

Investigation of Portable Event Based Monte Carlo Transport Using the Nvidia Thrust Library

Ryan C. Bleile^{*,†}, Patrick S. Brantley^{*}, Shawn A. Dawson^{*}, Matthew J. O'Brien^{*}, Hank Childs[†]

^{*}Lawrence Livermore National Laboratory, P.O. Box 808, Livermore, CA 94551

[†]Department of Computer and Information Science, University of Oregon, Eugene, OR 97403
bleile1@llnl.gov, brantley1@llnl.gov, dawson6@llnl.gov, obrien20@llnl.gov, hank@uoregon.edu

INTRODUCTION

Power consumption considerations are driving future high performance computing platforms toward many-core computing architectures. The Trinity machine to become available at Los Alamos National Laboratory in 2016 will use both Intel Xeon Haswell processors and Intel Xeon Phi Knights Landing many integrated core (MIC) architecture coprocessors. The Sierra machine to be available at Lawrence Livermore National Laboratory beginning in 2018 will use an IBM PowerPC architecture along with Nvidia graphics processing unit (GPU) architecture accelerators. As a result of these different advanced architectures, the computing landscape for the upcoming years is complex.

Traditional approaches to Monte Carlo transport do not work efficiently on these new computing platforms. MIC architectures require vectorization to operate efficiently, vectorization is difficult to achieve in Monte Carlo transport. GPU architectures require additional code to explicitly use the hardware, requiring significant code changes or hardware specific branches in the source code. A significant challenge for Monte Carlo transport projects is to simultaneously support efficient versions of simulation codes for both the current generation and the different advanced computing architectures within a single source code base.

In order to address these issues, two important changes are typically made: a new algorithmic approach to solving Monte Carlo transport, and explicit use of the GPU hardware in software. In this paper, we describe initial research investigations of an event based Monte Carlo transport algorithm [1] implemented using the Nvidia Thrust library [2] on a GPU for a Monte Carlo test code. The event based algorithm targets many-core architectures by increasing SIMD (single instruction multiple data) parallelism, while Thrust provides portable performance by allowing one source code base to compile code targeted for both CPUs and GPUs

HISTORY BASED APPROACH

For this research, we began with the ALPSMC Monte Carlo test code [3] that models particle transport in a one-dimensional planar geometry binary stochastic medium. The ALPSMC code was originally implemented in C++ using a standard history based Monte Carlo transport algorithm as shown in Alg. 1. This approach follows a single particle at a time from creation until it is absorbed or leaked. Parallelism is easily achieved with parallelism over particle histories, (main foreach loop), and each thread is left to work independently on one particle at a time. This approach does not use vectorization and uses a MIMD (multiple instruction multiple data) parallelism scheme.

Algorithm 1: History based Monte Carlo algorithm

```
1 foreach particle history do
2   generate particle from boundary condition or
   source
3   while particle not escaped or absorbed do
4     sample distance to collision in material
5     sample distance to material interface
6     compute distance to cell boundary
7     select minimum distance, move particle, and
   perform event
8     if particle escaped spatial domain then
9       update leakage tally
10      end particle history
11    if particle absorbed then
12      update absorption tally
13    end particle history
```

EVENT BASED APPROACH

Previous researchers [4, 5, 6] have noted that the use of an event based Monte Carlo particle transport algorithm [1] may be beneficial for GPU or vector-based architectures. We investigated this idea through the event based algorithm shown in Alg. 2 as a way to potentially optimize performance on GPU and vector-type architectures. In event based particle tracking, the individual events can be treated by a series of data parallel operations. The data parallel model matches the vector and GPU hardware with an emphasis on performing the same operations on many pieces of data at one time through SIMD parallelism.

Thrust

Thrust is a C++ header library using a STL-like template interface [2]. Thrust provides a number of parallel algorithms and data structures designed to provide access to GPU computing without needing to write in CUDA [7]. Additionally, Thrust provides backend capabilities allowing these algorithms and data structures to target different devices, even CPUs with OpenMP threads. This design was used for studying portable performance techniques with Thrust, providing a method of maintaining only one source code.

Thrust algorithms are used for implementing procedures across all particles in a batch. These algorithms perform operations such as the data parallel map, reduce, gather, scatter, or scan operations defined in reference [8]. Each of these operations can be performed in a data parallel way.

Thrust also provides data types that can be used to man-

Algorithm 2: Event based Monte Carlo algorithm

```
1 foreach batch of particle histories (fits in memory  
   constraint) do  
2   generate all particles in batch from boundary  
   condition or source  
3   determine next event for all particles (collision,  
   material interface crossing, cell boundary  
   crossing)  
4   while particles remaining in batch do  
5     foreach event E in (collision, material  
       interface crossing, cell boundary crossing)  
       do  
6       identify all particles whose next event is  
       E  
7       perform event E for identified particles  
       and determine next event for these  
       particles  
8     if particle escaped spatial domain then  
9       update leakage tally  
10    if particle absorbed then  
11      update absorption tally  
12    delete particles absorbed or leaked
```

age memory for GPU devices. The `thrust::device_vector` and `thrust::host_vector` data structures operate similarly to a C++ `std::vector` but with automatic memory copying between host and devices whenever necessary. These data types allow for simple memory management schemes that work on both GPU and CPU based architectures.

Algorithm Detail

An event based algorithm focuses on performing data parallel operations across all particles performing the same action. Additional overhead is needed to find the grouping of particles that will be operated on and to decide an access pattern for particles. This reorganization stage can be costly and is not directly related to solving the transport problem.

Thrust provides permutation iterators that allow for the unaligned access of data elements according to an index map. Using this iterator scheme, data elements do not need to be copied into new locations for each operation. This is done at the cost of performing non-contiguous memory accesses for reading and writing the information.

In order to perform an event operation on particles using this scheme, a series of data parallel operations is used to establish the correct index mapping for the permutation iterator. This scheme is defined as follows and describes in detail lines six and seven of algorithm 2:

- Step 1: `thrust::transform` — Fill out a stencil map of 1's and 0's of all particles doing event E (where each particle whose next event is E will get a 1 in the stencil map at its index location)
- Step 2: `thrust::reduce` — Count the number of elements labeled 1 in the stencil (determines the number of particles

that will perform event E)

- Step 3: Check if the number of elements is greater than 0 (check if any particles are performing event E)
- Step 4: `thrust::exclusive_scan` — generate indices for index mapping from stencil map (indices for each particle performing event E)
- Step 5: Allocate a new map of appropriate size (map to hold indices for all particles performing event E)
- Step 6: Scatter indexes from scan into new index map (reduces the `exclusive_scan` generated indices into the map that holds only enough for particles performing event E)
- Step 7: Use new index map in `permutation_iterator` loops over all particles (combining the index map with the permutation iterator allows loops over all particles to operate only on the ones selected in the index map)

Related Work

Other researchers have performed related work implementing event based algorithms on GPUs [4, 5, 6]. Comparing and contrasting our event based method to theirs shows differences both in our approaches, target problem, and the way chosen to create data parallelism. While their research first studied the idea of event based Monte Carlo on GPUs, our research studied the event based Monte Carlo in the light of the portable performance platform of Thrust, as well as an algorithm designed around performing only data parallel operations.

The research performed by Bergmann et. al [6] with WARP most closely matches our work. We both implemented an algorithm that does a series of operations on the GPU in a single loop over particles still active. Additionally, we both implemented a remapping vector to identify particles for a task without copying them. While in WARP the remap vector is then sorted so that particles will most likely be close to other particles doing the same operations, in our implementation only particles who will do the same event are called in one kernel. Additionally, our remap vector does no sorting and only identifies whether particles will undergo a specific event. In this way, we implemented our event based model as a series of data parallel calls with Thrust. Each of our calls has almost no divergence except when deciding if a particle leaked or was absorbed.

The work done by Liu et al [5] within ARCHER to test an event based model, differs from our work algorithmically. Liu et. al did use Thrust for some operations but chose to make the important kernels CUDA only. Additionally, Liu's algorithm used a double while loop focusing all attention on one event and then switching to the other, their code having two events to consider.

The work done by Nelson et. al [4] is related in the attempt to use and optimize performance on a GPU, but event based methods were not considered in that work.

IMPLEMENTATIONS

We implemented the event based version of ALPSMC using both the Nvidia CUDA programming model [7] and

the Nvidia C++ Thrust library [2]. Our Thrust implementation of ALPSMC utilizes many data parallel operations and utilizes Thrust data types for managing memory. In the native CUDA implementation of ALPSMC, we found it useful to continue to use Thrust algorithms in building various maps. ALPSMC is implemented using double precision floating point numbers throughout. The Thrust and CUDA implementations of ALPSMC give physics results identical to the original history based implementation.

Multiple versions of ALPSMC were implemented for study. Each version was implemented in Thrust, (compiled with both OMP and CUDA), and in CUDA directly. Both the history based and event based algorithms were implemented. Additionally a version with and without a zonal flux tally was implemented. This produced four thrust versions compiled both ways for a total of eight versions. Additionally there were four CUDA versions, and the original serial C version with and without the zonal flux tally.

The CUDA implementations of this study matched the steps taken in the Thrust implementations. The differences in performance come from the few areas that native CUDA programming provides that cannot be accomplished with Thrust. Using CUDA directly enables more fine grained control at the kernel level and enables important access to different memory spaces such as GPU shared memory. This includes a scheduling algorithm to optimize the number of active threads on the GPU for each kernel call. Additionally, this includes the use of the different available memory spaces, such as constant and shared memory. For example, Monte Carlo particles were initially allocated in GPU global memory and then copied to shared memory for all operations within a kernel. All problem constants such as cross sections and mean chord length values were placed in GPU constant memory. These optimizations under certain conditions can have a significant impact on the performance of a GPU kernel.

NUMERICAL RESULTS

We performed scaling studies in which we varied the number of Monte Carlo particle histories (problem size) and the implementation methodology (Thrust or CUDA). (The results presented are for Case 1a [3] with a spatial domain of 10 cm.) We also examined the differences in performance on three different computer platforms. The RZGPU platform has Intel Xeon Westmere-EP 2.8 GHz host cores with Nvidia Tesla M2070 GPU device accelerators. The Max platform has Intel Sandy Bridge 2.6 GHz host cores with Nvidia Tesla K20X GPU device accelerators. The Tesla K20X GPU has improved double precision performance over the Tesla M2070. The RZHASGPU platform has Intel Xeon Haswell 3.2 GHz host cores with Nvidia Tesla K80 GPU device accelerators. We did not use a multiple GPU implementation and are therefore only able to utilize around half of the computational power of the Nvidia Tesla K80s.

Our first study aimed to identify the speedups of our event based algorithm when compared to the initial serial implementation. We computed speedups over a serial calculation by dividing the wallclock time of a serial run of the history based version of ALPSMC on the host core of the given machine by the wallclock time of the event based

version of ALPSMC running on both the host and GPU device. The speedups obtained on the RZGPU and Max platforms are shown in Table I.

For this test the speedups obtained using the CUDA implementation of the event based algorithm are significantly larger than those obtained using the Thrust implementation by up to over a factor of four. We attribute this improved performance to the fact that CUDA offers more control over the memory spaces available on the GPU (e.g. shared memory). Thrust does not offer such flexibility and manages the memory allocation internally. We conclude based on these preliminary investigations that a direct CUDA implementation is more efficient than a Thrust implementation for event based Monte Carlo. Also, the speedups on the Tesla K20X GPU are larger than on the Tesla M2070 GPU by up to a factor of approximately two, presumably a result of the improved double precision performance of the K20X.

TABLE I: ALPSMC Event Based Monte Carlo GPU Speedups

	Number Particle Histories		
	10^6	10^7	10^8
CUDA (K20X)	5.90	11.88	11.91
CUDA (M2070)	3.96	6.05	6.05
Thrust (K20X)	2.11	2.60	2.60
Thrust (M2070)	1.42	1.64	1.63

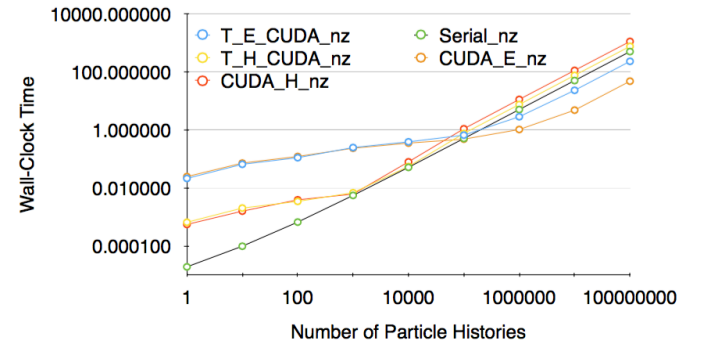


Fig. 1: Weak Scaling Study of ALPSMC under all versions tested

In a second test we performed a full scaling study of each version of the code on the RZHASGPU platform. Figure ?? shows the results of the versions of the code that did not have a zonal flux tally. In this diagram we can see that the event based methods have a significantly higher overhead. But at higher scale the overhead is less than the performance gains and we begin to see scaling improvements. At approximately 10^5 particle histories the event based versions of the code begin outperforming the history based versions. The serial plot also shows that it takes around 10^5 particle histories for the overhead of accessing the GPU to outperform the CPU in serial. Additionally we can see that the performance gains of the CUDA version over the Thrust version start to become significant at higher numbers of particle histories.

The Linear behavior at the right side of the plot is based

on the batching scheme we used to not overflow GPU device memory. Once we have entered the batching cycles we are no longer gaining any additional performance increases. At this stage our only improvements would be to process a greater number of particle histories and that number is hardware dependent.

CONCLUSIONS

We describe preliminary investigations of event based Monte Carlo algorithms for GPU architectures using a research Monte Carlo test code. We found that a CUDA implementation of an event based Monte Carlo algorithm performed significantly better than a Thrust implementation, most likely a result of additional flexibility in access to different memory spaces on the GPU. Additionally, we showed that on GPU platforms and at large enough problem sizes the event based implementations perform better than the history based implementations.

While investigating this problem we also discovered that the performance of the algorithm is effected by what tallies are being accounted for by each particle. The zonal flux tally included atomic operations that significantly impacted the performance of the code. We decided that since we wished to consider the effectiveness of the event based algorithm we should try removing that tally and see how that effected performance. In the future we would like to research better ways of handling tallies like this one. Additionally, we would like to consider new ways for optimizing both the Thrust and CUDA versions, in order to see how much performance we have yet to achieve.

ACKNOWLEDGMENTS

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. Funding was provided by the LLNL Livermore Graduate Scholar Program.

REFERENCES

1. F. B. BROWN and W. R. MARTIN, "Monte Carlo Methods for Radiation Transport Analysis on Vector Computers," *Progress in Nuclear Energy*, **14**, 269–299 (1984).
2. "Thrust Web Site," (2014), <https://developer.nvidia.com/Thrust>.
3. P. S. BRANTLEY, "A Benchmark Comparison of Monte Carlo Particle Transport Algorithms for Binary Stochastic Mixtures," *Journal of Quantitative Spectroscopy and Radiative Transfer*, **112**, 599–618 (2011).
4. A. G. NELSON, *Monte Carlo Methods for Neutron Transport on Graphics Processing Units Using CUDA*, M.S. Thesis, The Pennsylvania State University (2009).
5. T. LIU, X. DU, W. JI, X. G. XU, and F. B. BROWN, "A Comparative Study of History-Based Versus Vectorized Monte Carlo Methods in the GPU/CUDA Environment for a Simple Neutron Eigenvalue Problem," in "Proceedings of Supercomputing in Nuclear Applications and Monte Carlo (SNA+MC)," Paris France (October 27-31, 2013 2013).
6. R. M. BERGMANN and J. L. VUJIC, "Algorithmic Choices in WARP - A Framework for Continuous Energy Monte Carlo Neutron Transport in General 3D Geometries on GPUs," *Annals of Nuclear Energy*, **77**, 176–193 (2015).
7. "CUDA Web Site," (2014), http://www.nvidia.com/object/cuda_home_new.html.
8. G. E. BLELLOCH, *Vector models for data-parallel computing*, vol. 356, MIT press Cambridge (1990).