Integrated GPUs for Green High-Performance Computing

<u>Tom Scogland</u> Heshan Lin Wu-chun Feng





GPUs





Powerful





Powerful: Graphics







Powerful: Graphics







Powerful: Graphics



6



















OpenCL

9





Speedups Using GPU vs CPU





18X

Transcoding HD video stream to H.264 for portable video³

20X



Simulation in Matlab using .mex file CUDA function⁴

24X



Astrophysics Nbody simulation⁵









Rank	Site	Computer
1	Oak Ridge National Laboratory United States	Jaguar - Cray XT5-HE Opteron Six Core 2.6 GHz Cray Inc.
2	National Supercomputing Centre in Shenzhen (NSCS) China	Nebulae - Dawning TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU Dawning
3	DOE/NNSA/LANL United States	Roadrunner - BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, Voltaire Infiniband IBM
4	National Institute for Computational Sciences/University of Tennessee United States	Kraken XT5 - Cray XT5-HE Opteron Six Core 2.6 GHz Cray Inc.







Rank	Site	Computer
1	Oak Ridge National Laboratory United States	Jaguar - Cray XT5-HE Opteron Six Core 2.6 GHz Cray Inc.
2	National Supercomputing Centre in Shenzhen (NSCS) China	Nebulae - Dawning TC3600 Blade, Intel X5650, NVidia Tesla C2050 GPU Dawning
3	DOE/NNSA/LANL United States	Roadrunner - BladeCenter QS22/LS21 Cluster, PowerXCell 8i 3.2 Ghz / Opteron DC 1.8 GHz, Voltaire Infiniband IBM
4	National Institute for Computational Sciences/University of Tennessee United States	Kraken XT5 - Cray XT5-HE Opteron Six Core 2.6 GHz Cray Inc.











Green?





What is Green?











Energy Efficiency







Energy Efficiency

Power Efficiency







Green: Energy Efficiency

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
1	773.38	Forschungszentrum Juelich (FZJ)	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D- Torus	57.54
1	773.38	Universitaet Regensburg	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D- Torus	57.54
1	773.38	Universitaet Wuppertal	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D- Torus	57.54
4	492.64	National Supercomputing Centre in Shenzhen (NSCS)	Dawning Nebulae, TC3600 blade CB60-G2 cluster, Intel Xeon 5650/ nVidia C2050, Infiniband	2580





Green: Energy Efficiency

Green500 Rank	MFLOPS/W	Site*	Computer*	Total Power (kW)
1	773.38	Forschungszentrum Juelich (FZJ)	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D- Torus	57.54
1	773.38	Universitaet Regensburg	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D- Torus	57.54
1	773.38	Universitaet Wuppertal	QPACE SFB TR Cluster, PowerXCell 8i, 3.2 GHz, 3D- Torus	57.54
4	492.64	National Supercomputing Centre in Shenzhen (NSCS)	Dawning Nebulae, TC3600 blade CB60-G2 cluster, Intel Xeon 5650/ nVidia C2050, Infiniband	2580





Related Work

16



Wirginia Tech Invent the Future

Friday, September 17, 2010

Related Work

 Huang S, Xiao S, Feng W (2009) "On the Energy Efficiency of Graphics Processing Units for Scientific Computing." In Proceedings of the 5th IEEE Workshop on High-Performance, Power-Aware Computing (in conjunction with the 23rd International Parallel & Distributed Processing Symposium), Rome, Italy, May 2009.





Related Work

- Huang S, Xiao S, Feng W (2009) "On the Energy Efficiency of Graphics Processing Units for Scientific Computing." In Proceedings of the 5th IEEE Workshop on High-Performance, Power-Aware Computing (in conjunction with the 23rd International Parallel & Distributed Processing Symposium), Rome, Italy, May 2009.
- Rofouei M, Stathopoulos T, Ryffel S, Kaiser W, Sarrafzadeh M (2008) "Energy-Aware High Performance Computing with Graphic Processing Units." In: Workshop on Power Aware Computing and System





Power of processing unit only

17







Friday, September 17, 2010





Green: Power Efficiency

- Systems which are
 - Embedded
 - Portable
 - Power constrained
 - Cooling constrained







VirginiaTech

Invent the Future

What We Need

19





What We Need









Integrated GPUs





Integrated GPUs: What?





Discrete GPU



Discrete GPU







Discrete GPU










Common Perception



VirginiaTech Invent the Future

Discrete GPU = Integrated GPU =





Discrete GPU = Fast Integrated GPU =





Unvent the Future

Discrete GPU = Fast Integrated GPU = Slow







Discrete GPU = Fast Integrated GPU = Slow?







Tested GPUs

29







Tested GPUs

	NVIDIA GTX 280	NVIDIA ION (GeForce 9400)	
Multiprocessors	30	2	
Cores	240	16	
Clock Rate	I.3 GHz	I.I GHz	
Memory	IGB (device)	256MB (system)	



WirginiaTech

Invent the Future

1000

Tested Systems

			GTX 280		ION	
	CPU Cores Clock Rate		Xeon e5405 x2	Atom 230		
			8	2 I.6 GHz		
			2 GHz			Z
	Memory		4GB		3GB	



Invent the Future

WirginiaTech

Software

- OS: Linux
- Distribution: Ubuntu Hardy
- CUDA SDK 3.0
- GCC 4.2: all optimizations enabled





Tested Applications

- Tested 4 dwarfs
 - N-Body
 - Structured grid
 - Dynamic programming
 - Dense linear algebra





N-Body: GEM

- Application
 - Molecular dynamics
 - Computes electrostatic potential along the surface of macromolecules
- Signature
 - Compute intensive







GEM: Signature All-Pairs Vertices (φi) <-> Atoms (qi)

 $\mathbf{r} = \mathbf{A} + \mathbf{p}$

q_i

34

 ϵ_{in}

Eout

p.

di



Structured Grid: SRAD

- Application
 - Image cleanup
 - Used to smooth sonic and radar imagery without losing important detail
- Signature
 - Kernel launch intensive





SRAD: Signature



Dynamic Programming: Needleman Wunsch (NW)

- Application
 - Sequence alignment
 - NW performs a global alignment between two sequences, usually DNA
- Signature
 - Synchronization intensive





NW: Signature







Dense Linear Algebra: K-Means

- Application
 - Clustering
 - Uses iterative refinement to cluster items in a space
- Signature
 - PCI-E Memory transfer intensive







40



Uirginia Tech

Invent the Future

Results



Performance





42

PerformanceXeon SMPGTX 280











Friday, September 17, 2010







Friday, September 17, 2010







Friday, September 17, 2010



Power

44





Power: GPU

Single Xeon Double Xeon GTX 280 ION






Power: GPU



Power: GPU



Power: GPU



Power: System













Power: System



Power: System









Unvent the Future

Metric: Energy Delay Product (EDP)

Common circuit design metric





Metric: Energy Delay Product (EDP)
Common circuit design metric
EDP = (joules) * (seconds)





Metric: Energy Delay Product (EDP)
Common circuit design metric
EDP = (joules) * (seconds)
EDP = (Watts * seconds) * (seconds)





- Metric: Energy Delay Product (EDP)
 - Common circuit design metric
- EDP = (joules) * (seconds)
- EDP = (Watts * seconds) * (seconds)
- Performance-centric energy efficiency





Xeon SMP GTX 280 ION























Xeon SMP GTX 280 ION









Friday, September 17, 2010





Friday, September 17, 2010









Friday, September 17, 2010



		Xeon SMP	GTX 280	ION
Pe	rformance			
	Power			
	EDP			
	Energy			

Best Average Worst

51



WirginiaTech

Invent the Future

Conclusions



WirginiaTech
• GPUs are "Green" for some applications





- GPUs are "Green" for some applications
- Discrete GPUs are highly efficient according to EDP





- GPUs are "Green" for some applications
- Discrete GPUs are highly efficient according to EDP
- Integrated GPUs are highly efficient according to total energy consumption



- GPUs are "Green" for some applications
- Discrete GPUs are highly efficient according to EDP
- Integrated GPUs are highly efficient according to total energy consumption
- Integrated and low-power GPUs are viable options for power constrained systems









• Evaluation of further platforms:

- AMD GPUs
- Embedded GPUs (e.g. tegra)
- Traditional embedded accelerators





• Evaluation of further platforms:

- AMD GPUs
- Embedded GPUs (e.g. tegra)
- Traditional embedded accelerators
- Expand into OpenCL testing





• Evaluation of further platforms:

- AMD GPUs
- Embedded GPUs (e.g. tegra)
- Traditional embedded accelerators
- Expand into OpenCL testing
- Investigate clustering low-power GPUs





Questions?

Contact Email: <u>tom.scogland@vt.edu</u>





Questions?



54





Friday, September 17, 2010