
- $T \leq 0.1\mu$ ($\Rightarrow \text{Proba}(\text{Poisson}(\frac{T}{\mu}) \geq 2) \leq 0.005$)

① Protocol Overhead

Coordinated checkpointing

Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

Applications

Platforms

④ Experimental results

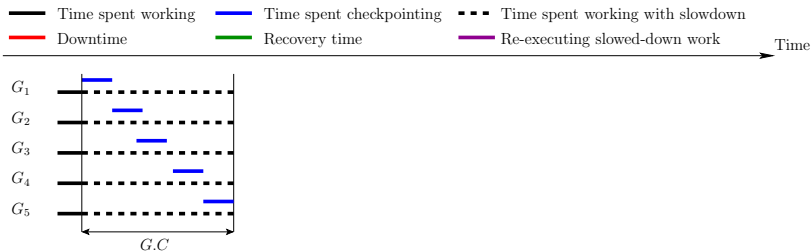
Plotting formulas

Simulations

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

Impact of checkpointing

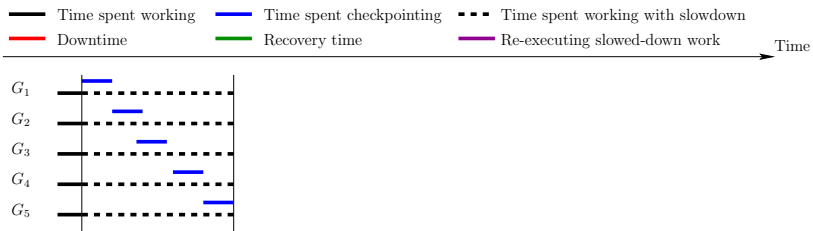


When a group checkpoints, its own computation speed is slowed-down

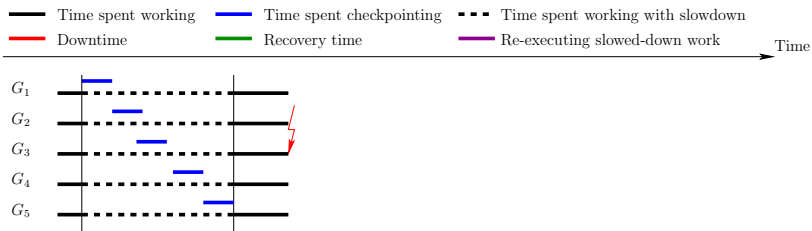
This holds for all groups because of the tightly-coupled assumption

$$\text{WASTE} = \frac{T - \text{WORK}}{T} \text{ where } \text{WORK} = T - (1 - \alpha)GC(q)$$

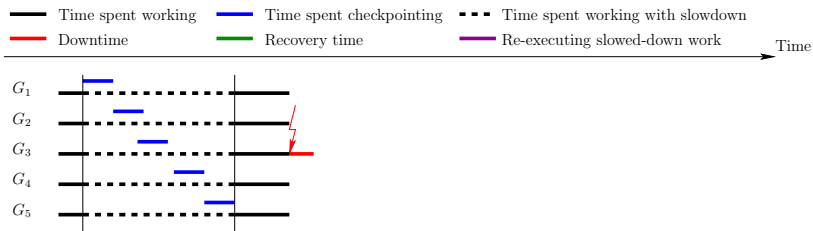
Impact of checkpointing



Impact of checkpointing

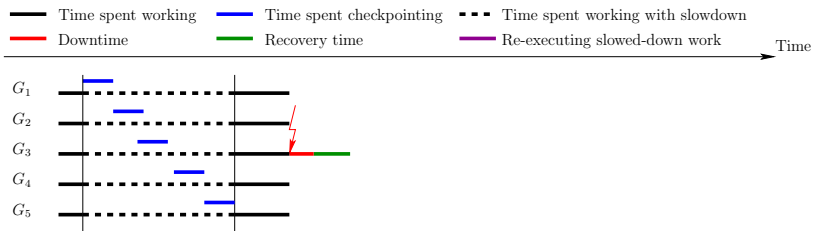


Impact of checkpointing



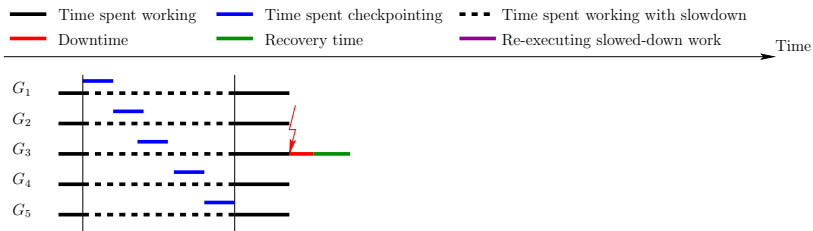
Tightly-coupled model: while one group is in downtime, none can work

Failure during computation phase



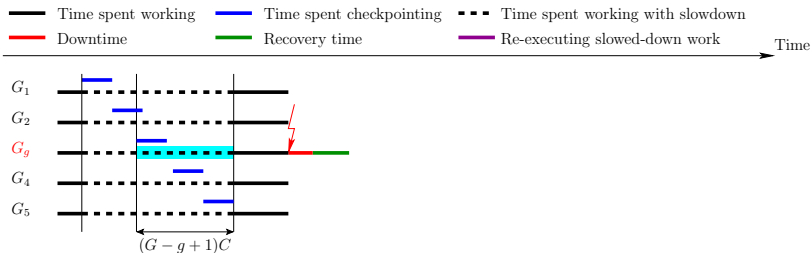
Tightly-coupled model: while one group is in recovery, none can work

Failure during computation phase



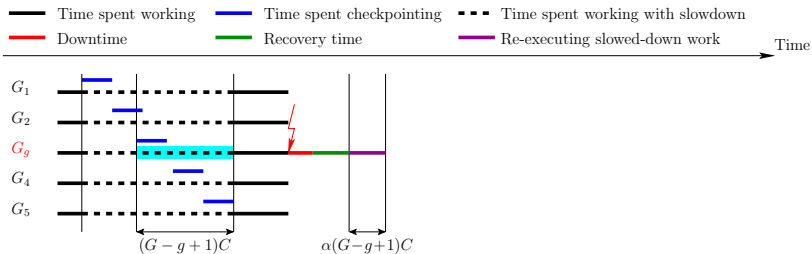
Groups must have completed the same amount of work in between two consecutive checkpoints, independently of the fact that a failure may or may not have happened on the platform in between these checkpoints. Hence, no checkpointing is possible during the rollback.

Failure during computation phase



Redo work done during previous checkpointing phase and that was destroyed by the failure

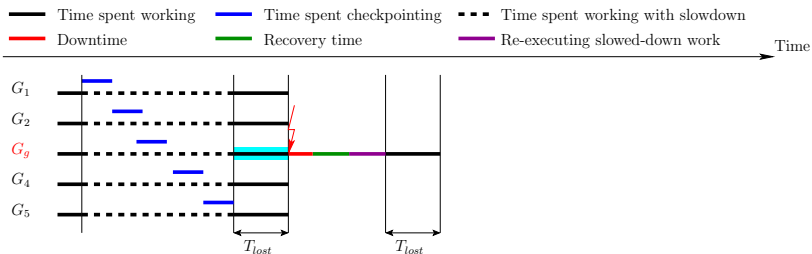
Failure during computation phase



Redo work done during previous checkpointing phase and that was destroyed by the failure

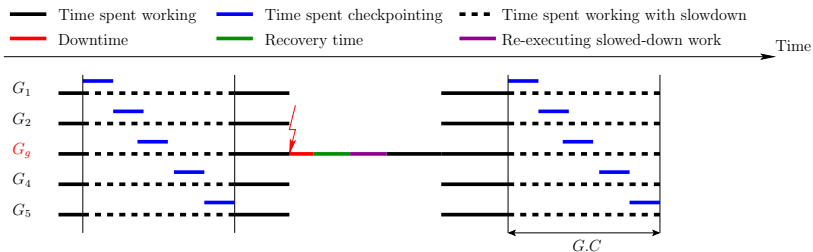
But no checkpoint is taken in parallel, hence this re-computation is faster than the original computation

Failure during computation phase



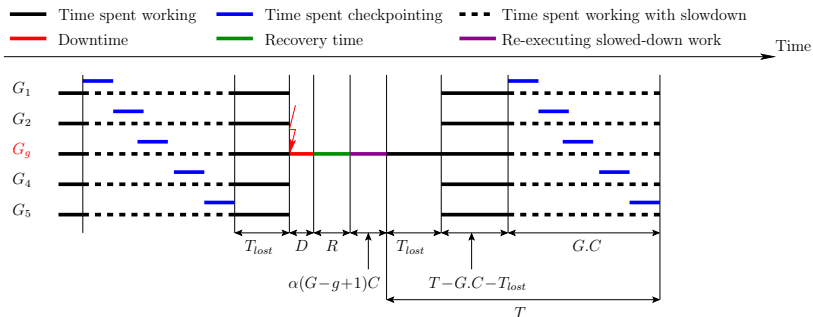
Redo work done in computation phase and that was destroyed by the failure

Failure during computation phase



Finally, perform checkpointing phase

Failure during computation phase

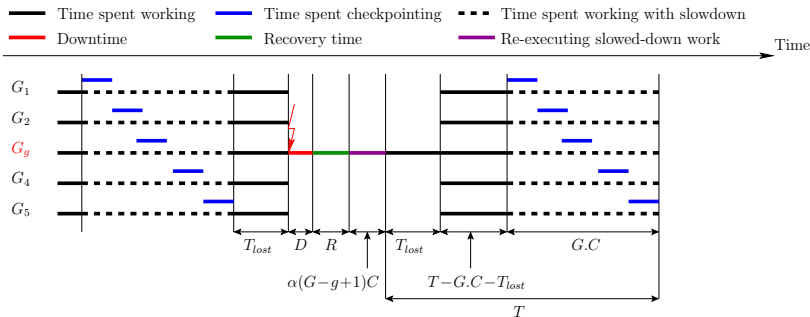


$$\text{RE-EXEC: } T_{lost} + \alpha(G - g + 1)C$$

$$\text{Expectation: } T_{lost} = \frac{1}{2}(T - G.C)$$

$$\text{Approximated RE-EXEC: } \frac{T - G.C}{2} + \alpha(G - g + 1)C$$

Failure during computation phase

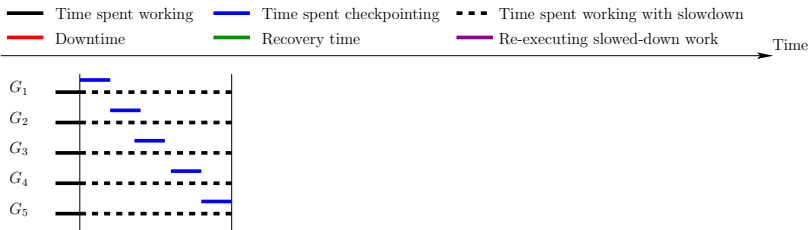


$$\text{Approximated RE-EXEC: } \frac{T - G.C}{2} + \alpha(G - g + 1)C$$

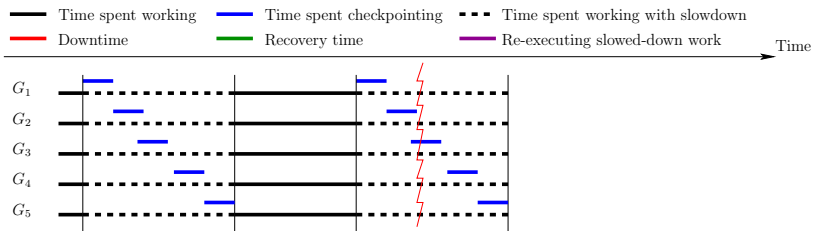
Average approximated RE-EXEC:

$$\begin{aligned} & \frac{1}{G} \sum_{g=1}^G \left[\frac{T - G.C(q)}{2} + \alpha(G - g + 1)C(q) \right] \\ &= \frac{T - G.C(q)}{2} + \alpha \frac{G+1}{2} C \end{aligned}$$

Failure during checkpointing phase



Failure during checkpointing phase



When does the failing group fail?

- ① Before starting its own checkpoint
- ② While taking its own checkpoint
- ③ After completing its own checkpoint

Average waste for failures during checkpointing phase

Average RE-EXEC when the failing-group g fails

Overall average RE-EXEC: $\text{RE-EXEC}_{ckpt} =$

$$\frac{1}{G} \left((g-1) \cdot \text{RE-EXEC}_{before_ckpt} + 1 \cdot \text{RE-EXEC}_{during_ckpt} + (G-g) \cdot \text{RE-EXEC}_{after_ckpt} \right)$$

Average over all groups:

$$\text{AVG_RE-EXEC}_{ckpt} = \frac{G+1}{2G} T + \frac{\alpha C(q)(G+3)}{2} + \frac{C(q)(1-2\alpha)}{2G} - \frac{C(q)(G+1)}{2}$$

Average waste

$$\begin{aligned} \text{WASTE}_{\text{hierarch}} &= \frac{T - \text{WORK}}{T} + \frac{1}{\mu} \left(D(q) + R(q) + \text{RE-EXEC} \right) \\ &= \frac{1}{2\mu T} \times \left(\begin{array}{l} T^2 \\ +GC(q) [(1 - \alpha)(2\mu - T) + (2\alpha - 1)C(q)] \\ +T [2(D(q) + R(q)) + (\alpha + 1)C(q)] \\ +(1 - 2\alpha)C(q)^2 \end{array} \right) \end{aligned}$$

Minimize $\text{WASTE}_{\text{hierarch}}$ subject to:

- $GC(q) \leq T$ (by construction)
- $T \leq 0.1\mu$ ($\Rightarrow \text{Proba}(\text{Poisson}(\frac{T}{\mu}) \geq 2) \leq 0.005$)

Outline

- ① Protocol Overhead
- ② Accounting for message logging
- ③ Instantiating the model
- ④ Experimental results

Impact on work

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

$$\text{WASTE}_{\text{hierarch}} = \frac{T - \lambda \text{WORK}}{T} + \frac{1}{\mu} \left(D(q) + R(q) + \frac{\text{RE-EXEC}}{\rho} \right)$$

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

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

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

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

- Constraint $GC(q) \leq T$ translates into

$$GC_0(q)\beta\lambda\alpha \leq 1 \text{ and } T \geq \frac{GC_0(q)}{1 - GC_0(q)\beta\lambda\alpha}$$

Outline

- ① Protocol Overhead
- ② Accounting for message logging
- ③ Instantiating the model
- ④ Experimental results

Three case studies

Coord-IO

Coordinated approach: $C = C_{\text{Mem}} = \frac{\text{Mem}}{b_{io}}$

where Mem is the memory footprint of the application

Hierarch-IO

Several (large) groups, *I/O-saturated*

⇒ groups checkpoint sequentially

$$C_0(q) = \frac{C_{\text{Mem}}}{G} = \frac{\text{Mem}}{Gb_{io}}$$

Hierarch-Port

Very large number of smaller groups, *port-saturated*

⇒ some groups checkpoint in parallel

Groups of q_{\min} processors, where $q_{\min} b_{port} \geq b_{io}$

① Protocol Overhead

Coordinated checkpointing

Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

Applications

Platforms

④ Experimental results

Plotting formulas

Simulations

Three applications

- 1 2D-stencil
- 2 3D-Stencil
 - Plane
 - Line
- 3 Matrix product

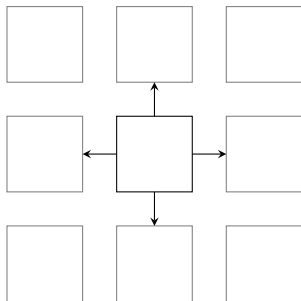
Computing β for Stencil-2D

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

Real $n \times n$ matrix and $p \times p$ grid

$$\text{Work} = \frac{9b^2}{s_p}, \quad b = n/p$$

Each process sends a block to its 4 neighbors



Computing β for Stencil-2D

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

Real $n \times n$ matrix and $p \times p$ grid

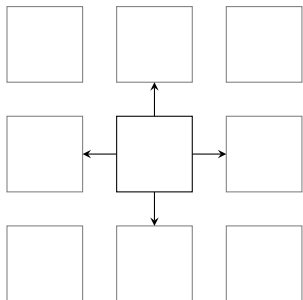
$$\text{Work} = \frac{9b^2}{s_p}, \quad b = n/p$$

Each process sends a block to its 4 neighbors

HIERARCH-IO:

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

- $\beta = \frac{2s_p}{9b^3}$



Computing β for Stencil-2D

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

Real $n \times n$ matrix and $p \times p$ grid

$$\text{Work} = \frac{9b^2}{s_p}, \quad b = n/p$$

Each process sends a block to its 4 neighbors

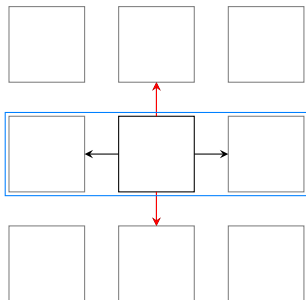
HIERARCH-IO:

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

- $\beta = \frac{2s_p}{9b^3}$

HIERARCH-PORT:

- β doubles



Three applications: 2) 3D-stencil

- Real matrix of size $n \times n \times n$ partitioned across a $p \times p \times p$ processor grid
- Each processor holds a cube of size $b = n/p$
- At each iteration:
 - average each matrix element with its 27 closest neighbors
 - exchange the six faces of its cube
- (Parallel) work for one iteration is $WORK = \frac{27b^3}{s_p}$

Three hierarchical variants

- 1 HIERARCH-IO-PLANE: group = horizontal plane of size p^2 :

$$\beta = \frac{2s_p}{27b^3}$$
- 2 HIERARCH-IO-LINE: group = horizontal line of size p :

$$\beta = \frac{4s_p}{27b^3}$$
- 3 HIERARCH-PORT: groups of size q_{min} : $\beta = \frac{6s_p}{27b^3}$

① Protocol Overhead

- Coordinated checkpointing
- Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

- Applications
- Platforms

④ Experimental results

- Plotting formulas
- Simulations

Four platforms: basic characteristics

Name	Number of cores	Number of processors p_{total}	Number of cores per processor	Memory per processor	I/O Network Bandwidth (b_{io})	
					Read	Write
Titan	299,008	16,688	16	32GB	300GB/s	300GB/s
K-Computer	705,024	88,128	8	16GB	150GB/s	96GB/s
Exascale-Slim	1,000,000,000	1,000,000	1,000	64GB	1TB/s	1TB/s
Exascale-Fat	1,000,000,000	100,000	10,000	640GB	1TB/s	1TB/s

Four platforms: 2D-STENCIL and MATRIX-PRODUCT

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

Four platforms: 3D-STENCIL

Name	Scenario	G	β for 3D-STENCIL
Titan	COORD-IO	1	/
	HIERARCH-IO-PLANE	26	0.001476
	HIERARCH-IO-LINE	675	0.002952
	HIERARCH-PORT	1,246	0.004428
K-Computer	COORD-IO	1	/
	HIERARCH-IO-PLANE	44	0.003422
	HIERARCH-IO-LINE	1,936	0.006844
	HIERARCH-PORT	17,626	0.010266
Exascale-Slim	COORD-IO	1	/
	HIERARCH-IO-PLANE	100	0.003952
	HIERARCH-IO-LINE	10,000	0.007904
	HIERARCH-PORT	200,000	0.011856
Exascale-Fat	COORD-IO	1	/
	HIERARCH-IO-PLANE	46	0.001834
	HIERARCH-IO-LINE	2,116	0.003668
	HIERARCH-PORT	33,333	0.005502

Outline

- ① Protocol Overhead
- ② Accounting for message logging
- ③ Instantiating the model
- ④ Experimental results

① Protocol Overhead

Coordinated checkpointing

Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

Applications

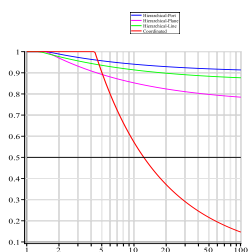
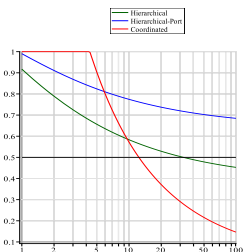
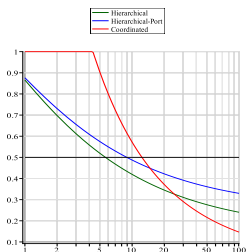
Platforms

④ Experimental results

Plotting formulas

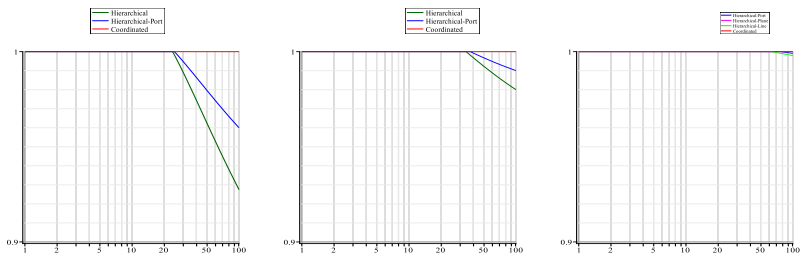
Simulations

Platform: Titan



Waste as a function of processor MTBF μ

Platform: K-Computer



Waste as a function of processor MTBF μ

Platform: Exascale

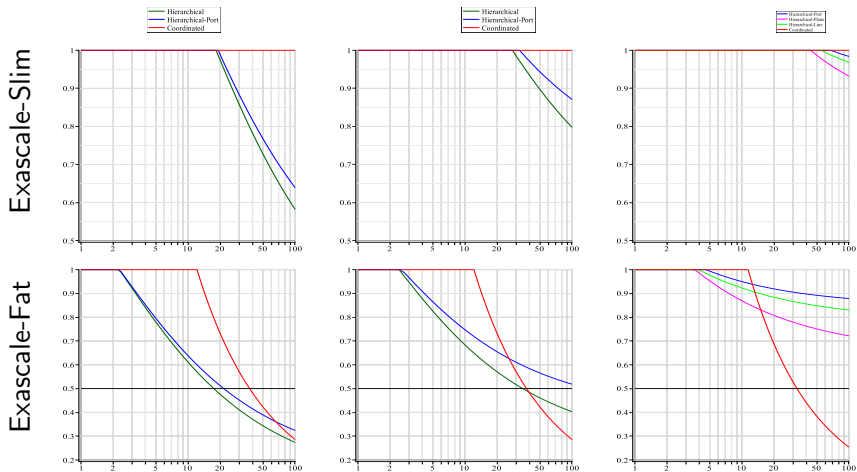
WASTE = 1 for all scenarios!!!

Platform: Exascale

WASTE = 1 for all scenarios!!!

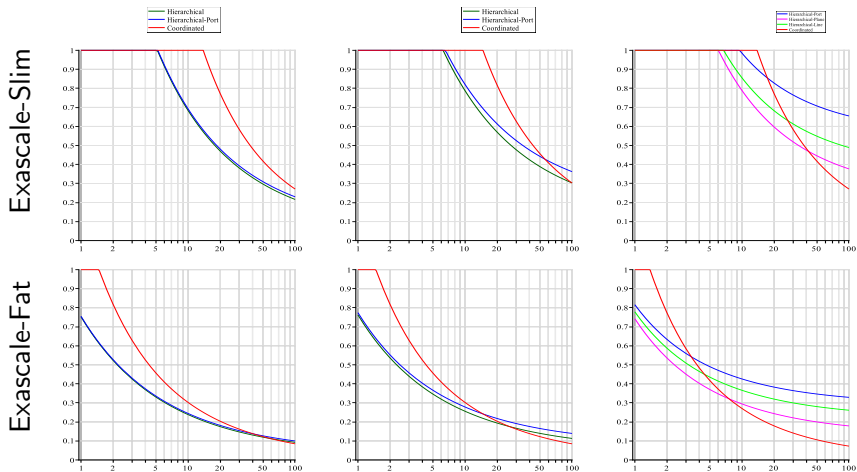
Goodbye Exascale?!

Platform: Exascale with $C = 1,000$



Waste as a function of processor MTBF μ , $C = 1,000$

Platform: Exascale with $C = 100$



Waste as a function of processor MTBF μ , $C = 100$

① Protocol Overhead

Coordinated checkpointing

Hierarchical checkpointing

② Accounting for message logging

③ Instantiating the model

Applications

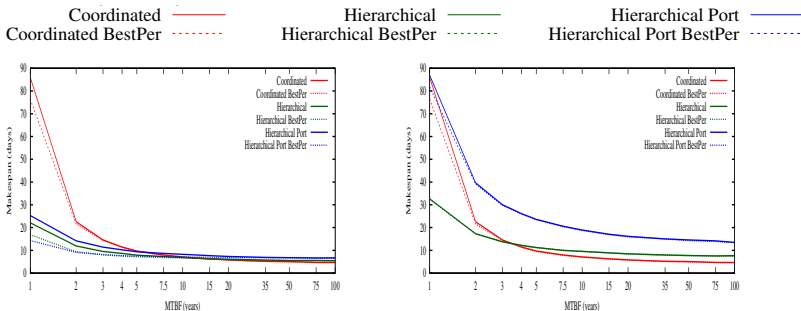
Platforms

④ Experimental results

Plotting formulas

Simulations

Platform: Titan



Makespan (in days) as a function of processor MTBF μ

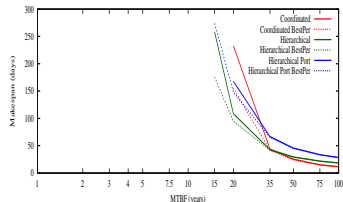
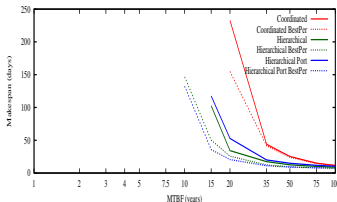
Platform: Exascale with $C = 1,000$

Exascale-Slim

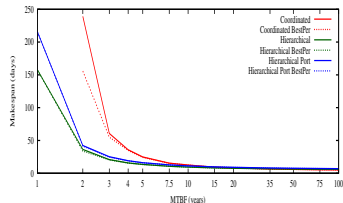
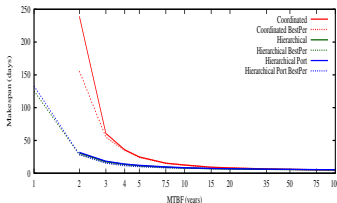
Coordinated ———
Coordinated BestPer - - - - -

Hierarchical ———
Hierarchical BestPer - - - - -

Hierarchical Port ———
Hierarchical Port BestPer - - - - -



Exascale-Fat



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

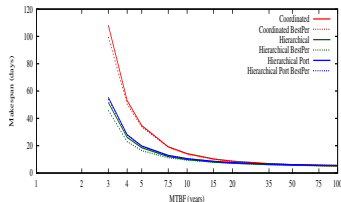
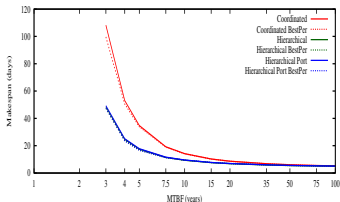
Platform: Exascale with $C = 100$

Exascale-Slim

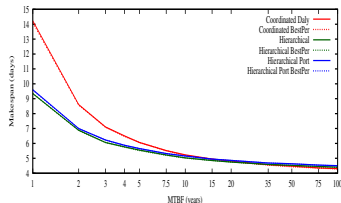
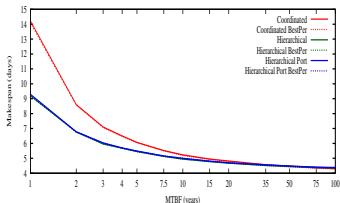
Coordinated —
Coordinated BestPer - - -

Hierarchical —
Hierarchical BestPer - - -

Hierarchical Port —
Hierarchical Port BestPer - - -



Exascale-Fat



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

Conclusion

- First attempt at analytical comparison of coordinated and hierarchical checkpointing protocols
- Classical models (Young, Daly) extended
 - Several new parameters (α , λ , ρ)
 - Message logging impact (β)
- Instantiation
 - Scenarios: COORD-IO, HIERARCH-IO, HIERARCH-PORT
 - Realistic application/platform combinations
- Future work: (partial) replication, prediction, energy