Towards Statistical Guarantees in Controlling Quality Tradeoffs for Approximate Acceleration

Divya Mahajan

Amir Yazdanbakhsh Jongse Park Bradley Thwaites Hadi Esmaeilzadeh

Alternative Computing Technologies (ACT) Lab Georgia Institute of Technology





Approximate Computing

Relax the abstraction of "near perfect" accuracy in



Data Processing

Storage



Accept imprecision to improve performance energy efficiency





Approximate Acceleration







Approximate Acceleration



Application



















Motivation

Challenges in devising MITHRA

A hardware software solution

Detailing the components of MITHRA





Exploiting Accelerator Characteristics









Factors Influencing Error



Accelerator Error = | Output_{accelerator} - Output_{original}|

Output_{accelerator} = f (accelerator inputs, accelerator configuration)





Factors Influencing Error



Accelerator Error = | Output_{accelerator} - Output_{original} |

Output_{accelerator} = f (accelerator inputs, accelerator configuration)





Factors Influencing Error



Accelerator Error = | Output_{accelerator} - Output_{original}|

Output_{accelerator} = f (accelerator inputs, accelerator configuration) Constant Output_{accelerator} = f (accelerator inputs)







How to eliminate anomalous invocations?



Look at the accelerator inputs



















What guarantees?













What algorithm at runtime?

Application

Application

Application

Application

Application

Application

Binomial Proportion Confidence Interval

Less than the desired programmer quality loss → nsuccess

Greater than the desired programmer quality loss

Binomial Proportion Confidence Interval

with a confidence level

$(n_{trials}, n_{success}) \Rightarrow r < SuccessRate$ with a confidence level

E.g., (100, 80) **•** 72% < SuccessRate

with 95% confidence level

$(n_{trials}, n_{success}) \Rightarrow r < SuccessRate$ with a confidence level

E.g., (100, 80) **•** 72% < SuccessRate

with 95% confidence level

Final Quality Level, Success Rate and Confidence Interval programmer specified

Statistical Optimization

if desired metrics are not met:

 $th_{t+1} = th_t - \Delta$

else if desired metrics are met: $th_{t+1} = th_t + \Delta$

Reiterate; till the tht meets the metrics but tht+1 doesn't

Statistical Optimization

if desired metrics are not met:

$$th_{t+1} = th_t - \Delta$$

else if desired metrics are met: th_{t+1}= th_t + Δ

Reiterate; till the tht meets the metrics but tht+1 doesn't

Tighter threshold better Final Quality Loss Level, Success Rate and Confidence Interval but lower benefits from approximation

Training the Classifiers

The training data used to generate classifier topology

Hardware Classifiers

Simple algorithm that can be easily implemented in hardware.

We use two techniques for this work:

- 1. Table Based
- 2. Neural Network Based

Table-based Classifiers

Table-based Classifiers

A small **ensemble** of table-based classifiers achieve better accuracy and performance

Neural Network Based Classifiers

Benchmarks

	Accelerator Topology	Baseline Error	Neural Classifier Topology
Blackscholes	$6 \rightarrow 8 \rightarrow 8 \rightarrow 1$	6.02%	$6 \rightarrow 4 \rightarrow 2$
		-	
FFT	$1 \rightarrow 4 \rightarrow 4 \rightarrow 2$	7.22%	$1 \rightarrow 4 \rightarrow 2$
		-	
Inversek2j	$2 \rightarrow 8 \rightarrow 2$	7.50%	$2 \rightarrow 4 \rightarrow 2$
JMEINT	$18 \rightarrow 32 \rightarrow 8 \rightarrow 2$	17.69%	18 → 16 → 2
JPEG Encoding	$64 \rightarrow 16 \rightarrow 64$	7.00%	64 → 2 → 2
		-	
Sobel	9 → 8 → 1	9.96%	9 → 4 →2

Table Classifier Topology8 tables each of size 0.5 KB

Energy and Performance Benefits

Invocation Rate

Conclusion

Thank you

Thank you

False Positive and False Negative Results

Binomial Proportion Confidence Interval

Evaluation of MITHRA

False Positive and False Negatives

Multi-Input Hashing and Input Signature Generation

Multiple Tables for Improved Performance

A small ensemble of table-based Classifiers achieve better accuracy and performance

Neural Network Based Mechanism

Compiler Support for MITHRA

- Programmer provides:
 - Application specific quality requirement
 - Quality Metric
 - Set of representative application inputs
- Algorithm maps final output requirement to a threshold (th) on accelerator error.

Accelerator error Accelerator error > th : learn to fall back to the original function

Quality vs Benefits Tradeoffs

Output Quality

