# Algorithms for coping with silent errors

Anne Benoit[1], Aurélien Cavelan[1], Yves Robert[1,2] and Hongyang Sun[1]

1. ENS Lyon
2. University of Tennessee Knoxville

yves.robert@inria.fr

ICL Lunch – October 24, 2014

## Definitions

- Instantaneous error detection $\Rightarrow$ fail-stop failures, e.g. resource crash
- Silent errors (data corruption) $\Rightarrow$ detection latency

**Silent error detected only when corrupt data is activated and modifies application behavior**

- Includes some software faults, some hardware errors (soft errors in L1 cache, ALU), double bit flip
- Cannot always be corrected by ECC memory

## Probability distributions for silent errors

**?**

**Theorem:** $\mu_p = \dfrac{\mu_{\text{ind}}}{p}$ for arbitrary distributions

(a.k.a, scale is the enemy)

# Probability distributions for silent errors



**Theorem:** $\mu_p = \dfrac{\mu_{\text{ind}}}{p}$ for arbitrary distributions

(a.k.a, scale is the enemy)

## Lesson learnt for fail-stop failures

**(Not so) Secret data**

- Tsubame 2: 962 failures during last 18 months so $\mu = 13$ hrs
- Blue Waters: 2-3 node failures per day
- Titan: a few failures per day
- Tianhe 2: wouldn't say

$$T_{\mathrm{opt}} = \sqrt{2\mu C} \quad \Rightarrow \quad \mathrm{WASTE}_{\mathrm{opt}} \approx \sqrt{\frac{2C}{\mu}}$$

| | | | |
|---|---|---|---|
| Petascale: | $C = 20$ min | $\mu = 24$ hrs | $\Rightarrow \mathrm{WASTE}_{\mathrm{opt}} = 17\%$ |
| Scale by 10: | $C = 20$ min | $\mu = 2.4$ hrs | $\Rightarrow \mathrm{WASTE}_{\mathrm{opt}} = 53\%$ |
| Scale by 100: | $C = 20$ min | $\mu = 0.24$ hrs | $\Rightarrow \mathrm{WASTE}_{\mathrm{opt}} = 100\%$ |

## Lesson learnt for fail-stop failures

**(Not so) Secret data**
- Tsubame ~ 962 failures during last 18 months so ~ 13 hrs
- Blue Waters: 2-3 node failures per day
- Titan: a few failures per day
- Tianhe

Exascale $\neq$ Petascale $\times 1000$
Need more reliable components
Need to checkpoint faster

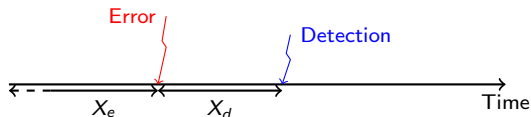| | | | |
|---|---|---|---|
| Petascale | $C = 20$ min | $\mu = 24$ hrs | $\Rightarrow \text{WASTE}_{opt} = 17\%$ |
| Scale by 10: | $C = 20$ min | $\mu = 2.4$ hrs | $\Rightarrow \text{WASTE}_{opt} = 53\%$ |
| Scale by 100: | $C = 20$ min | $\mu = 0.24$ hrs | $\Rightarrow \text{WASTE}_{opt} = 100\%$ |

## Lesson learnt for fail-stop failures

**(Not so) Secret data**
- Tsubame 2: 962 failures during last 18 months so $\mu = 13$ hrs
- Blue Waters: 2-3 node failures per day
- Titan: a few failures per day
- Tianhe 2: wouldn't say

> Silent errors:
> detection latency $\Rightarrow$ additional problems

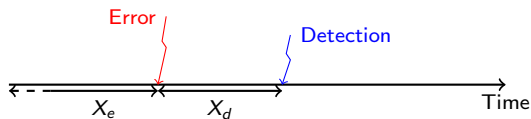| | | | |
|---|---|---|---|
| Petascale: | $C = 20$ min | $\mu = 24$ hrs | $\Rightarrow \mathrm{WASTE}_{\mathsf{opt}} = 17\%$ |
| Scale by 10: | $C = 20$ min | $\mu = 2.4$ hrs | $\Rightarrow \mathrm{WASTE}_{\mathsf{opt}} = 53\%$ |
| Scale by 100: | $C = 20$ min | $\mu = 0.24$ hrs | $\Rightarrow \mathrm{WASTE}_{\mathsf{opt}} = 100\%$ |

## Outline

## Outline

## General-purpose approach



Error and detection latency

- Last checkpoint may have saved an already corrupted state

- Saving $k$ checkpoints (Lu, Zheng and Chien):
  ① Critical failure when all live checkpoints are invalid
  ② Which checkpoint to roll back to?

## General-purpose approach



Error and detection latency

- Last checkpoint may have saved an already corrupted state

- Saving $k$ checkpoints (Lu, Zheng and Chien):
  ① Critical failure when all live checkpoints are invalid
     Assume unlimited storage resources
  ② Which checkpoint to roll back to?
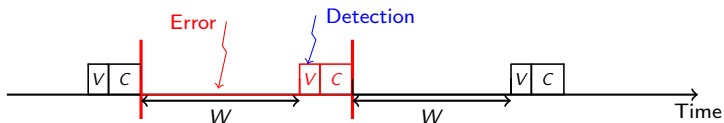     Need a verification mechanism ☹ ☹ ☹

# Outline

# Coupling checkpointing and verification

- Verification mechanism of cost $V$
- Silent errors detected only when verification is executed
- Approach agnostic of the nature of verification mechanism (checksum, error correcting code, coherence tests, triple modular redundancy, etc)
- Fully general-purpose (application-specific information, if available, can always be used to decrease $V$)

# Base pattern (and revisiting Young/Daly)



| | Fail-stop (classical) | Silent errors |
|---|---|---|
| Pattern | $T = W + C$ | $S = W + V + C$ |
| $\text{WASTE}_{\text{FF}}$ | $\frac{C}{T}$ | $\frac{V+C}{S}$ |
| $\text{WASTE}_{\text{fail}}$ | $\frac{1}{\mu}(D + R + \frac{W}{2})$ | $\frac{1}{\mu}(R + W + V)$ |
| Optimal | $T_{\text{opt}} = \sqrt{2C\mu}$ | $S_{\text{opt}} = \sqrt{(C + V)\mu}$ |
| $\text{WASTE}_{\text{opt}}$ | $\sqrt{\frac{2C}{\mu}}$ | $2\sqrt{\frac{C+V}{\mu}}$ |

# On-line ABFT scheme for PCG



1 : Compute $r^{(0)} = b - Ax^{(0)}, z^{(0)} = M^{-1}r^{(0)}, p^{(0)} = z^{(0)}$,
     and $\rho_0 = r^{(0)T}z^{(0)}$ for some initial guess $x^{(0)}$
2 : checkpoint: $A$, $M$, and $b$
3 : for $i = 0, 1, \ldots$
4 :      if ( (i>0) and (i%d = 0) )
5 :          if ( $\frac{p^{(i+1)T}q^{(i)}}{||p^{(i+1)}||.||q^{(i)}||} > 10^{-10}$
         or $\frac{||r^{(i+1)} + Ax^{(i+1)} - b||}{||b||.||A||} > 10^{-10}$ )
6 :             recover: $A$, $M$, $b$, $i$, $\rho_i$,
            $p^{(i)}$, $x^{(i)}$, and $r^{(i)}$.
7 :          else if ( i%(cd) = 0 )
8 :             checkpoint: $i$, $\rho_i$, $p^{(i)}$, and $x^{(i)}$
9:          endif
10:      endif
11:      $q^{(i)} = Ap^{(i)}$
12:      $\alpha_i = \rho_i/p^{(i)T}q^{(i)}$
13:      $x^{(i+1)} = x^{(i)} + \alpha_i p^{(i)}$
14:      $r^{(i+1)} = r^{(i)} - \alpha_i q^{(i)}$
15:      solve $Mz^{(i+1)} = r^{(i+1)}$, where $M = M^T$
16:      $\rho_{i+1} = r^{(i+1)T}z^{(i+1)}$
17:      $\beta_i = \rho_{i+1}/\rho_i$
10:      $p^{(i+1)} = z^{(i+1)} + \beta_i p^{(i)}$
19:      check convergence; continue if necessary
20: end

Zizhong Chen, PPoPP'13

- Iterate PCG
  Cost: SpMV, preconditioner solve, 5 linear kernels
- Detect soft errors by checking orthogonality and residual

- Verification every $d$ iterations
  Cost: scalar product+SpMV
- Checkpoint every $c$ iterations
  Cost: three vectors, or two vectors + SpMV at recovery

- Experimental method to choose $c$ and $d$

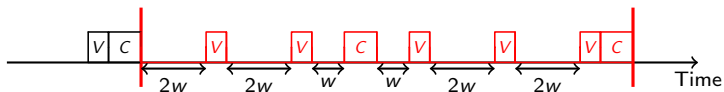## Outline

## Problems

- Given a pattern with $p$ checkpoints and $q$ verifications:
  where to position them?
  what is the optimal period (pattern length)?

  PCG example:
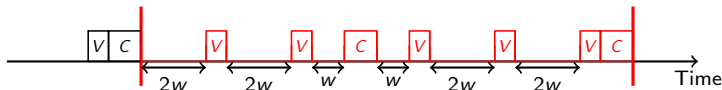  1 checkpoint, $c$ verifications, length $= cd$ iterations

- What is the best pattern?

# BALANCEDALGORITHM



- $p$ checkpoints and $q$ verifications, $p \leq q$
- $p = 2$, $q = 5$, $S = 2C + 5V + W$
- $W = 10w$, six chunks of size $w$ or $2w$
- May store invalid checkpoint (error during third chunk)
- After successful verification in fourth chunk, preceding checkpoint is valid
- Keep only two checkpoints in memory and avoid any fatal failure

# BALANCEDALGORITHM



① ( proba $2w/W$) $T_{\text{lost}} = R + 2w + V$

② ( proba $2w/W$) $T_{\text{lost}} = R + 4w + 2V$

③ ( proba $w/W$) $T_{\text{lost}} = 2R + 6w + C + 4V$

④ ( proba $w/W$) $T_{\text{lost}} = R + w + 2V$

⑤ ( proba $2w/W$) $T_{\text{lost}} = R + 3w + 2V$

⑥ ( proba $2w/W$) $T_{\text{lost}} = R + 5w + 3V$

$$\text{WASTE}_{\text{opt}} \approx 2\sqrt{\frac{7(2C + 5V)}{20\mu}}$$

# BALANCEDALGORITHM



## Theorem

① *Given pattern with p checkpoints and q verifications,* BALANCEDALGORITHM *is optimal when MTBF is large in front of resilience parameters $C, R, V$*

② *Given pattern, can compute optimal period length*

③ *Given $C/V$ ratio, can compute optimal pattern (best values of p and q)*

## Outline

## Linear chain

- $\{T_1, T_2, \ldots, T_n\}$ : linear chain of $n$ tasks
- Each task $T_i$ fully parametrized:
    - $w_i$ computational weight
    - $C_i, R_i, V_i$ : checkpoint, recovery, verification
- Error rates:
    - $\lambda^F$ rate of fail-stop errors
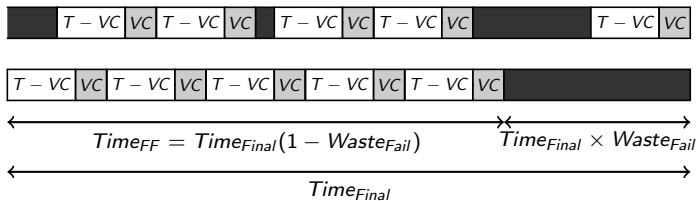    - $\lambda^S$ rate of silent errors

# VC-only



$$\min_{0 \le k < n} Time_C^{rec}(n, k)$$

$$Time_C^{rec}(j, k) = \min_{k \le i < j} \{ Time_C^{rec}(i, k-1) + T_C^{SF}(i+1, j) \}$$

$$T_C^{SF}(i, j) = p_{i,j}^F \left( T_{lost_{i,j}} + R_{i-1} + T_C^{SF}(i, j) \right)$$
$$+ \left( 1 - p_{i,j}^F \right) \left( \sum_{\ell=i}^{j} w_\ell + V_j + p_{i,j}^S \left( R_{i-1} + T_C^{SF}(i, j) \right) + \left( 1 - p_{i,j}^S \right) C_j \right)$$

# Young/Daly



$$\text{Waste} = \text{Waste}_{ef} + \text{Waste}_{fail}$$

$$\text{Waste} = \frac{V + C}{T} + \lambda^F(s)(R + \frac{T}{2}) + \lambda^S(s)(R + T)$$

$$T_{\text{opt}} = \sqrt{\frac{2(V + C)}{\lambda^F(s) + 2\lambda^S(s)}}$$

## Extensions

- VC-ONLY and VC+V
- Different speeds with DVFS, different error rates
- Different execution modes
- Optimize for time or for energy consumption

☺ ☺ ☺

- Use verification to correct some errors (ABFT)
- Same analysis (smaller error rate but higher verification cost)

- VC-ONLY and VC+V
- Different speeds with DVFS, different error rates
- Different execution modes
- Optimize for time or for energy consumption

☺ ☺ ☺

- Use verification to correct some errors (ABFT)
- Same analysis (smaller error rate but higher verification cost)
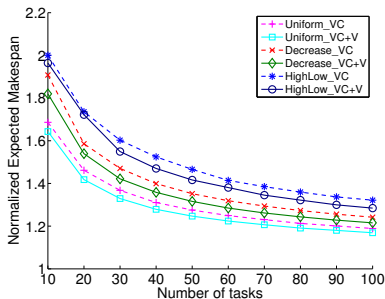
## Outline

## Settings

- Linear chains with $n$ tasks
  - Total work $W \approx 14$ hours
  - Patterns: (1) Uniform; (2) Decrease; (3) HighLow
- Set of speeds from Intel Xscale processor
  - Normalized speeds $\{0.15, 0.4, 0.6, 0.8, 1\}$
  - Fitted power function $P(s) = 1550s^3 + 60$
  - $\lambda^F(s) = \lambda^F_{\mathrm{ref}} \cdot 10^{\frac{d \cdot |s_{\mathrm{ref}} - s|}{s_{\max} - s_{\min}}}$
  - Reference speed $s_{\mathrm{ref}} = 0.6$ and $\lambda^F_{\mathrm{ref}} = 10^{-5}$ for fail-stop errors
  - Sensitivity parameter $d = 3$
  - Corresponds to $0.83 \sim 129$ errors over entire chain
  - Silent errors: $\lambda^S(s) = \eta \cdot \lambda^F(s)$
- Checkpoint and verification costs for a task
  - $cr$ ratio of checkpointing cost over computational cost
  - $vr$ ratio of verification cost over computational cost
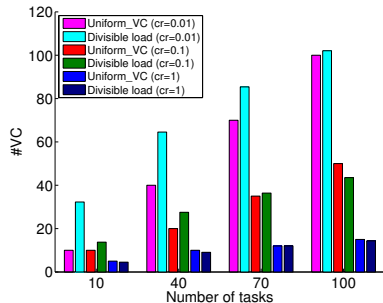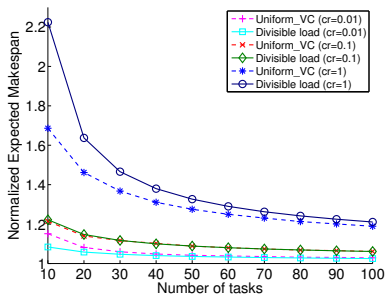  - Default: checkpoint cost $\gg$ verification cost

## Outline

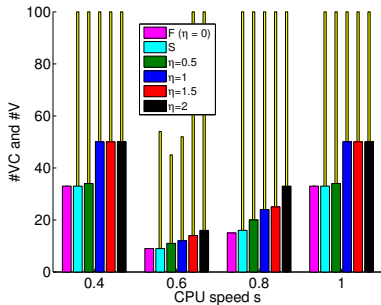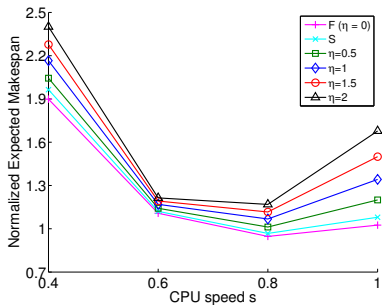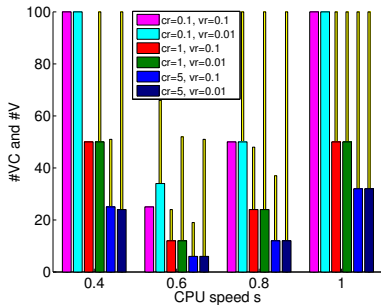## Impact of $n$ and cost distribution

# Comparison with a divisible load application

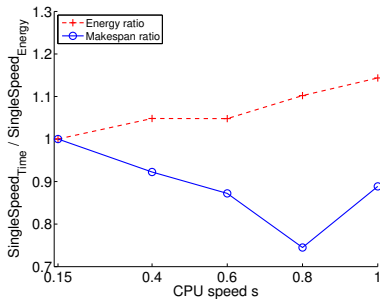# Impact of $\eta$ (TIME-VC+V, $n = 100$, Underline{Uniform})

# Impact of $cr$ and $vr$ ($\textsc{Time-VC+V}$, $n = 100$, <u>Uniform</u>)

# Outline

# Impact of CPU speed $s$

# Impact of $P_{idle}$ and $P_{io}$

## Outline

# TIME-VC+V (HighLow)

# Energy-VC+V (HighLow)

# Conclusion

- Soft errors difficult to cope with, even for divisible workloads or linear chains

- Investigate general task graphs

- Combine checkpointing, replication and application-specific techniques

- Multi-criteria optimization problem
  execution time/energy/reliability
  best resource usage (performance trade-offs)

Several challenging algorithmic/scheduling problems ☺

## Back to task graphs?

### Framework

- You're given a (very big) task graph
- Each task produces files that you can save (checkpoint) or not
- Each task can choose from different execution speeds, with different error probabilities
- You can replicate some tasks, either for verification or for faster execution of successor tasks
- You may also be able to verify results by some application-specific mechanism

### Problem

- Given energy budget or power cap, minimize execution time
- For each task, many things to be decided by schedule ☹

## Back to task graphs?

### Framework

- You're given
- Each task pro                      checkpoint) or not
- Each task car                      on speeds, with different erro
- You can repli                      ification or for faster executi
- You may also                      ne application-sp

### Problem

- Given energy                      execution time
- For each task                      schedule ☹

# Back to task graphs?

## A few questions

### Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

## A few questions

### Silent errors

- Error rate? MTBE?

- Selective reliability?

- New algorithms beyond iterative? matrix-product, FFT, ...

- Resilient research?

## A few questions

### Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

## A few questions

### Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

## A few questions

#### Silent errors

- Error rate? MTBE?

- Selective reliability?

- New algorithms beyond iterative? matrix-product, FFT, ...

- Resilient research?

$$\text{Model} + \text{Theorem} \rightarrow \text{Paper}$$

## A few questions

### Silent errors

- Error rate? MTBE?

- Selective reliability?

- New algorithms beyond iterative? matrix-product, FFT, ...

- Resilient research?

$$\text{Model} + \text{Theorem} \nrightarrow \text{Paper} \ \odot$$
$$\text{Model} + \text{Theorem} + \text{Matlab code} \rightarrow \text{Paper}$$

## A few questions

Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

$$\text{Model} + \text{Theorem} \nrightarrow \text{Paper} \;\; \ddot{\frown}$$
$$\text{Model} + \text{Theorem} + \text{Matlab code} \nrightarrow \text{Paper} \;\; \ddot{\frown}$$
$$\text{Model} + \text{Theorem} + \text{Matlab code} + \text{Simulations} \rightarrow \text{Paper}$$

## A few questions

Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

Model + Theorem $\not\rightarrow$ Paper ☹
Model + Theorem + Matlab code $\not\rightarrow$ Paper ☹
Model + Theorem + Matlab code + Simulations $\not\rightarrow$ Paper ☹
Model + Theorem + Matlab code + Simulations +
Experiments on large platforms ➜ Paper

## A few questions

### Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

<div align="center">

Model + Theorem $\not\rightarrow$ Paper ☹

Model + Theorem + Matlab code $\not\rightarrow$ Paper ☹

Model + Theorem + Matlab code + Simulations $\not\rightarrow$ Paper ☹

~~Model + Theorem + Matlab code + Simulations~~ +
Experiments on large platforms $\rightarrow$ Paper

</div>

## A few questions

Silent errors

- Error rate? MTBE?
- Selective reliability?
- New algorithms beyond iterative? matrix-product, FFT, ...
- Resilient research?

### Models needed to assess techniques at scale
### without bias ☺