

# New-Sum: A Novel Online ABFT Scheme For General Iterative Methods

**Dingwen Tao (University of California, Riverside)**

Shuaiwen Leon Song (Pacific Northwest National Laboratory)

Sriram Krishnamoorthy (Pacific Northwest National Laboratory)

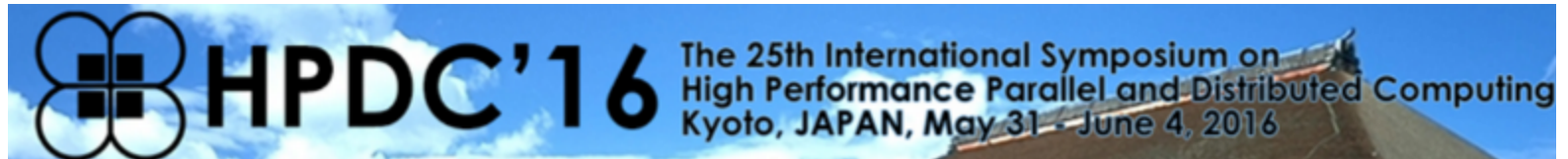
Panruo Wu (University of California, Riverside)

Xin Liang (University of California, Riverside)

Eddy Z. Zhang (Rutgers University)

Darren Kerbyson (Pacific Northwest National Laboratory)

Zizhong Chen (University of California, Riverside)



# Introduction

- ◆ Supercomputers are increasingly susceptible to Silent Data Corruption (SDC)
  - Large number of complex architecture components (e.g., novel memory designs)
  - Each component has growing on-chip transistor density
  - Limited power and energy consumption
  
- ◆ This phenomena has been already observed on several real-world leadership-class supercomputers
  
- ◆ Titan/Jaguar case study<sup>†</sup>
  - Constant stream of single bit flips
  - Double-bit flip every 24 hrs.
  - 20 faults per hour
    - Heartbeat fault every 3 minutes
    - 12 kernel panics in 3 days



<sup>†</sup>AI Geist, How To Kill A Supercomputer, 2016

# Outline

- ◆ Algorithm-Based Fault Tolerance (ABFT)
- ◆ Limitations of Traditional ABFT for matrix-vector multiplication (MVM)
- ◆ Limitations of Existing Techniques for FT-Iterative Methods

# Outline

- ◆ Algorithm-Based Fault Tolerance (ABFT)
- ◆ Limitations of Traditional ABFT for matrix-vector multiplication (MVM)
- ◆ Limitations of Existing Techniques for FT-Iterative Methods
- ◆ Our Designs
  - Novel Error-Preserving Checksum for MVM
  - New Online ABFT Schemes

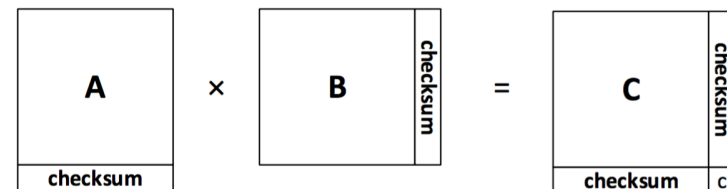


# Outline

- ◆ Algorithm-Based Fault Tolerance (ABFT)
- ◆ Limitations of Traditional ABFT for matrix-vector multiplication (MVM)
- ◆ Limitations of Existing Techniques for FT-Iterative Methods
- ◆ Our Designs
  - Novel Error-Preserving Checksum for MVM
  - New Online ABFT Schemes
- ◆ Theoretical Comparison
- ◆ Empirical Evaluation
- ◆ Conclusions

# Algorithm-Based Fault Tolerance (ABFT)

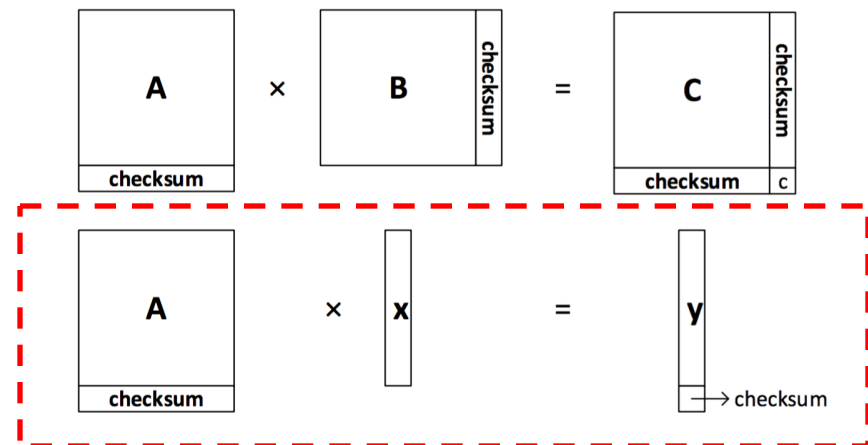
- ◆ Algorithm-Based Fault Tolerance (ABFT) is a checksum-based approach to detect and correct/locate SDCs at a low cost
- ◆ Traditional checksum encoding scheme
  - **Encoding:** augment the input matrices with a checksum computed from the rows or columns
  - **Computation:** perform a matrix operation on the augmented matrices and compute a checksum for the output matrix automatically
  - **Verification:** any error in the computation will break the encoding relationship between the output matrix and its checksum



# Algorithm-Based Fault Tolerance (ABFT)

- ◆ Algorithm-Based Fault Tolerance (ABFT) is a checksum-based approach to detect and correct/locate SDCs at a low cost
- ◆ Traditional checksum encoding scheme
  - **Encoding:** augment the input matrices with a checksum computed from the rows or columns
  - **Computation:** perform a matrix operation on the augmented matrices and compute a checksum for the output matrix automatically
  - **Verification:** any error in the computation will break the encoding relationship between the output matrix and its checksum

- ◆ Example: applying traditional checksum encoding scheme [Huang and Abraham] to **matrix-vector multiplication (MVM)**



# Algorithm-Based Fault Tolerance for MVM

◆ **Encoding:**  $A \rightarrow A^* = \begin{pmatrix} A \\ \bar{c}^T A \end{pmatrix}$

➤  $\bar{c}$  is a predefined vector, normally,  $\bar{c} = (1, 1, \dots, 1)^T$

◆ **Computation on encoded matrix and vectors:**

$$\bar{y}^* = A^* \bar{x} = \begin{pmatrix} A \\ \bar{c}^T A \end{pmatrix} \bar{x} = \begin{pmatrix} A\bar{x} \longrightarrow \text{computed vector } \mathbf{y} \\ \bar{c}^T A\bar{x} \longrightarrow \text{checksum}(\mathbf{y}) \end{pmatrix}$$

◆ **Verification:**  $\text{checksum}(\mathbf{y}) = \bar{c}^T \bar{y}$

- Checksum relationship
- Satisfied: **no error** in the computation
- Not Satisfied: **error(s) occurred** in the computation

# Limitations of Traditional ABFT for MVM

- ◆ Consider an error in  $\mathbf{x}$  before the MVM, resulting in an erroneous vector  $\mathbf{x}'$

- ◆ **Computation:**

$$\bar{\mathbf{y}}'^* = A * \bar{\mathbf{x}}' = \begin{pmatrix} A \\ \bar{\mathbf{c}}^T A \end{pmatrix} \bar{\mathbf{x}}' = \begin{pmatrix} A\bar{\mathbf{x}}' \\ \bar{\mathbf{c}}^T A\bar{\mathbf{x}}' \end{pmatrix} \begin{matrix} \longrightarrow \text{erroneous computed } \mathbf{y} \\ \longrightarrow \text{erroneous checksum}(\mathbf{y}) \end{matrix}$$

- ◆ **Verification:**

- Checksum relationship for erroneous computed vector  $\mathbf{y}$  still holds
- Checking the output vector's checksum relationship cannot identify all SDCs

Problems With the Traditional ABFT Scheme:

- (1) Traditional encoding scheme may **FAIL** for **MVM**
- (2) Traditional encoding scheme can **NOT** address **cache/register bit-flips** and protect **preconditioners**

# ABFT for Iterative Methods

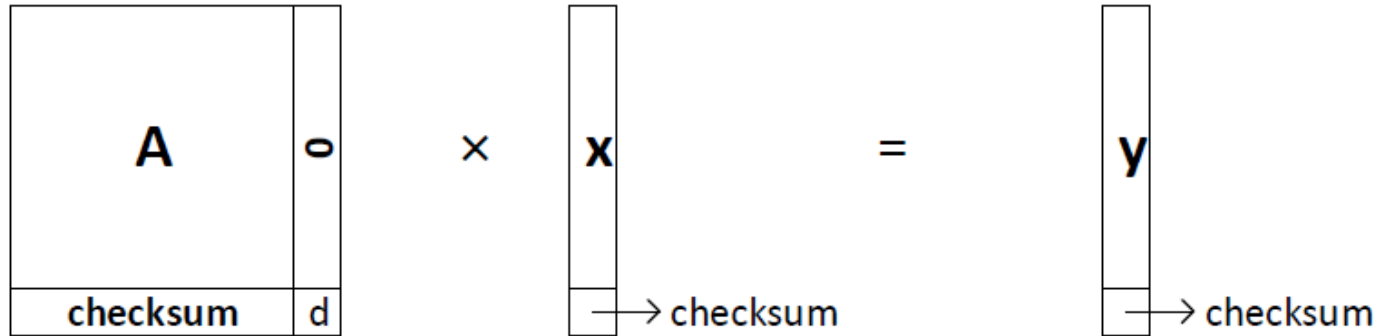
- ◆ Iterative Methods are widely used for solving systems of equations or computing eigenvalues of large sparse matrices
- ◆ Existing fault tolerant techniques for iterative methods
  - [Shantharam et al., ICS'12]: requires matrix  $A$  to be ***strictly diagonally dominant*** – **NOT GENERAL**
  - [Chen, PPOPP'13]: only covers a subset of Krylov methods that can offer ***orthogonality*** – **NOT GENERAL**
  - [Sloan et al., DSN'13]: uses the ***traditional checksum***-encoding mechanism – **MAY FAIL**
- ◆ Developing a ***general*** and ***low-cost*** ABFT scheme with ***good coverage*** for iterative methods is in high demand

# Outline

- ◆ Algorithm-Based Fault Tolerance (ABFT)
- ◆ Limitations of Traditional ABFT for matrix-vector multiplication (MVM)
- ◆ Limitations of Existing Techniques for FT-Iterative Methods
- ◆ Our Designs
  - **Novel Error-Preserving Checksum for MVM**
  - New Online ABFT Schemes
- ◆ Theoretical Comparison
- ◆ Empirical Evaluation
- ◆ Conclusions

# Novel Error-Preserving Checksum for MVM

## ◆ New Encoding Mechanism



$$A \rightarrow A^* = \begin{pmatrix} A & \bar{\mathbf{0}} \\ \bar{c}^T A - d\bar{c}^T & d \end{pmatrix} \quad \bar{x} \rightarrow \bar{x}^* = \begin{pmatrix} \bar{x} \\ \bar{c}^T \bar{x} \end{pmatrix}$$

$$\bar{y}^* = A^* \bar{x}^*$$

$$= \begin{pmatrix} A & \bar{\mathbf{0}} \\ \bar{c}^T A - d\bar{c}^T & d \end{pmatrix} \begin{pmatrix} \bar{x} \\ \bar{c}^T \bar{x} \end{pmatrix} = \begin{pmatrix} A\bar{x} \\ \bar{c}^T A\bar{x} \end{pmatrix} \rightarrow \text{checksum}(y) = \bar{c}^T \bar{y}$$

**Checksum Relationship Maintained**

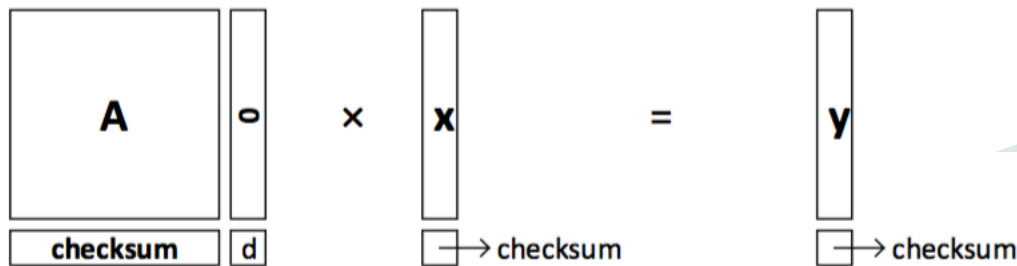


# Novel Error-Preserving Checksum for MVM (cont.)

## ◆ Computation

$$\bar{y}^* = A * \bar{x}^*$$

$$= \begin{pmatrix} A & \bar{0} \\ checksum(A) & d \end{pmatrix} \begin{pmatrix} \bar{x} \\ checksum(\bar{x}) \end{pmatrix} = \begin{pmatrix} A\bar{x} \\ checksum(A)\bar{x} + d \cdot checksum(\bar{x}) \end{pmatrix}$$



**Separation**

## ◆ Checksum can be updated separately after **MVM** $y = Ax$

$$checksum(\bar{y}) = checksum(A)\bar{x} + d \cdot checksum(\bar{x})$$

# Novel Error-Preserving Checksum for MVM (cont.)

## Jacobi Method

```

2: for k = 1, 2, ... do
3:   for i = 1, 2, ..., n do
4:     xi = 0
5:     for j = 1, 2, ..., i - 1, i + 1, ..., n do
6:       xi = xi + ai,jxj(k-1) ○
7:     end for
8:     xi = (bi - xi)/ai,i ○
9:   end for
10:  x(k) = x
11:  check convergence; continue if necessary
12: end for
  
```

## Preconditioned CG

```

2: for i = 0, 1, ... do
3:   q(i) = Ap(i) ●
4:   αi = ρi/p(i)T q(i) ■
5:   x(i+1) = x(i) + αip(i) ○
6:   r(i+1) = r(i) - αiq(i) ○
7:   solve Mz(i+1) = r(i+1) ◀
8:   ρi+1 = r(i+1)T z(i+1) ■
9:   βi = ρi+1/ρi ■
10:  p(i+1) = z(i+1) + βip(i) ○
11:  check convergence; continue if necessary
12: end for
  
```

## Preconditioned Chebyshev

```

3: for i = 1, 2, ... do
4:   solve Mz(i-1) = r(i) ◀
5:   if (i = 1) then
6:     p(1) = z(0)
7:     α1 = 2/θ
8:   else
9:     βi-1 = αi-1(ε/2)2 ■
10:    αi = 1/(θ - βi-1) ■
11:    p(i) = z(i) + βi-1p(i-1) ○
12:   end if
13:   x(i) = x(i-1) + αip(i) ○
14:   r(i) = b - Ax(i) (= r(i-1) - αiAp(i)) ●
15:   check convergence; continue if necessary
16: end for
  
```

- Represents MVM
- Represents VLO
- ◀ Represents PCO
- Represents Vector Inner-Product (VDP) or scalar computation

- ◆ Vector Linear Operation (VLO):  $\mathbf{z} = \alpha\mathbf{x} + \beta\mathbf{y}$  and Solving Preconditioned System (PCO):  $M\mathbf{z} = \mathbf{r}$

$$checksum(\bar{z}) = \alpha \cdot checksum(\bar{x}) + \beta \cdot checksum(\bar{y})$$

$$checksum(\bar{z}) = (checksum(M)\bar{z} - checksum(\bar{r})) / d$$

# Novel Error-Preserving Checksum for MVM (cont.)

- ◆ **Theorem:** For any matrix-vector multiplication (MVM), vector linear operation (VLO), preconditioning operation (PCO), the **checksum relationship** of the output vector is **preserved** if and only if there are **no soft errors before or during** the operation. (key theorem)

Proof: see our paper for the details.

- ◆ **Verification:**

- **No error** occurred **before or during** the computation  $\Leftrightarrow$   
checksum( $\mathbf{y}$ ) =  $\mathbf{c}^T \mathbf{y}$

- Any errors  $\Rightarrow$  Checksum relationship will **preserve to be BROKEN** in the **subset iterations**

# Outline

- ◆ Algorithm-Based Fault Tolerance (ABFT)
- ◆ Limitations of Traditional ABFT for matrix-vector multiplication (MVM)
- ◆ Limitations of Existing Techniques for FT-Iterative Methods
- ◆ Our Designs
  - Novel Error-Preserving Checksum for MVM
  - **New Online ABFT Schemes**
- ◆ Theoretical Comparison
- ◆ Empirical Evaluation
- ◆ Conclusions

## New Online ABFT Schemes

- ◆ Based on our new designed checksum encoding scheme, we design efficient online ABFT solutions for iterative methods
  - ◆ **“Lazy” Detection:** Low-Cost Online ABFT Algorithm
  - ◆ **“Eager” Recovery for MVM:** Triple Checksums
  - ◆ **“Hybrid” Detection:** Two-Level Online ABFT Algorithm

# “Lazy” Detection: Low-Cost Online ABFT Algorithm

- ◆ Step 1: update checksum after each MVM, VLO, PCO (red in algorithm 1)
- ◆ Step 2: low-cost error detection
  - Verify checksum relationship after **every MVM, VLO, PCO – high detection cost**
  - *Do we need to verify the checksum relationship every iteration?*
    - ✧ Observation 1: Soft errors in **p, x, r** will propagate to the subsequent iterations
    - ✧ **Verify checksum relationship every several iterations (blue in algorithm 1)**
  - *Do we need to verify all the operations?*
    - ✧ Observation 2: Soft errors in **z, p, q** will eventually propagate to **x** and **r**
    - ✧ **Only verify 2 checksum relationships, x and r (blue in algorithm 1)**

**Algorithm 1** Online-ABFT Preconditioned Conjugate Gradient Algorithm Based on New Checksum with checkpoint/restart Technique

```

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: Compute  $checksum(A) = c^T A - dc^T$ ,  $checksum(M) = c^T M - dc^T$ ,  $checksum(b) = c^T b^{(0)}$ ,  $checksum(x) = c^T x^{(0)}$ ,  $checksum(r) = c^T r^{(0)}$ ,  $checksum(z) = [(checksum(r) - checksum(M)z^{(0)})]/d$ ,  $checksum(p) = checksum(z)$ 
3: Checkpoint  $A, M$ 
4: for  $i = 0, 1, \dots$  do
5:   if  $((i > 0) \text{ and } (i \% d = 0))$  then
6:     if  $(\|checksum(r) - c^T r^{(i)}\|/n > \theta)$  or  $(\|checksum(x) - c^T x^{(i)}\|/n > \theta)$  then
7:       Rollback: recover  $A, M, i, \rho_i, p^{(i)}, x^{(i)}, r^{(i)}$ 
8:     else if  $(i \% cd) = 0$  then
9:       Checkpoint:  $i, \rho_i, p^{(i)}, x^{(i)}$ 
10:    end if
11:  end if
12:   $q^{(i)} = Ap^{(i)}$ 
13:   $checksum(q) = checksum(A)p^{(i)} + d \cdot checksum(p)$ 
14:   $\alpha_i = \rho_i / p^{(i)T} q^{(i)}$ 
15:   $x^{(i+1)} = x^{(i)} + \alpha_i p^{(i)}$ 
16:   $checksum(x) = checksum(x) + \alpha_i \cdot checksum(p)$ 
17:   $r^{(i+1)} = r^{(i)} - \alpha_i q^{(i)}$ 
18:   $checksum(r) = checksum(r) - \alpha_i \cdot checksum(q)$ 
19:  solve  $Mz^{(i+1)} = r^{(i+1)}$ 
20:   $checksum(z) = \frac{checksum(r) - checksum(M)z^{(i+1)}}{d}$ 
21:   $\rho_{i+1} = r^{(i+1)T} z^{(i+1)}$ 
22:   $\beta_i = \rho_{i+1} / \rho_i$ 
23:   $p^{(i+1)} = z^{(i+1)} + \beta_i p^{(i)}$ 
24:   $checksum(p) = checksum(z) + \beta_i \cdot checksum(p)$ 
25:  check convergence; continue if necessary
26: end for
  
```

# “Lazy” Detection: Low-Cost Online ABFT Algorithm (cont.)

**Algorithm 1** Online-ABFT Preconditioned Conjugate Gradient Algorithm Based on New Checksum with checkpoint/restart Technique

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: Compute $checksum(A) = c^T A - dc^T, checksum(M) = c^T M - dc^T, checksum(b) = c^T b^{(0)}, checksum(x) = c^T x^{(0)}, checksum(r) = c^T r^{(0)}, checksum(z) = [(checksum(r) - checksum(M)z^{(0)})]/d, checksum(p) = checksum(z)$	
3: Checkpoint $A, M$	
4: for $i = 0, 1, \dots$ do	
5: if $((i > 0) \text{ and } (i\%d = 0))$ then	
6: if $(\ checksum(r) - c^T r^{(i)}\ /n > \theta)$ or $(\ checksum(x) - c^T x^{(i)}\ /n > \theta)$ then	2 VDP (4n FLOPS)
7: Rollback: recover $A, M, i, \rho_i, p^{(i)}, x^{(i)}, r^{(i)}$	2 Mcpy + 4 Vcpy
8: else if $(i\%(cd) = 0)$ then	
9: Checkpoint: $i, \rho_i, p^{(i)}, x^{(i)}$	2 Vcpy
10: end if	
11: end if	
12: $q^{(i)} = Ap^{(i)}$	1 MVM (nnz FLOPS)
13: $checksum(q) = checksum(A)p^{(i)} + d \cdot checksum(p)$	1 VDP (2n FLOPS)
14: $\alpha_i = \rho_i/p^{(i)T}q^{(i)}$	1 VDP (2n FLOPS)
15: $x^{(i+1)} = x^{(i)} + \alpha_i p^{(i)}$	1 VLO (2n FLOPS)
16: $checksum(x) = checksum(x) + \alpha_i \cdot checksum(p)$	2 FLOPS
17: $r^{(i+1)} = r^{(i)} - \alpha_i q^{(i)}$	1 VLO (2n FLOPS)
18: $checksum(r) = checksum(r) - \alpha_i \cdot checksum(q)$	2 FLOPS
19: solve $Mz^{(i+1)} = r^{(i+1)}$	1 PCO
20: $checksum(z) = \frac{checksum(r) - checksum(M)z^{(i+1)}}{d}$	1 VDP (2n FLOPS)
21: $\rho_{i+1} = r^{(i+1)T}z^{(i+1)}$	1 VDP (2n FLOPS)
22: $\beta_i = \rho_{i+1}/\rho_i$	1 FLOPS
23: $p^{(i+1)} = z^{(i+1)} + \beta_i p^{(i)}$	1 VLO (2n FLOPS)
24: $checksum(p) = checksum(z) + \beta_i \cdot checksum(p)$	2 FLOPS
25: check convergence; continue if necessary	
26: end for	

## ◆ Step 3: low-cost error recovery

- Error detect every several iterations => Checkpoint/Rollback
- Do we need to checkpoint every vector?
  - ✧ Observation 3: Using  $p$  and  $x$  can compute the other 3 vectors
  - ✧ **Only checkpoint 2 vectors,  $p$  and  $x$  (in purple)**

## ◆ Overhead Summary

- Checksum update: 2 vector dot-products every iteration
- Error detection: 2 vector dot-products every  $d$  iterations
- Rollback recovery: 2 matrix copies every  $cd$  iterations
- Checkpoint: 2 vector copies every  $cd$  iterations



# “Lazy” Detection: Low-Cost Online ABFT Algorithm (cont.)

**Algorithm 1** Online-ABFT Preconditioned Conjugate Gradient Algorithm Based on New Checksum with checkpoint/restart Technique

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: Compute $checksum(A) = c^T A - dc^T, checksum(M) = c^T M - dc^T, checksum(b) = c^T b^{(0)}, checksum(x) = c^T x^{(0)}, checksum(r) = c^T r^{(0)}, checksum(z) = [(checksum(r) - checksum(M)z^{(0)})]/d, checksum(p) = checksum(z)$	
3: Checkpoint $A, M$	
4: for $i = 0, 1, \dots$ do	
5: if $((i > 0) \text{ and } (i\%d = 0))$ then	
6: if $(\ checksum(r) - c^T r^{(i)}\ /n > \theta)$ or $(\ checksum(x) - c^T x^{(i)}\ /n > \theta)$ then	2 VDP (4n FLOPS)
7: Rollback: recover $A, M, i, \rho_i, p^{(i)}, x^{(i)}, r^{(i)}$	2 Mcpy + 4 Vcpy
8: else if $(i\%(cd) = 0)$ then	
9: Checkpoint: $i, \rho_i, p^{(i)}, x^{(i)}$	2 Vcpy
10: end if	
11: end if	
12: $q^{(i)} = Ap^{(i)}$	1 MVM (nnz FLOPS)
13: $checksum(q) = checksum(A)p^{(i)} + d \cdot checksum(p)$	1 VDP (2n FLOPS)
14: $\alpha_i = \rho_i / p^{(i)T} q^{(i)}$	1 VDP (2n FLOPS)
15: $x^{(i+1)} = x^{(i)} + \alpha_i p^{(i)}$	1 VLO (2n FLOPS)
16: $checksum(x) = checksum(x) + \alpha_i \cdot checksum(p)$	2 FLOPS
17: $r^{(i+1)} = r^{(i)} - \alpha_i q^{(i)}$	1 VLO (2n FLOPS)
18: $checksum(r) = checksum(r) - \alpha_i \cdot checksum(q)$	2 FLOPS
19: solve $Mz^{(i+1)} = r^{(i+1)}$	1 PCO
20: $checksum(z) = \frac{checksum(r) - checksum(M)z^{(i+1)}}{d}$	1 VDP (2n FLOPS)
21: $\rho_{i+1} = r^{(i+1)T} z^{(i+1)}$	1 VDP (2n FLOPS)
22: $\beta_i = \rho_{i+1} / \rho_i$	1 FLOPS
23: $p^{(i+1)} = z^{(i+1)} + \beta_i p^{(i)}$	1 VLO (2n FLOPS)
24: $checksum(p) = checksum(z) + \beta_i \cdot checksum(p)$	2 FLOPS
25: check convergence; continue if necessary	
26: end for	

## ◆ Step 3: low-cost error recovery

- Error detect every several iterations => Checkpoint/Rollback
- Do we need to checkpoint every vector?
  - ❖ Observation 3: Using  $p$  and  $x$  can compute the other 3 vectors
  - ❖ **Only checkpoint 2 vectors,  $p$  and  $x$  (in purple)**

## ◆ Overhead Summary

- Checksum update: 2 vector dot-products every iteration
- Error detection: 2 vector dot-products every  $d$  iterations
- Rollback recovery: 2 matrix copies every  $cd$  iterations
- Checkpoint: 2 vector copies every  $cd$  iterations
- Original algorithm: 1 MVM and 1 PCO
- MVM and PCO >> vector dot-product, matrix and vector copy



# “Eager” Recovery for MVM: Triple Checksums

- ◆ **MVM:** computation-intensive, vulnerable => faster recovery under a high error rate can be beneficial

## ◆ Encoding:

$$A \rightarrow A^* = \begin{pmatrix} A & \bar{0} & \bar{0} & \bar{0} \\ \bar{c}_1^T A - d_1 \bar{c}_1^T - d_2 \bar{c}_2^T - d_3 \bar{c}_3^T & d_1 & d_2 & d_3 \\ \bar{c}_2^T A - d_2 \bar{c}_1^T - d_3 \bar{c}_2^T - d_1 \bar{c}_3^T & d_2 & d_3 & d_1 \\ \bar{c}_3^T A - d_3 \bar{c}_1^T - d_1 \bar{c}_2^T - d_2 \bar{c}_3^T & d_3 & d_1 & d_2 \end{pmatrix} \quad \bar{x} \rightarrow \bar{x}^* = \begin{pmatrix} \bar{x} \\ \bar{c}_1^T \bar{x} \\ \bar{c}_2^T \bar{x} \\ \bar{c}_3^T \bar{x} \end{pmatrix} \begin{matrix} \rightarrow \text{checksum}_1(\mathbf{x}) \\ \rightarrow \text{checksum}_2(\mathbf{x}) \\ \rightarrow \text{checksum}_3(\mathbf{x}) \end{matrix}$$

- ◆ **Checksum update:** similar to 1 checksum (see in the paper)

## ◆ Verification:

- Detect if there is **any error**:  $\text{checksum}_1(y) \stackrel{?}{=} \bar{c}_1^T \bar{y}$
- Identify whether there is **more than one error**:
  - Choose  $\mathbf{c}_1 = (1, 1, \dots, 1)^T$ ,  $\mathbf{c}_2 = (1, 2, \dots, n)^T$ ,  $\mathbf{c}_3 = (1, 1/2, \dots, 1/n)^T$   
 $(\text{checksum}_1(y) - \bar{c}_1^T \bar{y})^2 \stackrel{?}{=} (\text{checksum}_2(y) - \bar{c}_2^T \bar{y})(\text{checksum}_3(y) - \bar{c}_3^T \bar{y})$
- If there is one error, **locate and correct**:
  - Erroneous locate:  $(\text{checksum}_2(y) - \bar{c}_2^T \bar{y}) / (\text{checksum}_1(y) - \bar{c}_1^T \bar{y})$

# “Hybrid” Detection: Two-Level Online ABFT Algorithm

**Algorithm 2** Two-Level Online-ABFT Preconditioned Conjugate Gradient Algorithm

```

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: Compute  $checksum_1, checksum_2, checksum_3$  of  $A, M, x, r, z, p$ 
3: Checkpoint  $A, M$ 
4: for  $i = 0, 1, \dots$  do
5:   if  $((i > 0) \text{ and } (i \% d = 0))$  then
6:     if  $(\|checksum_1(r) - c_1^T r^{(i)}\|/n^2 > \theta)$  or
        $(\|checksum_1(x) - c_1^T x^{(i)}\|/n^2 > \theta)$  then
7:       Rollback: recover  $A, M, i, \rho_i, p^{(i)}, x^{(i)}, r^{(i)}$ 
8:     else if  $(i \% cd = 0)$  then
9:       Checkpoint:  $i, \rho_i, p^{(i)}, x^{(i)}$ 
10:    end if
11:   end if
12:    $q^{(i)} = Ap^{(i)}$ 
13:    $checksum_1(q) = checksum_1(A)p^{(i)} + d_1 \cdot checksum_1(p)$ 
      $+ d_2 \cdot checksum_2(p) + d_3 \cdot checksum_3(p)$ 
14:    $checksum_2(q) = checksum_2(A)p^{(i)} + d_2 \cdot checksum_1(p)$ 
      $+ d_3 \cdot checksum_2(p) + d_1 \cdot checksum_3(p)$ 
15:    $checksum_3(q) = checksum_3(A)p^{(i)} + d_3 \cdot checksum_1(p)$ 
      $+ d_1 \cdot checksum_2(p) + d_2 \cdot checksum_3(p)$ 
16:    $\delta_1 = c_1^T q^{(i)} - checksum_1(q)$ 
17:    $\delta_2 = c_2^T q^{(i)} - checksum_2(q)$ 
18:    $\delta_3 = c_3^T q^{(i)} - checksum_3(q)$ 
19:   if  $(|\delta_1| > \theta)$  then
20:     if  $(|1 - \delta_1^2 / \delta_2 \delta_3| < \theta)$  then
21:       Locate and correct now:
22:        $j = \delta_2 / \delta_1$ 
23:        $q_j^{(i)} = q_j^{(i)} - \delta_1$ 
24:     else
25:       Rollback: recover  $A, M, b, i, \rho_i, p^{(i)}, x^{(i)}, r^{(i)}$ 
26:     end if
27:   end if
28:    $\alpha_i = \rho_i / p^{(i)T} q^{(i)}$ 

```

Outer Protection

Inner Protection

```

29:  $x^{(i+1)} = x^{(i)} + \alpha p^{(i)}$ 
30:  $checksum_1(x) = checksum_1(x) + \alpha_i \cdot checksum_1(p)$ 
31:  $checksum_2(x) = checksum_2(x) + \alpha_i \cdot checksum_2(p)$ 
32:  $checksum_3(x) = checksum_3(x) + \alpha_i \cdot checksum_3(p)$ 
33:  $r^{(i+1)} = r^{(i)} - \alpha q^{(i)}$ 
34:  $checksum_1(r) = checksum_1(r) - \alpha_i \cdot checksum_1(q)$ 
35:  $checksum_2(r) = checksum_2(r) - \alpha_i \cdot checksum_2(q)$ 
36:  $checksum_3(r) = checksum_3(r) - \alpha_i \cdot checksum_3(q)$ 
37: solve  $Mz^{(i+1)} = r^{(i+1)}$ 
38: Compute  $checksum_1(z), checksum_2(z), checksum_3(z)$ 
39:  $\rho_{i+1} = r^{(i+1)T} z^{(i+1)}$ 
40:  $\beta_i = \rho_{i+1} / \rho_i$ 
41:  $p^{(i+1)} = z^{(i+1)} + \beta_i p^{(i)}$ 
42:  $checksum_1(p) = checksum_1(z) + \beta_i \cdot checksum_1(p)$ 
43:  $checksum_2(p) = checksum_2(z) + \beta_i \cdot checksum_2(p)$ 
44:  $checksum_3(p) = checksum_3(z) + \beta_i \cdot checksum_3(p)$ 
45: check convergence; continue if necessary
46: end for

```

## ◆ Inner protection

- Detect a single error in MVM: locate and correct
- Detect multiple errors in MVM: rollback immediately

## ◆ Outer Protection

- Detect error(s) in VLOs and multiple errors can not be recovered in inner protection

# “Hybrid” Detection: Two-Level Online ABFT Algorithm (cont.)



General procedure to construct a Two-Level ABFT:

1. Encode matrices and vectors
2. Compute checksum updates after MVM, VLO, PCO
3. Analyze dependency relationship between vectors
4. Every  $d$  iterations, invoke the outer-level protection
5. Every  $cd$  iterations, checkpoint the minimum number of vectors
6. After each MVM, add the inner-level protection

# Outline

- ◆ Algorithm-Based Fault Tolerance (ABFT)
- ◆ Limitations of Traditional ABFT for matrix-vector multiplication (MVM)
- ◆ Limitations of Existing Techniques for FT-Iterative Methods
- ◆ Our Designs
  - Novel Error-Preserving Checksum for MVM
  - New Online ABFT Schemes
- ◆ Theoretical Comparison
- ◆ Empirical Evaluation
- ◆ Conclusions

# Theoretical Comparison - Error Coverage and Feature

	Offline residual	Online MV	Online orthogonality	Basic/two-level online ABFT method
Can protect arithmetic error	Yes	Yes	Yes	Yes
Can protect memory bit flips	Yes	Yes	Yes	Yes
Can protect cache or register bit flips	Yes	No	No	Yes
Can be applied to all iterative methods	Yes	Yes	No	Yes
Not necessary to check every iteration	Yes	No	Yes	Yes
Not necessary to check every operation	Yes	No	Yes	Yes

**Offline residual:** verify residual at the end of computation and recompute

**Online MV:** online MVM scheme using the traditional checksum proposed by [Sloan et al., DSN'13]

**Online Orthogonality:** online orthogonally checking proposed by [Chen, PPOPP'13]

**Basic Online ABFT:** our proposed “lazy” online ABFT using checksum updates and C/R

**Two-Level ABFT:** our proposed two-level online ABFT using triple-checksum mechanism

# Theoretical Comparison - Performance

- ◆ 3 designed scenarios
  - One error in MVM during the **entire execution** – **low error rate**
  - One error in MVM **every cd iterations** – **medium/high error rate**
  - One error in MVM **every iteration** – **extremely high error rate**

	Basic online ABFT ( $O_1$ )	Two-level online ABFT ( $O_2$ )	Online MV ( $O_3$ )
Scenario 1	$(2/d+2)VDP+2VLO/cd$	$(2/d+9)VDP+2VLO/cd$	$1PCO+2VDP+3VLO$
Scenario 2	$0.5MVM+(2/d+5)VDP+0.5PCO+(6(1+c_0)/cd + 1.5)VLO$	$(2/d+9)VDP+2VLO/cd$	$1PCO+(5/cd+2)VDP+3VLO$
Scenario 3	$+\infty$	$(2/d+9)VDP+2VLO/cd$	$1PCO+7VDP+3VLO$

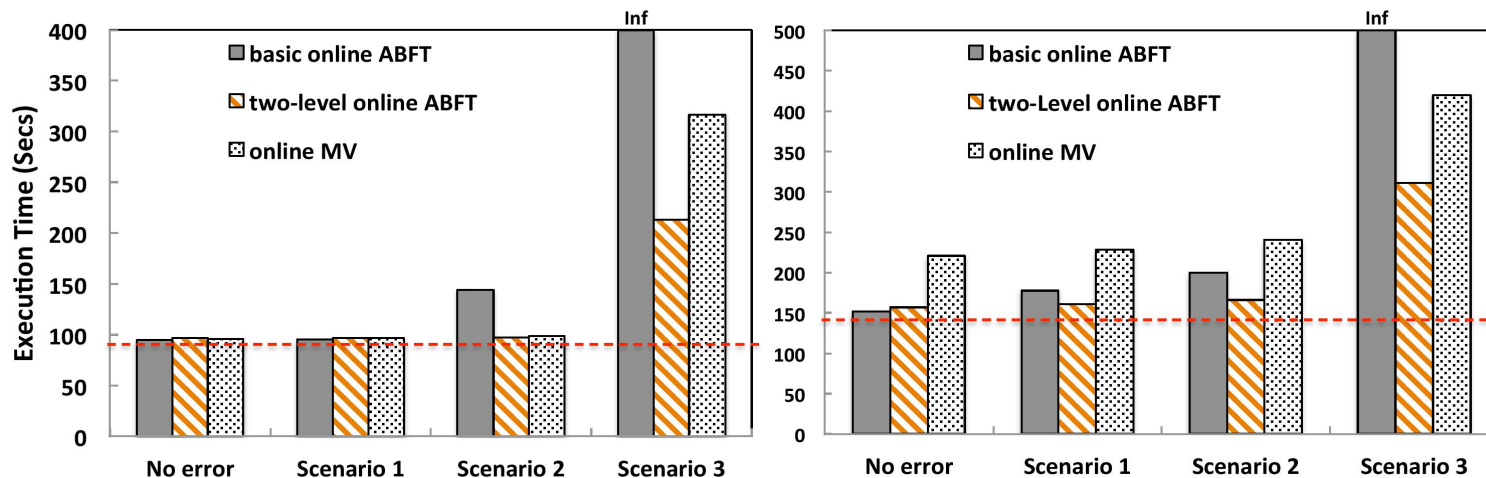
- ◆ Theoretical Comparison on PCG
  - ◆ Low error rate: **basic online ABFT** has the lowest overhead
  - ◆ Medium/high error rate: **two-level online ABFT** has the lowest overhead
  - ◆ Extremely high error rate: **two-level online ABFT** has the lowest overhead
  
- ◆ Overall, no matter the error rate, one of our approaches will outperform the online MV for PCG.

# Empirical Evaluation – Configurations

- ◆ Platforms
  - **Stampede** supercomputer at TACC, each node with 2 Intel Xeon E5-2680 processors
  - **Tianhe-2** supercomputer (No.1 in Top 500), each node with 2 Intel Xeon E5-2692 processors
  
- ◆ Implemented our proposed online ABFT schemes in **PETSc**
  
- ◆ Evaluated solvers
  - Preconditioned Conjugate Gradient (**PCG**): has orthogonality relation
  - Preconditioned Biconjugate Gradient Stabilized (**PBiCGSTAB**): no orthogonality relation (against [Chen, PPOPP'13])
  
- ◆ Input Matrix: G3\_circuit
  - The **largest** SPD matrix from the University of Florida Sparse Matrix Collection
  - 1,585,478 rows and columns with 7,660,826 nonzero elements



# Empirical Evaluation – Results



(a) Comparison with **PCG** on **Stampede**      (b) Comparison with **PBiCGSTAB** on **Stampede**

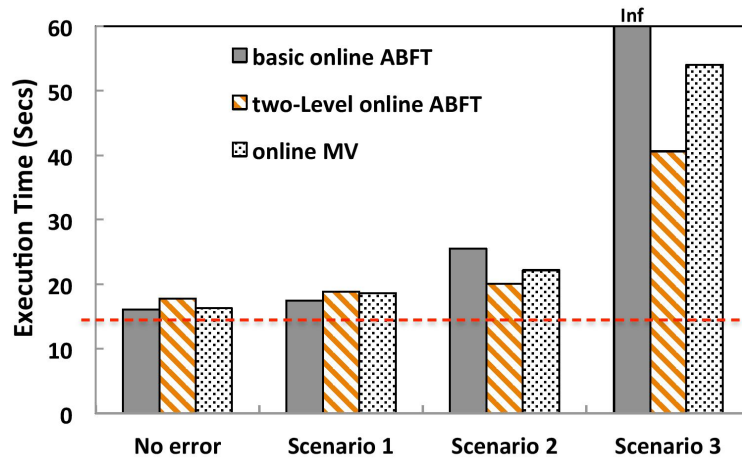
## ◆ Empirical comparison

- Failure-free: overhead is low for both proposed online ABFT (**0.4%** and **1.3%** for PCG, **1.0%** and **4.0%** for PBiCGSTAB)
- Low error rate: **basic online ABFT** has the lowest overhead
- Medium/high error rate: **two-level online ABFT** has the lowest overhead
- Extremely high error rate: **two-level online ABFT** has the lowest overhead (online MV is **48%** higher than two-level)

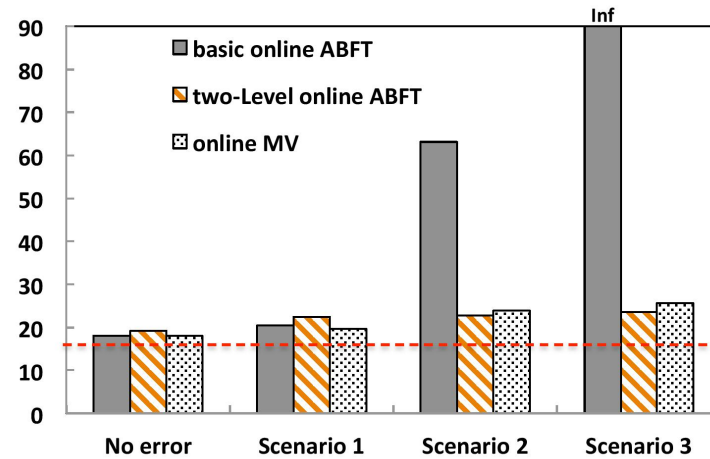
*Consistent with theoretical analysis*



# Empirical Evaluation – Results (cont.)



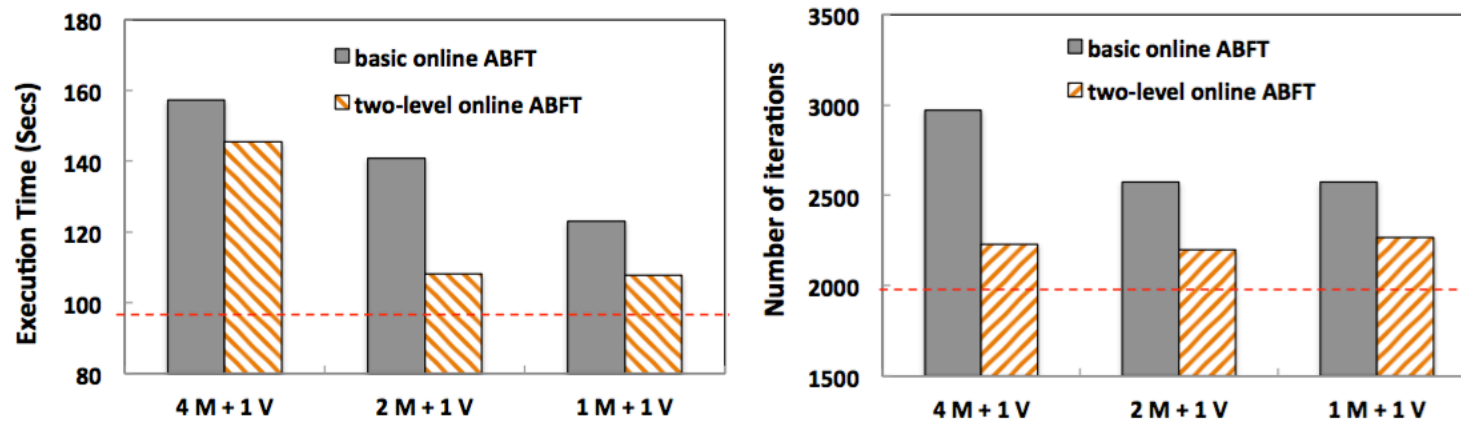
(c) Comparison with **PCG** on *Tianhe-2*



(d) Comparison with **PBiCGSTAB** on *Tianhe-2*

**Similar conclusions as on *Stampede*.**

# Empirical Evaluation – Results (cont.)



(e) Comparison with **PCG** under a **high error-rate** scenario on **Stampede**

**Two-level** online ABFT outperforms **basic** online ABFT by **32.1%** on average under the high error rate scenario.

## Conclusions

- ◆ HPC platforms are anticipated to be more susceptible to soft errors in both logic circuits and memory subsystems
- ◆ Proposed a new checksum encoding mechanism
- ◆ Developed two online ABFT algorithms for general iterative methods – “basic” and “two-level” – that allow errors to be detected eagerly and lazily
- ◆ Experimental results demonstrate our designs are efficient and effective to detect and recover soft errors for general iterative methods

# Thank you !

Contact:

**Dingwen Tao** ([dtao001@cs.ucr.edu](mailto:dtao001@cs.ucr.edu))

**Shuaiwen Leon Song** ([Shuaiwen.Song@pnl.gov](mailto:Shuaiwen.Song@pnl.gov))

Acknowledgement: DOE CESAR Project, DOE ADEM Project and NSF