

Energy-aware Scheduling and Fault Tolerance Techniques for the Exascale Era

IV Congresso REDU

December 1st, 2016

Laércio Lima Pilla

laercio.pilla@ufsc.br

Federal University of Santa Catarina, Brazil



Agenda

EnergySFE Project

Global Scheduling

problems + research + interests

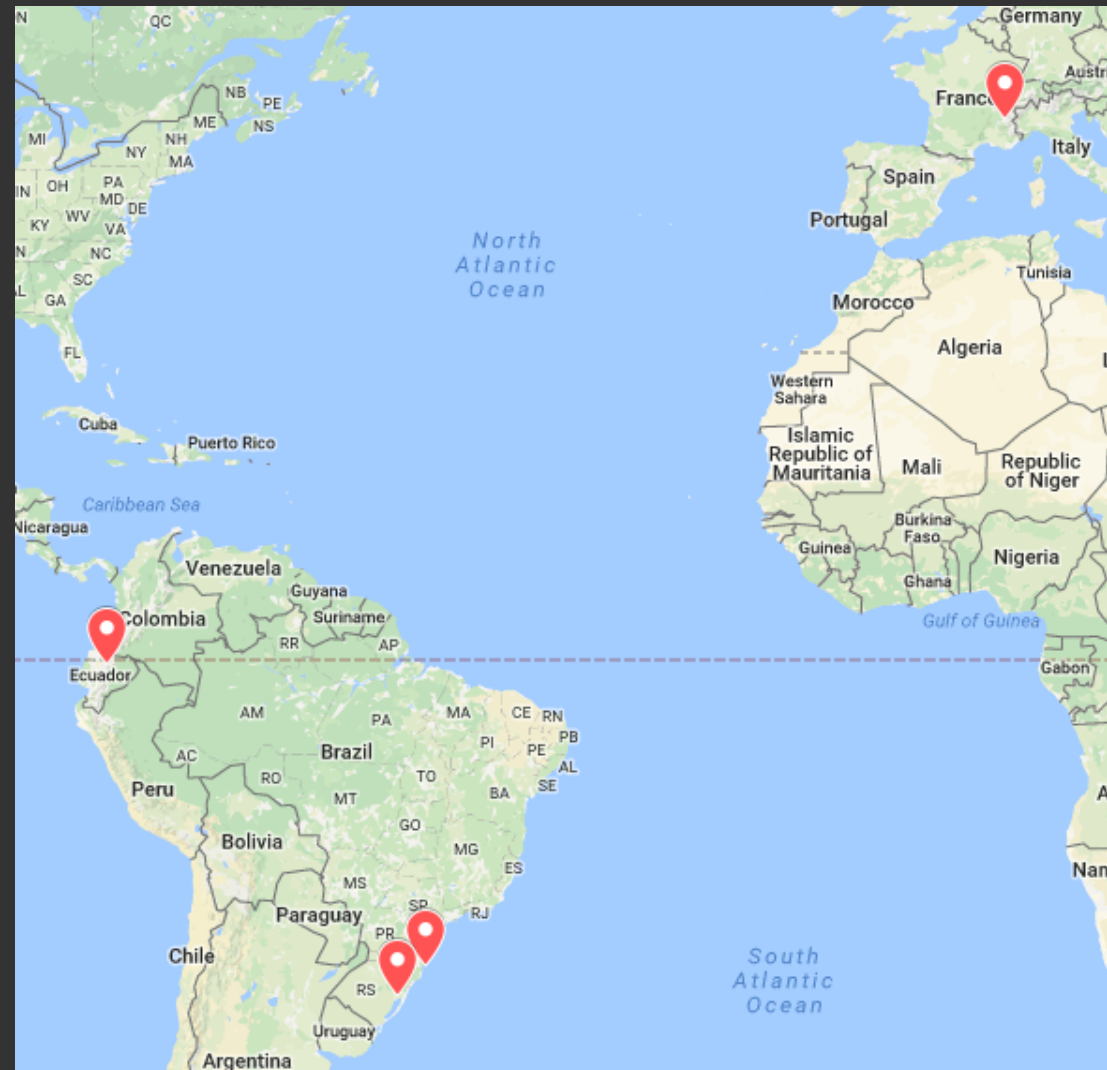
Fault Tolerance

problems + research + interests

EnergySFE Project

EnergySFE Project

STIC-AmSud Project Grant 99999.007556/2015-02



EnergySFE Project

Members

Enrique Vinicio Carrera - **EC**

Márcio Bastos Castro - **BR**

François Broquedis - **FR**

Mario Antônio Dantas - **BR**

Frédéric Desprez - **FR**

Pablo Francisco Ramos - **EC**

Jean-François Méhaut - **FR**

Paolo Rech - **BR**

Laércio Lima Pilla - **BR**

Philippe O. A. Navaux - **BR**

Lucas Mello Schnorr - **BR**

Vanessa C. Vargas - **EC**



<http://energysfe.ufsc.br/>



EnergySFE Project

Energy-aware Scheduling and Fault Tolerance Techniques for the **Exascale Era**

Increasing performance with limited energy



EnergySFE Project

Research Questions

1. How to **schedule tasks and threads** that compete for resources with different constraints while considering the complex hierarchical organization of future Exascale supercomputers?

EnergySFE Project

Research Questions

2. How to **tolerate faults** without incurring in too much overhead in future Exascale supercomputers?

EnergySFE Project

Research Questions

3. How scheduling and fault tolerance approaches can be adapted to be **energy-aware**?

Global Scheduling

problems + research + interests

Global Scheduling Problems

objectives

shortest execution time

highest throughput

highest utilization rate

fairness

...

Global Scheduling Problems

load imbalance

poor initial mapping

load dynamicity

platform sharing

heterogeneous platforms

DVFS changes

...

Global Scheduling Problems

communication slowdown

poor initial mapping

complex communication patterns

hierarchical machine topologies

platform sharing

heterogeneous platforms

...

Global Scheduling Problems

energy wastage

long computing times

low resource usage

excessive data movement

slow scheduling algorithms

...

Global Scheduling Research

What we want with scheduling

balanced loads

optimized communications

fast scheduling algorithms

least migrations possible

Global Scheduling Research

Periodic load balancing

principle of persistence

topology-aware load balancing

energy-aware load balancing

Global Scheduling Research

topology-aware load balancing

NucoLB

greedy algorithm + NUMA node latency

HwTopoLB

refinement algorithm + latency & bw + convergence

HierarchicalLB

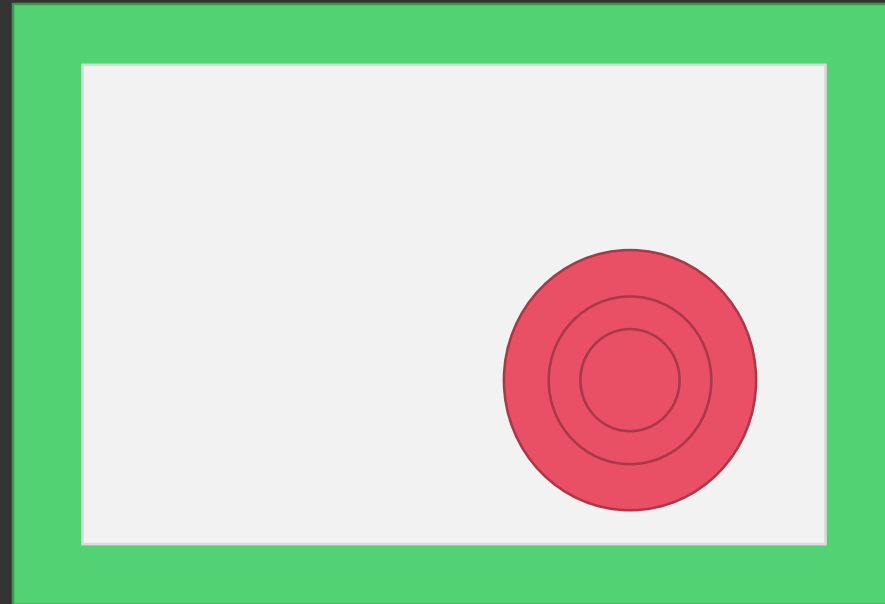
LB at node level and at whole machine level

Global Scheduling Research

Application: Ondes3D

irregular loads

dynamic loads

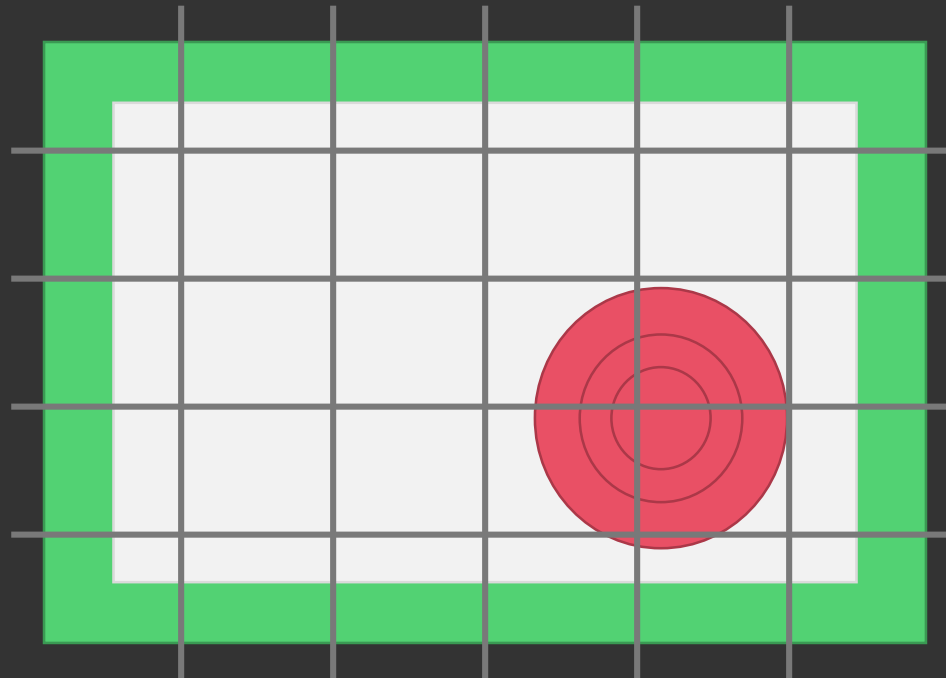


Global Scheduling Research

Application: Ondes3D

irregular loads

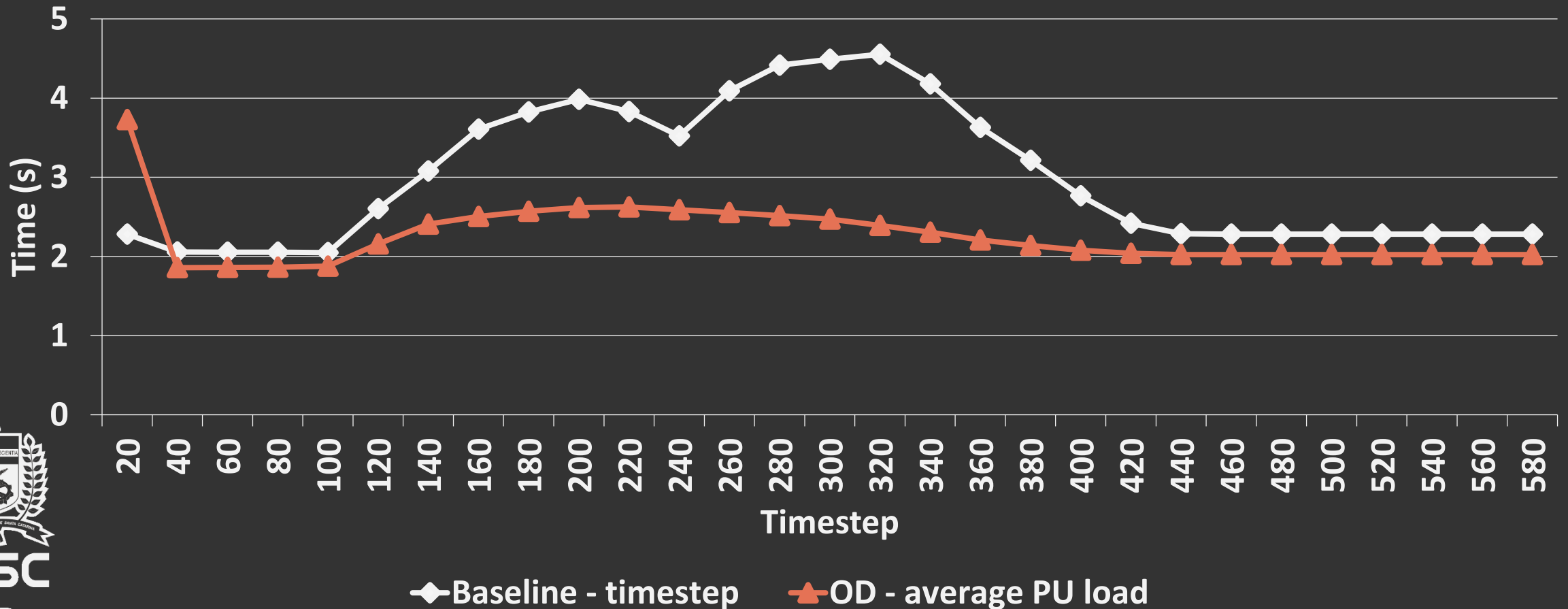
dynamic loads



Global Scheduling Research

Application: Ondes3D

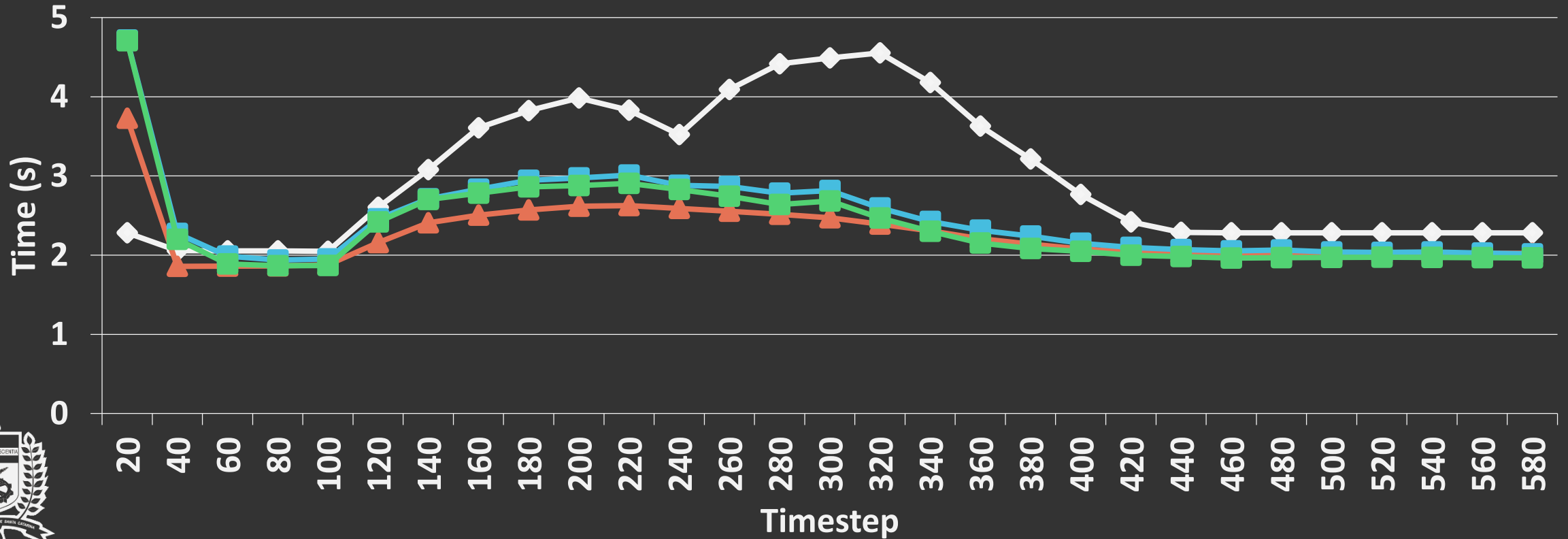
OD = overdecomposition (16 tasks / core)



Global Scheduling Research

Application: Ondes3D

OD = overdecomposition (16 tasks / core)

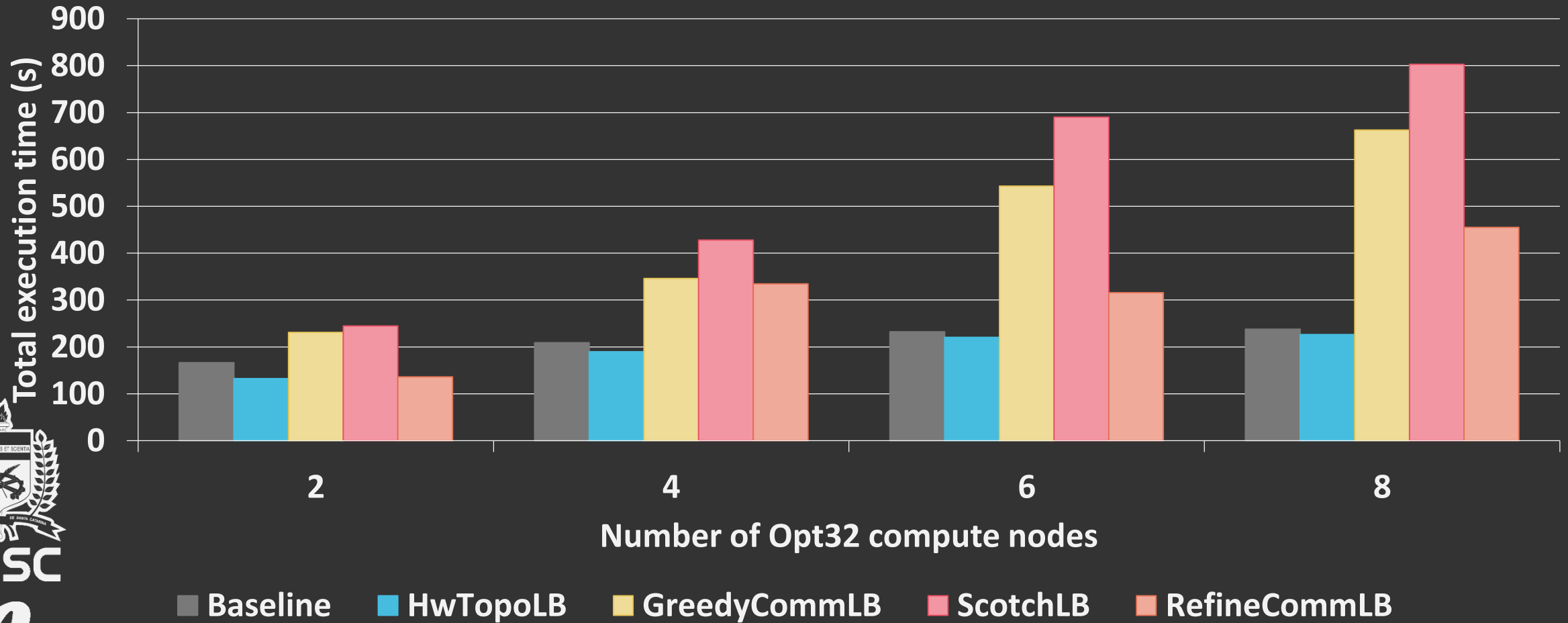


◆ Baseline - timestep ▲ OD - average PU load ■ HwTopoLB - timestep ■ NucoLB - timestep

Global Scheduling Research

Application: **LeanMD**

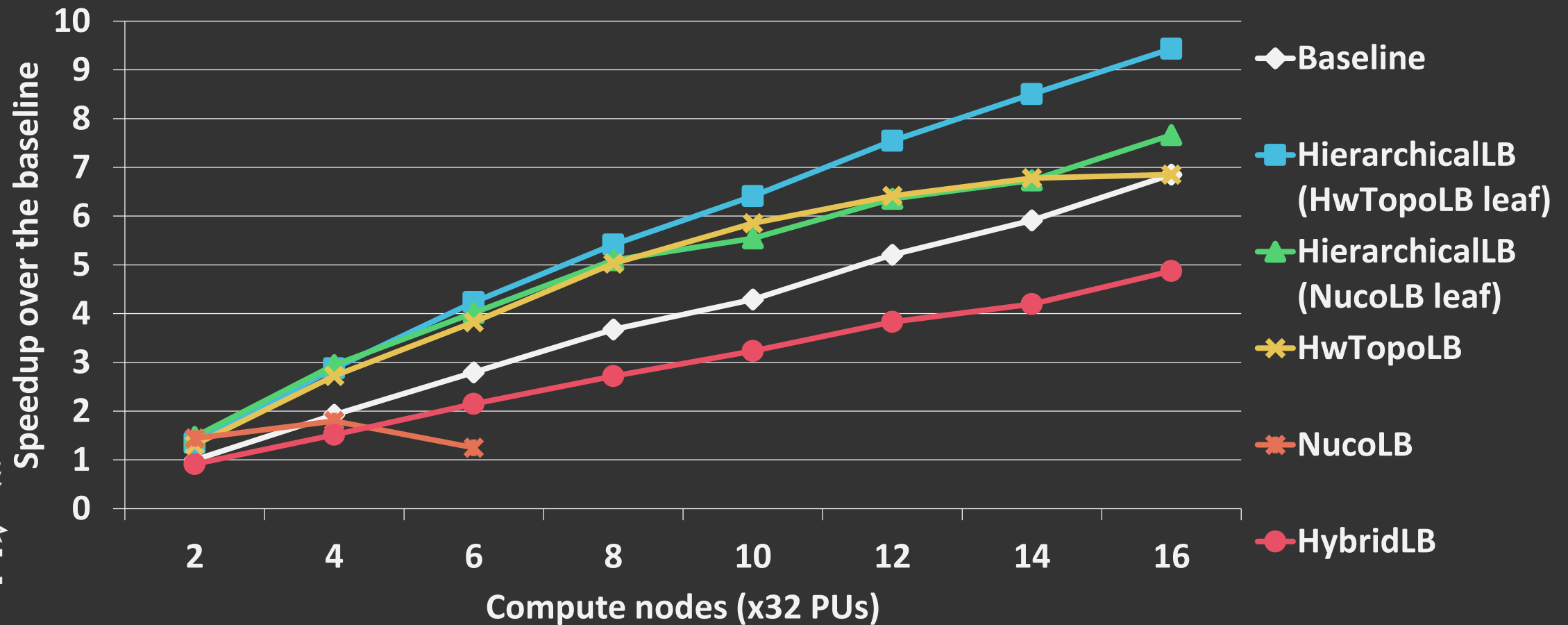
Weak scalability: Total execution time



Global Scheduling Research

Strong scalability: LeanMD

Speedup over the baseline on 2 CNs



Global Scheduling Research

energy-aware load balancing

LB + DVFS on underloaded resources

Fine-Grained EnergyLB

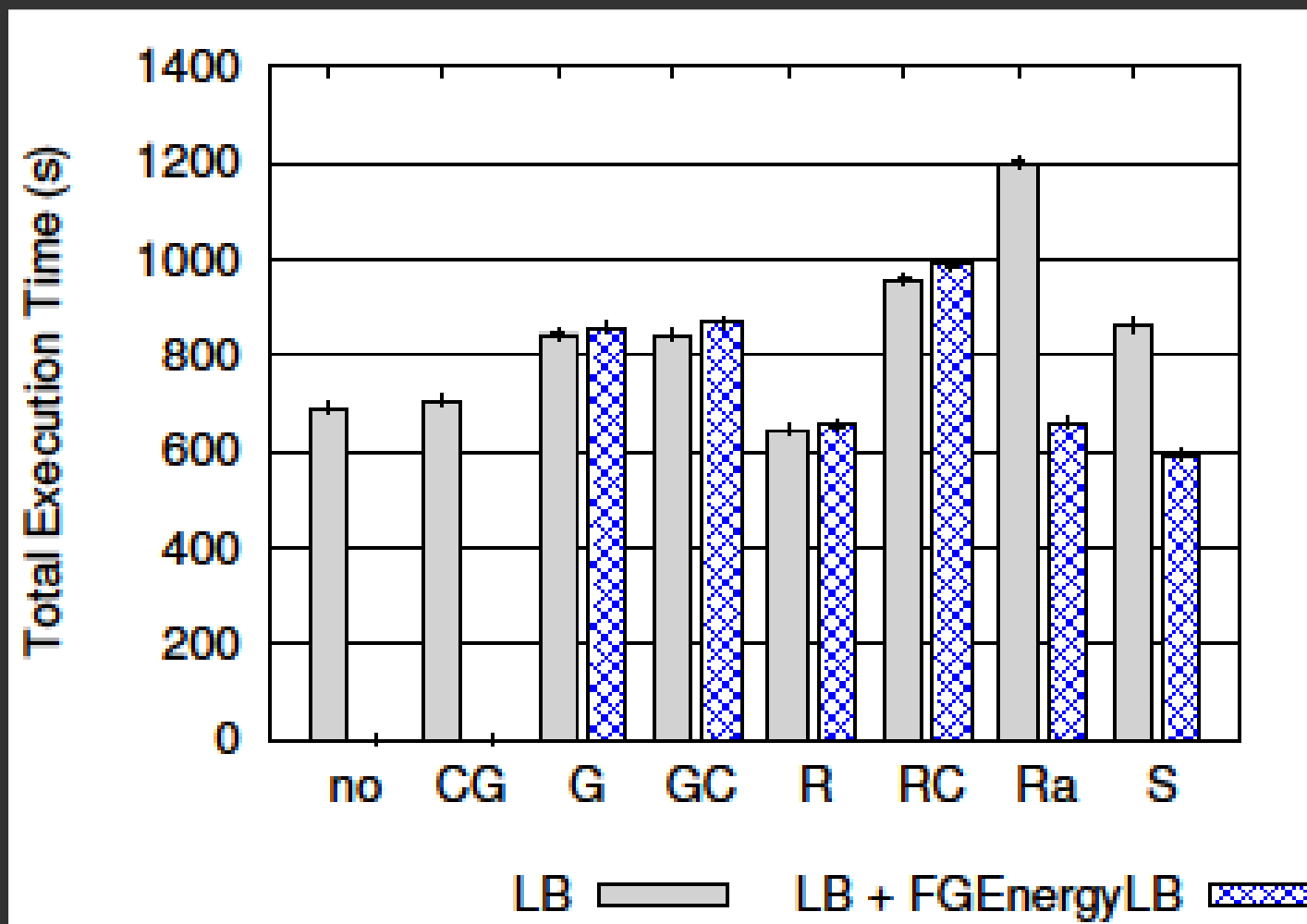
per core

Coarse-Grained EnergyLB

per socket + hierarchical algorithm

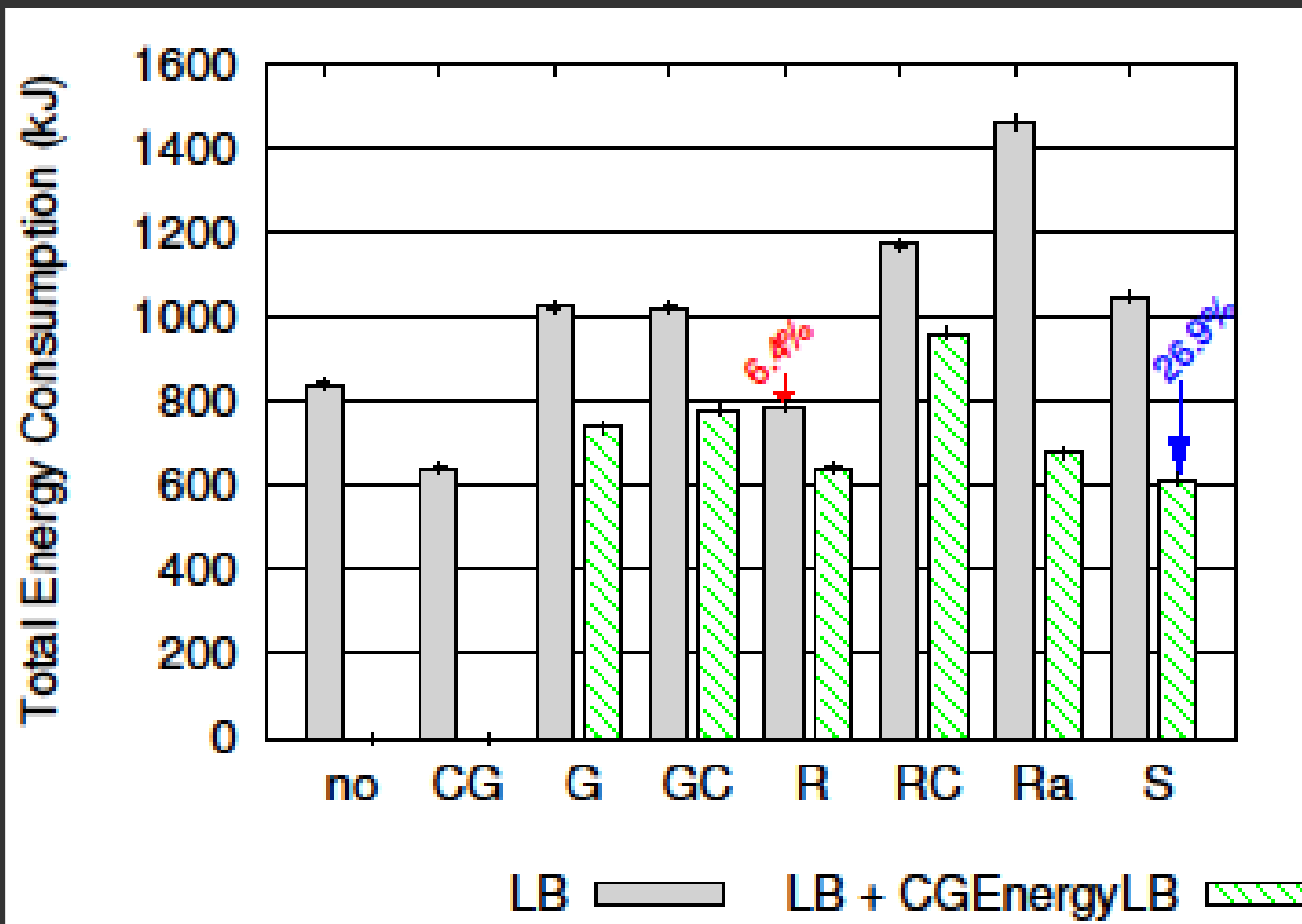
Global Scheduling Research

Example: Lulesh + FG EnergyLB



Global Scheduling Research

Example: Lulesh + FG EnergyLB



Global Scheduling Interests

- + hierarchical algorithms [+ energy]
- + distributed algorithms [+ energy]
- + platforms for experiments [+ energy]
- + applications

Fault Tolerance

problems + research + interests

Fault Tolerance Problems

failure: systems stops working as expected

fault → error → failure

Fault Tolerance

use redundancy to stop failures from occurring

Fault Tolerance Problems

usual fault tolerance schemes

checkpoint & restart

cannot notice corrupted data

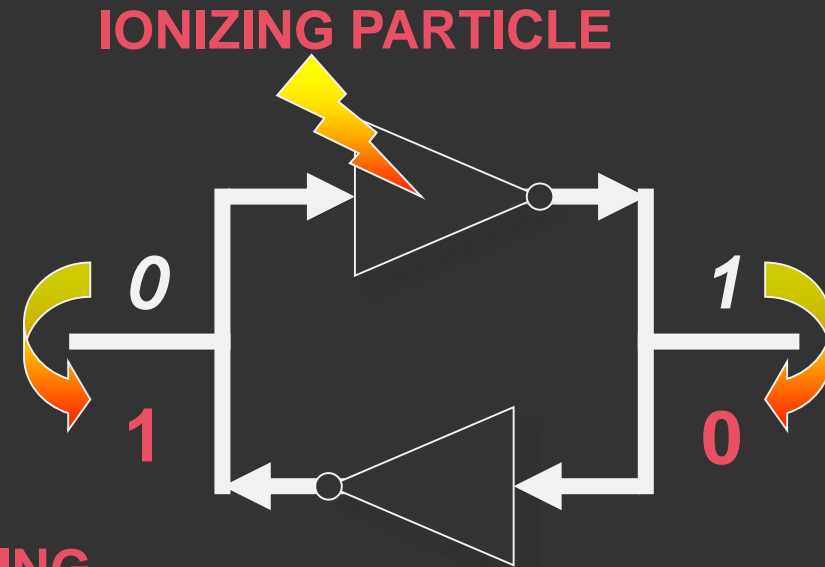
duplication with comparison

expensive in resource (time & energy)

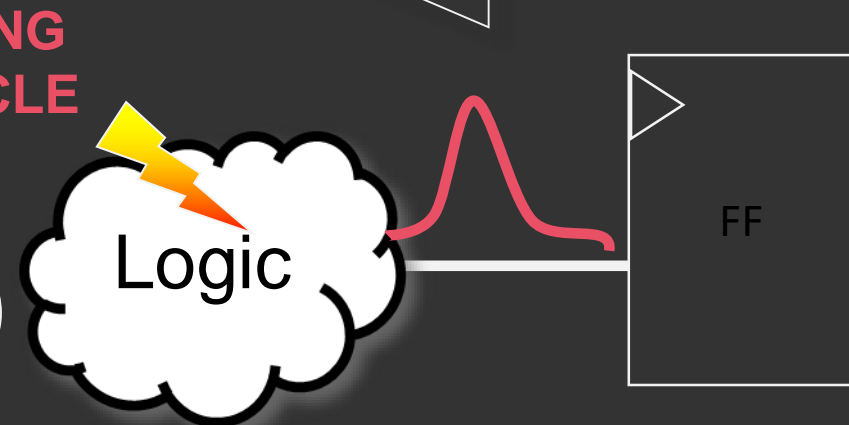
Fault Tolerance Problems

Soft Errors: the device is not permanently damaged, but a particle may generate:

One or more bit-flips
Single Event Upset (**SEU**)
Multiple Bit Upset (**MBU**)



Transient voltage pulse
Single Event Transient (**SET**)



Fault Tolerance Problems

Silent Data Corruption

data caches

register files

ALU

scheduler

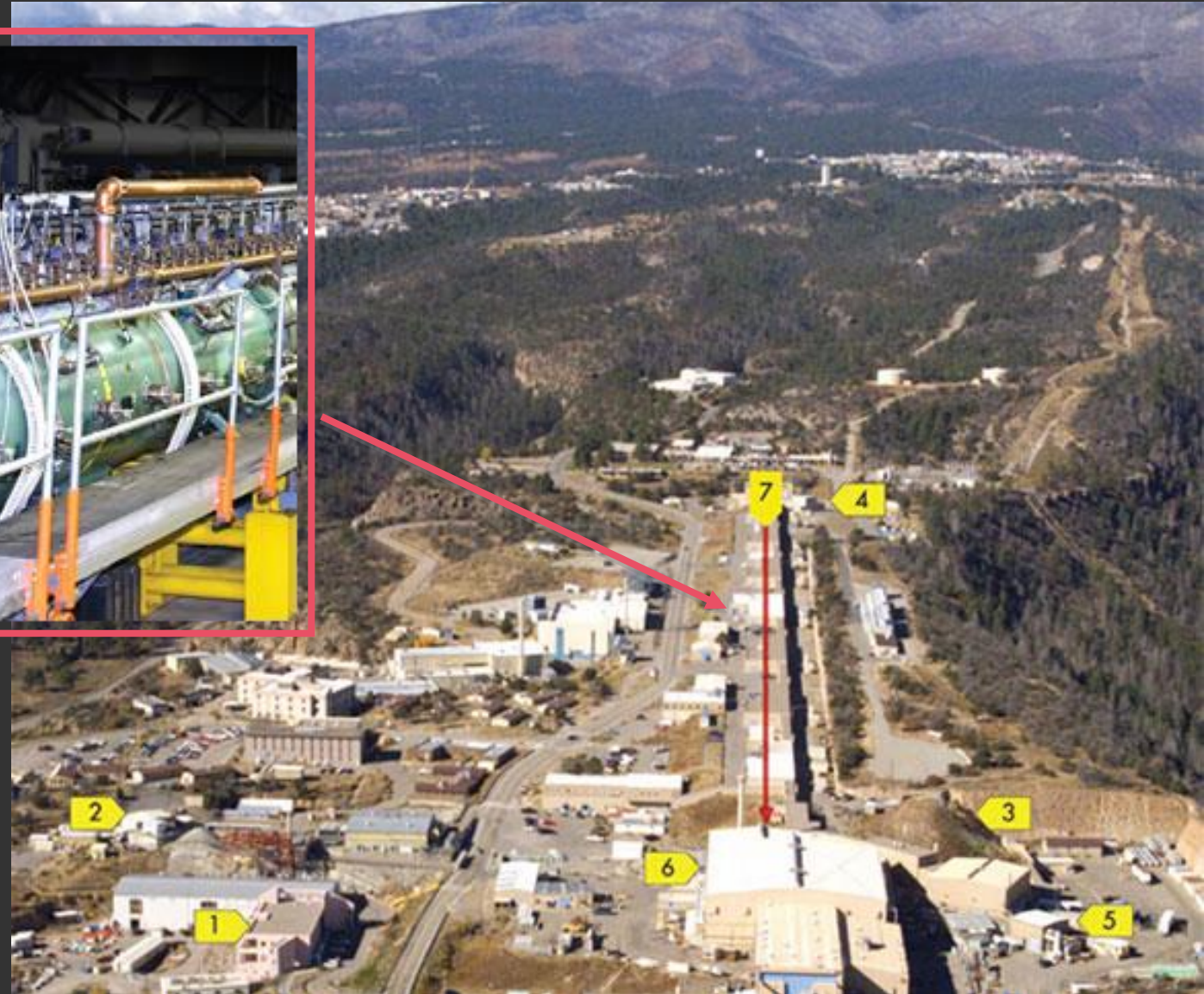
Crash

instruction cache

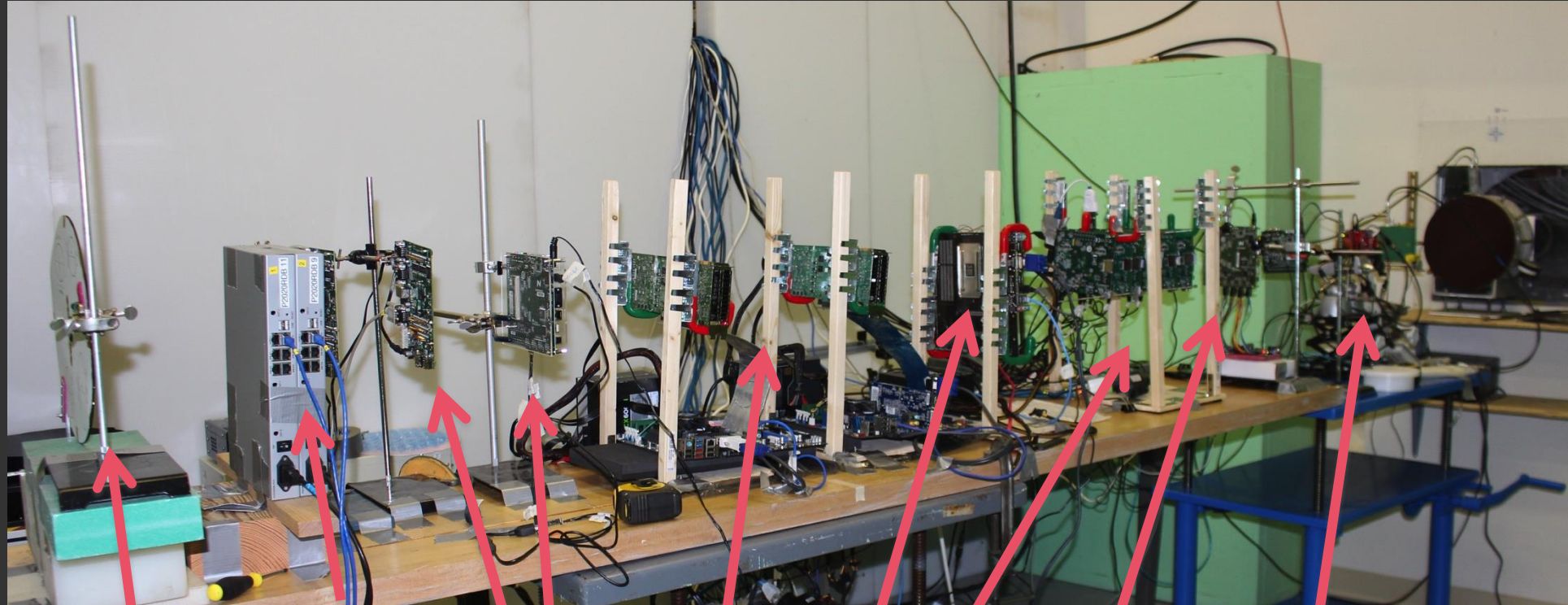
scheduler

PCIe bus

Fault Tolerance Research



Fault Tolerance Research



Flash

SoC

FPGA

GPU

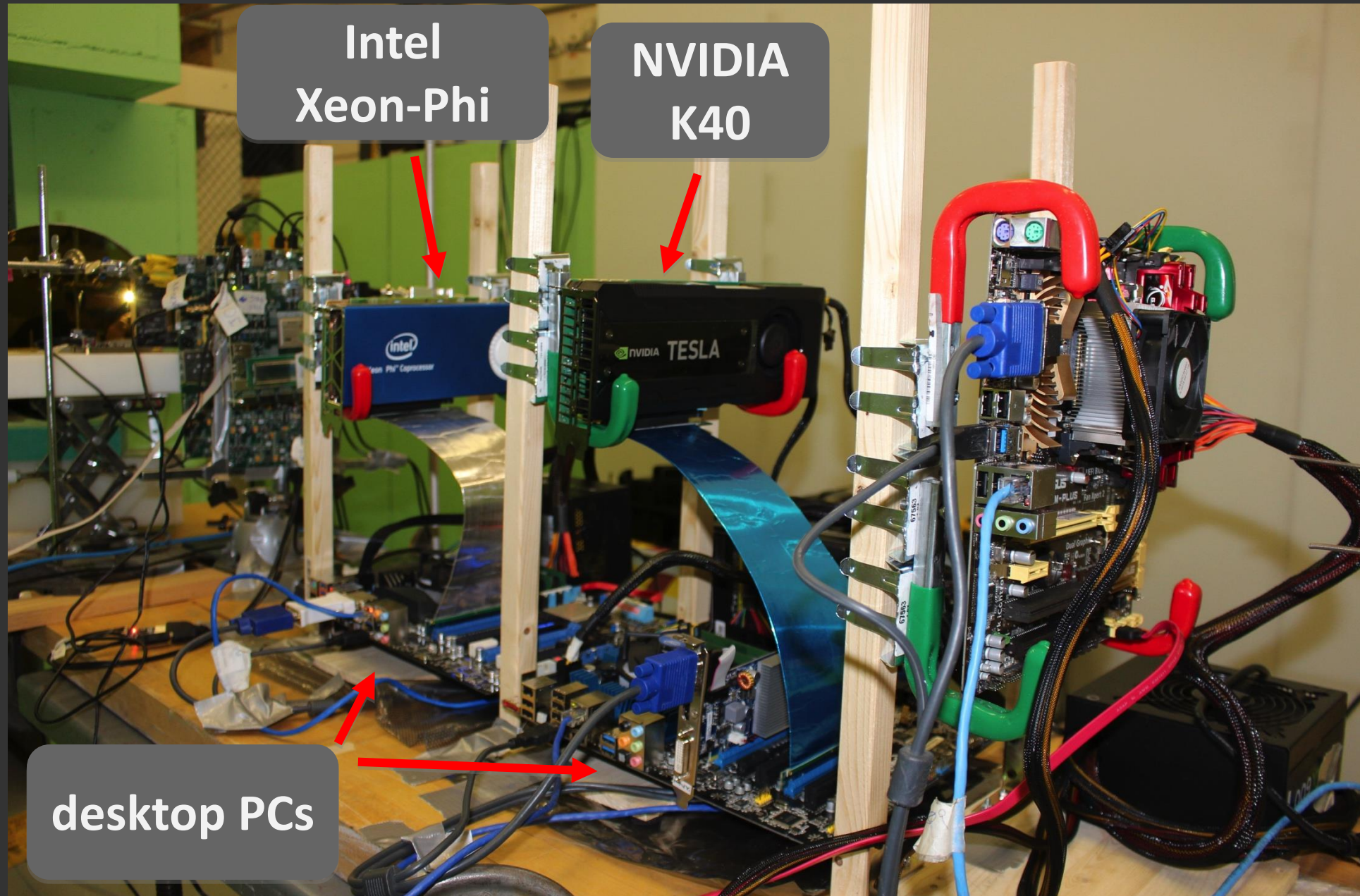
APU

SoC

FPGA

microcontrollers

Fault Tolerance Research



Intel
Xeon-Phi

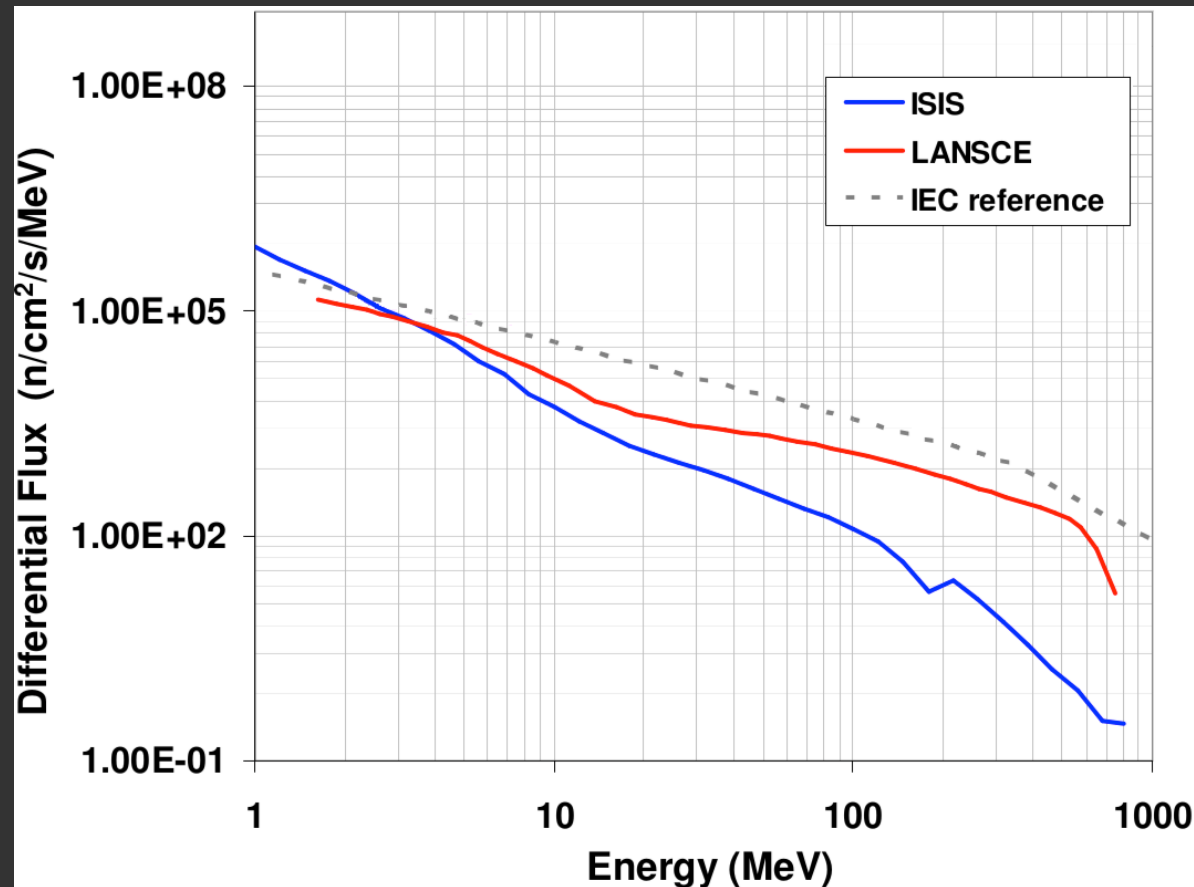
NVIDIA
K40

desktop PCs

Fault Tolerance Research

@LANSCE 1.8×10^6 n/(cm² h)

@NYC 13 n/(cm² h)



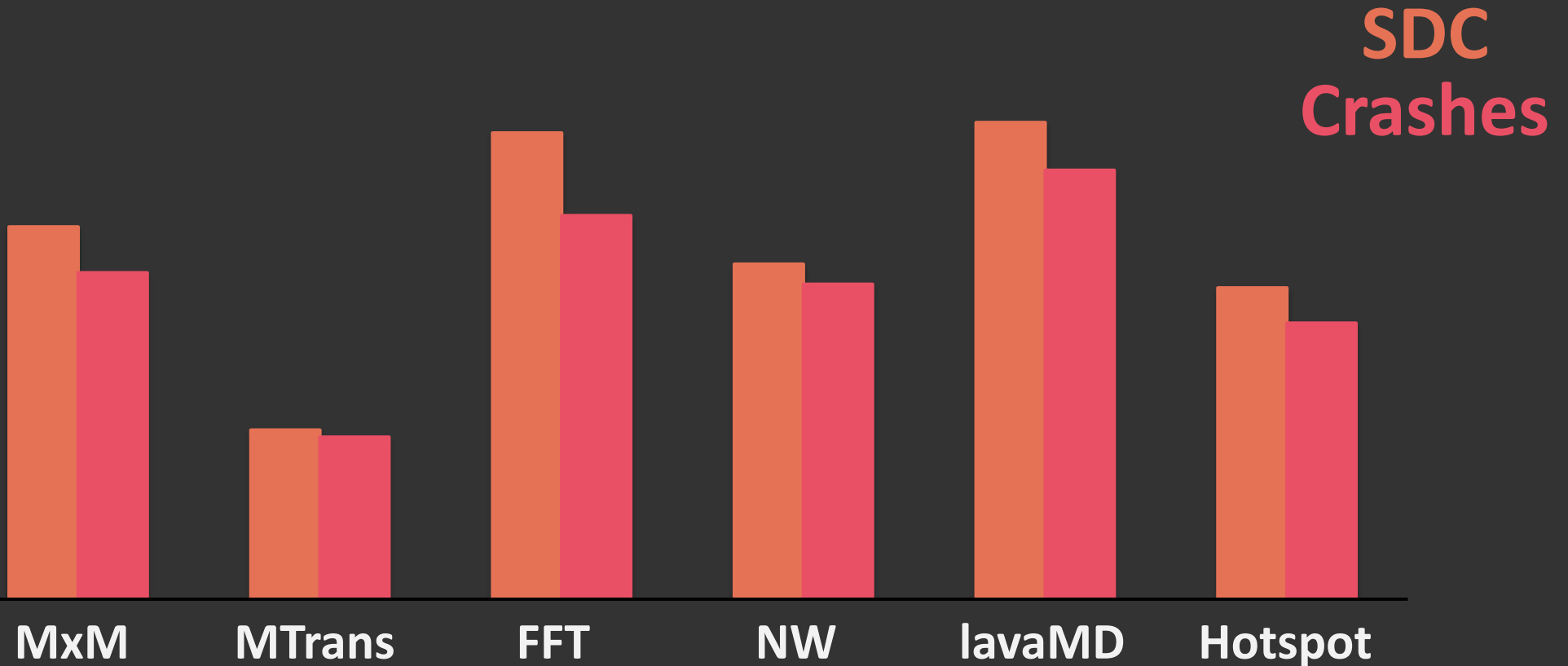
Fault Tolerance Research

Parallel algorithms' reliability

SDC rates vary ~3 orders of magnitude

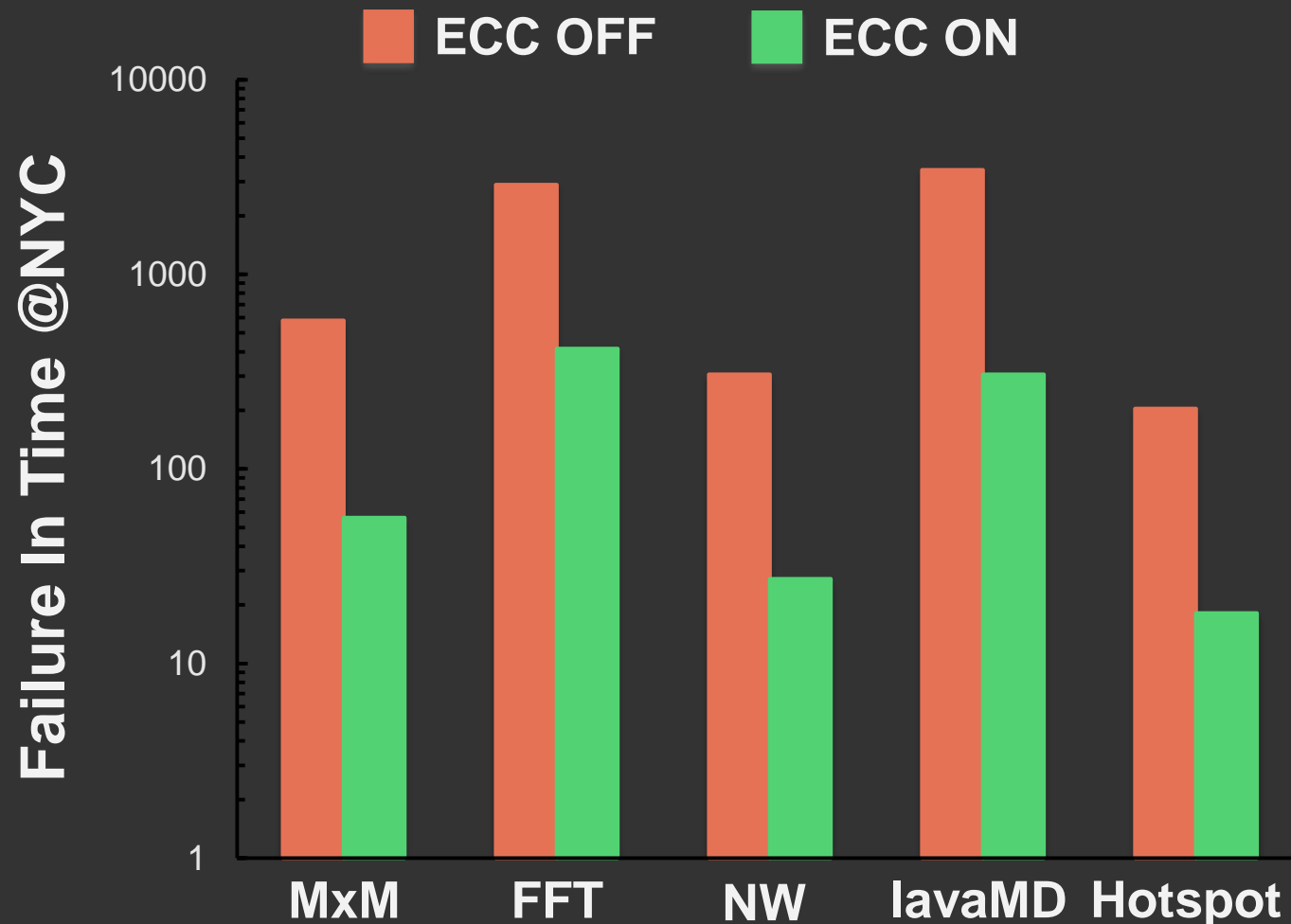
Failure In Time @NYC

10000
1000
100
10
1



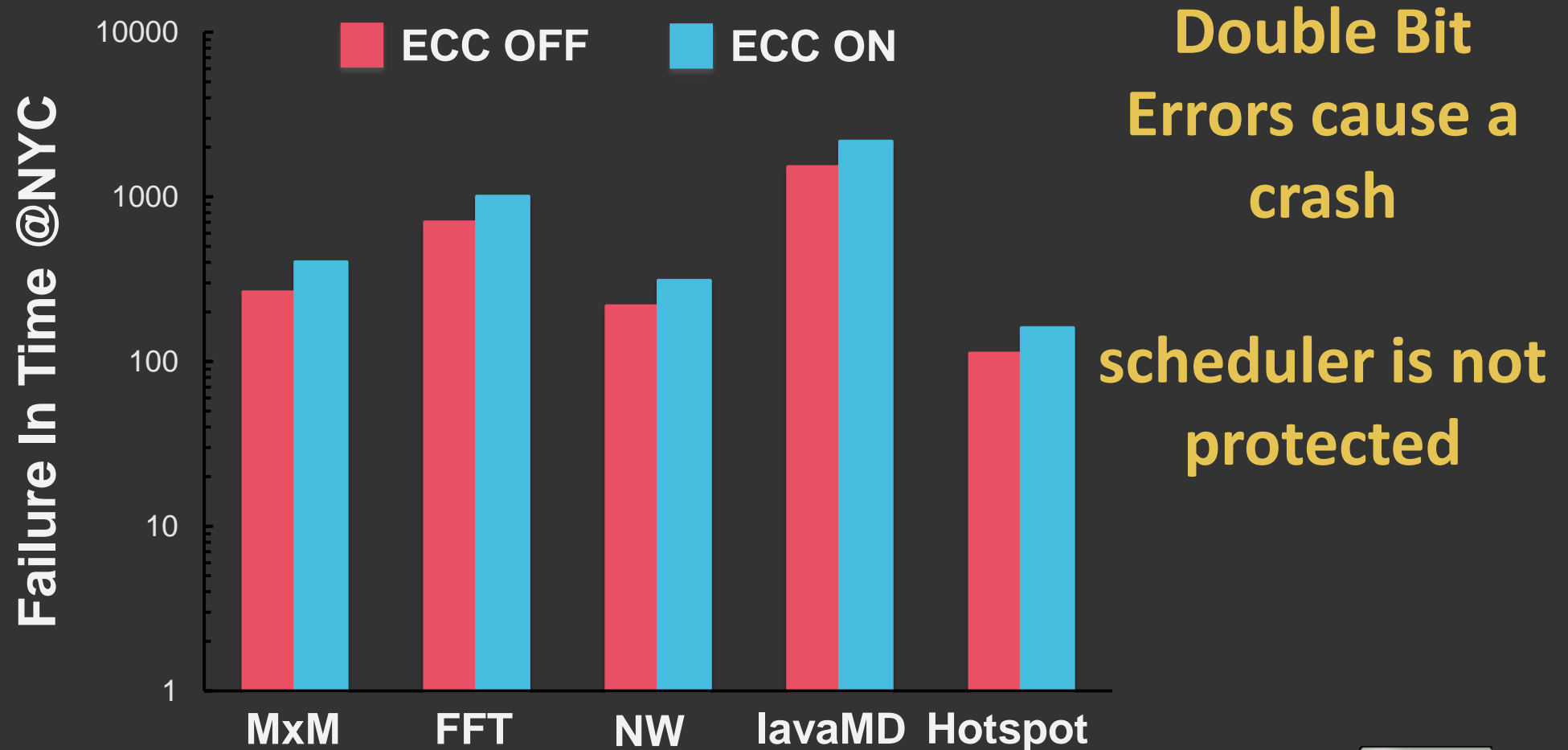
Fault Tolerance Research

ECC reduces the **SDC FIT** of ~1 order of magnitude
(there is almost no code dependence)



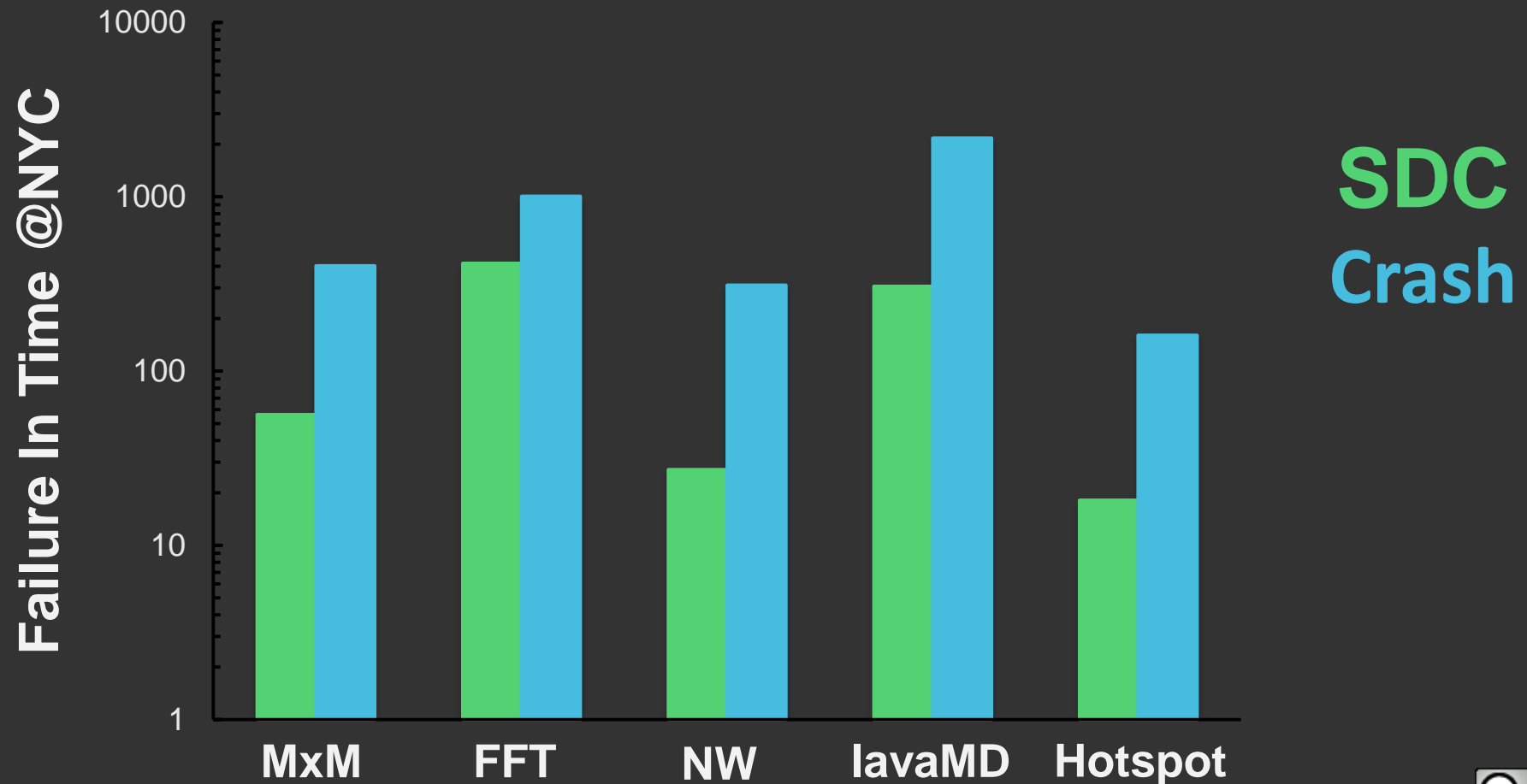
Fault Tolerance Research

ECC increases the **Crash FIT** of about 50%
(there is almost no code dependence)



Fault Tolerance Research

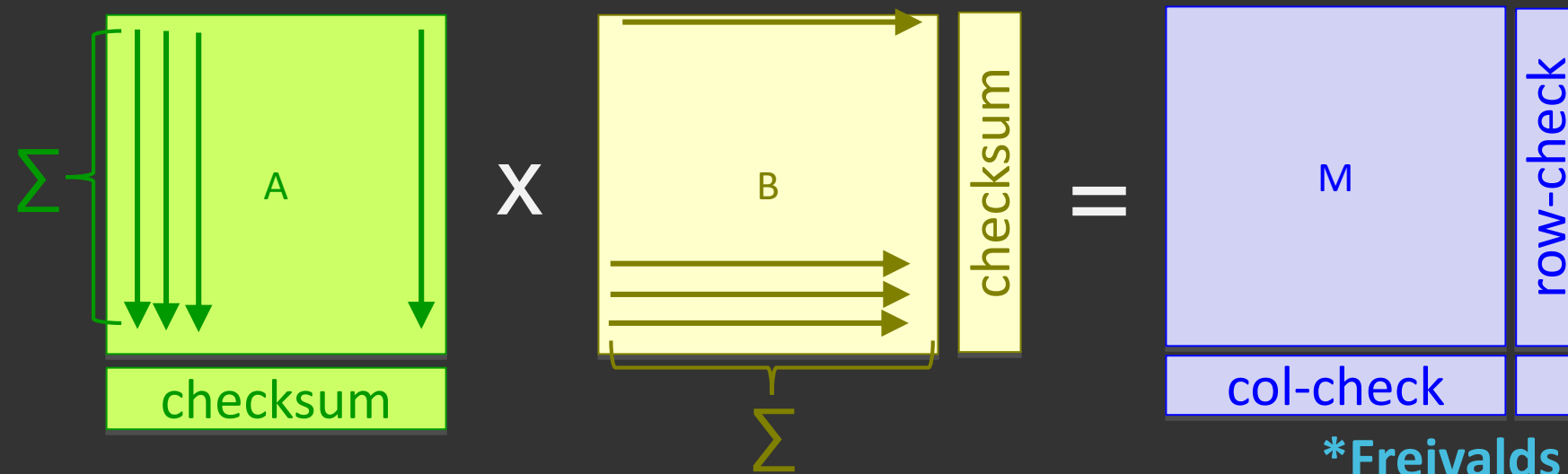
When the ECC is ON **Crashes** are more likely to occur than **SDCs** (this is GOOD for HPC centers!)



Fault Tolerance Research

ABFT: technique designed specifically for an algorithm.

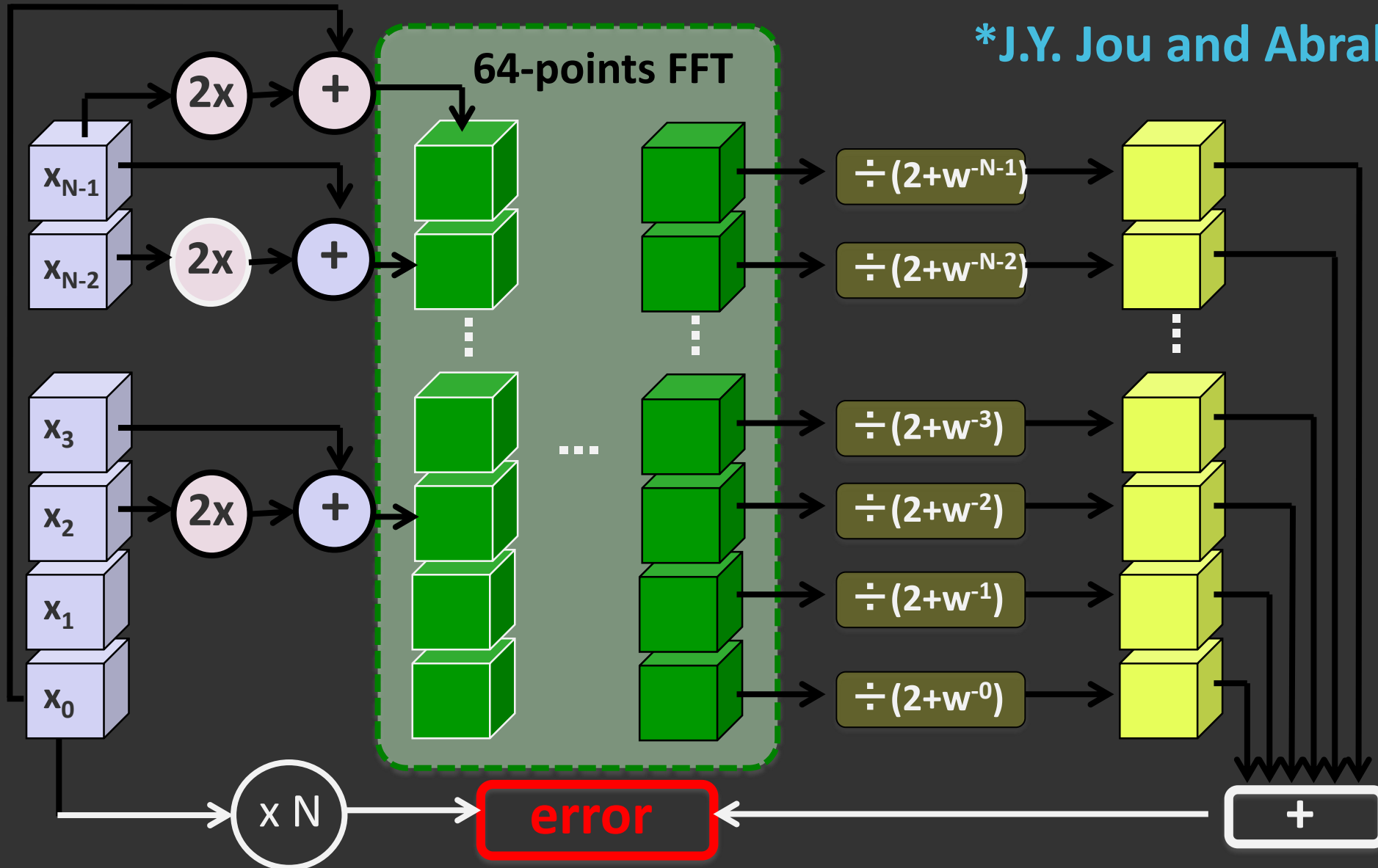
Usually ABFT requires input coding, algorithm modification, and output decoding with error detection/correction



*Freivalds '79

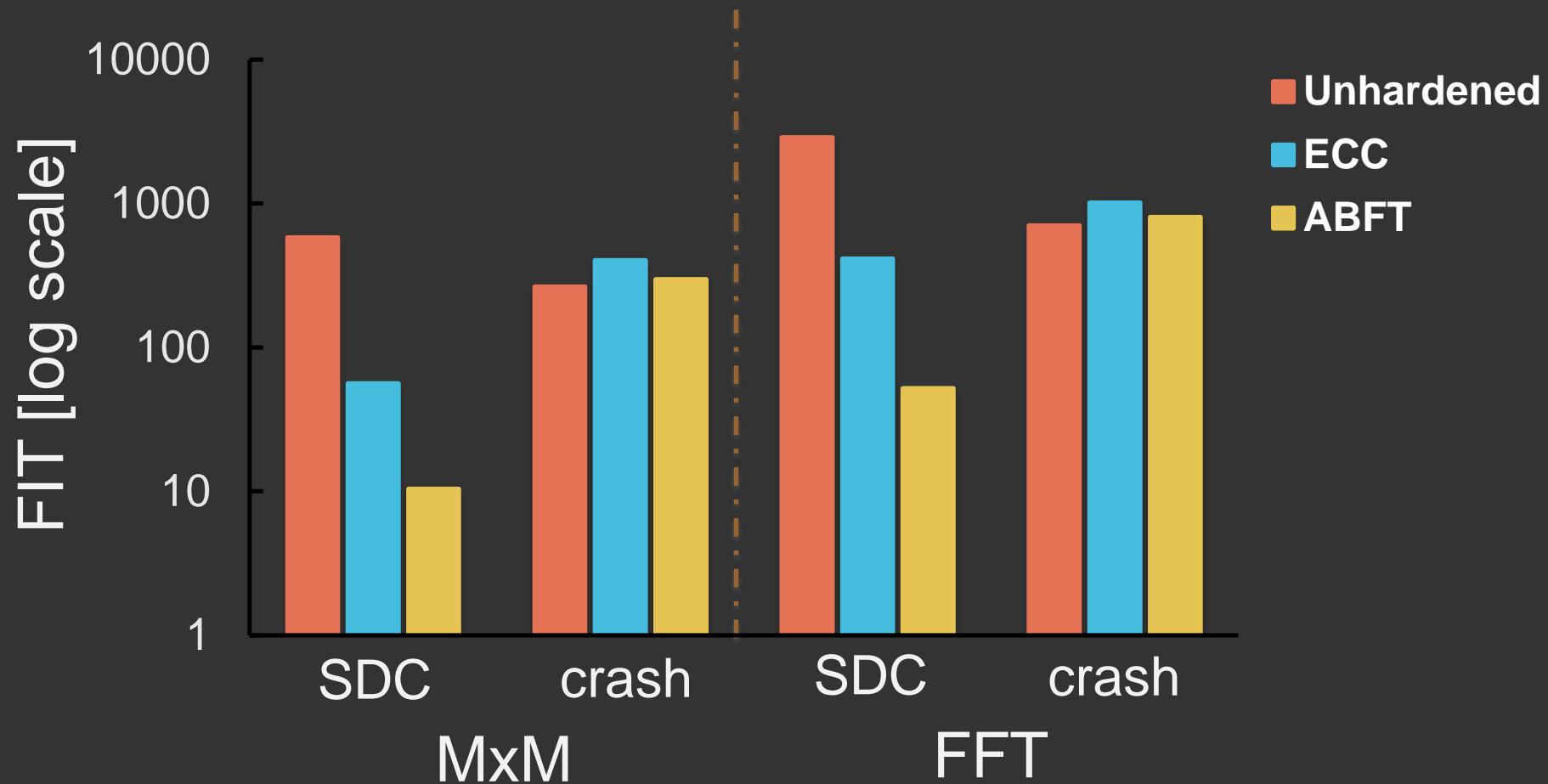
Fault Tolerance Research – ABFT for FFT

*J.Y. Jou and Abraham '88



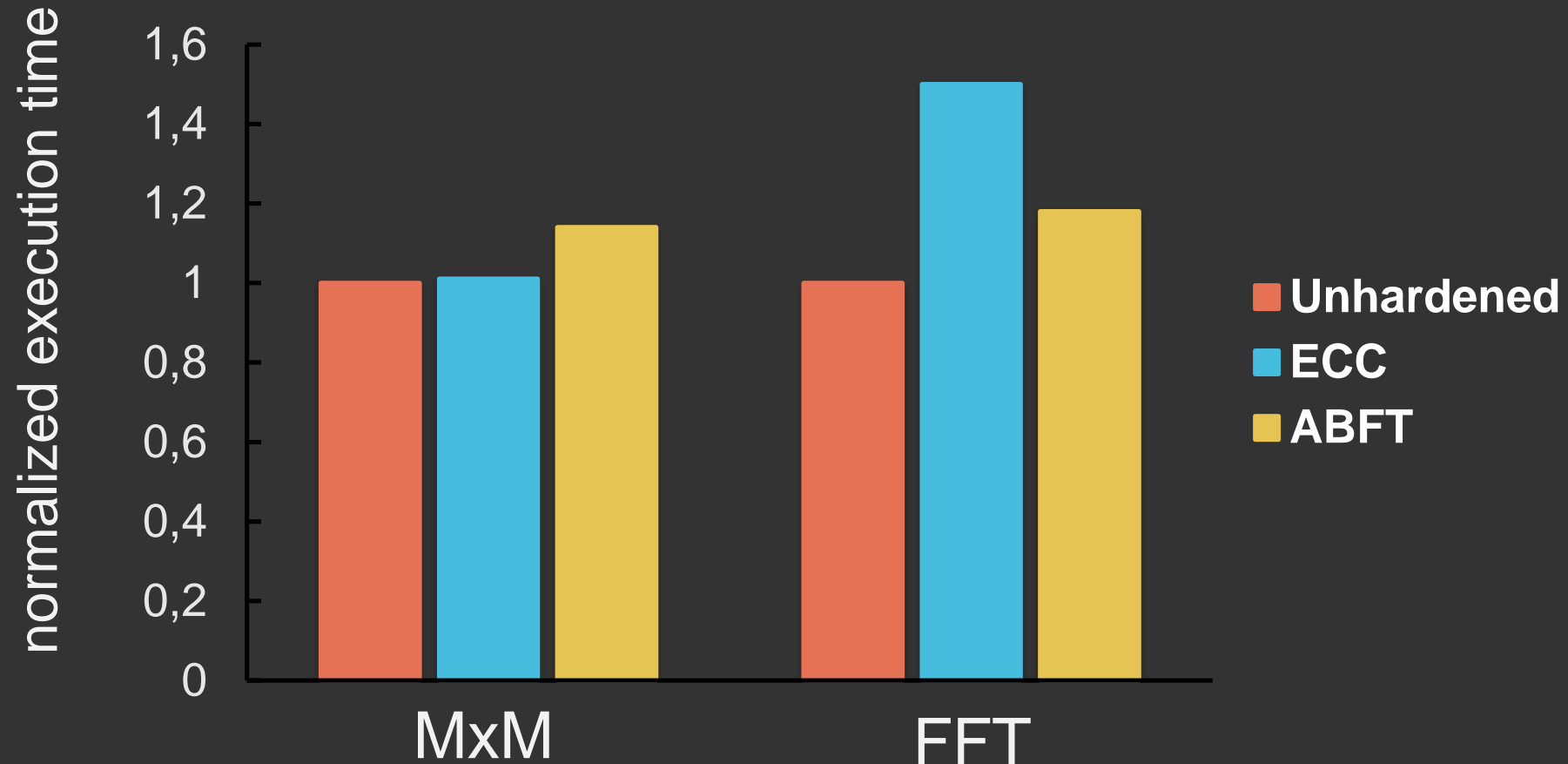
Fault Tolerance Research

ECC x ABFT



Fault Tolerance Research

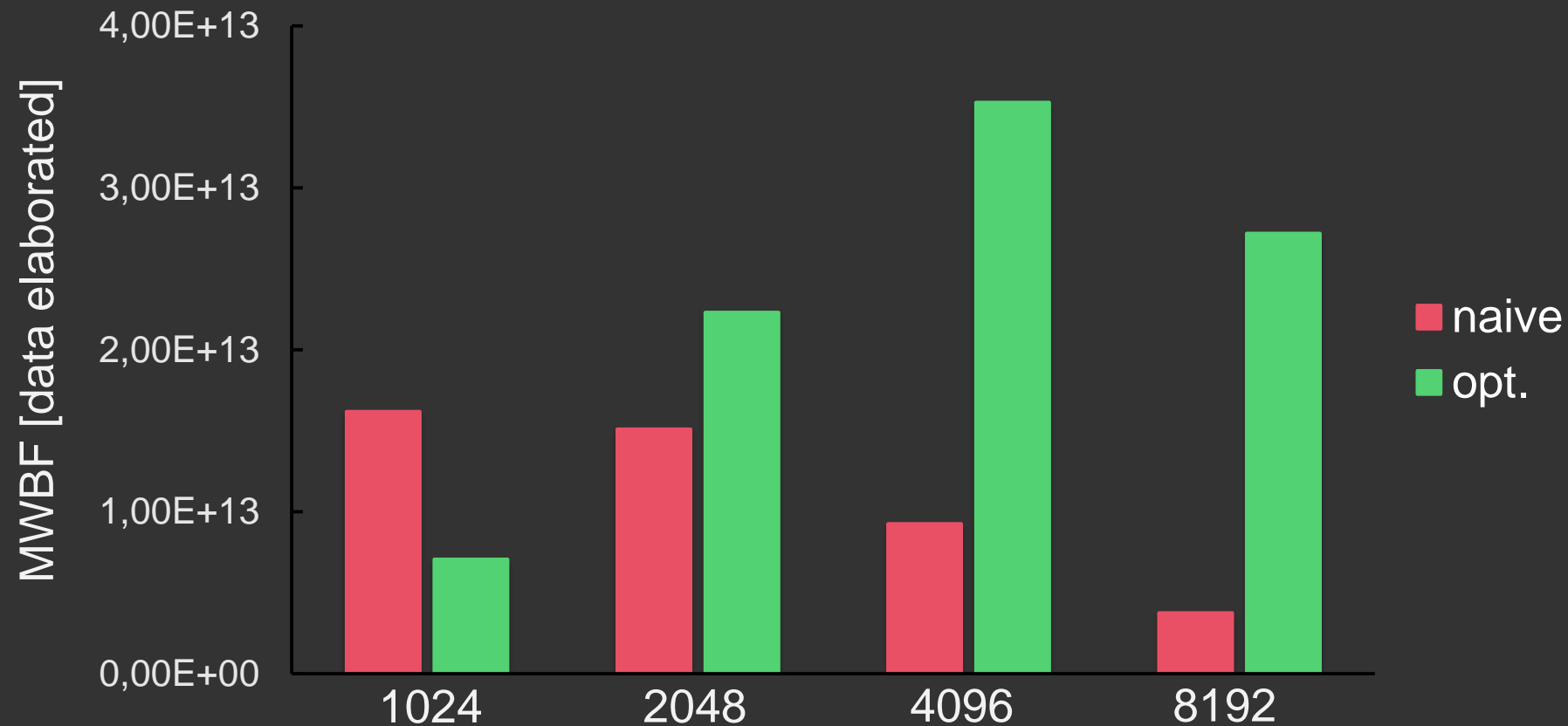
ECC x ABFT



Fault Tolerance Research

Effects of optimizations: MxM

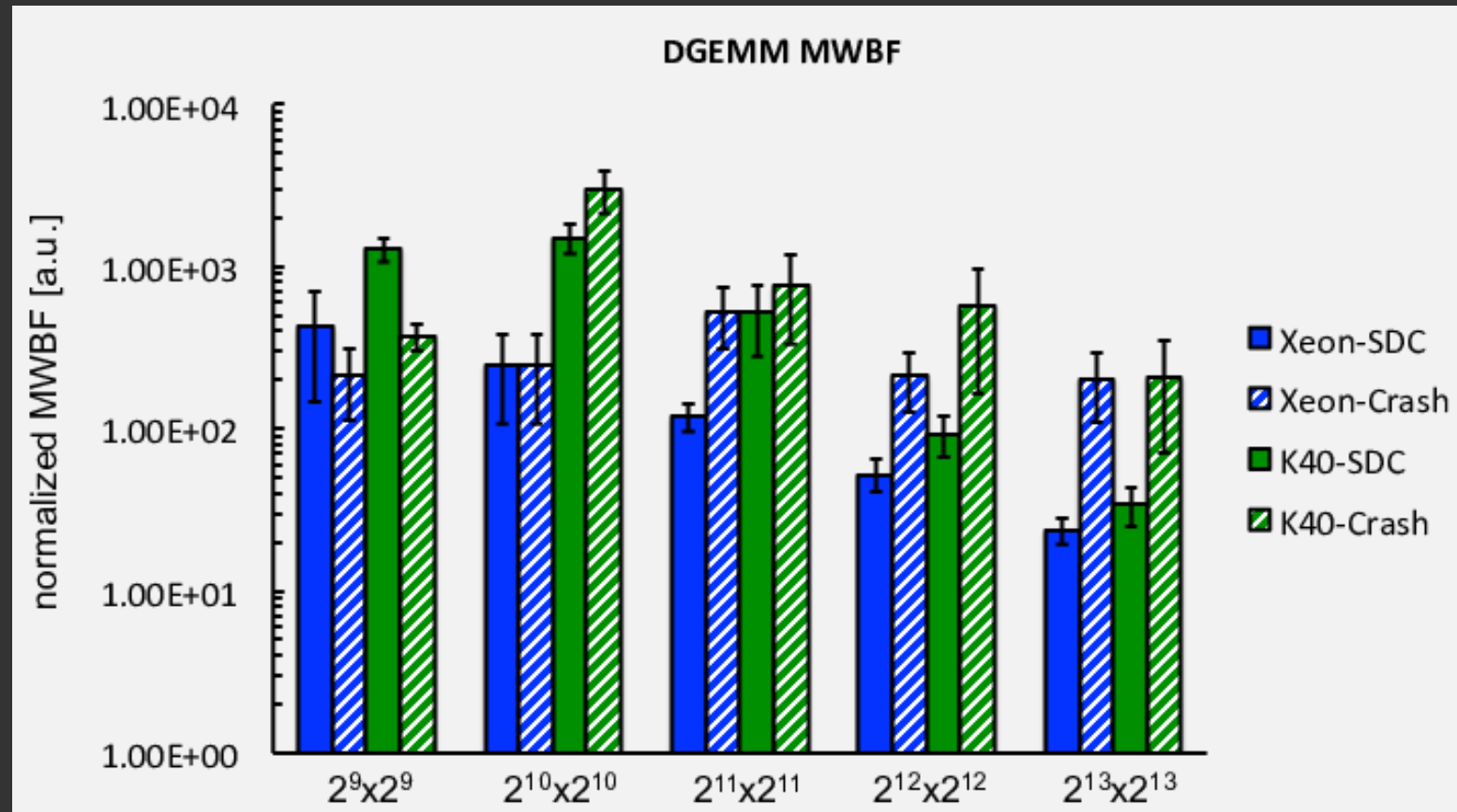
MWBF: Mean Workload Between Failures



Fault Tolerance Research

Xeon-Phi MWBF decreases significantly with input size.

Even if more prone to be corrupted, Kepler produces more correct data (if parallelism is exploited)



Fault Tolerance Interests

- + applications & kernels for CPU & GPU
- + platforms for energy measurement
- + experiments with radiation

Conclusions

Conclusions

EnergySFE

international collaboration

three research questions

lots of work in scheduling and fault tolerance

lots more to be done

Energy-aware Scheduling and Fault Tolerance Techniques for the Exascale Era

Thank you.

Laércio Lima Pilla

laercio.pilla@ufsc.br

Federal University of Santa Catarina, Brazil

