

Accelerating the BLAST code with hybrid MPI + OpenMP + CUDA programming

Tingxing(Tim) Dong

Lawrence Livermore National Laboratory

September 7th, 2012

Mentors: Tzanio Kolev, Robert Rieben



- Introduction of BLAST
- Motivation
- Implementation
- Results
- Conclusion

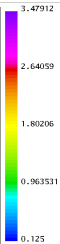
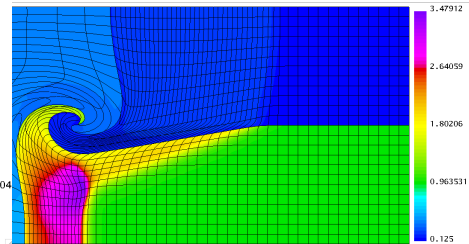
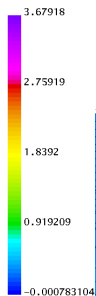
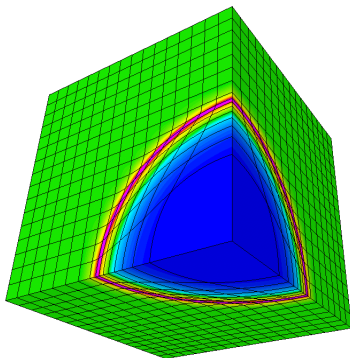
- Introduction of BLAST
- Motivation
- Implementation
- Results
- Conclusion

BLAST

- Solve equations of compressible hydrodynamics with Finite Element Method(FEM)
- Based on Lagrangian frame (moving mesh)
- C++ code, parallelized by MPI

BLAST's features

- Curvilinear zone geometries
- Higher order field representations
- Exact discrete energy conservation by construction
- Reduces to classical SGH under simplifying assumptions
- Support for 2D/3D meshes
- Multiple options for basis functions / quadrature order



- Introduction of BLAST
- **Motivation**
- Implementation
- Results
- Conclusion

Target Cluster: Edge

CPU	GPU	Memory/Node	Switch	Nodes
2x 6Core X5660	2 M2050	96GB	IB QDR	216

Table: Overview of Edge Cluster in LLNL

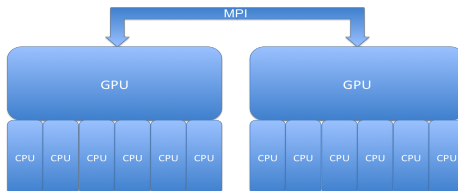


Figure: Architecture hierarchy

Two Strategies

MPI + CUDA

- 1 GPU is dedicated to 1 MPI task(before Kepler)
- 6 MPI task run on 6 cores
- 6 cores + 1 GPU : 1 MPI runs with GPU, 5 do not
- Load balance problem between MPIs
- We have to modify METIS, which takes care of load balance of BLAST

MPI + OpenMP + CUDA

- 6 OpenMP threads run on 6 cores
- 1 MPI task calls 6 OpenMP threads and 1 GPU

- Introduction of BLAST
- Motivation
- **Implementation**
- Results
- Conclusion

Euler's Equations in a Lagrangian Frame

Euler's Equations

Momentum Conservation: $\rho \frac{d\vec{v}}{dt} = \nabla \cdot \sigma$

Mass Conservation: $\frac{1}{\rho} \frac{d\rho}{dt} = -\nabla \cdot \vec{v}$

Energy Conservation: $\rho \frac{de}{dt} = \sigma : \nabla \vec{v}$

Equation of State: $p = EOS(e, \rho)$

Equation of Motion: $\frac{d\vec{x}}{dt} = \vec{v}$

Semi-discrete finite element method in BLAST

Momentum Conservation: $M_v \cdot \frac{d\mathbf{v}}{dt} = -\mathbf{F} \cdot \mathbf{1}$

Energy Conservation: $\frac{de}{dt} = M_e^{-1} \mathbf{F}^T \cdot \mathbf{v}$

Equation of Motion: $\frac{d\mathbf{x}}{dt} = \mathbf{v}$

Generalized corner forces on the GPU

Semi-discrete finite element method in BLAST

Momentum Conservation: $\mathbf{M}_v \cdot \frac{d\mathbf{v}}{dt} = -\mathbf{F} \cdot \mathbf{1}$

Energy Conservation: $\frac{de}{dt} = \mathbf{M}_e^{-1} \mathbf{F}^T \cdot \mathbf{v}$

Equation of Motion: $\frac{d\mathbf{x}}{dt} = \mathbf{v}$

Matrix \mathbf{F} is highly floating point operation intensive and thread independent

\mathbf{F} is constructed by two loops:

- Loop over zones in the domain (in each processor)
 - Loop over quadrature points in this zone
- Compute hydro forces associated with this quadrature point

on each point we compute this value absolutely independently

$$(\mathbf{F}_z)_{ij} = \int_{\Omega_z(t)} (\sigma : \nabla \vec{w}_i) \phi_j = \sum_k \alpha_k \hat{\sigma}(\hat{q}_k) : \mathbf{J}_z^{-1}(\hat{q}_k) \hat{\nabla} \hat{w}_i(\hat{q}_k) \hat{\phi}_j(\hat{q}_k) |\mathbf{J}_z(\hat{q}_k)|$$

Pre-computed constant values that
can go in constant memory

have to solve eigen vectors/values, SVD.

Method	Corner Force Matrix	CG Solver	Total time
Q4Q3	198.6	53.6	262.7
Q3Q2	72.6	26.2	103.7
Q2Q1 3D	90	56.7	164

Table: Corner force tasks 55%-75% of total time. CG solver takes 20%-34%. Measurement is based on three hundred iterations. Time is in seconds.

Hybrid programming model

two layers of parallelism

- MPI-based parallel domain-partition and communication between CPUs
- CUDA and OpenMP based parallel corner force inside each MPI task
- CUDA OpenMP can be turned on off

OpenMP & CUDA to accelerate Corner Force

- Host thread distributes work on GPU and return immediately
- Spawn 6 threads running on 6 cores
- No communication during computation of corner force
- OpenMP and GPU are done; Synchronization of CPU and GPU.
- host thread resumes
- Auto balance to maintain load balance between OpenMP and GPU

Semi-discrete finite element method in BLAST

Momentum Conservation: $\mathbf{M}_v \cdot \frac{d\mathbf{v}}{dt} = -\mathbf{F} \cdot \mathbf{1}$

Energy Conservation: $\frac{de}{dt} = \mathbf{M}_e^{-1} \mathbf{F}^T \cdot \mathbf{v}$

Equation of Motion: $\frac{d\mathbf{x}}{dt} = \mathbf{v}$

CUDA kernel 1: Loop over quadrature points. Compute part of \mathbf{F} based on \mathbf{v} , \mathbf{e} , \mathbf{x} (transferred from CPU), and need to solve eigen vec/values and SVD.

CUDA kernel 2: Loop over zones. each zone does a DGEMM and assemble the \mathbf{F}

CUDA kernel 3 : Compute $\mathbf{F} \cdot \mathbf{1}$ and either return result to the CPU or keep on the GPU depending on the CG solver settings.

CUDA kernel 4 : Compute $\mathbf{F}^T \cdot \mathbf{v}$

Mass matrix solve on the GPU

Semi-discrete finite element method in BLAST

Momentum Conservation: $\frac{d\mathbf{v}}{dt} = -\mathbf{M}_v^{-1} \mathbf{F} \cdot \mathbf{1}$

Energy Conservation: $\frac{d\mathbf{e}}{dt} = \mathbf{M}_e^{-1} \mathbf{F}^T \cdot \mathbf{v}$

Equation of Motion: $\frac{d\mathbf{x}}{dt} = \mathbf{v}$

CUDA kernel 5 : CUDA CG solver for $\mathbf{M}_v^{-1} \mathbf{F} \cdot \mathbf{1}$ based on CUBLAS/CUSPARSE/MAGMA, with a diagonal preconditioner.

CUDA kernel 6 : SpMV to solve $\mathbf{M}_e^{-1} \mathbf{F}^T \cdot \mathbf{v}$ by calling CUSPARSE

Notice:

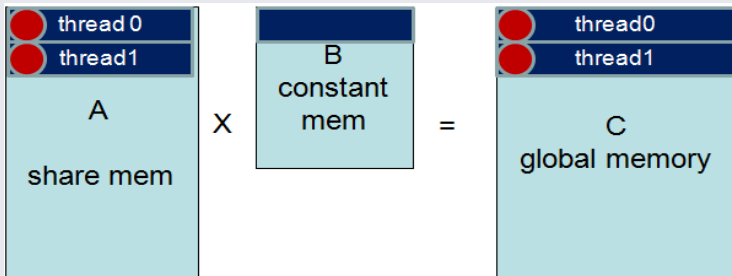
\mathbf{M}_v and \mathbf{M}_e^{-1} are computed once and read only thereafter (stay on GPU)

\mathbf{M}_v^{-1} is dense, so we did not use it directly

\mathbf{M}_e is a diagonal local dense matrix, so \mathbf{M}_e^{-1} is a sparse one, can be used directly

Kernel2: $AB^t = C$

- Each thread block(zone) does a $AB^t=C$
- A is updated in each time step; B is read only and common to every blocks
- Accessing global memory is coalesced



Problem	Kernel 2.1	Kernel 2.2	Kernel 2.3	CUBLAS	Kernel 1
Sedov 2D	0.63	0.42	0.34	27.2	2.06
Triple-pt	0.14	0.095	0.08	5.7	0.48

Table: Running time of different versions of kernel 2. In kernel 2.3 the bandwidth is highly improved by taking advantage of the memory hierarchy. For comparison, time of Kernel 1 is also presented. Kernel time is measured with CUDA event in milliseconds.

Kernel 2.1 only use global memory to read A and B. Kernel 2.2 use shared memory to read A. Kernel 2.3 used shared memory to read A and constant memory to read B.

- In corner force, each zone computes independently;
- Some zones go to GPU; The others go to CPU (how to find the optimal ratio is the key) .
- Because GPU runs asynchronously with CPU, control will return to host thread prior to GPU completing work.
- After launch CUDA kernels, host thread will spawns OpenMP threads and distributes the zones (loop) among threads.
- Each OpenMP thread executes like normal serial code.
- Synchronization between CPU and GPU.

- Repeated iterations of the same computational component;
- First, work load(zones) is distributed among CPU and GPU in an arbitrary ratio, say half to half.
- In each iteration, GPU and CPU are timed separately.
- If the ratio of their running time exceed the interval say $[0.9 - 1.1]$, move zones will go to the one who finished earlier.
- Just a few sampling periods.
- This idea can be extend to other tuning parameters.

Problem	Optimal ratio	Convergence period
Sedov 3D	0.45	3
Sedov 2D	0.75	14
Triple-pt	0.77	12

Table: The optimal ratio refers to the percentage of zones distributed on GPU out of total zones. The starting tentative ratio is 0.5, and target interval is $[0.9 - 1.1]$. One sampling period consists of forty iterations.

- Introduction of BLAST
- Motivation
- Implementation
- **Results**
- Conclusion

Validation: Engery preservation

Platform	Time out	Kinetic	Internal	Total	Total Change
CPU/GPU	Initial	0.0000000e+00	1.0050000e+01	1.005e+01	N/A
CPU	0.6	5.0423596e-01	9.5457640e+00	1.005e+01	-9.21929199648e-13
GPU	0.6	5.0418618e-01	9.5458131e+00	1.005e+01	-4.93827201353e-13

Table: Results of CPU and GPU for triple-pt problem with Q3Q2 method; Total Energy includes kinetic energy and internal energy. Both CPU and GPU results preserve energy conservation very well.

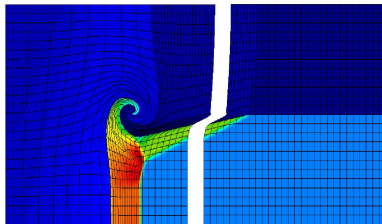


Figure: 2MPI tasks each with 1 M2050 and 6 Xeon Cores

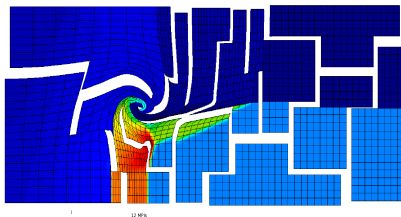


Figure: 12 MPI CPU tasks

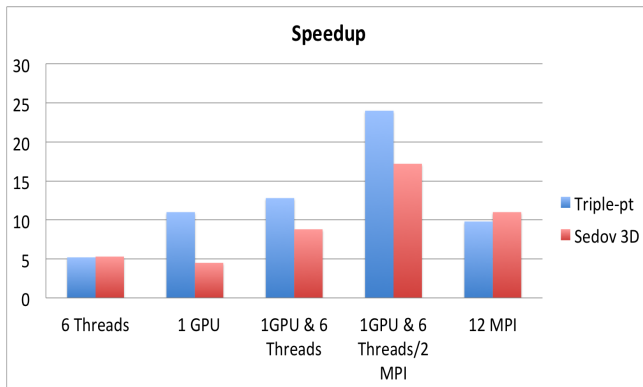


Figure: Speedup of corner force compared to serial code

Performance of CUDA-CG: 1 GPU vs 1 core

Problem	Method	MFEM PCG	CUDA PCG	Speedup
Triple-pt	Q3Q2	90.18	20.8	4.3
Sedov 3D	Q2Q1	27.81	10.55	2.6

Table: Speedup of CUDA-PCG compared to PCG in BLAST. Memory transfer overhead is counted in CUDA-PCG. Measurement is based on 1000 iterations. Time is in seconds.

Limitation: CUDA-CG only runs on 1 GPU until now.

Overall speedup

Overall testing consider every component, including CUDA corner force and CUDA-PCG and the unparallelized part. Tests are performed on a single GPU compared to a single CPU core.

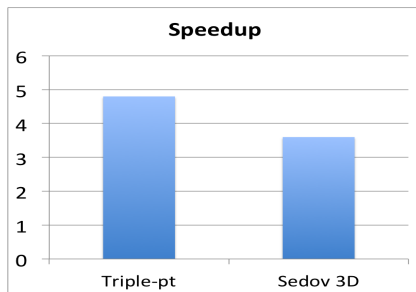


Figure: Overall speedup: including CUDA corner force, CUDA-PCG and unparallelized part

Method	Corner Force Matrix	CG Solver	Maximum Speedup
Triple-pt Q3Q2	70%	20%	10
Sedov Q2Q1 3D	54%	34%	8

- Introduction of BLAST
- Motivation
- Implementation
- Results
- **Conclusion**

Conclusion

- GPU computing is a good choice for CFD problems.
- We only parallelize part of code on GPU/OpenMP, the overall speedup is limited by Amdahl's law.
- On Kepler, MPI tasks "truly" share one GPU by Hyper-Q.

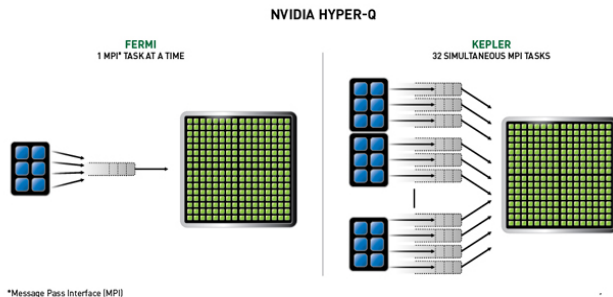


Figure: Kepler: Hyper-Q