



PORTABLE PARALLEL HALO AND CENTER FINDERS FOR HACC

Christopher Sewell, LANL Katrin Heitmann, ANL Ollie Lo, LANL Salman Habib, ANL Jim Ahrens, LANL

> OLCF User Meeting July 23, 2014

Outline

Background: SDAV, PISTON, VTK-m, dataparallel abstractions, and halo analysis
Algorithms: Halo and center finding
Results: Moonlight (LANL), Stampede (TACC), and Titan (ORNL)
Extensions: in-situ integration and Poisson center

finder

Future needs: current bottlenecks, expected trends, possible hardware and software solutions



Background

SDAV, PISTON, VTK-m, Data-Parallel Abstractions, and Halo Analysis





SDAV VTK-m Frameworks

Objective: Enhance existing multi/many-core technologies in anticipation of in situ analysis use cases with LCF codes

Benefit to scientists: These frameworks will make it easier for domain scientists' simulation codes to take advantage of the parallelism available on a wide range of current and nextgeneration hardware architectures, especially with regards to visualization and analysis tasks

Projects

EAVL, Oak Ridge National Laboratory Dax, Sandia National Laboratory DIY, Argonne National Laboratory PISTON, Los Alamos National Laboratory

Work on integrating these projects with VTK is on-going, in collaboration with Kitware



Office of ENERGY Office of Science

Algorithms for Science Applications Using VTK-m

- The PISTON component of VTK-m focuses on developing data-parallel algorithms that are portable across multi-core and many-core architectures for use by LCF codes of interest
- PISTON consists of a library of visualization and analysis algorithms implemented using Thrust, as well as a set of extensions to Thrust
- PISTON algorithms are integrated into LCF codes in-situ either directly or though integration with ParaView Catalyst



PISTON isosurface with curvilinear coordinates





Ocean temperature isosurface generated across four GPUs using distributed PISTON



PISTON integration with VTK and ParaView



Brief Introduction to Data-Parallelism and Thrust

What algorithms does Thrust provide?

Sorts	input	4	5	2	1	3
Transforms	transform(+1)	5	6	3	2	
Reductions	<pre>inclusive_scan(+) exclusive_scan(+)</pre>	4 0	9 4	11 9	12 11	15 12
Scans	<pre>exclusive_scan(max) transform_inscan(*2,+)</pre>	0 8	4 18	5 22	5 24	5 30
Binary searches	for_each(-1) sort	3 1	4 2	1 3	0 4	2 5
Stream compactions	copy_if(n % 2 == 1) reduce(+)	5	1	3		15
Scatters / gathers	input1	0	0	2	4	8
hallenge: Write operators in terms	input2	3	4 	1	0 2	2
these primitives only	permutation iterator	с 4	4	0	0	2

Reward: Efficient, portable code



C

0



Halo Analysis

- An important and time-consuming analysis function within HACC is finding halos and the centers of those halos
- HACC is written as an MPI+X code, designed to be easy to port and optimize for different architectures
- It would be very time-consuming to optimize the full cross-product of all analysis operators and all architectures
- Data-parallelism, using a framework such as PISTON, allows a single implementation of an analysis operator to take advantage of all supported architectures
- Total implementation work is O((# of analysis operators) + (# of architectures)) rather than O((# of analysis operators) * (# of architectures))





Images of the matter density field, produced by the HACC simulation on Titan Image credits: Joe Insley and Silvio Rizzi





Definitions

- Friend-of-friends (FOF) halo: connect each particle to all "friends", i.e., all other particles within a specified "linking length" of it; two particles will end up in the same "halo" if there exists any chain of "friends" between them
- Most connected particle (MCP) center: the particle within a halo with the most "friends"
- Most bound particle (MBP) center: the particle within a halo with the lowest potential, where the potential for a given particle is computed as the sum over all other particles of the negative of mass divided by distance



MBP center:

$$\arg\min_{i\in H} \left(-\sum_{j\in H, j\neq i} \frac{m_j}{d_{ij}}\right)$$





Algorithms

Halo and Center Finding





Distributed Parallel Halo Finder

- Particles are distributed among processors according to a decomposition of the physical space
- Overload zones (where particles are assigned to two processors) are defined such that every halo will be fully contained within at least one processor
- Each processor finds halos within its domain
- At the end, the parallel halo finder performs a merge step to handle "mixed" halos (shared between two processors), such that a unique set of halos is reported globally
- Developed primarily by Pat Fasel at LANL, with the serial halo finder running on each individual processor using a KD-tree based algorithm designed by C.H. Hsu





Distributed Parallel Halo Finder

- Particles are distributed among processors according to a decomposition of the physical space
- Overload zones (where particles are assigned to two processors) are defined such that every halo will be fully contained within at least one processor
- Each processor finds halos within its domain: Drop in PISTON multi-/many-core accelerated algorithms
- At the end, the parallel halo finder performs a merge step to handle "mixed" halos (shared between two processors), such that a unique set of halos is reported globally





Halo Finder Algorithm: Connected Components



// Pointer jumping on each vertex

for all i pardo set D(i):=D(D(i))

All rooted stars

Reference: An Introduction to Parallel Algorithms, Joseph JáJá, 1992





Halo Finder Algorithm: Computing Edges on the Fly

- If we define an edge to exist between two particles if and only if their distance is less than the linking length, the connected components solution is the FOF halos
- However, it could take O(n²) time and O(n²) memory to directly compute and store all edges
- Instead, partition the domain into bins with edge length equal to the linking length
 - Friends can only exist in its own bin or one of 27 neighbor bins
 - Compute edges on the fly in the algorithm, comparing each particle to each other particle in its bin and its neighbor bins to see which edges exist
 - The number of bins (most empty) may be too large to store pointer to each, so instead store only pointers to neighbor bins for each particle
 - Neighbor bins will be located in 1D vector in groups of three, so 9 not 27 pointers needed to neighbors for each particle





Halo Finding: Binning Example



Two vectors are computed, N1 and N2, which contain the beginning and ending of three ranges (nine in 3D) in the sorted particle vectors for which each particle will need to search for its potential friends.

Scalable Data Management, Analysis, and Visualization



Center Finding Algorithms

- Most connected particle can be found easily by counting the number of friends for each particle during any iteration through the virtual edge list, and then using a max-reduce to find the maximum
- Similarly, an approximation of DB scan can exclude any particle with too few friends from all edges
- Most bound particle can be found by computing the potential for each particle in a highly parallel brute force approach
- Centers for all halos can be computed simultaneously using segmented vectors





MBP Center Finding Example



The inputs are the particle ids (I), coordinates (X, Y), and halo id (D, found using the halo finding algorithm). The output is a vector C containing for each particle the id of the MBP center for its halo.

Scalable Data Management, Analysis, and Visualization



Results

Moonlight, Stampede, and Titan





Results: Moonlight



FOF + MBP: PISTON ~4.9x faster than original with 16 rpn FOF + MCP: PISTON ~2.5x faster than original with 16 rpn





Results: Visual comparison of halos



Original Algorithm VTK-m







Results: Xeon Phi (MIC) on Stampede



• To demonstrate the portability of our algorithms, the same code was compiled to the Thrust OpenMP backend (including our own OpenMP implementation of scan) and run on a 256³ particle data set on an Intel Xeon Phi SE10P (MIC) Coprocessor on a single node of Stampede at TACC

• PISTON version scales to more cores than running the existing serial algorithms with multiple MPI processes





Results: Titan



• This test problem has ~90 million particles per process.

• Due to memory constraints on the GPUs, we utilize a hybrid approach, in which the halos are computed on the CPU but the centers on the GPU.

• The PISTON MBP center finding algorithm requires much less memory than the halo finding algorithm but provides the large majority of the speed-up, since MBP center finding takes much longer than FOF halo finding with the original CPU code.





Results: Large Run on Titan



• Because the memory requirements increase with the number of MPI processes due to overload regions, the HACC simulation with this data can only use one MPI process per node, along with the associated GPU, given the memory available on Titan.

• Thus, there is an even greater potential for speed-up by utilizing the GPUs.

• The performance improvements using PISTON on GPUs allowed halo analysis to be performed on a very large 8192³ particle data set across 16,384 nodes on Titan for which analysis using the existing CPU algorithms was not feasible.





Extensions

In-situ Integration and Poisson Center Finder





In-situ Integration in HACC

- Successfully ran 500 time-step, 512³ particle simulation on Moonlight using our GPU halo and center finders integrated with HACC in-situ
- Integrated with CosmoTools





Extension: Potential Field Algorithm

- Optimization for MBP center finder: for each halo...
 - Superimpose a grid over the (extended) extents of the halo
 - Estimate particle density on the grid using binning
 - Solve Poisson equation for the potential on the grid using FFT (actually DST for zero boundary conditions approximation)
 - Find the grid point with minimum potential
 - Search around the neighborhood of the minimum potential grid for the particle with minimum potential
 - Return the position of such particle as the halo center







Future Needs

Current Bottlenecks, Expected Trends, and Possible Hardware and Software Solutions





Future Needs: Memory

- Current bottlenecks
 - Main memory: 8192³ simulation run with only 1 rpn because number of "ghost" particles increases with total number of ranks, and even 2 rpn exceeded total system memory
 - GPU memory: Halo finding performed on CPU because of GPU memory limitations
- Expected future trends
 - Growth in computational rates will outpace growth in memory capacity and bandwidth
- Possible solutions
 - Hardware: More memory
 - Software: Streaming and external memory algorithms





Future Needs: Resiliency

- Current bottlenecks (16k node simulation)
 - Node failures on start-up
 - Node failures while running
- Expected future trends
 - Systems will have more nodes with more cores; more things that could fail
- Possible solutions
 - Hardware: Components with longer MTBF
 - Software: Improved checkpointing (perhaps using burst buffers?); more robust run-time schedulers





Future Needs: Portability

- Current bottlenecks
 - Simulation code has to be rewritten to run, or at least to run efficiently, on new architectures
 - Some potential optimizations have been foregone because platform-specific optimizations would inhibit portability
- Expected future trends
 - More types of accelerators will come to market; individual systems will become more heterogeneous
- Possible solutions
 - Hardware: Standardize on an architecture (not likely)
 - Software: Write programs using higher-level abstractions





Future Needs: In-Situ

- Due to power and I/O constraints, we expect to need to perform analysis in-situ rather than in post-processing
- Traditional options: efficient but highly constrained insitu, or expensive but highly interactive post-processing
- Recent work explores using image databases to allow trade-offs in the space between these extremes





LA-UR-14-25437



Jim Ahrens, Patrick

O'Leary, et. al.