

Petascale Application Improvement Discovery at Blue Waters: Best Practices in Accelerating Gather/Scatter Operations in Particle-in-Cell Algorithms on GPU

Bei Wang * William Tang *

November 9, 2015

Driven by chip- and system power limitations, the HPC community is moving to the era of multi- and many-core architectures with greatly increased thread and vector parallelism on shared memory processors. However, achieving high performance and architectural efficiency is challenging for gather and scatter operations as modern computer architectures rely on single-instruction multiple-data execution to provide high computing capabilities. In this report, we describe optimization techniques to improve the performance of the gather and scatter operations from particle-in-cell (PIC) algorithms on GPU. This includes considerations in choosing appropriate data layouts and memory storage for particle and mesh data, and techniques to enable coalesced access and spatial and temporal locality. For a representative PIC application, we see a performance improvement up to 3.4x between the un-optimized and optimized versions of the code. We hope that this report can serve as a “cook book” for those who plan to port their applications (featuring gather and scatter operations) on GPU architectures.

1 Gather/Scatter Patterns in Particle-in-Cell Algorithms

The basic particle method has long been a well established approach for simulating the behavior of charged particles interacting with each other through pair-wise electromagnetic forces. Unfortunately, the $O(N^2)$ complexity makes a particle method impractical for plasma simulations using millions of particles per process. Rather than calculating $O(N^2)$ forces, the particle-in-cell (PIC) method employs a grid to calculate the long range electromagnetic forces. This reduces the complexity from $O(N^2)$ to $O(N + M \log M)$, where M is the number of grid points ($M \ll N$). Figure 1 shows the particle interpolation to the nearest neighbors (linear interpolation) on a 2D poloidal grid in a PIC algorithm. Similar data access patterns occur in a wide range of applications, such as astrophysics, molecular dynamics and statistical analysis.

In general, for each particle, the neighboring grid points are unknown until run time. For two particles stored at contiguous memory addresses, their neighboring grid data are likely to be stored at non-contiguous addresses far way from each other. Non-contiguous memory access will cause more cache/page misses and less efficiently utilize the single-instruction multiple-data execution units. In multithreading environment, simultaneously writing to the same memory address will lead to data hazard, which must be guarded by synchronization mechanisms in order to guarantee correct results. It is well known that synchronization is expensive as it serializes the instructions.

*Princeton Institute of Computational Science and Engineering, Princeton University, Princeton, NJ 08540 USA

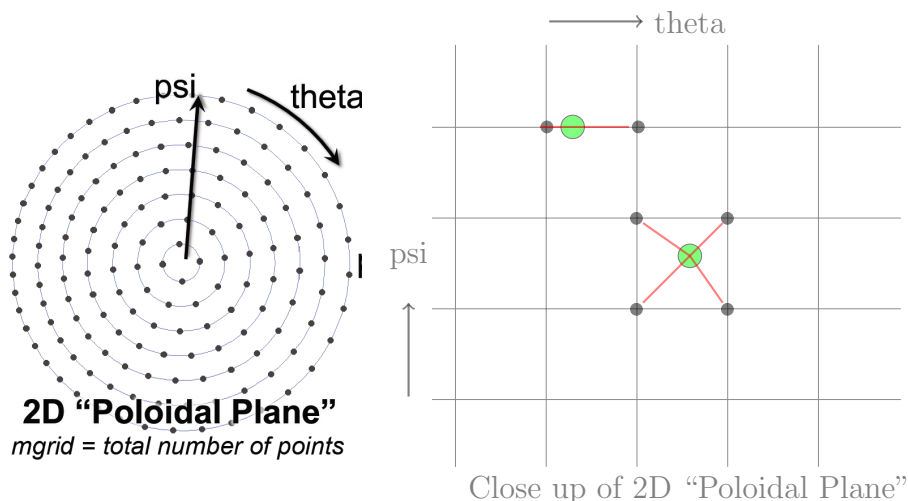


Figure 1: Particle mesh interpolation on a 2D poloidal grid.

```

1| #if USE_TEXTURE
2| texture<int2, 1, cudaReadModeElementType> evectorTexRef;
3| int size = mgrid*3*(mzeta+1); // mgrid*(mzeta+1) is the number of grid points in a 3D torus segment
4| cudaBindTexture(0, evectorTexRef, (void *)(&evector), size*sizeof(real));
5| static __inline__ __device__ real fetch_evector(int i)
6| {
7|     int2 e = tex1Dfetch(evectorTexRef,i);
8|     return __hiloint2double(e.y,e.x);
9| }
10| #define EVECTOR(i) (fetch_evector(i))
11| #else
12| #define EVECTOR(i) (evector[i])
13| #endif

```

Figure 2: Example: bind double precision *evector* array to texture memory

GPU architectures exhibit multiple distinguishing computational features including high-bandwidth memory, limited coherence support, low-overhead scheduling and high-throughput vector computing. In this report, we share some of the practices we found effective to improve the GPU performance of a PIC code. We start with choosing the appropriate data layouts for grid-based data and particle-based data. Then we use techniques such as *cooperative threading* and *binning* to exploit locality. Finally, we found it is more efficient to put the read only grid-based fields on GPU’s texture memory.

2 Choosing Data Layout and Storage

Two of the most common used data layouts in high performance applications are array-of-structures (AoS) and structure-of-arrays (SoA). AoS is generally good for better cache behavior and is also more meaningful to represent physical variables, e.g., particles. However, AoS does not lead to streaming (unit) data access and SIMD instructions. In addition, AoS may cause less efficient utilization of the cached data. For example, reading a single variable from a particle array will result in unnecessary memory transaction. SoA can often lead to streaming access and more efficient utilization of the cached data. However, SoA can lead to poor cache behavior because of the memory distance in accessing multiple struct members of a particle element.

A particle-in-cell code involves two types of data structures, grid-based and particle-based. The grid-based data includes a set of one-dimensional and multi-dimensional arrays, for example, *density* (1D array to store charge density) and *evector* (3D array to store electric field). The particle-based data includes a multi-dimensional array, *zion*, which stores particle’s phase space position and weight. Associated with the grid-based and particle-based data structures are two types of operations in a PIC algorithm, grid-based one such as the kernel to solve Poisson’s equation, and particle-based one such as charge deposition (scatter) and field interpolation (gather). The computational workload is proportional to the number of grid points and the number of particles for these two operations. Since the number of particles is several order of magnitude larger than the number of grid points, the particle-based kernels are usually the hot spot and consists of 80%- 90% of the total wall-clock time.

Before choosing the appropriate data layouts for field and particles, we will first examine memory access patterns of these two arrays. In the representative PIC code, GTC-P [5, 3, 4], we offload the particle-based kernels such as charge deposition and field interpolation to GPU, while leave the grid-based kernels on CPU. In this case, the multidimensional field variable *evector* will be read only, but features random access pattern (gather). The multidimensional particle array will be read and written, but in a pattern of unit or stride access depending on the choice of data layout.

GPU has a special cache called texture memory. The read-only texture memory is optimized for 2D spatial locality. We thus choose AoS for *evector* and bind it to texture memory. Note the 2D *evector* data is fattened in 1D. This will significantly improve the performance of gather operation in field interpolation. Figure 2 shows a code example of how to bind variable *evector* to texture memory. To enable coalesced data access to particles, *zion* is organized as SoA data layout. On GPU with compute capabilities 3.5 (or newer) architecture, we can utilize texture memory simpler with “_ldg or const restrict”.

3 Optimizing Gather/Scatter Patterns

In multithreading environment, two particles operating by two different threads may access the same grid point and cause a read-after-write data hazard. Thus, the data access must be either guarded with a synchronization mechanism or redirected to a private copy depending on the availability of memory resources.

3.1 Private Copy Strategy

One often used approach for GPU based PIC algorithm based on Cartesian grid is to group particles to supercells [2]. Each supercell consists of multiple number of grid points. Then we assign one GPU block to a specified supercell, within which the thread-based data conflict is resolved through fast shared memory atomic operations. The charge in shared memory is added to global memory at once at the end. Given the extensive use of shared memory, we can configure the GPU to favor shared memory.

However, applying the private copy strategy to PIC algorithm based on unstructured grid in magnetic field-line coordinate is challenging. Keeping in mind that particle grouping is based on the particle position. In magnetic field-line coordinate, the final interpolation grid points is not necessary within the associated supercell due to field-line twisting along the toroidal coordinate (see line 35 and 36 in Figure 3a). As such, large number of ghost cells may needed for each supercell. Additionally, it is difficult to applying the private copy strategy to gyrokinetic PIC algorithm, where each particle is not a single particle,

Processor	Freq. (GHz)	D\$/core (KB)	Cores/CPU	D\$/CPU (MB)	CPUs/Node	DP Gflop/s	STREAM (GB/s)	Memory GB
AMD Opteron	2.3	16+2048	8	8	2	140.8	31	32
NVIDIA	0.733	64	14	1.5	1	1314	171	6

Table 1: Blue Waters XK7 nodes configuration

Kernels	Times (ms)			
	baseline	binning	binning+cooperative threading	binning+texture memory
Charge	51.0	51.0	21.5	-
Push	126.9	83.4	-	37.3

Table 2: Performance comparison of two GPU accelerated kernels with different optimization options.

but represents a gyro-ring of varying radius. Since the interpolation is not from a particle to the nearest neighboring grid points, but from multiple points on the gyro-ring of the particle to their neighboring grid points, up to 8 ghost cells in each dimension is required.

The introduction of the NVIDIA Kepler architecture (employed in Blue Waters) enables fast double-precision atomic increment (implemented via a compare-and-swap operation). As such, we develop a cooperative threading algorithm that uses global atomic write for scatter operations as described in details below. The cooperative threading algorithm can be considered as an alternative for dealing with data hazard on unstructured grid with less organized data.

3.2 Synchronization Strategy

Figure 3a shows the cooperative threading algorithm (highlighted) that uses shared memory for data transpose. Specifically, computation is organized so that CUDA threads access successive particles when reading, but is reorganized so that threads work together on one charge deposition (of up to 32 neighboring grid points) when writing. The reorganization is through transposing the data write in shared memory. Given the limited use of shared memory in this case, the GPU is configured to favor the L1 cache.

As a comparison, the implementation that used global atomic write only is also shown in Figure 3a. For gather operation, the analysis indicates that a simple decomposition of one particle per CUDA thread and 64 threads per block perform well. The key optimization for gather operation is to place the grid-based array, such *evector* on texture memory. Additionally, some variables are also judiciously placed in CUDA’s constant memory.

3.3 Data Locality with Binning

Particle binning is often used for PIC algorithm to improve performance by providing better data locality. Here we rely on the CUDA Thrust library to bin the particles according to their radial dimension positions.

4 Performance Evaluation

Blue Waters [1] is a Cray XE6/XK7 MPP at the the National Center for Supercomputing Applications. It contains 22,500 XE6 Interlagos-based nodes augmented by more than 4200 XK7 hybrid nodes containing NVIDIA GK110 Kepler accelerator and connected by Cray’s Gemini network into a 3D torus. Each XK7 node includes two 8-core processors running at 2.3GHz in which each core has a private 16KB L1 data cache, each pair of cores shares a 256b-wide SIMD unit, each pair of cores shares a 2MB L2 cache, and all cores share a 32MB L3 cache. Each processor is connected to 16GB of DDR3-1333 DRAM via a 128-bit memory controller. As such, each compute node contains two NUMA nodes and programmers are motivated to run 2 processes of 8 threads per node. The Kepler GK110 includes 14 streaming multiprocessor units (SMX). Each SMX has 64KB on-chip memory that can be configure as 48/16KB shared memory and 16/48KB L1 cache and 1536KB L2 on-chip cache memory. Currently, the NVIDIA accelerator is not directly interact with the Gemini interconnect, so data has to be moved to GPU or accessed as mapped memory. Table 1 lists the system configuration of blue waters XK7 nodes.

We have evaluated the performance of charge deposition and field interpolation operations on GPU with different optimization techniques on one Kepler processor. The total number of grid points in 2D poloidal plane and the total number of particles are 39,693 and 3,959,300, respectively. Table 2 shows the wall-clock time spending on charge deposition (charge) and field interpolation (push) for one time step. Comparing with the baseline performance, optimization techniques such as binning, cooperative threading and texture memory lead to 2.4x and 3.4x speed up for charge and push kernels, respectively.

Acknowledgments

This report is generated as supported by the Petascale Application Improvement Discovery (PAID) program at Blue Waters.

```

1 const int tid = threadIdx.x; const int bid = blockIdx.x;
2 const int nblocks = gridDim.x; const int nthreads = blockDim.x;
3 const int gid = tid+bid*nthreads; const int np = nblocks * nthreads;
4
5 #if COOP
6 extern __shared__ int shared_buffer[];
7 int *update_idx = shared_buffer;
8 double *update_val = (real *)shared_buffer[nthreads*4];
9 #endif
10 ...
11
12 for (int m=gid; m<me; m+=np){
13     double psitmp = z0[m];
14     double thetatmp = z1[m];
15     double zetatmp = z2[m];
16     double weight = z4[m];
17
18     double wzt = (zetatmp-zetamin)*delz;
19     int kk = abs_min_int(mzeta-1,(int)wzt);
20
21     double wz1 = weight*(wzt-(double)kk);
22     double wz0 = weight-wz1;
23
24     double r_diff = psitmp-a0;
25     double rdum = delr*abs_min_real(a_diff,r_diff);
26     int ii = abs_min_int(mpsi_max,(int)rdum);
27
28     double wp1 = rdum-(real)ii;
29     double wp0 = 1.0-wp1;
30
31     double tflr = thetatmp;
32     int im = ii;
33     int im2 = ii + 1;
34
35     double tdumtmp = pi2_inv*(tflr-zetatmp*qtinv[im]);
36     double tdumtmp2 = pi2_inv*(tflr-zetatmp*qtinv[im2]);
37     double tdum = (tdumtmp-(int)tdumtmp)*delt[im];
38     double tdum2 = (tdumtmp2-(int)tdumtmp2)*delt[im2];
39
40     int j00 = abs_min_int(mtheta[im]-1,(int)tdum);
41     int j01 = abs_min_int(mtheta[im2]-1,(int)tdum2);
42     int jtelectron0tmp = igrid[im]+j00;
43     int jtelectronitmp = igrid[im2]+j01;
44
45     double wtelectron0tmp = tdum-(double)j00;
46     double wtelectronitmp = tdum2-(double)j01;
47
48     double wt10 = wp0*wtelectron0tmp;
49     double wt00 = wp0-wt10;
50     double wt11 = wp1*wtelectronitmp;
51     double wt01 = wp1-wt11;
52
53     int ij1 = kk+(mzeta+1)*jtelectron0tmp;
54     int ij2 = kk+(mzeta+1)*jtelectronitmp;
55
56     /* interpolate from particle to grid - scatter operation */
57     #if COOP // use shared memory for optimizing scatter operation
58     update_idx[4*tid]=ij1;
59     update_idx[4*tid+1]=ij1+1;
60     update_idx[4*tid+2]=ij1+mzeta+1;
61     update_idx[4*tid+3]=ij1+mzeta+2;
62     update_val[4*tid]= wz0*wt00;
63     update_val[4*tid+1]= wz1*wt00;
64     update_val[4*tid+2]= wz0*wt10;
65     update_val[4*tid+3]= wz1*wt10;
66     __syncthreads();
67
68     atomicDPupdate(density+update_idx[tid], update_val[tid]);
69     atomicDPupdate(density+update_idx[stride*tid], update_val[stride*tid]);
70     atomicDPupdate(density+update_idx[2*stride*tid], update_val[2*stride*tid]);
71     atomicDPupdate(density+update_idx[3*stride*tid], update_val[3*stride*tid]);
72
73     update_idx[4*tid]=ij2;
74     update_idx[4*tid+1]=ij2+1;
75     update_idx[4*tid+2]=ij2+mzeta+1;
76     update_idx[4*tid+3]=ij2+mzeta+2;
77     update_val[4*tid]= wz0*wt01;
78     update_val[4*tid+1]= wz1*wt01;
79     update_val[4*tid+2]= wz0*wt11;
80     update_val[4*tid+3]= wz1*wt11;
81     __syncthreads();
82
83     atomicDPupdate(density+update_idx[tid], update_val[tid]);
84     atomicDPupdate(density+update_idx[stride*tid], update_val[stride*tid]);
85     atomicDPupdate(density+update_idx[2*stride*tid], update_val[2*stride*tid]);
86     atomicDPupdate(density+update_idx[3*stride*tid], update_val[3*stride*tid]);
87 #else // use global atomic for scatter operation
88     atomicDPupdate(density+ij1, wz0*wt00);
89     atomicDPupdate(density+ij1+1, wz1*wt00);
90     atomicDPupdate(density+ij1+mzeta+1, wz0*wt10);
91     atomicDPupdate(density+ij1+mzeta+2, wz1*wt10);
92
93     atomicDPupdate(density+ij2, wz0*wt01);
94     atomicDPupdate(density+ij2+1, wz1*wt01);
95     atomicDPupdate(density+ij2+mzeta+1, wz0*wt11);
96     atomicDPupdate(density+ij2+mzeta+2, wz1*wt11);
97 #endif
98 }

```

(a) GPU charge deposition (scatter operation)

```

1 const int tid = threadIdx.x; const int bid = blockIdx.x;
2 const int nblocks = gridDim.x; const int nthreads = blockDim.x;
3 const int gid = tid + bid*nthreads; const int np = nblocks * nthreads;
4 ...
5
6 for (int m=gid; m<me; m+=np){
7     double psitmp = z0[m];
8     double thetatmp = z1[m];
9     double zetatmp = z2[m];
10
11     double wzt = (zetatmp-zetamin)*delz;
12     int kk = d_abs_min_int(mzeta-1, (int)wzt);
13
14     double wz1 = wzt - (real) kk;
15     double wz0 = 1.0 - wzt;
16
17     double r_diff = psitmp-a0;
18     double rdum = delr * d_abs_min_real(a_diff, r_diff);
19     int ii = d_abs_min_int(mpsi_max, (int) rdum);
20
21     double wp1 = rdum - (real) ii;
22     double wp0 = 1.0 - wp1;
23
24     double tflr = thetatmp;
25     int im = ii;
26     int im2 = ii + 1;
27
28     double tdumtmp = pi2_inv * (tflr - zetatmp * qtinv[im]);
29     double tdumtmp2 = pi2_inv * (tflr - zetatmp * qtinv[im2]);
30     double tdum = (tdumtmp - (int) tdumtmp) * delt[im];
31     double tdum2 = (tdumtmp2 - (int) tdumtmp2) * delt[im2];
32
33     int j00 = d_abs_min_int(mtheta[im]-1, (int) tdum);
34     int j01 = d_abs_min_int(mtheta[im2]-1, (int) tdum2);
35     int jtelectron0tmp = igrid[im] + j00;
36     int jtelectronitmp = igrid[im2] + j01;
37
38     double wtelectron0tmp = tdum - (real) j00;
39     double wtelectronitmp = tdum2 - (real) j01;
40
41     double wt10 = wtelectron0tmp;
42     double wt11 = wtelectronitmp;
43     double wt01 = 1.0 - wt11;
44     double wt00 = 1.0 - wt10;
45
46     int ij1 = kk+(mzeta+1)*jtelectron0tmp;
47     int ij2 = kk+(mzeta+1)*jtelectronitmp;
48
49     int idx1 = 6*ij1;
50     int idx2 = 6*ij2;
51
52     double e1=0.0;
53     double e2=0.0;
54     double e3=0.0;
55
56     /* interpolate from grid to particle - gather operation */
57     e1 =e1+wp0*wt00*(wz0*EVECTOR(idx1+0)+wz1*EVECTOR(idx1+3));
58     e2 =e2+wp0*wt00*(wz0*EVECTOR(idx1+1)+wz1*EVECTOR(idx1+4));
59     e3 =e3+wp0*wt00*(wz0*EVECTOR(idx1+2)+wz1*EVECTOR(idx1+5));
60
61     e1 =e1+wp0*wt10*(wz0*EVECTOR(idx1+6+0)+wz1*EVECTOR(idx1+6+3));
62     e2 =e2+wp0*wt10*(wz0*EVECTOR(idx1+6+1)+wz1*EVECTOR(idx1+6+4));
63     e3 =e3+wp0*wt10*(wz0*EVECTOR(idx1+6+2)+wz1*EVECTOR(idx1+6+5));
64
65     e1 =e1+wp1*wt01*(wz0*EVECTOR(idx2+0)+wz1*EVECTOR(idx2+3));
66     e2 =e2+wp1*wt01*(wz0*EVECTOR(idx2+1)+wz1*EVECTOR(idx2+4));
67     e3 =e3+wp1*wt01*(wz0*EVECTOR(idx2+2)+wz1*EVECTOR(idx2+5));
68
69     e1 =e1+wp1*wt11*(wz0*EVECTOR(idx2+6+0)+wz1*EVECTOR(idx2+6+3));
70     e2 =e2+wp1*wt11*(wz0*EVECTOR(idx2+6+1)+wz1*EVECTOR(idx2+6+4));
71     e3 =e3+wp1*wt11*(wz0*EVECTOR(idx2+6+2)+wz1*EVECTOR(idx2+6+5));
72
73     wpi0[m] = e1;
74     wpi1[m] = e2;
75     wpi2[m] = e3;
76
77 }
78

```

(b) GPU field interpolation (gather operation). The code example of using texture memory for vector is given in Figure 2.

References

- [1] Blue Waters Cray XE6. <http://www.ncsa.illinois.edu/enabling/bluewaters>.
- [2] V. K. Decyk and T. V. Singh. Particle-in-cell algorithms for emerging computer architectures. *Computer Physics Communications*, 185(3):708 – 719, 2014.
- [3] K. Z. Ibrahim, K. Madduri, S. Williams, B. Wang, S. Ethier, and L. Oliker. Analysis and optimization of gyrokinetic toroidal simulations on homogenous and heterogenous platforms. *International Journal of High Performance Computing Applications*, 2013.
- [4] K. Madduri, K. Z. Ibrahim, S. Williams, E.-J. Im, S. Ethier, J. Shalf, and L. Oliker. Gyrokinetic toroidal simulations on leading multi- and manycore HPC systems. In *Proc. Int'l. Conf. for High Performance Computing, Networking, Storage and Analysis (SC '11)*, pages 23:1–23:12, New York, NY, USA, 2011. ACM.
- [5] B. Wang, S. Ethier, W. Tang, T. Williams, K. Z. Ibrahim, K. Madduri, S. Williams, and L. Oliker. Kinetic turbulence simulations at extreme scale on leadership-class systems. In *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC '13*, pages 82:1–82:12, New York, NY, USA, 2013. ACM.