Reliability Evaluation of Mixed-Precision Architectures

Fernando F. Santos, Caio Lunardi, Daniel Oliveira, <u>Fabiano Libano</u>, and Paolo Rech





Double (64 bit)

Single (32 bit)

floating point operations increased graphics quality and user experience...



Double (64 bit)

Single (32 bit)

floating point operations increased graphics quality and user experience...

as well as physical simulation precision and accuracy





Double (64 bit)

Single (32 bit)

Mixed-precision architectures can improve performance, efficiency and reliability.



Approximate computing has shown that some operations in some algorithms can be approximated without affecting significantly the final result







Reduced precision operations are particularly interesting for **Neural Networks training and execution**



Reduced precision operations are particularly interesting for **Neural Networks training and execution**



Reduced precision operations are particularly interesting for **Neural Networks training and execution**

Neural Nets accuracy*

Network	FP32 Baseline	Mixed precision	Half (16 bit)
AlexNet	56.8%	56.9%	
VGG-D	65.4%	65.4%	Short INT (8 bit)
GoogLeNet	68.3%	68.4%	
Inception v2	70.0%	70.0%	(1 bit)
Inception v3	73.9%	74.1%	
Resnet 50	75.9%	76.0%	
ResNeXt 50	77.3%	77.5%	
		*data fro	om nvidia nvidia

Reduced precision operations are particularly interesting for **Neural Networks training and execution**



Reliability of Mixed-Precision

Mixed-Precision delivers:

- -Higher Performance
- -Smaller Area



-Lower Power Consumption

But how does it affect overall system Reliability?

Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions







Reliability Evaluation of Mixed-Precision Architectures – INF, UFRGS



Reliability Evaluation of Mixed-Precision Architectures – INF, UFRGS

Terrestrial Radiation Environment



Galactic Cosmic rays interact with atmosphere and produce a shower of energetic particles:

- Muons
- Pions
- Protons
- Gamma Rays
- Neutrons

13 n/(cm²x h) @sea level* *JEDEC JESD89A Standard

Terrestrial Radiation Environment



Galactic Cosmic rays interact with atmosphere and produce a shower of energetic particles:

- Muons
- Pions
- Protons
- Gamma Rays
- Neutrons

13 n/(cm²x h) @sea level* *JEDEC JESD89A Standard

Soft Errors: the device is not permanently damaged, but the particle may generate **bit flip(s)** in memory or **logic error(s)**



Silent Data Corruption vs Crash

Soft Errors in: -data cache -register files -logic gates (ALU) -scheduler

Silent Data Corruption

Silent Data Corruption vs Crash

Soft Errors in: -data cache -register files -logic gates (ALU) -scheduler

Soft Errors in: -instruction cache -scheduler / dispatcher -PCI-e bus controller Silent Data Corruption

DUE (Crash)

Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions

Experimental Methodologies (Devices)

FPGA: resources utilization is tailored by the user higher precision => higher area

Experimental Methodologies (Devices)

FPGA: resources utilization is tailored by the user higher precision => higher area

X86 (Xeon Phi): the Vector Processing Unit (VPU) executes operations in 64 or 32 bits, on the same HW.Compiler decides how to use the VPUs

Experimental Methodologies (Devices)

FPGA: resources utilization is tailored by the user higher precision => higher area

X86 (Xeon Phi): the Vector Processing Unit (VPU) executes operations in 64 or 32 bits, on the same HW.Compiler decides how to use the VPUs

GPU (Volta V100): dedicated HW for double and single/half precision

1 double operation, 1 single or 2 half operations

Fault Injection

Fault Injection

We purposely inject faults in registers and variables.

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Neutron Beam

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Neutron Beam

We irradiate the devices while running benchmarks.

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Neutron Beam

We irradiate the devices while running benchmarks.

Realistic estimation of the error rate of a given device executing a given application.

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Neutron Beam

We irradiate the devices while running benchmarks.

Realistic estimation of the error rate of a given device executing a given application.

> [FIT] Failure in Time

Fault Injection

We purposely inject faults in registers and variables.

Probability of faults propagating to the output.

[PVF/AVF] Program Vulnerability Factor Architectural Vulnerability Factor

Neutron Beam

We irradiate the devices while running benchmarks.

Realistic estimation of the error rate of a given device executing a given application.

> [FIT] Failure in Time





Radiation Experiments @ChipIR





Radiation Experiments @ChipIR




Radiation Experiments @ChipIR



Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions

Mixed-Precision Reliability



It is much more likely for a double value/operation to be corrupted than a half value/operation

Mixed-Precision Reliability

Double has a much larger area than Half Memory: 4x Functional Units: ~16x Double Half

It is much more likely for a double value/operation to be corrupted than a half value/operation

However, a fault in a double value is much less critical than a fault in half

Double 52/64 (81%) bits are mantissa Half 10/16 (60%) bits are mantissa

Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions



Matrix Multiplication (128x128) $\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \times \begin{bmatrix} 2 & 0 \\ 1 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 4 \\ 10 & 8 \end{bmatrix}$













As we lower the precision, we also decrease circuit area, ultimately reducing the overall error rate [FIT].

















Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions



Xeon Phi - Beam

One-size fits all: Xeon Phi does not have dedicated HW for double/single





Xeon Phi - Beam

One-size fits all: Xeon Phi does not have dedicated HW for double/single





Xeon Phi - Beam

One-size fits all: Xeon Phi does not have dedicated HW for double/single



We inject faults in variables during execution

DUE Masked SDC 100% **Program Vulnerability Factor** 80% 60% 40% 20% 0% Double Single Double Single Double Single LavaMD LUD MxM

We inject faults in variables during execution



We inject faults in variables during execution



We inject faults in variables during execution



Xeon Phi - Criticality



Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions

Volta GPU



Clock cycles depend only on data-type, not on operation. MUL, ADD, FMA have different complexity.

We test micro-benchmarks to investigate GPU reliability (realistic codes are shown in the paper)



Volta GPU - FIT rates

SDC Double SDC Single

SDC Half DUE

Double has higher FIT rate even if there are fewer cores than **Single** and **Half** (2,688 vs 5,376)

MUL hardware complexity (and, then, error rate) increases a lot when precision is increased







GPU - fault injection

We inject errors in registers (lowest possible level, yet)

SDC DUE Masked 100% Architectural Vulnerability Factor 80% 60% 40% 20% 0% **Double Single** Half **Double Single** Half **Double Single** Half Micro MUL Micro ADD Micro FMA



GPU - fault injection

We inject errors in registers (lowest possible level, yet)





GPU - fault injection

We inject errors in registers (lowest possible level, yet)





GPU - Error Criticality



GPU - Neural Networks





We tested YOLOv3 implemented in Double, Single, Half precision

GPU - Neural Networks





We tested YOLOv3 implemented in Double, Single, Half precision



We consider as critical faults that significantly modify detection/classification


GPU - Neural Networks





GPU - Neural Networks



Reliability Evaluation of Mixed-Precision Architectures – INF, UFRGS

Outline

- Radiation Effects Introduction
- Experimental Methodology
- Reliability of Mixed-Precision Architectures
 - FPGAs
 - x86
 - GPUs
- Conclusions

Conclusions and Future Work

-Mixed precision architectures significantly improve performance

-Reducing precision impacts the code error rate in a non-obvious way

-Low precision can improve reliability: more data can be processed before experiencing an error (details in the paper)

-Future work: duplication using mixed-precision to detect critical faults with low-overhead



Reliability Evaluation of Mixed-Precision Architectures – INF, UFRGS

Acknowledgments



Lucas Weigel Lucas Klein Pedro Pimenta Philippe Navaux Luigi Carro

Nathan DeBardeleben Sean Blanchard Los Alamos Heather Quinn Thomas Fairbanks Steve Wender



Chris Frost Carlo Cazzaniga



Timothy Tsai Siva Hari **NVIDIA**. Michael Sullivan **Steve Keckler**



Matteo Sonza Reorda Luca Sterpone





Reliability Evaluation of Mixed-Precision Architectures

Fernando F. Santos, Caio Lunardi, Daniel Oliveira, <u>Fabiano Libano</u>, and Paolo Rech



