

**ANALYZING OPTIMIZATION
TECHNIQUES FOR POWER
EFFICIENCY ON HETEROGENEOUS
PLATFORMS**

Yash Ukidave and David Kaeli
Department of Electrical and Computer Engineering
Northeastern University, Boston, USA

AsHES 2013
Boston, MA
20th May, 2013

Northeastern University 1 | AsHES 2013 | May 2013

WHAT IS THIS TALK ABOUT ?

- Study of power and energy profile of different optimization techniques used in heterogeneous applications
- Evaluation of power/performance of such optimization techniques on heterogeneous applications such as FFT

Northeastern University 2 | AsHES 2013 | May 2013

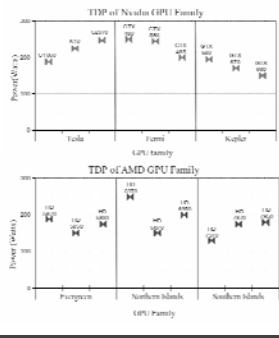
TOPICS

- Applications using OpenCL and power consumption of heterogeneous devices
- Fast Fourier Transforms (FFT) & evaluation methodology
- Optimization techniques used for analysis
- Results for power-performance of FFT implementations
- Analysis of power-performance of different optimization techniques
- Energy profile of different optimization techniques
- Conclusion
- Future work

Northeastern University 3 | AsHES 2013 | May 2013

MOTIVATION

- Increasing use of the CPU-GPU environment to accelerate data-parallel heterogeneous applications
- **Thermal Design Power (TDP)** of latest generation of GPUs used for heterogeneous compute
- Understanding the effects of software design methods contributing to power consumption
- Power and Energy aspect of different optimization techniques for heterogeneous platforms



Northeastern University 4 | AsHES 2013 | May 2013

FAST FOURIER TRANSFORM (FFT)

- FFT is an algorithm to compute Discrete Fourier Transforms (DFT)
 - Reduces time complexity to $O(n \log n)$ from $O(n^2)$
- FFT classified as Decimation in Time (DIT) or Decimation in Frequency (DIF)
 - DIT** : Operates on odd and even components of signal
 - DIF** : Operates on two halves of the signal
- Butterfly** structure is a major component of FFT of a given Radix
 - Each FFT work item performs computes butterfly on given data points

Northwestern University | ASRES 2013 | May 2013

FFT IMPLEMENTATIONS

- MR-SC FFT** : Multi-Radix single kernel call, based on the FFT implementation in AMD SDK.
- MR-MC FFT** : Based on Cooley-Tukey algorithm, uses Multiple kernel calls and Multiple Radix combinations for compute
- Stockham FFT** : Based on the Stockham algorithm for FFT. Single Radix and Single Kernel call computation
- Apple FFT** : A Multiple Kernel call based FFT provided by Apple Inc. using OpenCL

FFT Implementation	Kernel Calls	Memory Access Patterns	Twiddle factor Computation
MR-SC	Single	Global & Local Memory	Kernel Compile-Time
Stockham FFT	Multiple	Global Memory	Kernel Run-Time
Apple FFT			
MR-MC FFT			

Northwestern University | ASRES 2013 | May 2013

PLATFORMS FOR EVALUATION

- Shared Memory APUs

Device Features	Intel Core i7 3300	AMD Fusion A8 APU
Device Generation	Ivy Bridge	Evergreen
Compute Units(CU)	16	5
Processing Elements(PE)/CU	4	16
TDP(Watts)	70	100
Memory Bandwidth(GB/s)	26	17
Register File Per CU(KB)	64	64
- Discrete GPUs

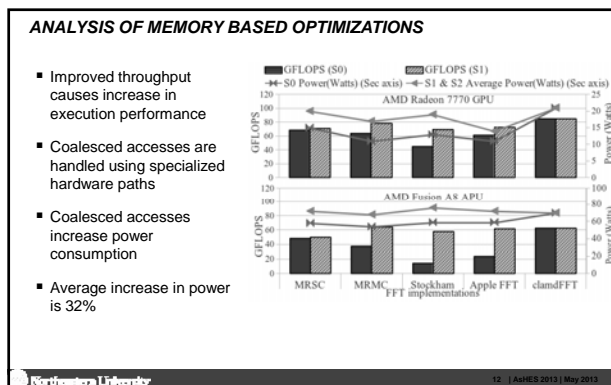
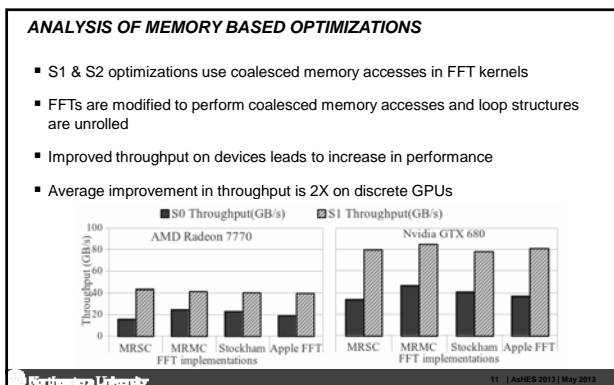
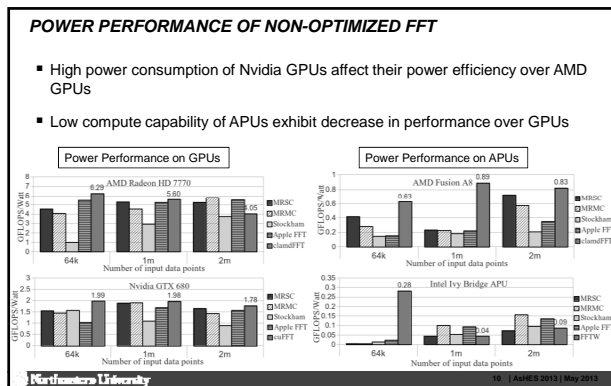
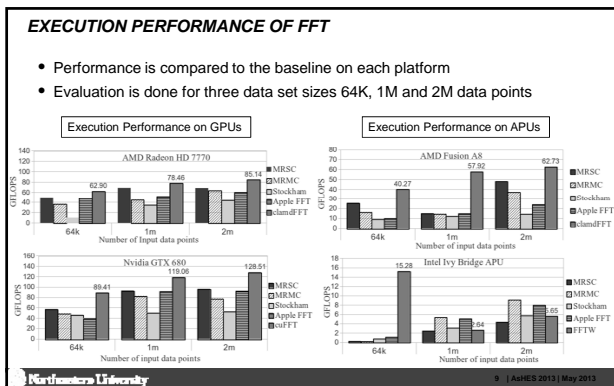
Device Features	Nvidia GTX 680	AMD Radeon HD 7770
Device Generation	Kepler	Southern Islands
Compute Units(CU)	8	10
Processing Elements(PE)/CU	192	64
TDP(Watts)	195	80
Memory Bandwidth(GB/s)	192	72
Register File Per CU(KB)	256	256

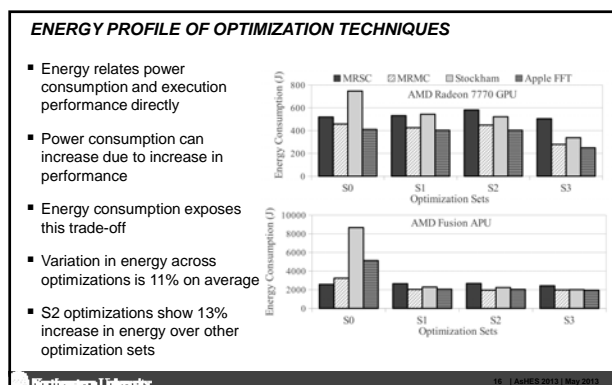
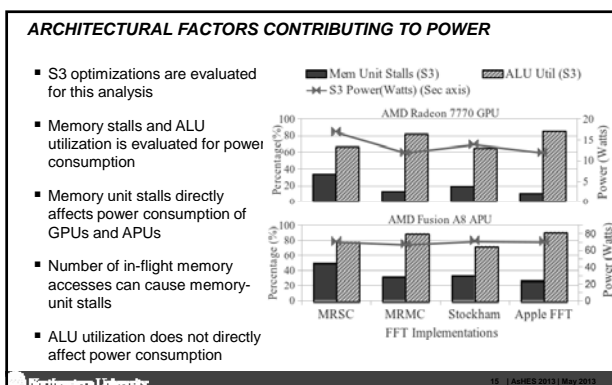
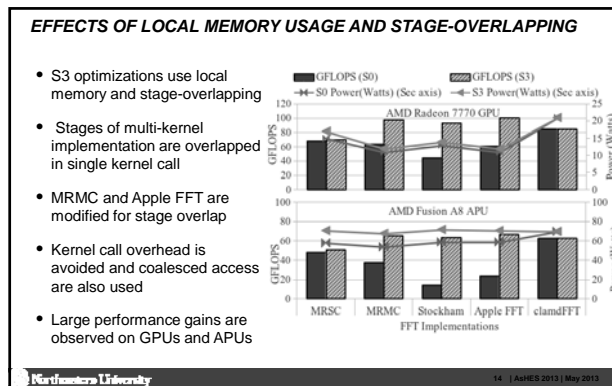
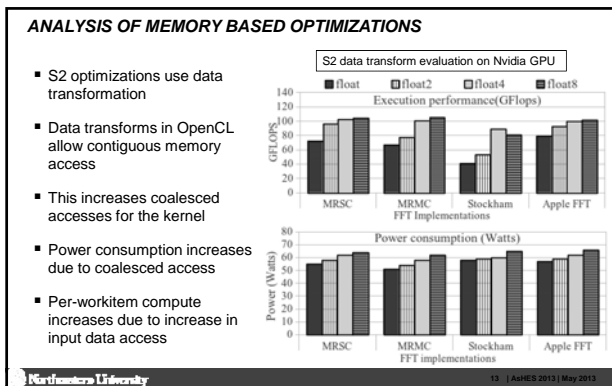
Northwestern University | ASRES 2013 | May 2013

OPTIMIZATION TECHNIQUES CLASSIFIED IN SETS

- Set S0** : No modifications ("out-of-box") performance
 - Set S1** : Coalesced Global Memory Accesses, Loop Unrolling
 - Set S2** : Data Transformation (float → float2, float → float4, float → float8)
 - Set S3** : Local Memory Usage, Stage overlapping for stage-based compute

Northwestern University | ASRES 2013 | May 2013





RESULTS SUMMARY

GPUs

- S1 & S2 improve performance efficiency with cost for power consumption
- Local memory in S3 improve power-performance
- S3 are compute and power efficient

APUs

- S1 & S2 are not power efficient
- Local memory improves compute efficiency and power efficiency
- Stage overlapping increases load on resources and increases power consumption

Northwestern University | 17 | ASRES 2013 | May 2013

RESULTS SUMMARY

✓ Highly Efficient	✗ Less Efficient	Ⓟ Moderately Efficient				
More than 40% improvement	Less than 10% improvement	10-40% improvement				
Optimization Techniques	Power Efficiency		Performance Efficiency		Power Performance (Gflops/Watts)	
	GPUs	APUs	GPUs	APUs	GPUs	APUs
Coalesced Memory Access	✗	✗	✓	✓	Ⓟ	Ⓟ
Loop Unrolling	✗	✗	✓	✓	Ⓟ	Ⓟ
Data Transformation	✗	✗	✓	✓	✗	✗
Local Memory Usage	✓	✓	✓	✓	✓	✓
Stage Overlapping	✓	Ⓟ	✓	✓	✓	Ⓟ

Northwestern University | 18 | ASRES 2013 | May 2013

CONCLUSION

- Analyzed different optimization techniques for their power-performance on once class of applications
- Study helps developer identify potential of power-aware kernel development
- Optimizations related to coalesced memory accesses exhibit increase in power consumption on GPUs and APUs
- Local Memory utilization is observed as the most power efficient optimization technique
- Power increment due to performance improvement is captured effectively in the energy profile

Northwestern University | 19 | ASRES 2013 | May 2013

FUTURE WORK

- Analyze power consumption on SoC (System-on-Chip) devices with GPUs, such as **TI OMAP4, Samsung Exynos, Qualcomm Snapdragon**
- Power-performance analysis to multi-GPU environments such as clusters
- Study of microarchitectural features responsible for power consumption on GPUs
- Extend study to different heterogeneous multi-core devices such as **Adapteva EpiPhany Processor, Tiler TilePro processors**

Northwestern University | 20 | ASRES 2013 | May 2013



THANK YOU !
QUESTIONS ? COMMENTS ?

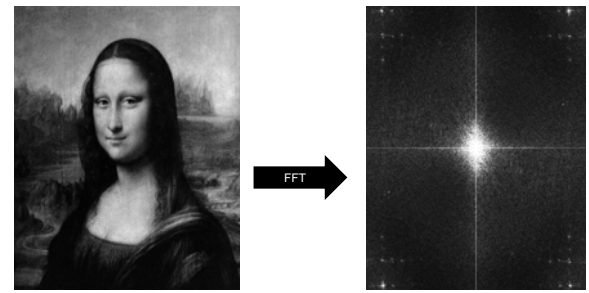
Yash Ukidave
 yukidave@ece.neu.edu

Northwestern University | 21 | ASRES 2013 | May 2013

EXTRA SLIDES

Northwestern University | 22 | ASRES 2013 | May 2013

TITLE PAGE PICTURE

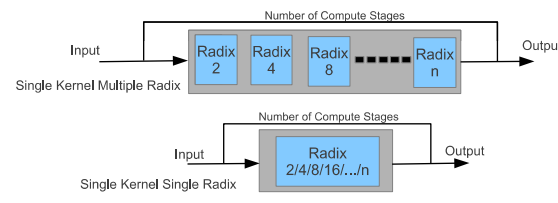


Source: http://commons.wikimedia.org/wiki/File:Mona_Lisa_bw_square.jpg

Northwestern University | 23 | ASRES 2013 | May 2013

FFT IMPLEMENTATIONS WITH SINGLE KERNEL CALL

- **MR-SC FFT** : FFT implementation based on AMD SDK. Multiple Radix sizes used
- **Stockham FFT** : Stockham algorithm for FFT computation. Single Radix computation



Number of Compute Stages

Input → [Radix 2 | Radix 4 | Radix 8 | ... | Radix n] → Output
 Single Kernel Multiple Radix

Number of Compute Stages

Input → [Radix 2/4/8/16/.../n] → Output
 Single Kernel Single Radix

Northwestern University | 24 | ASRES 2013 | May 2013

FFT IMPLEMENTATIONS WITH MULTIPLE KERNEL CALLS

- **MR-MC FFT** : Multiple kernel calls and Multiple Radix combinations for compute
- **Apple FFT** : Multiple Kernel call based FFT by Apple Inc. using OpenCL

The diagram illustrates two FFT implementation strategies. The top strategy, 'Multiple Kernel Single Radix', shows a sequence of three 'Radix 2/4/8/.../n' blocks connected by arrows, with 'Input' on the left and 'Output' on the right. Below this, 'Compute Stages' are indicated. The bottom strategy, 'Multiple Kernel Multiple Radix', shows three separate blocks: 'Radix-2 kernel', 'Radix-4 kernel', and 'Radix-n kernel', each with its own 'Input' and 'Output' and 'Compute Stages'.

Northwestern University | 25 | APRIL 2013 | May 2013

POWER MEASUREMENT SETUP

Power measurement on GPU:

- Discrete GPUs were isolated from the system for power measurement
- A dedicated external Power supply is used for Discrete GPUs
- Power consumption of GPU is measured by profiling the External PSU

The diagram shows a power measurement setup for a discrete GPU. It starts with 'Wall Power' entering a 'Power Meter'. The power then goes to a 'System Power supply'. A 'Main Bus Power Controller' is connected to the system power supply and the 'Main Board'. An 'External Power supply for GPU' is connected to the 'Main Board' and the 'GPU'. 'Memory' and 'CPU' are also connected to the 'Main Board'. 'Power Recording' is shown as an output from the power meter.

Power Measurement for APU:

- Power supply for APUs is measured using a similar power meter
- System power is measured to record power consumption on APUs
- Graphic processor of APUs cannot be completely isolated off the host CPU device

The diagram shows a power measurement setup for an APU. It starts with 'Wall Power' entering a 'Power Meter'. The power then goes to a 'System Power supply'. The 'System Power supply' is connected to the 'Main Board'. 'Memory' and 'CPU' are also connected to the 'Main Board'. 'Power Recording' is shown as an output from the power meter.

Northwestern University | 26 | APRIL 2013 | May 2013