

# Faster, Cheaper, Better: Biomolecular Simulation with NAMD, VMD, and CUDA

John Stone

Theoretical and Computational Biophysics Group  
Beckman Institute for Advanced Science and Technology  
University of Illinois at Urbana-Champaign

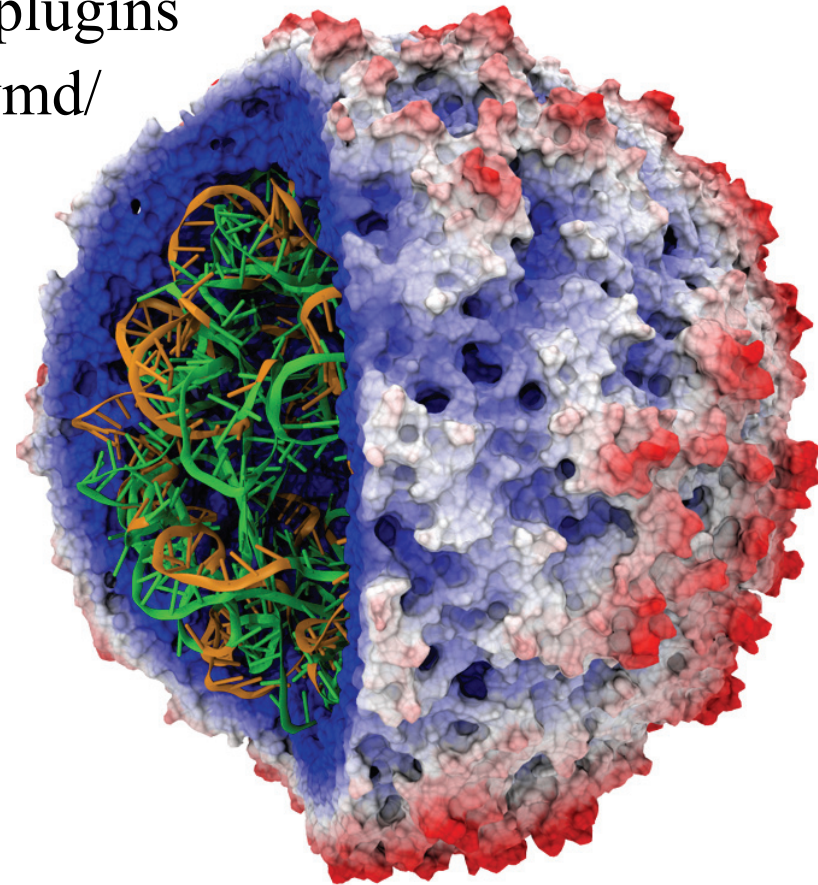
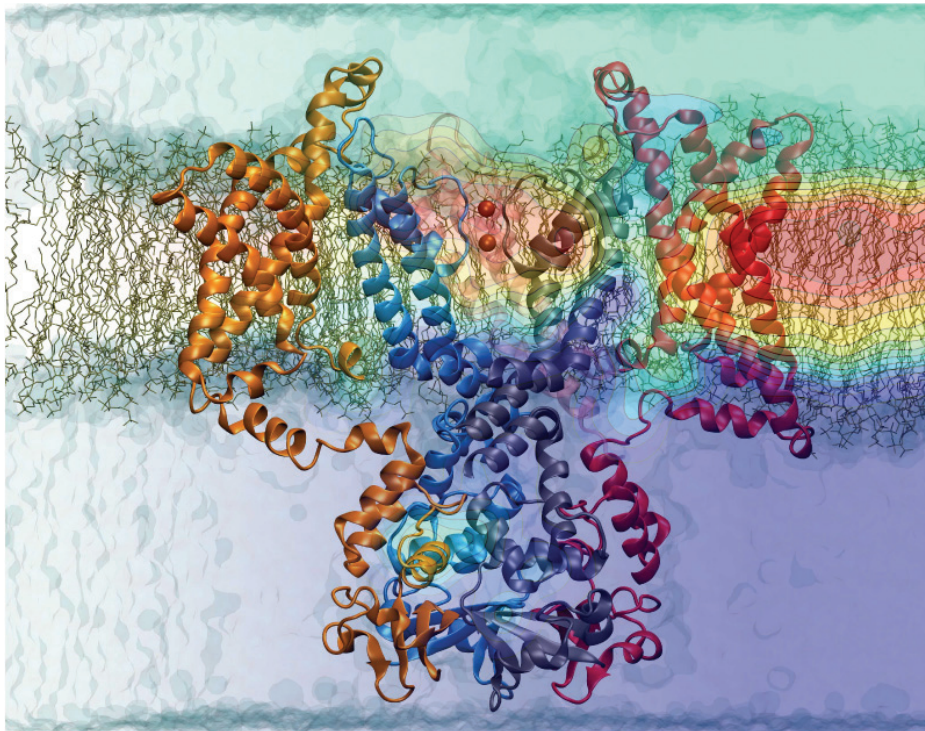
**<http://www.ks.uiuc.edu/Research/gpu/>**

Supercomputing 2010

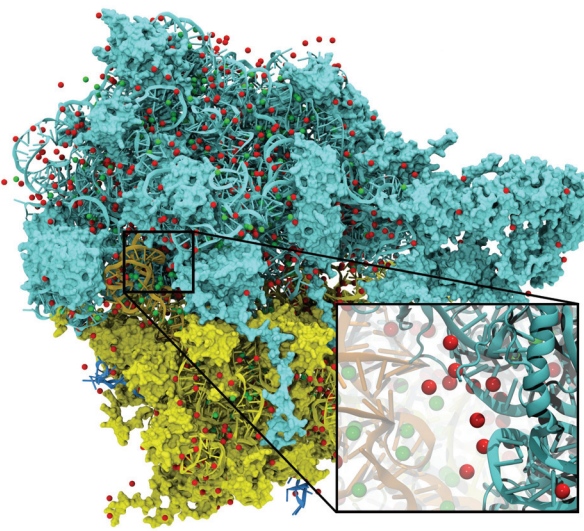
New Orleans, LA, Nov 16, 2010

# VMD – “Visual Molecular Dynamics”

- Visualization and analysis of molecular dynamics simulations, sequence data, volumetric data, quantum chemistry simulations, particle systems, ...
- User extensible with scripting and plugins
- <http://www.ks.uiuc.edu/Research/vmd/>

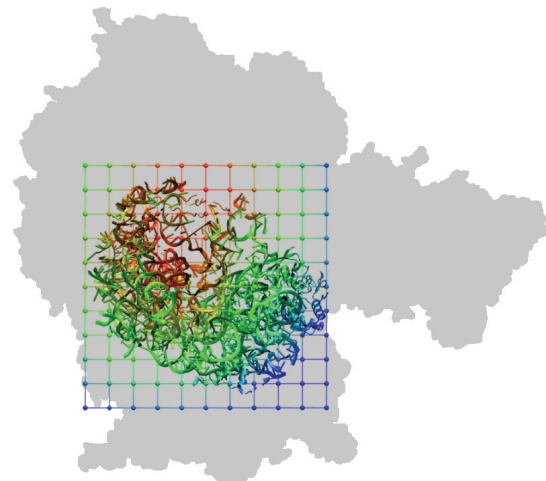


# CUDA Algorithms in VMD



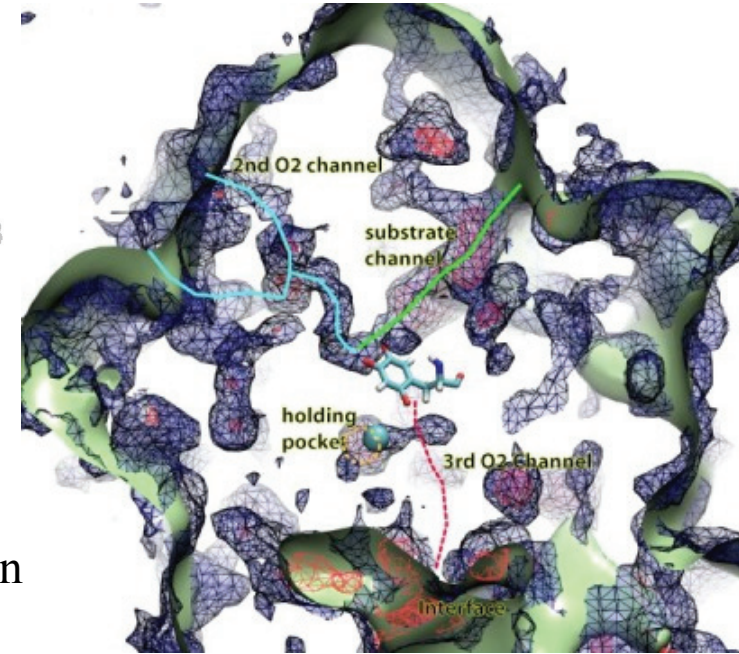
Ion placement

20x to 44x faster



Electrostatic field calculation

31x to 44x faster



Imaging of gas migration pathways in proteins with implicit ligand sampling

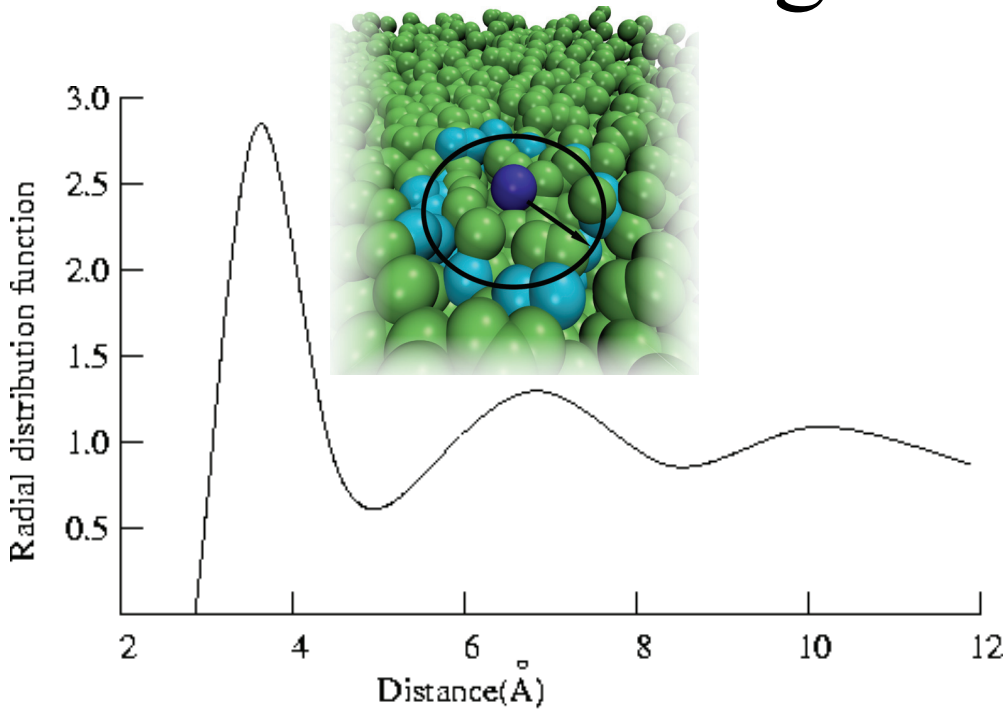
20x to 30x faster



GPU: massively parallel co-processor

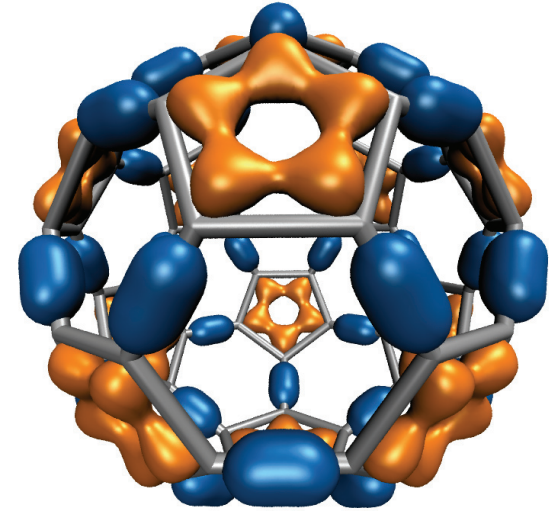


# CUDA Algorithms in VMD



Radial distribution functions

30x to 92x faster



Molecular orbital  
calculation and display

100x to 120x faster



GPU: massively parallel co-processor

# Quantifying GPU Performance and Energy Efficiency in HPC Clusters

- NCSA “AC” Cluster
- Power monitoring hardware on one node and its attached Tesla S1070 (4 GPUs)
- Power monitoring logs recorded separately for host node and attached GPUs
- Logs associated with batch job IDs



- 32 HP XW9400 nodes
- 128 cores, 128 Tesla C1060 GPUs
- QDR Infiniband

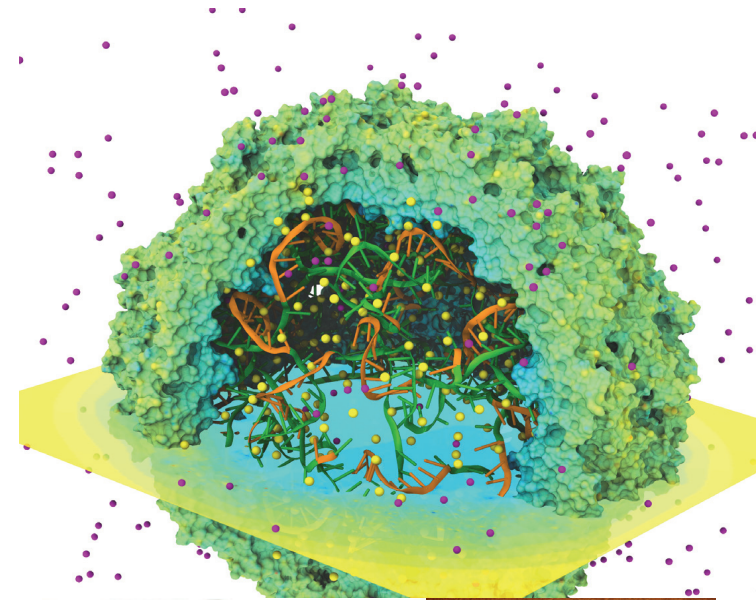
# NCSA GPU Cluster Power Measurements

State	Host Peak (Watt)	Tesla Peak (Watt)	Host power factor (pf)	Tesla power factor (pf)
power off	4	10	.19	.31
pre-GPU use idle	173	178	.98	.96
after NVIDIA driver module unload/reload	173	178	.98	.96
after deviceQuery (idle)	173	365	.99	.99
GPU memtest #10 (stress)	269	745	.99	.99
VMD Multiply-add	268	598	.99	.99
NAMD GPU STMV	321	521	.97-1.0	.85-1.0

**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.

# Energy Efficient GPU Computing of Time-Averaged Electrostatics

- **1.5 hour** job reduced to **3 min**
- Electrostatics of thousands of trajectory frames averaged
- Per-node power consumption on NCSA GPU cluster:
  - CPUs-only: 299 watts
  - CPUs+GPUs: 742 watts
- GPU Speedup: **25.5x**
- Power efficiency gain: **10.5x**



NCSA “AC” GPU cluster and Tweet-a-watt wireless power monitoring device

# AC Cluster GPU Performance and Power Efficiency Results

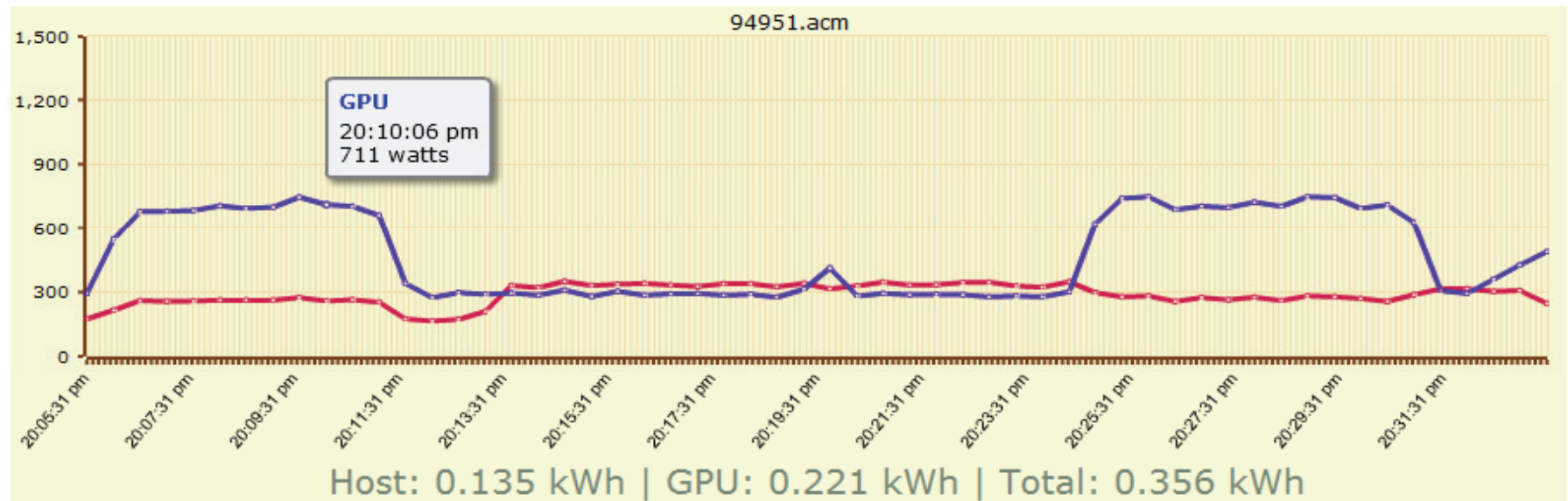
Application	GPU speedup	Host watts	Host+GPU watts	Perf/watt gain
NAMD	6	316	681	2.8
VMD	25	299	742	10.5
MILC	20	225	555	8.1
QMCPACK	61	314	853	22.6

**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*, 2010. In press.



# Power Profiling: Example Log

## AC Power Utilization

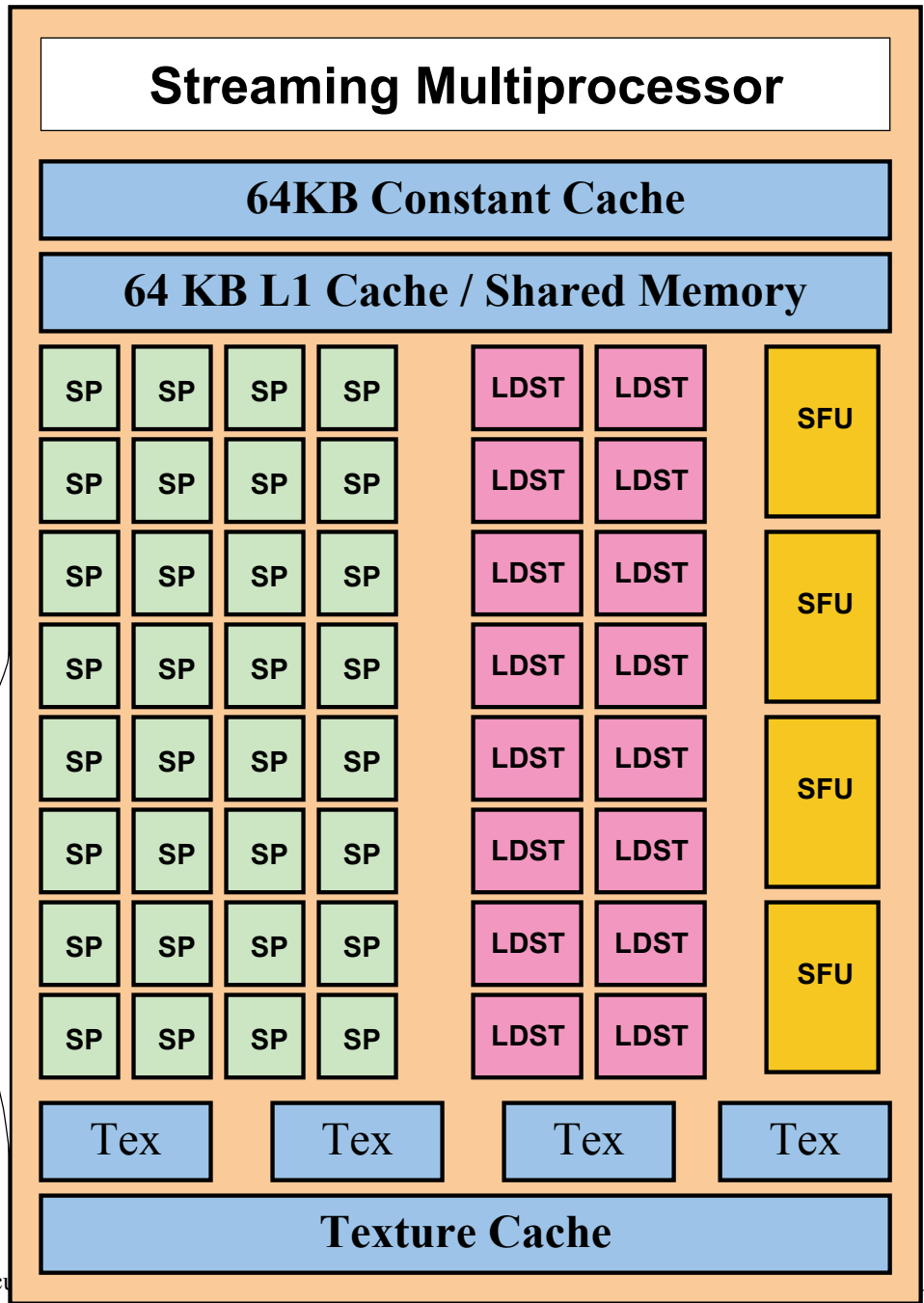
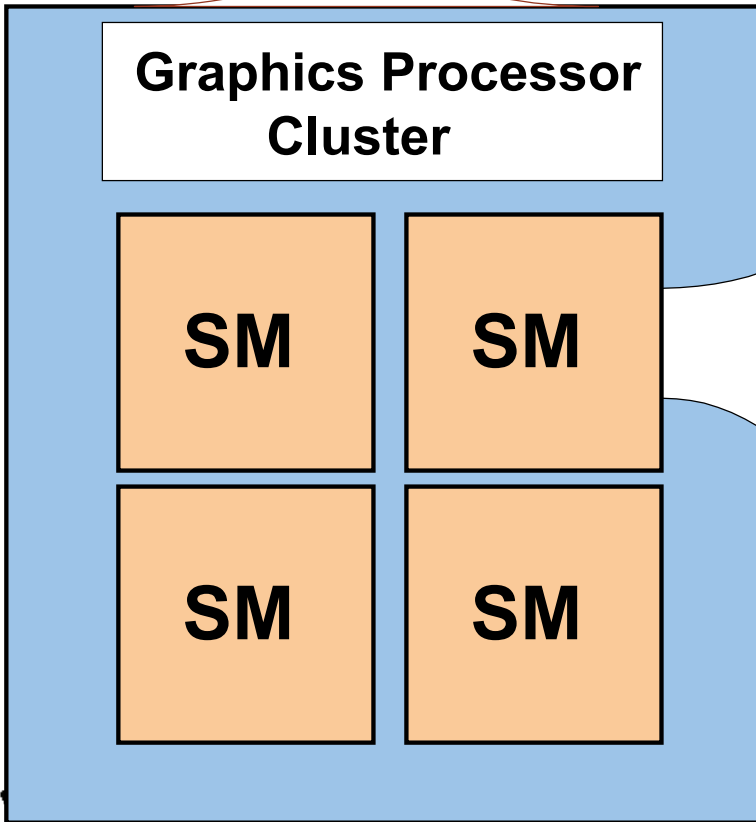
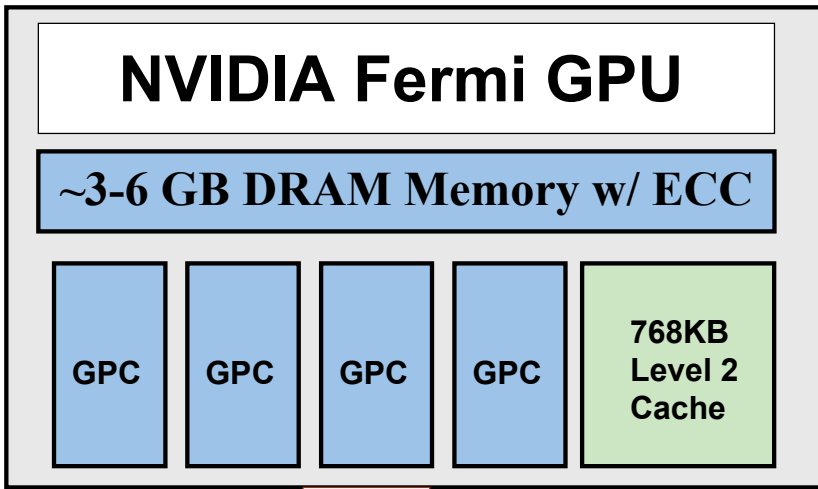


## JSON Data

- Mouse-over value displays
- Under curve totals displayed
- If there is user interest, we may support calls to add custom tags from application

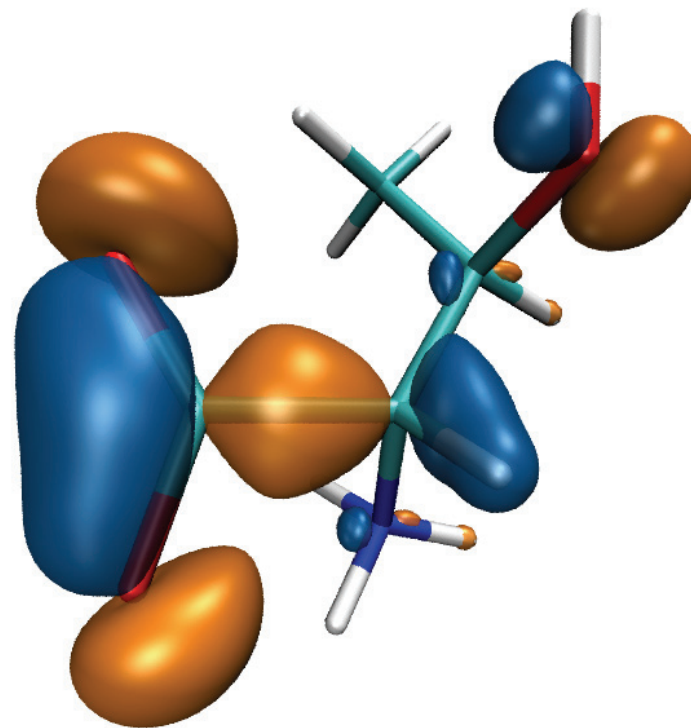
# Fermi GPUs Bring Higher Performance and Easier Programming

- NVIDIA's latest "Fermi" GPUs bring:
  - Greatly increased peak single- and double-precision arithmetic rates
  - Moderately increased global memory bandwidth
  - Increased capacity on-chip memory partitioned into shared memory and an L1 cache for global memory
  - Concurrent kernel execution
  - Bidirectional asynchronous host-device I/O
  - ECC memory, faster atomic ops, many others...



# Visualizing Molecular Orbitals

- Visualization of MOs aids in understanding the chemistry of molecular system
- Display of MOs can require **tens to hundreds of seconds** on multi-core CPUs, even with hand-coded SSE
- GPUs enable MOs to be computed and displayed in a **fraction of a second, fully interactively**





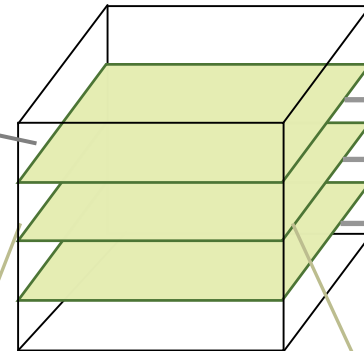
# MO GPU Parallel Decomposition

*MO 3-D lattice decomposes into 2-D slices (CUDA grids)*

*Small 8x8 thread blocks afford large per-thread register count, shared memory*

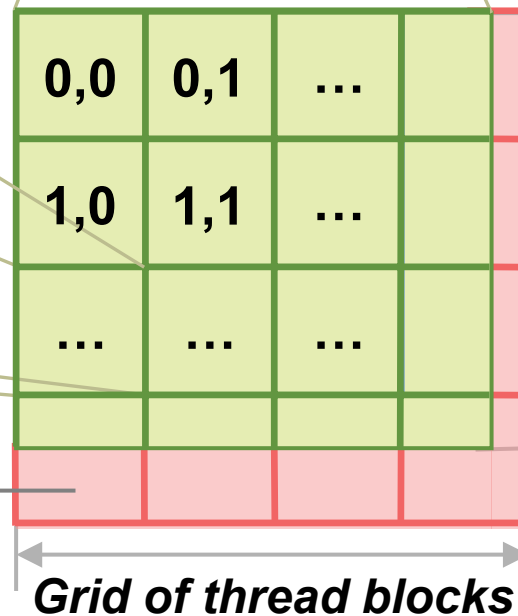
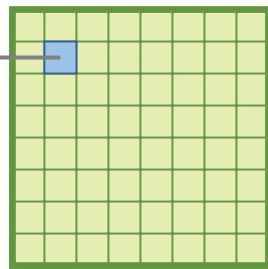
*Each thread computes one MO lattice point.*

*Padding optimizes global memory performance, guaranteeing coalesced global memory accesses*



...  
GPU 2  
GPU 1  
GPU 0

*Lattice can be computed using multiple GPUs*



*Threads producing results that are used*

*Threads producing results that are discarded*

# VMD MO GPU Kernel Snippet: Loading Tiles Into Shared Memory On-Demand

[... outer loop over atoms ...]

```
if ((prim_counter + (maxprim<<1)) >= SHAREDSIZE) {
    prim_counter += sblock_prim_counter;
    sblock_prim_counter = prim_counter & MEMCOAMASK;
    s_basis_array[sidx      ] = basis_array[sblock_prim_counter + sidx      ];
    s_basis_array[sidx + 64] = basis_array[sblock_prim_counter + sidx + 64];
    s_basis_array[sidx + 128] = basis_array[sblock_prim_counter + sidx + 128];
    s_basis_array[sidx + 192] = basis_array[sblock_prim_counter + sidx + 192];
    prim_counter -= sblock_prim_counter;
    __syncthreads();
}
```

```
for (prim=0; prim < maxprim; prim++) {
    float exponent      = s_basis_array[prim_counter      ];
    float contract_coeff = s_basis_array[prim_counter + 1];
    contracted_gto += contract_coeff * __expf(-exponent*dist2);
    prim_counter += 2;
}
```

[... continue on to angular momenta loop ...]

Shared memory tiles:

- Tiles are checked and loaded, if necessary, immediately prior to entering key arithmetic loops
- Adds additional control overhead to loops, even with optimized implementation

# VMD MO GPU Kernel Snippet:

## Fermi kernel based on L1 cache

[... outer loop over atoms ...]

```
// loop over the shells/basis funcs belonging to this atom
```

```
for (shell=0; shell < maxshell; shell++) {  
    float contracted_gto = 0.0f;  
    int maxprim = shellinfo[(shell_counter<<4)    ];  
    int shell_type = shellinfo[(shell_counter<<4) + 1];  
    for (prim=0; prim < maxprim; prim++) {  
        float exponent = basis_array[prim_counter    ];  
        float contract_coeff = basis_array[prim_counter + 1];  
        contracted_gto += contract_coeff * __expf(-  
            exponent*dist2);  
        prim_counter += 2;  
    }  
}
```

[... continue on to angular momenta loop ...]

L1 cache:

- Simplifies code!
- Reduces control overhead
- Gracefully handles arbitrary-sized problems
- Matches performance of constant memory

# VMD Single-GPU Molecular Orbital Performance Results for C<sub>60</sub> on Fermi

Intel X5550 CPU, GeForce GTX 480 GPU

Kernel	Cores/GPUs	Runtime (s)	Speedup
Xeon 5550 ICC-SSE	1	30.64	1.0
Xeon 5550 ICC-SSE	8	4.13	7.4
CUDA shared mem	1	0.37	83
<b>CUDA L1-cache (16KB)</b>	<b>1</b>	<b>0.27</b>	<b>113</b>
CUDA const-cache	1	0.26	117
CUDA const-cache, zero-copy	1	0.25	122

Fermi GPUs have caches: match perf. of hand-coded shared memory kernels. Zero-copy memory transfers improve overlap of computation and host-GPU I/Os.



# VMD Multi-GPU Molecular Orbital Performance Results for C<sub>60</sub>

Intel X5550 CPU, 4x GeForce GTX 480 GPUs,

Kernel	Cores/GPUs	Runtime (s)	Speedup
Intel X5550-SSE	1	30.64	1.0
Intel X5550-SSE	8	4.13	7.4
GeForce GTX 480	1	0.255	120
GeForce GTX 480	2	0.136	225
GeForce GTX 480	3	0.098	312
GeForce GTX 480	4	0.081	378

Uses persistent thread pool to avoid GPU init overhead,  
dynamic scheduler distributes work to GPUs

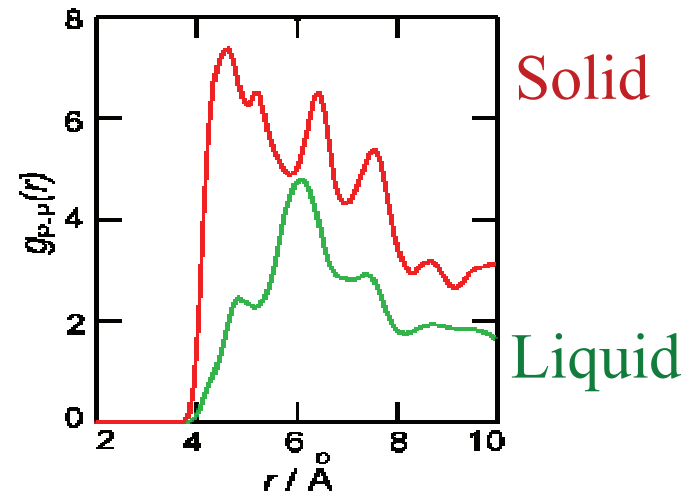
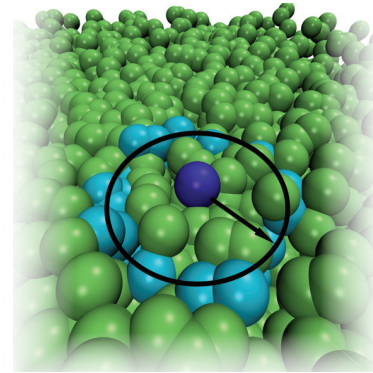
# Molecular Orbital Dynamic Scheduling Performance with Heterogeneous GPUs

Kernel	Cores/GPUs	Runtime (s)	Speedup
Intel X5550-SSE	1	30.64	1.0
Quadro 5800	1	0.384	79
Tesla C2050	1	0.325	94
GeForce GTX 480	1	0.255	120
GeForce GTX 480 + Tesla C2050 + Quadro 5800	3	0.114	268 (91% of ideal perf)

Dynamic load balancing enables mixture of GPU generations, SM counts, and clock rates to perform well.

# Radial Distribution Functions

- RDFs describes how atom density varies with distance
- Can be compared with experiments
- Shape indicates phase of matter: sharp peaks appear for solids, smoother for liquids
- Quadratic time complexity  $O(N^2)$



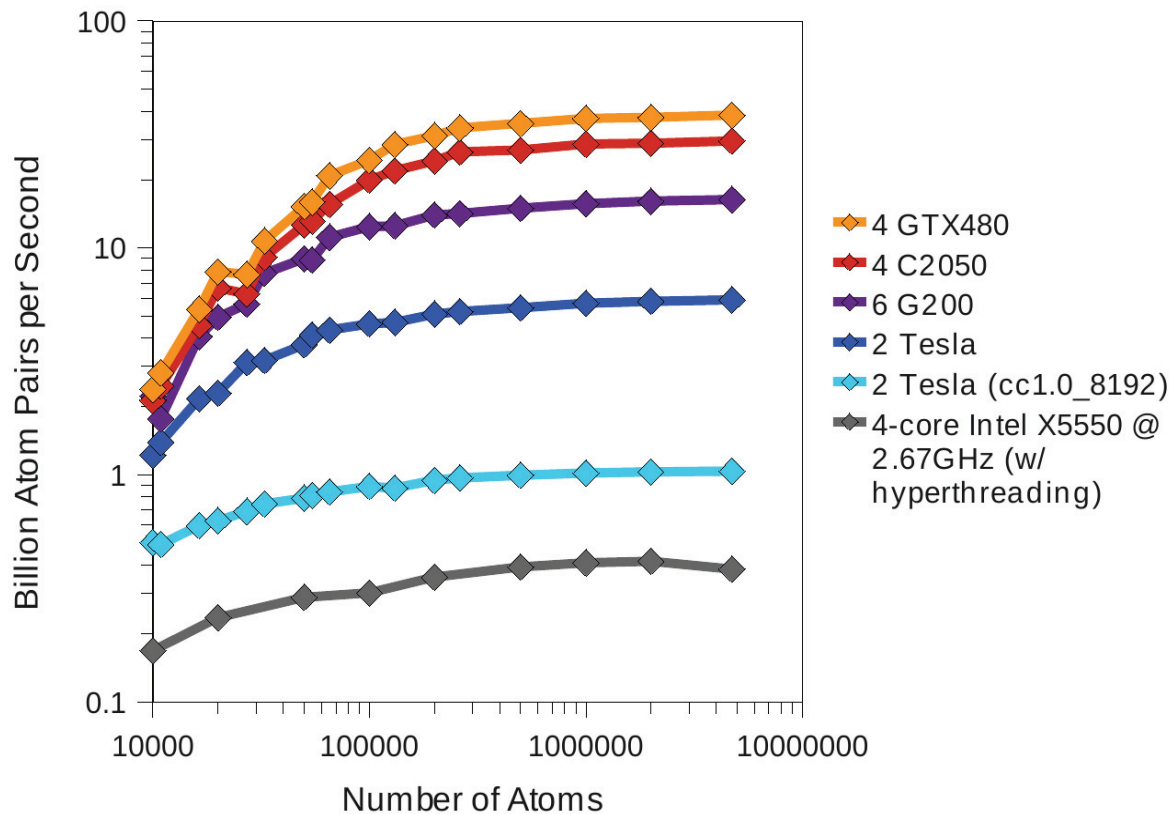
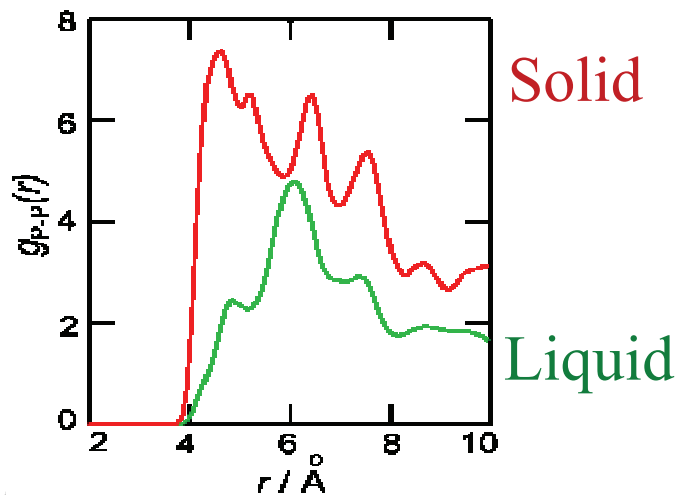
# Computing RDFs

- Compute distances for all pairs of atoms between two groups of atoms A and B
- A and B may be the same, or different
- Use nearest image convention for periodic systems
- Each pair distance is inserted into a histogram
- Histogram is normalized one of several ways depending on use, but usually according to the volume of the spherical shells associated with each histogram bin



# Multi-GPU RDF Performance

- 4 NVIDIA GTX480 GPUs 30 to 92x faster than 4-core Intel X5550 CPU
- Fermi GPUs ~3x faster than GT200 GPUs: larger on-chip shared memory



**Fast Analysis of Molecular Dynamics Trajectories with Graphics Processing Units – Radial Distribution Functions.** B. Levine, J. Stone, and A. Kohlmeier. 2010. (submitted)

# Acknowledgements

- Theoretical and Computational Biophysics Group, University of Illinois at Urbana-Champaign
- Ben Levine and Axel Kohlmeyer at Temple University
- NVIDIA CUDA Center of Excellence, University of Illinois at Urbana-Champaign
- NCSA Innovative Systems Lab
- The CUDA team at NVIDIA
- NIH support: P41-RR05969

# GPU Computing Publications

<http://www.ks.uiuc.edu/Research/gpu/>

- **Fast Analysis of Molecular Dynamics Trajectories with Graphics Processing Units – Radial Distribution Functions.** B. Levine, J. Stone, and A. Kohlmeier. 2010. (submitted)
- **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*, 2010. In press.
- **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 15<sup>th</sup> International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 347-358, 2010.

# GPU Computing Publications

<http://www.ks.uiuc.edu/Research/gpu/>

- **Probing Biomolecular Machines with Graphics Processors.** J. Phillips, J. Stone. *Communications of the ACM*, 52(10):34-41, 2009.
- **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 Processing Units (GPGPU-2)*, *ACM International Conference Proceeding Series*, volume 383, pp. 9-18, 2009.
- **Multilevel summation of electrostatic potentials using graphics processing units.** D. Hardy, J. Stone, K. Schulten. *J. Parallel Computing*, 35:164-177, 2009.



# GPU Computing Publications

<http://www.ks.uiuc.edu/Research/gpu/>

- **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.