Programming for Hybrid Architectures

John E. Stone

Theoretical and Computational Biophysics Group Beckman Institute for Advanced Science and Technology University of Illinois at Urbana-Champaign http://www.ks.uiuc.edu/Research/gpu/ GPGPU 2015: Advanced Methods for Computing with CUDA, University of Cape Town, April 2015



Solid Growth of GPU Accelerated Apps Top HPC Applications



Major Approaches For Programming Hybrid Architectures

- Use drop-in libraries in place of CPU-based libraries
 - Little or no code development
 - Speedups limited by Amdahl's Law and overheads associated with data movement between CPUs and GPU accelerators
 - Examples: MAGMA, BLAS-variants, FFT libraries, etc.
- Generate accelerator code as a variant of CPU source, e.g. using OpenMP w/ OpenACC, similar methods
- Write lower-level accelerator-specific code, e.g. using CUDA, OpenCL, other approaches



GPU Accelerated Libraries "Drop-in" Acceleration for your Applications



Courtesy NVIDIA

OpenACC: Open, Simple, Portable



Using the CPU to Optimize GPU Performance

- GPU performs best when the work evenly divides into the number of threads/processing units
- Optimization strategy:
 - Use the CPU to *"regularize"* the GPU workload
 - Use fixed size bin data structures, with "empty" slots skipped or producing zeroed out results
 - Handle exceptional or irregular work units on the CPU;
 GPU processes the bulk of the work concurrently
 - On average, the GPU is kept highly occupied, attaining a high fraction of peak performance



CUDA Grid/Block/Thread Decomposition



Avoiding Shared Memory Bank Conflicts: Array of Structures (AOS) vs. Structure of Arrays (SOA) • AOS: • SOA typedef struct { typedef struct { float x; float x[1024]; float y; float y[1024]; float z[1024]; float z; myvec; } myvecs; myvec aos[1024]; myvecs soa; aos[threadIdx.x].x = 0;soa.x[threadIdx.x] = 0;aos[threadIdx.x].y = 0;soa.y[threadIdx.x] = 0;



NIH BTRC for Macromolecular Modeling and Bioinformatics http://www.ks.uiuc.edu/ Beckman Institute, U. Illinois at Urbana-Champaign

Time-Averaged Electrostatics Analysis on Energy-Efficient GPU Cluster

- 1.5 hour job (CPUs) reduced to 3 min (CPUs+GPU)
- Electrostatics of thousands of trajectory frames averaged
- Per-node power consumption on NCSA "AC" GPU cluster:
 - CPUs-only: 448 Watt-hours
 - CPUs+GPUs: 43 Watt-hours
- GPU Speedup: 25.5x
- Power efficiency gain: **10.5**x

Quantifying the Impact of GPUs on Performance and Energy Efficiency in HPC Clusters. J. Enos, C. Steffen, J. Fullop, M. Showerman, G. Shi, K. Esler, V. Kindratenko, J. Stone, J. Phillips. *The Work in Progress in Green Computing*, pp. 317-324, 2010.



AC Cluster GPU Performance and Power Efficiency Results

Application	GPU speedup	Host watts	Host+GPU watts	Perf/watt gain
NAMD	6	316	681	2.8
VMD	25	299	742	10.5
MILC	20	225	555	8.1
QMCPACK	61	314	853	22.6

Quantifying the Impact of GPUs on Performance and Energy Efficiency in HPC Clusters. J. Enos, C. Steffen, J. Fullop, M. Showerman, G. Shi, K. Esler, V. Kindratenko, J. Stone, J. Phillips. *The Work in Progress in Green*

Computing, pp. 317-324, 2010.

NIH BTRC for Macromolecular Modeling and Bioinformatics http://www.ks.uiuc.edu/



Optimizing GPU Algorithms for Power Consumption

NVIDIA "Carma", "Kayla", "Jetson" single board computers

Tegra+GPU energy efficiency testbed





Time-Averaged Electrostatics Analysis on NCSA Blue Waters

NCSA Blue Waters Node Type	Seconds per trajectory frame for one compute node
Cray XE6 Compute Node: 32 CPU cores (2xAMD 6200 CPUs)	9.33
Cray XK6 GPU-accelerated Compute Node: 16 CPU cores + NVIDIA X2090 (Fermi) GPU	2.25
Speedup for GPU XK6 nodes vs. CPU XE6 nodes	XK6 nodes are 4.15x faster overall
Tests on XK7 nodes indicate MSM is CPU-bound with the Kepler K20X GPU. Performance is not much faster (yet) than Fermi X2090 Need to move spatial hashing, prolongation, interpolation onto the GPU	In progress XK7 nodes 4.3x faster overall

Preliminary performance for VMD time-averaged electrostatics w/ Multilevel Summation Method on the NCSA Blue Waters Early Science System





Multilevel summation of electrostatic potentials using graphics processing units.

D. Hardy, J. Stone, K. Schulten. J. Parallel Computing, 35:164-177, 2009.



NIH BTRC for Macromolecular Modeling and Bioinformatics http://www.ks.uiuc.edu/

Multi-GPU NUMA Architectures:

- Example of a "balanced" PCIe topology
- NUMA: Host threads should be pinned to the CPU that is "closest" to their target GPU
- GPUs on the same PCIe I/O Hub (IOH) can use CUDA peer-to-peer transfer APIs
- Intel: GPUs on different IOHs can't use peer-to-peer



Multi-GPU NUMA Architectures:

- Example of a very "unbalanced" PCIe topology
- CPU 2 will overhelm its QP/HT link with host-GPU DMAs
- Poor scalability as compared to a balanced PCIe topology



Multi-GPU NUMA Architectures:

- GPU-to-GPU peer DMA operations are much more performant than other approaches, particularly for moderate sized transfers
- Likely to perform even better in future multi-GPU cards with direct GPU links, e.g. announced "NVLink"





Simulation of reaction diffusion processes over biologically relevant size and time scales using multi-GPU workstations Michael J. Hallock, John E. Stone, Elijah Roberts, Corey Fry, and Zaida Luthey-Schulten. Journal of Parallel Computing, 40:86-99, 2014. http://dx.doi.org/10.1016/j.parco.2014.03.009



NIH BTRC for Macromolecular Modeling and Bioinformatics http://www.ks.uiuc.edu/ Beckman Institute, U. Illinois at Urbana-Champaign

Single CUDA Execution "Stream"

- Host CPU thread launches a CUDA "kernel", a memory copy, etc. on the GPU
- GPU action runs to completion
- Host synchronizes with completed GPU action





Multiple CUDA Streams: Overlapping Compute and DMA Operations



Acknowledgements

- Theoretical and Computational Biophysics Group, University of Illinois at Urbana-Champaign
- NVIDIA CUDA Center of Excellence, University of Illinois at Urbana-Champaign
- NVIDIA CUDA team
- NCSA Blue Waters Team
- Funding:
 - NSF OCI 07-25070
 - NSF PRAC "The Computational Microscope"
 - NIH support: 9P41GM104601, 5R01GM098243-02





NIH BTRC for Macromolecular Modeling and Bioinformatics Beckman Institute University of Illinois at Urbana-Champaign



NIH BTRC for Macromolecular Modeling and Bioinformatics http://www.ks.uiuc.edu/ Beckman Institute, U. Illinois at Urbana-Champaign

- Runtime and Architecture Support for Efficient Data Exchange in Multi-Accelerator Applications Javier Cabezas, Isaac Gelado, John E. Stone, Nacho Navarro, David B. Kirk, and Wen-mei Hwu. IEEE Transactions on Parallel and Distributed Systems, 26(5):1405-1418, 2015.
- Unlocking the Full Potential of the Cray XK7 Accelerator Mark Klein and John E. Stone. Cray Users Group, Lugano Switzerland, May 2014.
- Simulation of reaction diffusion processes over biologically relevant size and time scales using multi-GPU workstations Michael J. Hallock, John E. Stone, Elijah Roberts, Corey Fry, and Zaida Luthey-Schulten. Journal of Parallel Computing, 40:86-99 2014.
- **GPU-Accelerated Analysis and Visualization of Large Structures Solved by Molecular Dynamics Flexible Fitting** John E. Stone, Ryan McGreevy, Barry Isralewitz, and Klaus Schulten. Faraday Discussions, 169:265-283, 2014.
- **GPU-Accelerated Molecular Visualization on Petascale Supercomputing Platforms.** J. Stone, K. L. Vandivort, and K. Schulten. UltraVis'13: Proceedings of the 8th International Workshop on Ultrascale Visualization, pp. 6:1-6:8, 2013.
- Early Experiences Scaling VMD Molecular Visualization and Analysis Jobs on Blue Waters. J. E. Stone, B. Isralewitz, and K. Schulten. Extreme Scaling Workshop (XSW), pp. 43-50, 2013.
- Lattice Microbes: High-performance stochastic simulation method for the reaction-diffusion master equation. E. Roberts, J. E. Stone, and Z. Luthey-Schulten. J. Computational Chemistry 34 (3), 245-255, 2013.



- Fast Visualization of Gaussian Density Surfaces for Molecular Dynamics and Particle System Trajectories. M. Krone, J. E. Stone, T. Ertl, and K. Schulten. *EuroVis Short Papers*, pp. 67-71, 2012.
- Fast Analysis of Molecular Dynamics Trajectories with Graphics Processing Units Radial Distribution Functions. B. Levine, J. Stone, and A. Kohlmeyer. J. Comp. Physics, 230(9):3556-3569, 2011.
- Immersive Out-of-Core Visualization of Large-Size and Long-Timescale Molecular Dynamics Trajectories. J. Stone, K. Vandivort, and K. Schulten. G. Bebis et al. (Eds.): *7th International Symposium on Visual Computing (ISVC 2011)*, LNCS 6939, pp. 1-12, 2011.
- Quantifying the Impact of GPUs on Performance and Energy Efficiency in HPC Clusters. J. Enos, C. Steffen, J. Fullop, M. Showerman, G. Shi, K. Esler, V. Kindratenko, J. Stone, J Phillips. *International Conference on Green Computing*, pp. 317-324, 2010.
- GPU-accelerated molecular modeling coming of age. J. Stone, D. Hardy, I. Ufimtsev, K. Schulten. J. Molecular Graphics and Modeling, 29:116-125, 2010.
- OpenCL: A Parallel Programming Standard for Heterogeneous Computing. J. Stone, D. Gohara, G. Shi. *Computing in Science and Engineering*, 12(3):66-73, 2010.



- An Asymmetric Distributed Shared Memory Model for Heterogeneous Computing Systems. I. Gelado, J. Stone, J. Cabezas, S. Patel, N. Navarro, W. Hwu. *ASPLOS '10: Proceedings of the 15th International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 347-358, 2010.
- **GPU Clusters for High Performance Computing**. V. Kindratenko, J. Enos, G. Shi, M. Showerman, G. Arnold, J. Stone, J. Phillips, W. Hwu. *Workshop on Parallel Programming on Accelerator Clusters (PPAC)*, In Proceedings IEEE Cluster 2009, pp. 1-8, Aug. 2009.
- Long time-scale simulations of in vivo diffusion using GPU hardware. E. Roberts, J. Stone, L. Sepulveda, W. Hwu, Z. Luthey-Schulten. In *IPDPS'09: Proceedings of the 2009 IEEE International Symposium on Parallel & Distributed Computing*, pp. 1-8, 2009.
- High Performance Computation and Interactive Display of Molecular Orbitals on GPUs and Multi-core CPUs. J. Stone, J. Saam, D. Hardy, K. Vandivort, W. Hwu, K. Schulten, 2nd Workshop on General-Purpose Computation on Graphics Pricessing Units (GPGPU-2), ACM International Conference Proceeding Series, volume 383, pp. 9-18, 2009.
- **Probing Biomolecular Machines with Graphics Processors**. J. Phillips, J. Stone. *Communications of the ACM*, 52(10):34-41, 2009.
- Multilevel summation of electrostatic potentials using graphics processing units. D. Hardy, J. Stone, K. Schulten. *J. Parallel Computing*, 35:164-177, 2009.



- Adapting a message-driven parallel application to GPU-accelerated clusters. J. Phillips, J. Stone, K. Schulten. *Proceedings of the 2008 ACM/IEEE Conference on Supercomputing*, IEEE Press, 2008.
- GPU acceleration of cutoff pair potentials for molecular modeling applications.
 C. Rodrigues, D. Hardy, J. Stone, K. Schulten, and W. Hwu. *Proceedings of the 2008 Conference On Computing Frontiers*, pp. 273-282, 2008.
- **GPU computing**. J. Owens, M. Houston, D. Luebke, S. Green, J. Stone, J. Phillips. *Proceedings* of the IEEE, 96:879-899, 2008.
- Accelerating molecular modeling applications with graphics processors. J. Stone, J. Phillips, P. Freddolino, D. Hardy, L. Trabuco, K. Schulten. *J. Comp. Chem.*, 28:2618-2640, 2007.
- Continuous fluorescence microphotolysis and correlation spectroscopy. A. Arkhipov, J. Hüve, M. Kahms, R. Peters, K. Schulten. *Biophysical Journal*, 93:4006-4017, 2007.

