

Assessing the impact of ABFT & Checkpoint composite strategies

George Bosilca¹, Aurélien Bouteiller¹, Thomas Hérault¹,
Yves Robert^{1,2} and Jack Dongarra¹

1. University of Tennessee Knoxville, USA
2. École Normale Supérieure de Lyon & INRIA, France

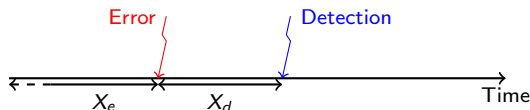
{bosilca,bouteill,herault,dongarra}@icl.utk.edu
yves.robert@ens-lyon.fr

November 26, 2013 - JLPC Workshop

- Fault prediction: checkpointing vs. migration (PPL)
- Model to assess checkpoint protocols (CCPE, online)
- Checkpointing and prediction (JPDC, online)
- In-memory checkpointing (APDCM'13)
- Multi-criteria: time vs resource utilization (Europar'13)
- Multi-criteria: time vs energy (PMBS'13)
- Silent errors, checkpoints & verifications (PRDC'13)

Detection latency

- Instantaneous error detection \Rightarrow fail-stop failures
- Silent errors (data corruption) \Rightarrow detection latency



Error and detection latency

- Last checkpoint may have saved an already corrupted state
- Even when saving k checkpoints: which one to roll back to?
- **Critical failure**: all checkpoints contain corrupted data

Coupling checkpointing and verification

- Verification mechanism of cost V
- Simplest idea: verify work before each checkpoint



V large compared to $w \Rightarrow$ large WASTE^{ff} , can we improve that?

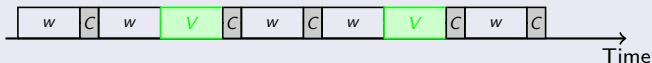
Coupling checkpointing and verification

- Verification mechanism of cost V
- Simplest idea: verify work before each checkpoint



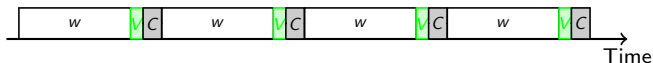
V large compared to $w \Rightarrow$ large WASTE^{ff} , can we improve that?

Is this better?



Coupling checkpointing and verification

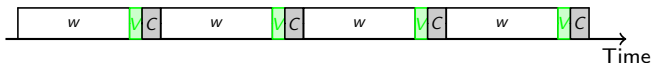
- Verification mechanism of cost V
- Simplest idea: verify work before each checkpoint



V small in front of $w \Rightarrow$ large $\text{WASTE}^{\text{fail}}$, can we improve that?

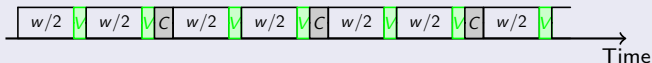
Coupling checkpointing and verification

- Verification mechanism of cost V
- Simplest idea: verify work before each checkpoint



V small in front of $w \Rightarrow$ large $\text{WASTE}^{\text{fail}}$, can we improve that?

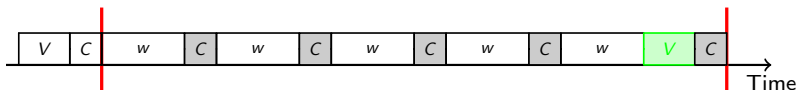
Is this better?



Coupling checkpointing and verification



Small cost V : 5 verifications for 1 checkpoint

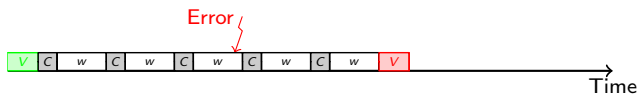


Large cost V : 5 checkpoints for 1 verification

More complicated periodic patterns? Different-size chunks?

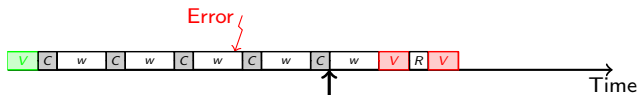
k checkpoints for 1 verification

Where did the error strike?



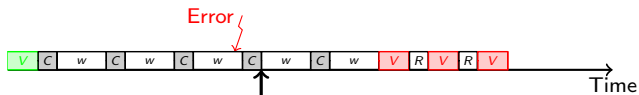
k checkpoints for 1 verification

Where did the error strike?



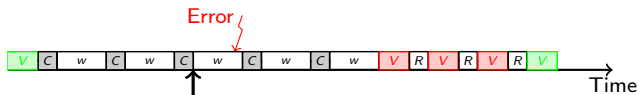
k checkpoints for 1 verification

Where did the error strike?



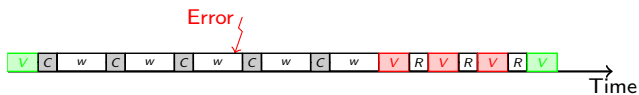
k checkpoints for 1 verification

Where did the error strike?



k checkpoints for 1 verification

Where did the error strike?



$$\text{RE-EXEC} = 2(w + C) + (w + V)$$

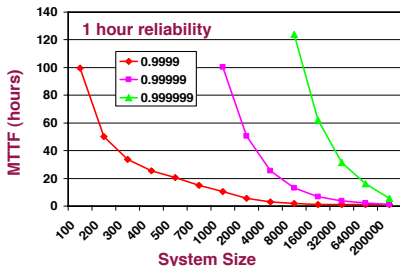
Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

Faults

- Assume independent failures
- Let N be the number of components (“System Size”)
- Let r be the probability of a component to operate for 1h
- Let R be the probability of the system to operate for 1h

$$R = r^N$$
$$R \approx \frac{1}{e^{\lambda N}}, \frac{1}{\lambda} = 1 - r$$

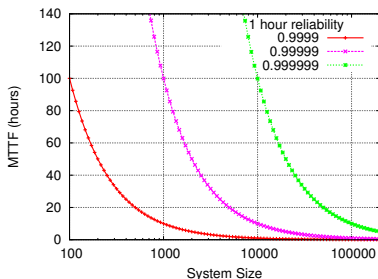


(Figure from Dan Reed “The Challenge of Complexity and Scale”)

Faults

- Assume independent failures
- Let N be the number of components (“System Size”)
- Let r be the probability of a component to operate for 1h
- Let R be the probability of the system to operate for 1h

$$R = r^N$$
$$R \approx \frac{1}{e^{\lambda N}}, \frac{1}{\lambda} = 1 - r$$

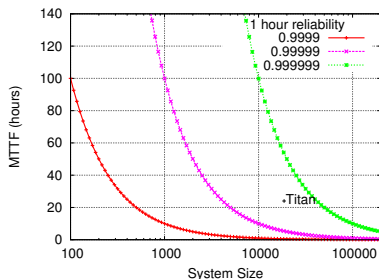


(Same figure with plotting software)

Faults

- Assume independent failures
- Let N be the number of components (“System Size”)
- Let r be the probability of a component to operate for 1h
- Let R be the probability of the system to operate for 1h

$$R = r^N$$
$$R \approx \frac{1}{e^{\lambda N}}, \frac{1}{\lambda} = 1 - r$$



(Same figure with plotting software)

Fault Tolerance Techniques

General Techniques

- Replication
- Rollback Recovery
 - Coordinated Checkpointing
 - Uncoordinated Checkpointing & Message Logging
 - Hierarchical Checkpointing

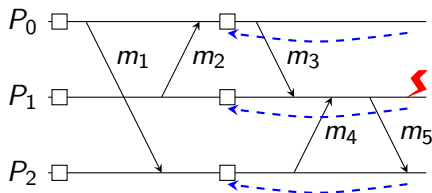
Application-Specific Techniques

- Algorithm Based Fault Tolerance (ABFT)
- Iterative Convergence



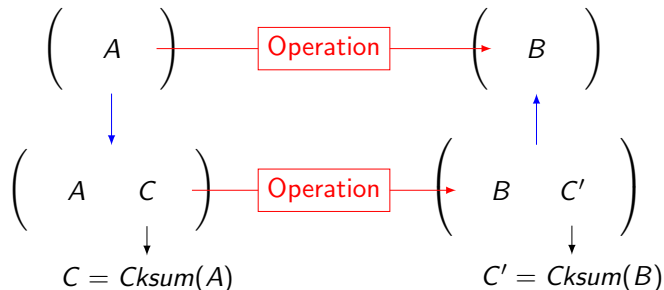
Coordinated Checkpointing and Rollback Recovery

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



- ☺ General technique (we assume preemptive checkpointing capability)
- ☹ All processors need to roll back
- ☹ All memory needs to be saved

Algorithm-Based Fault Tolerance



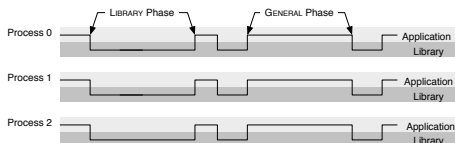
Principle of ABFT

- Input Data (A) and Result (B) are distributed
- *Operation* preserves *Checksum* properties
- Apply the operation on Data + Checksum (AC)
- In case of failure, recover the missing data by inversion of the checksum

Application

Typical Application

```
for( aninsanenummer ) {  
  /* Extract data from  
   * simulation, fill up  
   * matrix */  
  sim2mat();  
  
  /* Factorize matrix,  
   * Solve */  
  dgeqrf();  
  dsolve();  
  
  /* Update simulation  
   * with result vector */  
  vec2sim();  
}
```



Characteristics

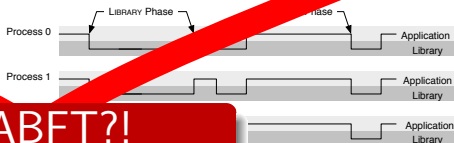
- ☺ Large part of (total) computation spent in factorization/solve
- Between LA operations:
 - ☹ use resulting vector / matrix with operations that do not preserve data checksums
 - ☹ modify data not covered by ABFT algorithms

Application

Typical Application

```
for( aninsanenum  
/* Extract data  
* simulation,  
* matrix */  
sim2mat();  
  
/* Factorize matrix,  
* Solve */  
dgeqrf();  
dsolve();  
  
/* Update simulation  
* with result vector */  
vec2sim();  
}
```

Goodbye ABFT?!



Characteristics

- 😊 Large part of (total) computation spent in factorization/solve
- Between LA operations:
 - ☹ use resulting vector / matrix with operations that do not preserve data checksums
 - ☹ modify data not covered by ABFT algorithms

Problem Statement

Typical

```
for ( z  
/* I  
* s  
* r  
sim2
```

How to use fault tolerant operations () within a non-fault tolerant (**) application? (***)*

```
/* I (*) ABFT, or other application-specific FT  
* s (**) Or within an application that does not have the same kind of FT  
dge (**) And keep the application globally fault tolerant...  
dsolve (,,
```

```
/* Update simulation  
* with result vector */  
vec2sim ();  
}
```

- ☹ use resulting vector / matrix with operations that do not preserve data checksums
- ☹ modify data not covered by ABFT algorithms

- Application
Library

- Application
Library

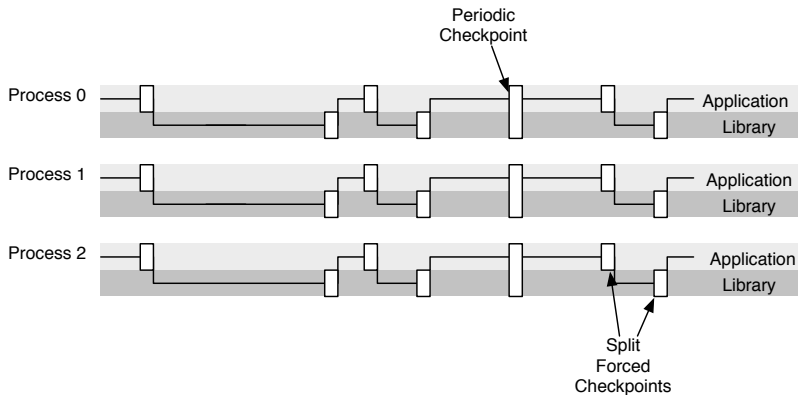
- Application
Library

Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

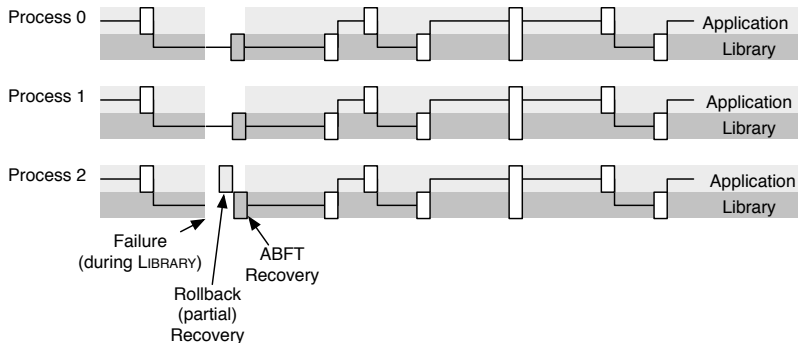
ABFT&PERIODICCKPT

ABFT&PERIODICCKPT: no failure



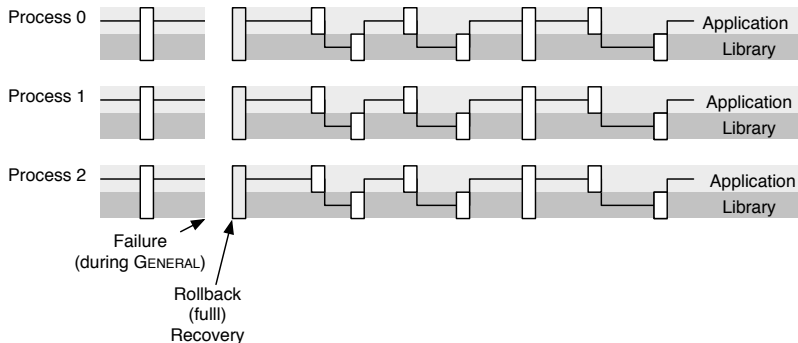
ABFT&PERIODICKPT

ABFT&PERIODICKPT: failure during LIBRARY phase

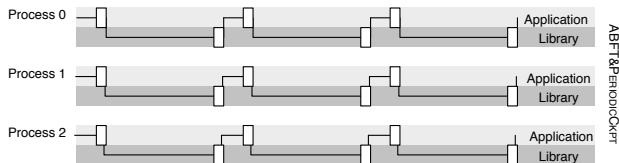


ABFT&PERIODICKPT

ABFT&PERIODICKPT: failure during GENERAL phase



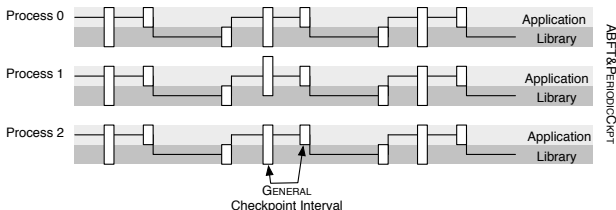
ABFT&PERIODICCKPT: Optimizations



ABFT&PERIODICCKPT: Optimizations

- If the duration of the `GENERAL` phase is too small: don't add checkpoints
- If the duration of the `LIBRARY` phase is too small: don't do ABFT recovery, remain in `GENERAL` mode
 - this assumes a performance model for the library call

ABFT&PERIODICCKPT: Optimizations



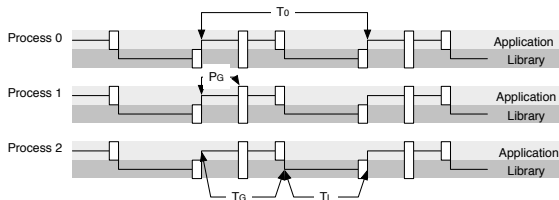
ABFT&PERIODICCKPT: Optimizations

- If the duration of the **GENERAL** phase is too small: don't add checkpoints
- If the duration of the **LIBRARY** phase is too small: don't do ABFT recovery, remain in **GENERAL** mode
 - this assumes a performance model for the library call

Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling**
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

A few notations



Times, Periods

T_0 : Duration of an Epoch (without FT)

$T_L = \alpha T_0$: Time spent in the LIBRARY phase

$T_G = (1 - \alpha) T_0$: Time spent in the GENERAL phase

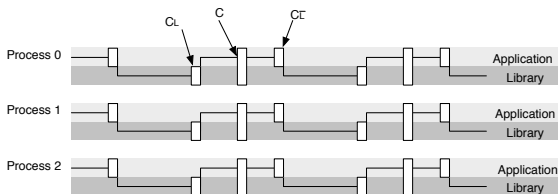
P_G : Periodic Checkpointing Period

$T_G^{ff}, T_G^{ff}, T_L^{ff}$: "Fault Free" times

t_G^{lost}, t_L^{lost} : Lost time (recovery overloads)

T_G^{final}, T_L^{final} : Total times (with faults)

A few notations



Costs

$C_L = \rho C$: time to take a checkpoint of the LIBRARY data set

$C_{\bar{L}} = (1 - \rho)C$: time to take a checkpoint of the GENERAL data set

$R, R_{\bar{L}}$: time to load a full / GENERAL data set checkpoint

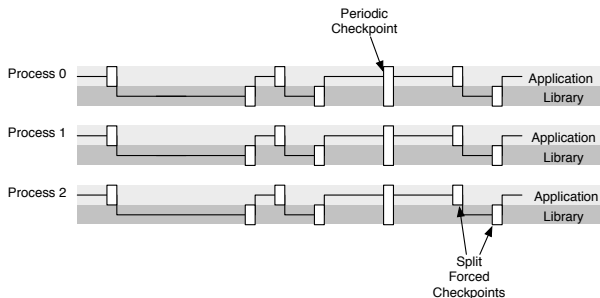
D : down time (time to allocate a new machine / reboot)

$\text{Recons}_{\text{ABFT}}$: time to apply the ABFT recovery

ϕ : Slowdown factor on the LIBRARY phase, when applying ABFT

GENERAL phase, fault free waste

GENERAL phase

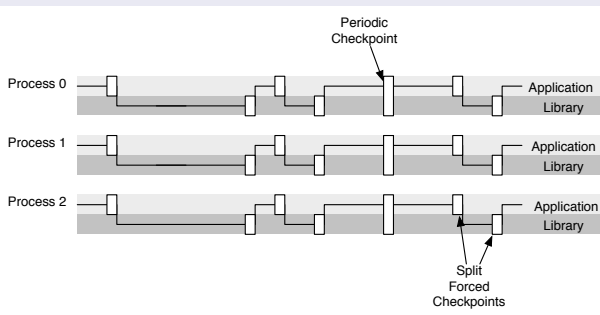


Without Failures

$$T_G^{\text{ff}} = \begin{cases} T_G + C_L & \text{if } T_G < P_G \\ \frac{T_G}{P_G - C} \times P_G & \text{if } T_G \geq P_G \end{cases}$$

LIBRARY phase, fault free waste

LIBRARY phase

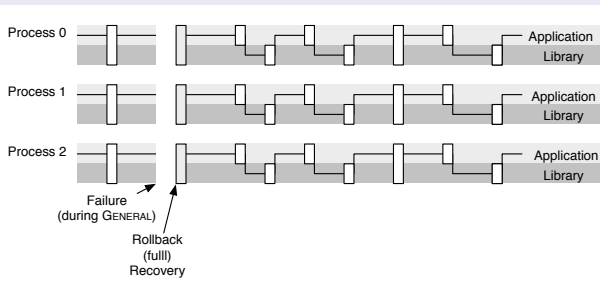


Without Failures

$$T_L^{ff} = \phi \times T_L + C_L$$

GENERAL phase, failure overhead

GENERAL phase

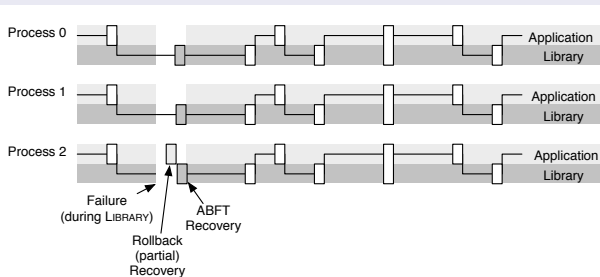


Failure Overhead

$$t_G^{\text{lost}} = \begin{cases} D + R + \frac{T_G^{\text{ff}}}{2} & \text{if } T_G < P_G \\ D + R + \frac{P_G}{2} & \text{if } T_G \geq P_G \end{cases}$$

LIBRARY phase, failure overhead

LIBRARY phase



Failure Overhead

$$t_L^{\text{lost}} = D + R_L + \text{Recons}_{\text{ABFT}}$$

Overall

Time (with overheads) of LIBRARY phase is constant (in P_G):

$$T_L^{\text{final}} = \frac{1}{1 - \frac{D+R_L+\text{Recons}_{\text{ABFT}}}{\mu}} \times (\alpha \times T_L + C_L)$$

Time (with overheads) of GENERAL phase accepts two cases:

$$T_G^{\text{final}} = \begin{cases} \frac{1}{1 - \frac{D+R+\frac{T_G+C_L}{2}}{\mu}} \times (T_G + C_L) & \text{if } T_G < P_G \\ \frac{T_G}{(1 - \frac{C}{P_G})(1 - \frac{D+R+\frac{P_G}{2}}{\mu})} & \text{if } T_G \geq P_G \end{cases}$$

Which is minimal in the second case, if

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

Waste

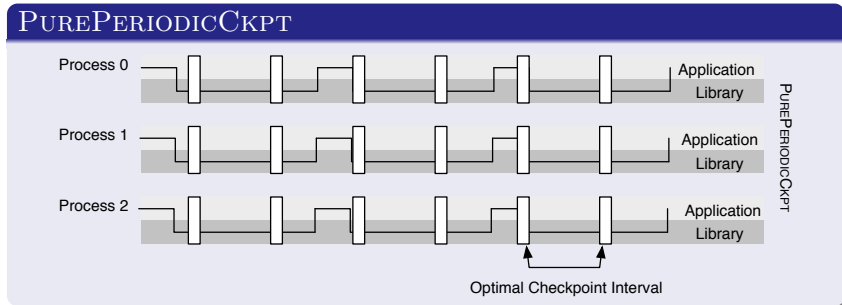
From the previous, we derive the waste, which is obtained by

$$\text{WASTE} = 1 - \frac{T_0}{T_G^{\text{final}} + T_L^{\text{final}}}$$

Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

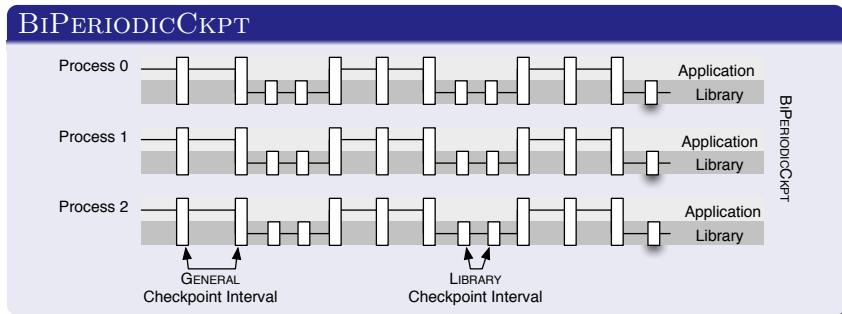
PUREPERIODICCKPT



Optimization

$$P_{PC}^{\text{opt}} = \sqrt{2C(\mu - D - R)}$$

BiPERIODICCKPT



Optimization

$$P_{BPC,G}^{\text{opt}} = \sqrt{2C(\mu - D - R)}$$

$$P_{BPC,L}^{\text{opt}} = \sqrt{2C_L(\mu - D - R)}$$

Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation**
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation**
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

Model & Simulations: PUREPERIODICKPT

PUREPERIODICKPT

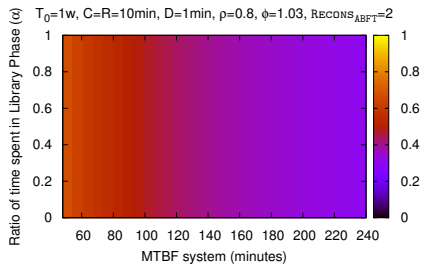
Model & Simulations: BiPERIODICkPT

BiPERIODICkPT

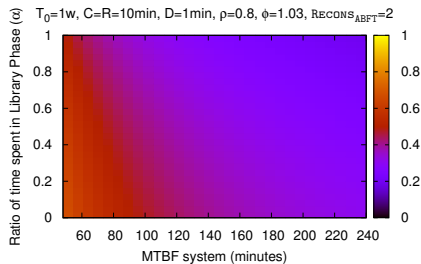
Model & Simulations: ABFT&PERIODICKPT

ABFT&PERIODICKPT

Model: PUREPERIODICKPT vs. BIAPERIODICKPT

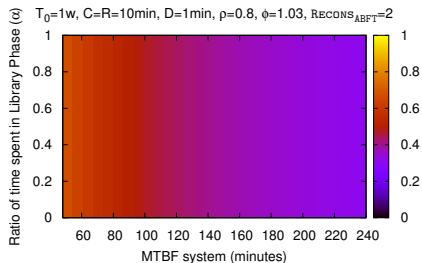


PUREPERIODICKPT

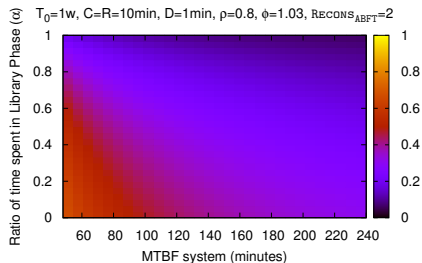


BIAPERIODICKPT

Model & Simulations: PUREPERIODICKPT vs. ABFT&PERIODICKPT

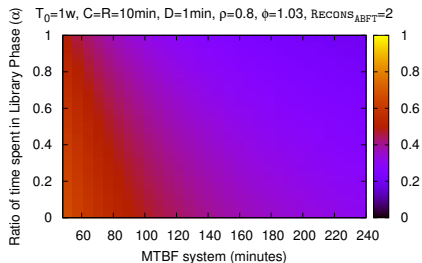


PUREPERIODICKPT

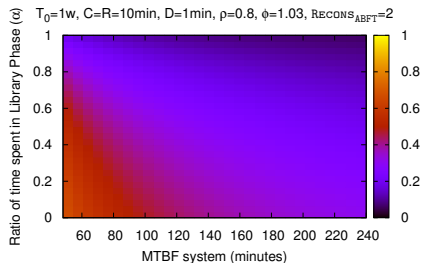


ABFT&PERIODICKPT

Model & Simulations: BiPERIODICkPT vs. ABFT&PERIODICkPT



BiPERIODICkPT



ABFT&PERIODICkPT

Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation**
 - As function of α and μ
 - **Weak Scaling**
- 6 Conclusion

Toward Exascale, and Beyond!

Let's think at scale

- Number of components $\nearrow \Rightarrow$ MTBF \searrow
 - Number of components $\nearrow \Rightarrow$ Problem Size \nearrow
 - Problem Size $\nearrow \Rightarrow$
Computation Time spent in LIBRARY phase \nearrow
- 😊 ABFT&PERIODICCKPT should perform better with scale
- 🤔 By how much?

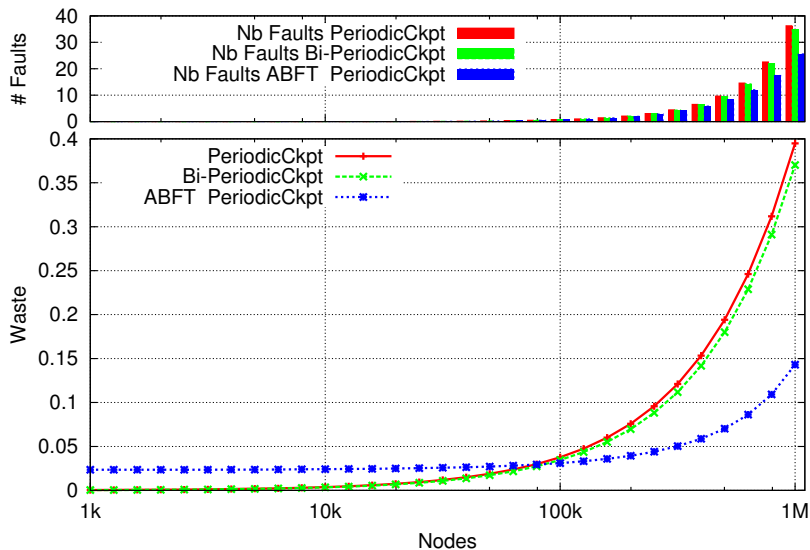
Weak Scale #1

Weak Scale Scenario #1

- Number of components, x , increases
- Memory per component M_{ind} remains constant
- PbSize n increases in $O(\sqrt{x})$ (e.g. matrix, $n^2 = xM_{ind}$)

- μ at $x = 10^5$: 1 day, is in $O(\frac{1}{x})$
- $C (=R)$ at $x = 10^5$, is 1 minute, is in $O(x)$
- α is constant at 0.8, as is ρ .
(both LIBRARY and GENERAL phase increase in time at the same speed)

Weak Scale #1



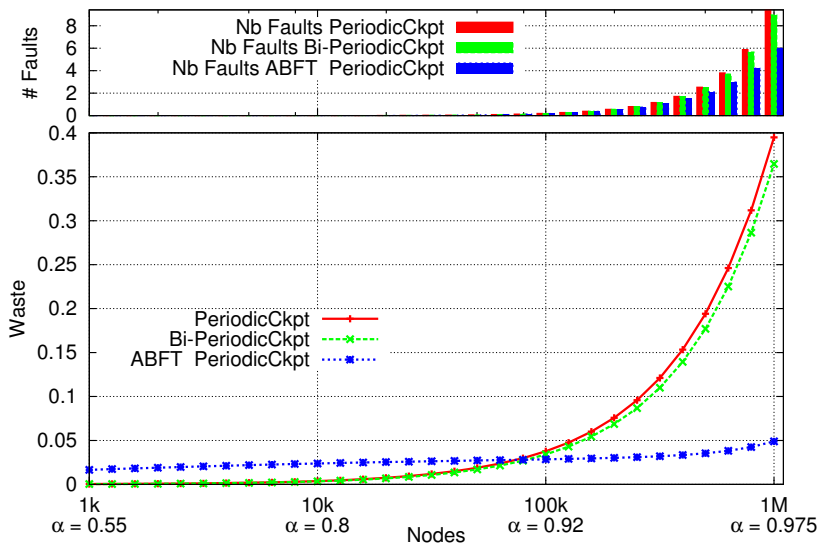
Weak Scale #2

Weak Scale Scenario #2

- Number of components, x , increases
- Memory per component M_{ind} remains constant
- PbSize n increases in $O(\sqrt{x})$ (e.g. matrix, $n^2 = xM_{ind}$)

- μ at $x = 10^5$: 1 day, is $O(\frac{1}{x})$
- $C (=R)$ at $x = 10^5$, is 1 minute, is in $O(x)$
- ρ remains constant at 0.8, but LIBRARY phase is $O(n^3)$ when GENERAL phases progresses in $O(n^2)$ (α is 0.8 at $x = 10^5$ nodes).

Weak Scale #2



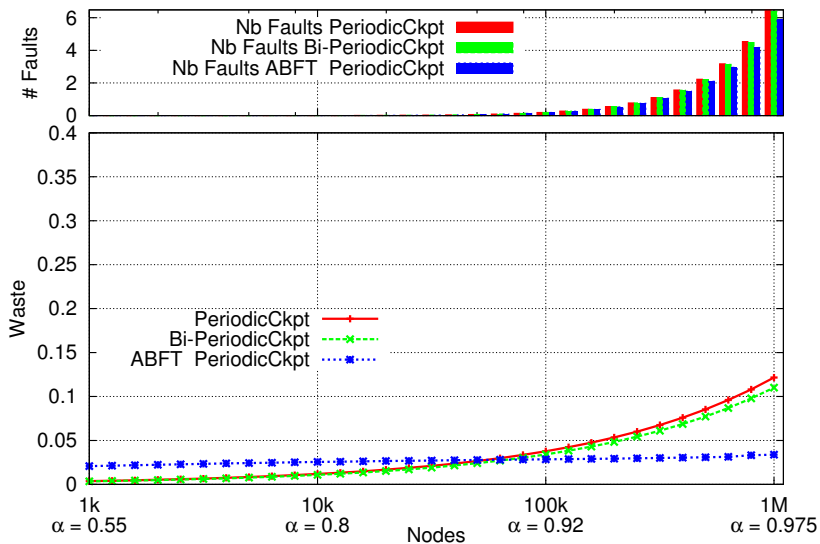
Weak Scale #3

Weak Scale Scenario #3

- Number of components, x , increases
- Memory per component M_{ind} remains constant
- PbSize increases in $O(\sqrt{x})$ (e.g. matrix, $n^2 = xM_{ind}$)

- μ at $x = 10^5$: 1 day, is $O(\frac{1}{x})$
- $C (=R)$ at $x = 10^5$, is 1 minute, **stays independent of x ($O(1)$)**
- ρ remains constant at 0.8, but LIBRARY phase is $O(n^3)$ when GENERAL phases progresses in $O(n^2)$ (α is 0.8 at $x = 10^5$ nodes).

Weak Scale #3



Outline

- 1 Motivation
- 2 ABFT&PERIODICCKPT
- 3 Performance Modeling
- 4 Periodic Checkpointing Protocols (for comparison)
- 5 Evaluation
 - As function of α and μ
 - Weak Scaling
- 6 Conclusion

Conclusion

- Method of composing fault tolerance approaches
 - applications that alternate between ABFT-aware and ABFT-unaware sections
 - each section is protected by its own mechanism
- Performance model shows good opportunity for scaling
 - even when checkpointing hypothesis is optimistic
 - composite approach benefits from checkpointing improvements too
- Energy Efficiency? Checkpointing on Buddies?
Checksumming? Better techniques to recover the ABFT-protected data in some cases.