

Performance and Accuracy of WARP - A Framework for Continuous Energy Monte Carlo Neutron Transport in General 3D Geometries on GPUs

Ryan M. Bergmann*, Kelly Rowland, Nikola Radnović, Jasmina L. Vujić

Department of Nuclear Engineering, 4155 Etcheverry Hall, University of California - Berkeley, Berkeley, CA 94703-1730

Abstract

In recent supercomputers, general purpose graphics processing units (GPGPUs) are a significant fraction of the supercomputer's total computational power. GPGPUs have different architectures compared to central processing units (CPUs), and for Monte Carlo neutron transport codes used in nuclear engineering to take advantage of these coprocessor cards, transport algorithms must be changed to execute efficiently on them. WARP is a continuous energy Monte Carlo neutron transport code that has been written to do this. The main thrust of WARP is to adapt previous event-based transport algorithms to the new GPU hardware; the algorithmic choices for all parts of which are presented in this paper. It is found that remapping history data references increases the GPU processing rate when histories start to complete. The main reason for this is that completed data are eliminated from the address space, threads are kept busy, and memory bandwidth is not wasted on checking completed data. Remapping also allows the interaction kernels to be launched concurrently, improving efficiency. The OptiX ray tracing framework and CUDPP library are used for geometry representation and parallel dataset-side operations, ensuring high performance and reliability.

Keywords: Monte Carlo, Neutron Transport, GPU, CUDA, CUDPP,

*Corresponding author. Tel.: +1 320 333 3089.

Email addresses: ryanbergmann@gmail.com (Ryan M. Bergmann), krowland@berkeley.edu (Kelly Rowland), radnovicn@gmail.com (Nikola Radnović), vujic@nuc.berkeley.edu (Jasmina L. Vujić)

1. Introduction

Graphics processing units, or GPUs, have gradually increased in computational power from the small, job-specific boards of the early 1990s to the programmable powerhouses of today. Compared to more common central processing units, or CPUs, GPUs have a higher aggregate memory bandwidth, much higher rate of floating-point operations per second (FLOPS), and lower energy consumption per FLOP. Because one of the main obstacles in exascale computing is power consumption, many new supercomputing platforms are gaining much of their computational capacity by incorporating GPUs into their compute nodes. Since CPU-optimized parallel algorithms are not directly portable to GPU architectures (or at least not without losing substantial performance), particle transport codes need to be rewritten to execute efficiently on GPUs. Unless this is done, reactor simulations cannot take full advantage of these new supercomputers. In this work, the algorithmic choices used in WARP and their consequences are presented. The ultimate accuracy of WARP, comparisons against other neutron transport codes, and benchmarking are not presented here. This work only concerns algorithmic decisions based on self-comparison. WARP benchmarking, accuracy and speedup comparison with other Monte Carlo codes will be presented in a companion paper.

WARP, which can stand for “Weaving All the Random Particles,” is a three-dimensional (3D) continuous energy Monte Carlo neutron transport code developed at UC Berkeley to efficiently execute on a CPU/GPU platform. WARP accelerates Monte Carlo simulations while preserving the benefits of using the Monte Carlo method, namely, that very few physical and geometrical simplifications are applied. WARP is able to calculate multiplication factors, flux tallies, and fission source distributions for time-independent problems, and can run in both criticality or fixed source modes. WARP can currently transport neutrons in unrestricted arrangements of parallelepipeds, hexagonal prisms, cylinders, and spheres.

The impetus for developing WARP was the research done by Martin and Brown in 1984 [?] for Monte Carlo and by Vujić and Martin in 1991 [?] for collision probability method. In the 1984 paper, a method for mapping the Monte Carlo problem onto SIMD (single instruction multiple data) vector

35 computers is described. SIMD is an execution model some processors use in
36 order to lower the number of instructions needed per amount of computation
37 done, which increases both power and computational efficiency [?]. SIMD
38 requires the same instructions to be carried out over every element in a
39 concurrently-processed data vector. The essential idea in the 1984 paper is
40 to bank the neutrons into vectors based on their required operation. If a
41 neutron is scattering, its data is placed in the scattering buffer. If a neutron
42 needs to do a surface crossing, it is put in the crossing buffer, and so on.
43 Once a buffer becomes full, it is processed in a SIMD fashion by the vector
44 computer. Processing the neutrons makes the buffers contain non-uniform
45 reactions, however, and a “shuffle” operation is done that actually moves
46 data back into contiguous blocks based on the reaction type.

47 This new approach was named “event-based” Monte Carlo, since the
48 neutron events are tracked and processed as a group [?]. This was a very
49 different way of performing a Monte Carlo simulation at the time. Almost all
50 computers were strictly serial, and SIMD lanes were only available in super-
51 computers. Therefore, the pervasive method was the “task-based” method
52 in which neutrons are tracked for their entire lifetime in series. Since GPUs
53 are massively parallel and rely on SIMD, an event-based algorithm seems to
54 be the appropriate approach to GPU-accelerated neutron transport.

55 WARP uses an event-based algorithm, but with some important differ-
56 ences. Vectorizing the Monte Carlo algorithm should allow for efficient GPU
57 execution, but a paper by Zhang [?] also shows that by remapping data ref-
58 erences, the thread divergence in GPU warps can be minimized. Moving data
59 is expensive, so WARP uses a remapping vector of pointer/index pairs to di-
60 rect GPU threads to the data they need to access. The remapping vector is
61 sorted by reaction type after every transport iteration using a high-efficiency
62 parallel radix sort, which serves to keep the reaction types as contiguous as
63 possible and removes completed histories from the transport cycle. Sorting
64 reduces the amount of divergence in GPU “thread blocks,” keeps the SIMD
65 units as full as possible, and eliminates using memory bandwidth to check
66 if a neutron in the batch has been terminated or not. Using a remapping
67 vector means the data access pattern is irregular, but this is mitigated by
68 using large batch sizes where the GPU can effectively eliminate the high cost
69 of irregular global memory access.

70 WARP sets itself apart from any previous endeavors in its breath of scope
71 and its novel adaption of the event-based transport algorithm. Previous codes
72 have also either used synthetic, simplified, and/or incomplete nuclear data [?

[? ?]. WARP loads standard data files and accurately simulate each reaction type specified in the data. WARP also uses a flexible, scalable, and optimized geometry representation where previous studies have used simplified and restricted geometry models [?]. Previous works have examined event-parallel algorithms, but have not parallelized them effectively and therefore did not see the benefits of adopting such an algorithm on a GPU [?]. WARP uses highly-parallelized algorithms and slightly modify the original vision of the event-based algorithm to better suit execution on the GPU. The previous event-based algorithms tried to implement a “shuffle” operation where neutron data was actually sorted into reaction-contiguous blocks [?], or used small, synthetic nuclear data and were not able to capture the effects of loading large nuclear datasets [?]. Other studies have been done that offload specific parts of the neutron transport routine to the GPU, such as that done on RMC at Tsinghua University [?]. They were able to obtain a 113x speedup over a CPU version without a flux tally and 36x speedup with a flux tally [?]. Again, one group cross sections were used.

WARP also changes the unionized energy grid data format to reduce the number of data loads needed to scan cross sections. Instead of storing a matrix of pointers indexed by reaction type and energy, WARP stores three matrices. The first contains cross section values, the second contains pointers to angular distributions, and a third contains pointers to energy distributions. This linked list type of layout increases memory usage, but lowers the number of data loads that are needed to determine a reaction by eliminating a pointer load to find a cross section value.

Optimized, high-performance GPU code libraries are also used by WARP wherever possible. The CUDA performance primitives (CUDPP) library is used to perform the parallel reductions, sorts and sums, the CURAND library is used to seed the linear congruential random number generators, and the OptiX ray tracing framework is used for geometry representation [? ? ?]. OptiX is a highly-optimized library developed by NVIDIA that automatically builds hierarchical acceleration structures around user-input geometry so only surfaces along a ray line need to be queried in ray tracing [?]. WARP also performs material and cell number queries with OptiX by using a point-in-polygon like algorithm. The reaction sampling routines have been coded in CUDA according to the appropriate ENDF laws.

WARP is designed to read ACE-formatted data [?] and perform all reaction types as prescribed by the data. ACE is an acronym for “A Compact ENDF,” and is the format which MCNP and Serpent both read [? ?]. ACE

111 data is loaded with the PyNE package [?] and passed to the main C++
 112 code via the Python C application programming interface (API). WARP uses
 113 a Serpent-like unionized energy grid to regularize data access (Serpent is a
 114 Finnish Monte Carlo reactor physics code). A unionized energy grid is one
 115 where the individual energy grids of every isotope used in a simulation have
 116 been combined into a single energy vector. The cross section vectors are
 117 then interpolated on this larger unionized energy grid [?]. The host-side
 118 code in WARP is written in C/C++. Single precision floating point num-
 119 bers are used throughout in order to realize the full computational capacity
 120 of the GPU and to allow simulations to be carried out on more affordable
 121 and higher clocked GeForce cards [? ?]. Using single precision numbers
 122 may lead to inaccuracies when there are very dilute isotopes or very rare
 123 reactions reactions present, as roundoff error may make their contributions
 124 zero. Buffer overflow and roundoff error in the tallies may also be a problem
 125 with single precision, but this can be mitigated by accumulating the tallies
 126 frequently in a double precision vector. Roundoff error may also be problem-
 127 atic in calculating the multiplication factor, but this can also be mitigated by
 128 accumulating the integer secondary neutron yield values into a large host-side
 129 variable. If double precision data is found to be needed, WARP can easily
 130 be changed, and doing so may be an interesting experiment in the future.

131 Section 2 discusses how OptiX is used to perform surface detection and
 132 query the current material a neutron is traveling through. The testing done
 133 to determine the optimal OptiX configuration is also discussed. Section 3
 134 deals with the way the nuclear data is loaded, reformatted to use a unionized
 135 energy grid, and the details about data layout and access patterns. Section 4
 136 details each individual CUDA kernel (i.e. GPU program) used by WARP to
 137 process the neutron histories. A Section 5 presents the results of the OptiX
 138 testing and performance improvements attributed to reference remapping,
 139 and conclusions based on the results are made in Section 6.

140 2. Ray Tracing and Geometry Representation with OptiX

141 2.1. *Instancing*

142 2.2. *Test Geometries*

143 3. Unionized Cross Section Data Layout

144 4. CUDA Kernels and Data-Parallel Tasks

145 4.1. *Criticality Source*

146 5. Results

147 5.1. *OptiX Scaling Study Results*

148 5.2. *The Effects of Reference Remapping*

149 6. Conclusions

150 WARP has shown that GPUs are an effective platform for performing
151 Monte Carlo neutron transport with continuous energy cross sections. Cur-
152 rently, WARP is the most detailed and feature-rich program in existence
153 for performing continuous energy Monte Carlo neutron transport in general
154 3D geometries on GPUs, but compared to production codes like Serpent
155 and MCNP, WARP has limited capabilities. Despite WARP’s lack of fea-
156 tures, its novel algorithm implementations show that high performance can
157 be achieved on a GPU despite the inherently divergent program flow and
158 sparse data access patterns. Remapping threads to active data is an effec-
159 tive way of raising the processing rate when the number of active neutrons
160 becomes small; this also reduces thread divergence in reaction kernels. Using
161 a radix sort to do the remapping is effective since it segregates reactions into
162 contiguous blocks, efficient since it can be done in place and in $O(kN)$ time,
163 and can eliminate completed data from being accessed if slight modifications
164 to the standard reaction number encodings are made.

165 Most of the performance gain in remapping data references comes from
166 being able to launch grids that are sized for only the active data rather than
167 the entire dataset for both global and reaction kernels. A non-remapping
168 algorithm does not keep track of where active data is, and therefore must
169 launch a grid that covers the entire dataset. When the number of active
170 neutrons drops below about 30% of the initial number, the overhead and
171 memory bandwidth cost of launching these extra threads, which only load a
172 “done” bit and return, is more than the cost of performing the radix sort and
173 edge detection. The majority of the transport iterations occur while there

174 are less than 30% of the initial neutrons left, and remapping references is
175 usually worthwhile.

176 Using the NVIDIA OptiX ray tracing framework was also shown to be an
177 effective way to handle the geometry representation in WARP. OptiX is flex-
178 ible, allows attachment of material and cell number to individual geometric
179 primitives, can perform surface detection with a randomly-distributed and
180 directed dataset, can incorporate the remapping vector created by a radix
181 sort, and is fast enough to be used in WARP. The acceleration structures
182 that OptiX can automatically build over the scene geometry was the initial
183 reason for using it, and it was determined that the BVH builder and tra-
184 verser provide the best performance as does using mesh primitive instancing
185 rather than a transform node approach. The number of objects present in
186 the scenes in reactors is small compared to many rendering scenes, and the
187 SBVH acceleration structure does not perform as well. This is presumably
188 due to some additional overhead related to traversing the objects that is not
189 offset when few objects (less than a few hundred thousand) are present in
190 the scene. Primitive instance provides better performance since using trans-
191 form nodes requires traversing a deeper geometry tree, which also has more
192 (redundant) data associated to it.

193 As mentioned in Section 1, a forthcoming companion paper is planned
194 that will benchmark WARP against Serpent and MCNP using identical cross
195 section libraries and problem geometries. In the paper, the relative accuracy
196 and speed of WARP will be determined by comparing it to these production-
197 level Monte Carlo neutron transport codes.

198 Acknowledgements

199 This research is based upon work partially supported by the U.S. De-
200 partment of Energy National Nuclear Security Administration under Award
201 Number DENA0000979 through the Nuclear Science and Security Consor-
202 tium: <http://nssc.berkeley.edu>.

203 Disclaimer

204 This report was prepared as an account of work sponsored by an agency
205 of the United States Government. Neither the United States Government
206 nor any agency thereof, nor any of their employees, makes any warranty,
207 express or limited, or assumes any legal liability or responsibility for the

208 accuracy, completeness, or usefulness of any information, apparatus, product,
209 or process disclosed, or represents that its use would not infringe privately
210 owned rights. Reference herein to any specific commercial product, process,
211 or service by trade name, trademark, manufacturer, or otherwise does not
212 necessarily constitute or imply its endorsement, recommendation, or favoring
213 by the United States Government or any agency thereof. The views and
214 opinions of authors expressed herein do not necessarily state or reflect those
215 of the United States Government or any agency thereof.