

A GPU-based Large-scale Monte Carlo Simulation Method for Systems with Long-range Interactions

Yihao Liang, Xiangjun Xing

*Institute of Natural Sciences and Department of Physics and Astronomy, Shanghai Jiao Tong
University, Shanghai, 200240 China*

Yaohang Li

Department of Computer Science, Old Dominion University, Norfolk, VA 23529, United States

Abstract

In this work we present an efficient implementation of Canonical Monte Carlo simulation for Coulomb many body systems on graphics processing units (GPU). Our method takes advantage of the GPU Single Instruction, Multiple Data (SIMD) architectures, and adopts the sequential updating scheme of Metropolis algorithm. It makes *no approximation* in the computation of energy, and reaches a remarkable 440-fold speedup, compared with the serial implementation on CPU. We further use this method to simulate primitive model electrolytes, and measure very precisely all ion-ion pair correlation functions at high concentrations. From these data, we extract the renormalized Debye length, renormalized valences of constituent ions, and renormalized dielectric constants. These results demonstrate unequivocally physics beyond the classical Poisson-Boltzmann theory.

Keywords: Monte Carlo, GPU, Parallel Computing, Coulomb Many Body Systems, Electrolytes, Charge Renormalization

1. Introduction

Molecular simulations generally fall into two categories: molecular dynamics (MD) and Monte Carlo (MC). In a MD simulation, one solves the Newtonian equation, from which both dynamical and static properties of studied systems can be extracted. In a Monte Carlo simulation, one carries out a Markovian stochastic process which converges to the equilibrium Gibbs distribution. The main advantage of MC is that it can often be accelerated substantially by performing unphysical moves that involve long displacement and/or large number of particles. Hence MC is usually more efficient than MD for simulation of equilibrium systems.

Long range interactions impose substantial difficulties on numerical simulations, because the computational complexity for one cycle (where every particle moves one step on average) scales as N^2 , in contrast with N for short range interactions, where N is the size of system being simulated. This severely limits the size of feasible simulations. There are three classes of methods to speed up the simulation for long range interacting systems: 1) Multipole expansion methods [1, 2, 3, 4], where interactions are computed approximately using truncated multipole expansions. This reduces the computational complexity to N or $N \log N$. 2) Fourier transforms augmented by Ewald-summations, which reduce the complexity to $N^{3/2}$. It can be further reduced to $N \log N$, by using Fast Fourier Transform. Unfortunately, the latter trick is not applicable for MC. Furthermore, additional artifacts arise due to periodic images. For a discussion of these artifacts, see reference [5]. 3) Multi-scale reaction-potential methods [5, 6, 7], whose idea is to simulate only a small portion of the system and use continuum theory to describe the remaining. These methods have difficulty to scale to system with large number of particles, since they do not use any speedup technique in the computation of energy.

Graphics processing units (GPU) offer a new possibility for speeding up large scale simulation of long range interacting systems *without sacrificing accuracy*. GPU is a powerful device which can process thousands of threads simultaneously with high memory bandwidth. Compared to CPU, GPU is designed with more transistors that are devoted to data processing rather than data caching and flow control [8]. It is suitable for computation-intensive and data-parallel computations such as graphics rendering - the original purpose of designing GPU. In recent years, the GPU is devoted to more and more general purpose fields, such as data mining, machine learning, finance, scientific computing and molecular simulation.

Many MD simulation methods have already been adapted to GPU in the past years [9, 10, 11, 12, 13, 14, 15, 16]. There are also many MD softwares and libraries that can be implemented on GPU, including AMBER [17, 18, 19, 20] FENZI [21, 22] LAMMPS [23, 24] NAMD [25, 26], HALMD [27], OpenMM [28] HOOMD-blue [29, 30], GROMACS [31], ACEMD[32]. Most of these packages have demonstrated efficiency in simulating long range interacting systems.

Implementation of Monte Carlo simulation on GPU turns out to be significantly harder. This is mainly because of the sequential nature of Monte Carlo stochastic dynamics, where particles are moved one by one. Nonetheless, there have been a few attempts to realize MC simulation of Ising model [33, 34, 35] and Hard disk fluid [36] on GPU. A. Yaseen and Y. Li used the remapping method to calculate the total energy on GPU for protein systems [37]. A group at Wayne State University [39, 40] realized a GPU code for Gibbs ensemble MC simulation of simple liquids. J. Kim *et. al.* developed an implementation with embarrassing parallel on GPU,

where each block performs an individual Monte Carlo simulation [41, 42, 43]. We have not found any previous realization of MC simulation for large-scale long range interacting systems on GPU.

In this work, we develop an efficient GPU approach to realize the canonical Monte Carlo of systems with long range interactions. The fundamental idea behind this method is to update every particle once during one invoking of GPU kernel, and to use one (or a few) thread(s) to control one particle in a synchronous mode with coalesced memory access, i.e., to calculate its energy change and to attempt to move it. The same idea has been used in MD simulations in [9], where each thread controls the evolution of one particle.¹ In our case of MC simulation, each thread also has to take care of Monte Carlo trials and decisions. In terms of the interaction table, this parallel metropolis scheme looks like a “brush”, hence we call it *the Brush Metropolis Algorithm*.

This approach enhances temporal locality and thereby improves cache performance. We benchmark this code on a Tesla K20 GPU and find a remarkable 440-fold speedup compared to sequential codes on Intel Xeon E5-2670. It is important to stress that this speedup is achieved *without* sacrifice of accuracy, since there is no approximation (such as truncated multipole expansion) used in the Brush Metropolis Algorithm. Using this program, we carry out large-scale Monte Carlo simulations of primitive model electrolytes containing as many as 10^6 ions and measure all pair correlation functions up to extremely high precision. The radius of the simulation box is hundreds of Debye length, so that all boundary artifacts are completely screened. Using this huge amount of data, we are able to measure precisely static linear response properties of the system, including the renormalized valences of constituent ions, the renormalized Debye lengths and renormalized dielectric constant. Comparison of these renormalized parameters with their bare values clearly demonstrates that the statistical physics of concentrated electrolytes is beyond the classical Poisson-Boltzmann theory.

The remaining of this paper is organized as follows: In section 2, we compare the feasibility of parallelism in random moving and sequential moving. In section 3, we discuss GPU implementation of sequential moving. We show the benchmark

¹However, there are also important differences between these two algorithms. In MD, the only dependency is that the particle coordinates updates and the force computations should be fairly separated in time. To satisfy this dependency, two basic kernels are needed: one is for position updating and the other one is for force computation. Moreover, the force computations are independent so that they can be performed without any inter-thread communication. In comparison, the Hamiltonian change in our MC method is more complicated, due to complex and strong dependency in updating coordinates and energy calculation. In our implementation for MC, the updates of coordinates and the corresponding energy computation are carried out within one kernel and the order of computations and updates together with inter-thread communications are deliberately designed.

results and present simulation results for linear response properties of electrolytes in sections 4 and 5, respectively. Section 6 summarizes our conclusions and future research directions.

2. Concurrency

In a Metropolis Monte Carlo simulation, a Monte Carlo step is the smallest unit of Markov chain where one particle is moved. A Monte Carlo cycle consists of N steps, where N is the total number of particles. A MC step contains three basic sub-steps:

- a) **Selection:** Select a particle k , either randomly or sequentially.
- b) **Trial:** Propose an unbiased perturbation of the selected particle. The new coordinate of the selected particle $\tilde{\mathbf{x}}_k$ is generated by a symmetric probability transition function $T(\mathbf{x}_k, \tilde{\mathbf{x}}_k) = T(\tilde{\mathbf{x}}_k, \mathbf{x}_k)$, which yields conditional probability density that $\tilde{\mathbf{x}}_k$ is selected as the new coordinate, given the current coordinate is \mathbf{x}_k . We then calculate the change of total energy ΔE_k due to this trial. For long range interacting system, the time cost for calculation of ΔE_k scales with N , and therefore is the most computationally expensive substep.
- c) **Acceptance/rejection:** Accept the tried state as the next state of the Markov chain with probability $\min\{1, e^{-\beta\Delta E}\}$. The new position is then given by

$$\mathbf{x}'_k = \begin{cases} \tilde{\mathbf{x}}_k, & \text{trial accepted;} \\ \mathbf{x}_k, & \text{trial rejected.} \end{cases} \quad (1)$$

Here $\beta = 1/(k_B T)$, k_B is the Boltzmann constant and T is the temperature of the system. Substep a) (selection) can be executed in two different ways: 1) Random Updating Scheme, where the particle is selected at random, and 2) Sequential Updating Scheme, where all particles are labeled and moved in ascending order. In Fig. 1, we schematically illustrate several consecutive MC steps in the random updating scheme (left) and the sequential updating orders (right). All particles are labeled by integers $0, 1, \dots, N - 1$. Each row of figure corresponds to a state in one MC cycle, where the current position of each particle is listed in ascending order. In the random scheme, every particle moves one step *on average* within one MC cycle. By contrast, in sequential updating scheme, every particle attempts to move exactly once within one cycle. In random updating scheme, the Markov process has a time-independent transition matrix, whereas in sequential scheme, the transition matrix is periodic with periodicity N . Note that detailed balance of the Markov process is guaranteed by the symmetry of the transition function $T(\mathbf{x}, \tilde{\mathbf{x}})$, together with the choice of acceptance probability, and therefore is valid for both updating schemes. Because of the Markovian nature of MC simulation, the acceptance ratio of a trial in the random updating scheme depends on all

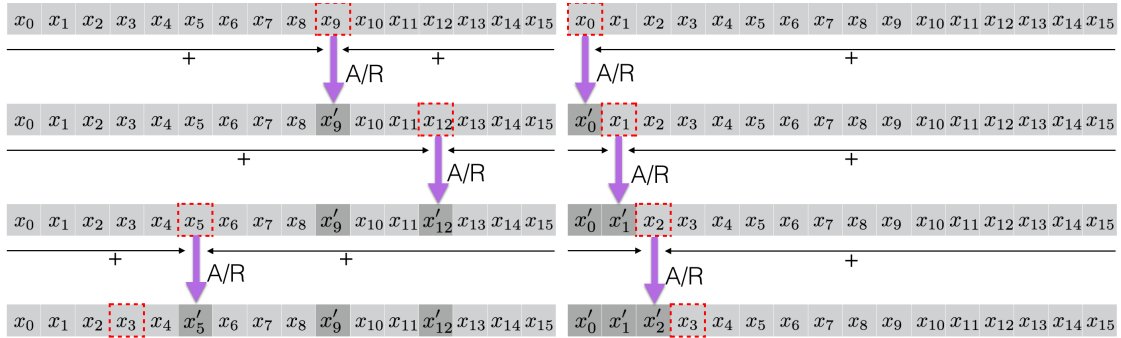


Figure 1: Schematic illustrations of several consecutive MC steps in the random updating scheme (left) and the sequential updating orders (right). Each row denotes a state of the system, and red dashed boxes denote the particles being updated, one in each step. The purple arrows denote the actions of updating, with A/R meaning accept and reject, respectively. In order to determine the results of updates (either acceptance or rejection), one needs to sum of the changes of all pairwise interactions. These are designated by black horizontal arrows with underlying plus signs.

details of all prior steps. This severely limits the potential of concurrency of the Random Updating Scheme. By contrast, the potential of concurrency of Sequential Updating Scheme is much higher, as we shall show in detail.

Starting from the initial state $\{\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{N-1}\}$, we can now propose N random trial new positions $\{\tilde{\mathbf{x}}_0, \tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_{N-1}\}$, one for each particle, so that the first several states in one MC cycle are as shown in Fig. 1. To illustrate the interdependency of tasks involved in one MC cycle of the sequential updating scheme, we introduce the graphic representations in Fig. 2 of all involved pairwise interactions that must be calculated during a MC cycle. The horizontal axis is the index of particles (in ascending order), whereas the vertical axis is the time line, with the numbers on the axis label particles selected in each step. Each square element (except the dark ones) denotes the change of a pairwise interaction ΔE_{ij} in step i where the i -th particle is being tried to move. It is defined as

$$\Delta E_{ij} \equiv \begin{cases} U(\tilde{\mathbf{x}}_i - \mathbf{x}_j) - U(\mathbf{x}_i - \mathbf{x}_j), & i < j; \\ U(\tilde{\mathbf{x}}_i - \mathbf{x}'_j) - U(\mathbf{x}_i - \mathbf{x}'_j), & i > j. \end{cases} \quad (2)$$

Now the important point is that for $i < j$, ΔE_{ij} does not depend on any of new positions \mathbf{x}'_k . (All these ΔE_{ij} are shown in the upper triangle in Fig. 2.) Therefore they can be calculated simultaneously before any particle is moved. Furthermore, all the squares in one row to the right of diagonal can be summed before any move. Both of these calculations can be done in parallel.

Now the 0-th particle is ready to take the trial move, and determines its new position according to Eq. (1). After this, all elements in the 0-th column below

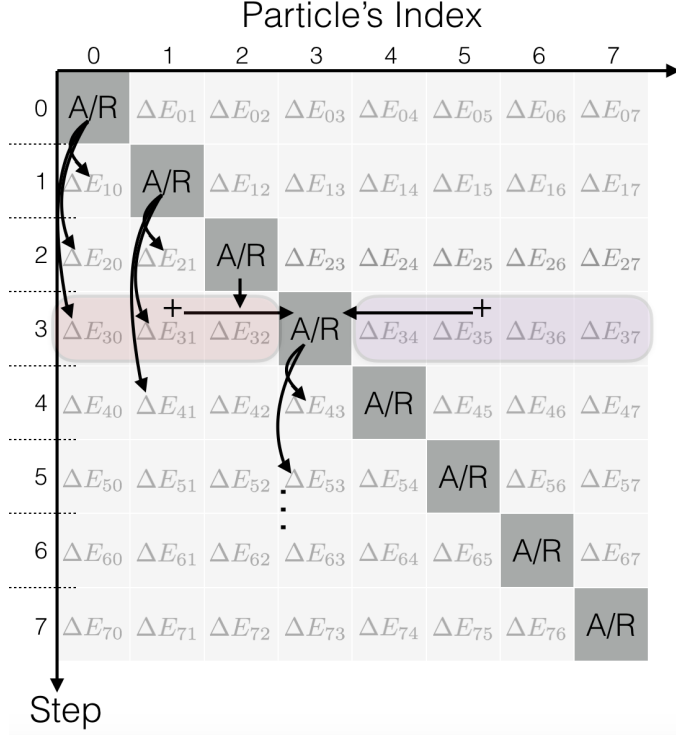


Figure 2: Table of pairwise interactions in sequential updating scheme. Each element (except the dark ones) denotes computation of a pairwise interaction. The dark square (one for each row) denotes the action of decision(acceptance or rejection) for the selected particle in each step. “+” denotes the summation of the pairwise interactions in selected region (red region or purple region). Decision must be done after the summations in each row. After the decision in one column is done, all squares below are to be updated.

the diagonal line can be calculated independently (because these elements depend only on the new position \mathbf{x}'_0 , but not on others). The inter-dependency of these operations is illustrated in Fig. 2. They can also be added to ΔE_i independently. After this, particle 1 is ready to take the trial move. The process keeps going until the last particle N is tried. This is the main idea behind our GPU implementation of the sequential updating scheme.

3. GPU implementation

We implement the sequential updating scheme with the NVIDIA CUDA programming model and perform simulations and benchmarks on Tesla K20 GPUs. CUDA is a programming platform with extensions of C/C++, which provides a convenient way for parallel programming on NVIDIA GPU. An Nvidia Tesla K20 card includes 13 stream multiprocessor(SMX), each of which has 192 CUDA

cores, result in totally 2496 CUDA cores and 1.17 Tflops double-precision peak performance or 3.52 Tflops single-precision peak performance. The size of global memory in a Tesla K20 GPU is 5GB, and in each SMX there is a 48KB shared memory. Details and terminologies of CUDA and GPU architecture can be found in the CUDA-C programming guide [8].

3.1. Data setting and random numbers

We use CUDA’s *float4* data type for particle parameters, where x, y, z components are positions and w stores valence of each particle. The information of particles makes up an N -elements array, which is stored in the device’s global memory. In the following text, we denote this array by X_{dev} . Using *float4* data type with aligned access allows coalesced memory access to the arrays of data in device memory, resulting in efficient memory requests and transfers [8]. All the tasks in MC moves are performed in GPU, except for generation of random numbers.

In our program, random numbers are generated by the code “ran4” described in the book *Numerical recipes* [44]. This code generates single-precision floating-point numbers which are uniformly distributed in the interval $[0, 1)$. We use two N -elements *float4* arrays to store the random numbers used in an MC cycle, one of which is on the host and the other is on the device. The array of random numbers on device is denoted by “rnum” in the following text and pseudo-code. The x, y, z components of an element in “rnum” are used to generate new coordinates whereas w is for decision(acceptance or rejection) in each trial. The computational cost of random number generations can be effectively masked by MC cycles carried out on GPU. This is due to the fact that during the GPU performing a MC cycle, the random numbers for the next MC cycle can be generated simultaneously on CPU supported by concurrent execution between CPU and GPU in CUDA. We use CUDA’s built-in function “cudaMemcpy” to upload the random numbers to the device. Function “cudaMemcpy” contains an intrinsic synchronization so that the system copies data only after the previous GPU job in the same stream finishes. Therefore there is no worry that the random number array changes while the GPU kernel is still using it.

Data sampling and analysis can be either on the device or on the host. If we want to put some analysis tasks in CPU cores, we can download this array by built-in functions such as “cudaMemcpy” from device as well.

3.2. Decomposition

When the host program invokes a kernel function of MC cycle, GPU constructs a one-dimensional grid with B thread blocks. Each thread block is one-dimensional and contains $S = \lceil N/B \rceil$ threads. In the following text, we use a two-element tuple

(bid, tid) to identify each thread, where bid shows which block the thread locates in and tid is the relative identity of this thread within the block bid .

Similarly, the array of particles is divided into B groups, each with S particles. The trial of p 'th particle in group g is assigned to the thread $tid = p$ in block $bid = g$. In this article, we call the assigned particle the host-particle of corresponding thread. The mission of each thread is to compute the change of energy due to the trial of its host-particle and then decide the acceptance of the new move. At the beginning of kernel grid, each thread loads the information of its host-particle and random number vector from the global memory to its own registers. Then it computes the new position of its host-particle. The total energy change ΔE_{tid} of the host-particle is stored as a double floating-point variable in thread's registers, which can be initialized with 0 or the change of its self-energy (if any). The pseudo-code of initialization and the statements of symbols are shown in Algorithm 1.

3.3. Energy computation

Despite the positions are represented by single floating-point data type, all the computations related to the energy should be performed by double floating-point operations. The reason is that the energy is a summation of billions of terms which can be either positive or negative in the charged system. Single float treatment in such summation amplifies the truncation error significantly.

If there are hard core repulsions between particles, energy may be infinite. To avoid this problem, we can set an individual variable to indicate overlap. All the updating scheme of this indicator is similar to the computation of energy. In the CPU-based sequential program, when the overlap is encountered, we can stop the energy computation and reject the trial directly. This pre-rejection strategy can improve the speed of simulation on CPU, especially in the high volume fraction system. However, this is not a GPU-friendly strategy, due to the fact that this strategy leads to a long divergence path among different MC trials, which offends the SIMD computing scheme of GPU [45] and degrades the parallel efficiency. In this work, we do not use the pre-rejection strategy for two reasons. Firstly, the purpose of this work is to give a general framework on how to parallelize the Metropolis Monte Carlo simulation, where the hard-core repulsion is just a specific choice to deal with the short-range interaction. Secondly, in almost all the systems we study in this work, the volume fraction of particles is low, so that the pre-rejection doesn't make significant reduction of the computation, and the probability that all threads within a warp be pre-rejected is very low. As a result, the naive pre-rejection strategy on the GPU implementation cannot improve the performance.

Algorithm 1 Statement and initialization

Require: bid : the index of block

Require: tid : the index of thread in its block

Require: B : number of groups and the number of blocks

Require: S : number of particles(threads) per group(block)

Require: X_{dev} : Particles' list on the global memory

Require: $rnum$: Random number sequence on the global memory

Require: D : Range of one moving

Require: ΔE : Total change of energy due to the trial of host-particle. It is double precision and located on each thread's registers.

Require: $BlockState$: An array with B elements on global memory which shows the state of each block. The initial value of each element is 0. The element is 1 if the corresponding block terminated and the coordinates of its host particle group updated.

Require: X : a float4 vector whose x,y,z components are old coordinates of the host-particle and w component is the valence

Require: X' : a float4 vector whose x,y,z components are new coordinates of the host-particle and w is the random number for decision

Require: Y : a float4 array with B elements on the shared memory, which stores the coordinates of guest particles

Require: λ : Bjerrum length

$$\begin{aligned} X &\leftarrow X_{dev}[bid*S+tid] \\ X' &\leftarrow rnum[bid*S+tid] \\ X'.x &\leftarrow (X'.x - 0.5) * D \\ X'.y &\leftarrow (X'.y - 0.5) * D \\ X'.z &\leftarrow (X'.z - 0.5) * D \\ \Delta E &\leftarrow \text{Change of self energy} \end{aligned}$$

3.4. Intergroup Calculation

When threads in a block compute interactions with group J where $J \neq bid$ (particles in J are called the **guest particles** and J is called the **guest group**), they perform a brush-like operation which is widely used in MD implementation [9]. Threads in this block firstly load the information of J from the global memory to an S -element float4 array Y on the shared memory. This loading procedure is one-to-one and aligned, so that we can enhance the cache hit rate and reduce the transaction of global memory.

After group J loaded (here a synchronize instruction “_syncthread” is taken to ensure that all the threads in this block finish loading data of guest-group J), all threads within this thread block take a loop to calculate the change of pairwise interactions and add them up, as showed in Fig. 3.

To avoid the bank conflict and take advantage of broadcasting feature of shared memory, all the loops in this block start at the same index. In this way, threads in the same warp access the same address simultaneously so that the shared memory can broadcast data. We choose 64-bit mode to make full use of the band-width of the shared memory.

At last, the block synchronize instruction “_syncthread” is taken again to ensure that all threads finish processing the group J before the next group are loaded. Fig. 3 is the schematic illustration of this part. The pseudo-code of this part is shown in algorithm 2.

Algorithm 2 Energy changes with group J ($J \neq bid$)

```
Y[tid] ←  $X_{dev}[J*S+tid]$ ;  
_syncthread();  
for  $l := 0$  to  $S - 1$  do  
     $\Delta E \leftarrow \Delta E + \Delta U(X, X', Y[l])$ ;  
end for  
_syncthread();
```

3.5. Self Calculation

When thread block computes the interactions within its host-group, the procedure involves two main tasks of computations: (1). Compute the upper triangle parts in Fig.4. (2). Decide the trial and calculate the lower triangle parts as in Fig.5. Here we describe the details of these two tasks respectively.

3.5.1. Upper triangle part

The upper triangle of the self-interaction table is independent of any other tasks. Therefore we can pre-calculate them by the brush scheme as before. Thread

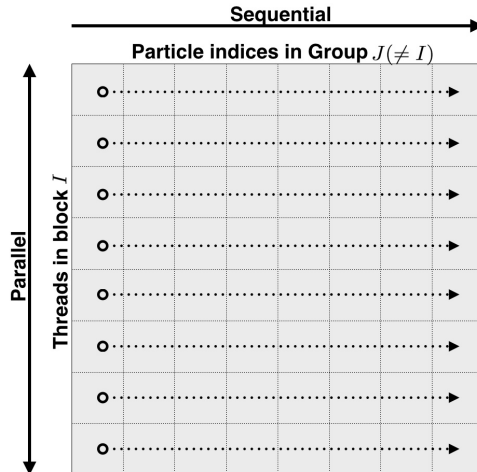


Figure 3: Schematic illustration of intergroup calculation. In this interaction table, the vertical indices indicate threads in thread block I and horizontal indices specify the particles in group $J(\neq I)$. Each element stands for a pairwise energy change of corresponding particle pairs. The dashed arrows show the direction of loop for summation. Threads in block I compute and accumulate these pairwise interactions along horizontal direction in sequential order, whereas the energy change on each particle in group I is evaluated in parallel.

$tid = p$ computes the change of energy with particles $p + 1, p + 2, \dots, S - 1$ in the host-group. To ensure threads within a warp accessing the same address, the for-loop of each thread should start at $\lfloor tid/w \rfloor \times w + 1$. Here w is the size of a warp. In most GPU, $w = 32$. Fig. 4 is an illustration of this procedure. The pseudo-code for this part is shown in algorithm 3.

Algorithm 3 Upper triangle part of self calculation

```

Y[tid] ← X;
__syncthreads();
for l := ⌊tid/w⌋ × w + 1 to S - 1 do
    if l > tid then
        ΔE ← ΔE + ΔU(X, X', Y[l]);
    end if
end for
__syncthreads();

```

In the first w loops, half of the threads on average within a warp are inactive. To reach the full warp efficiency we can use the **map algorithm** [37, 38, 39]. But this strategy increases usage of registers per thread that limits the number of concurrently executed blocks in one streaming multiprocessor(SMX) and thus decreases the SMX's efficiency. On the other hand, the computation of the first

w loops is light compared with the whole tasks. So we do not implement this strategy in our code for Tesla K20.

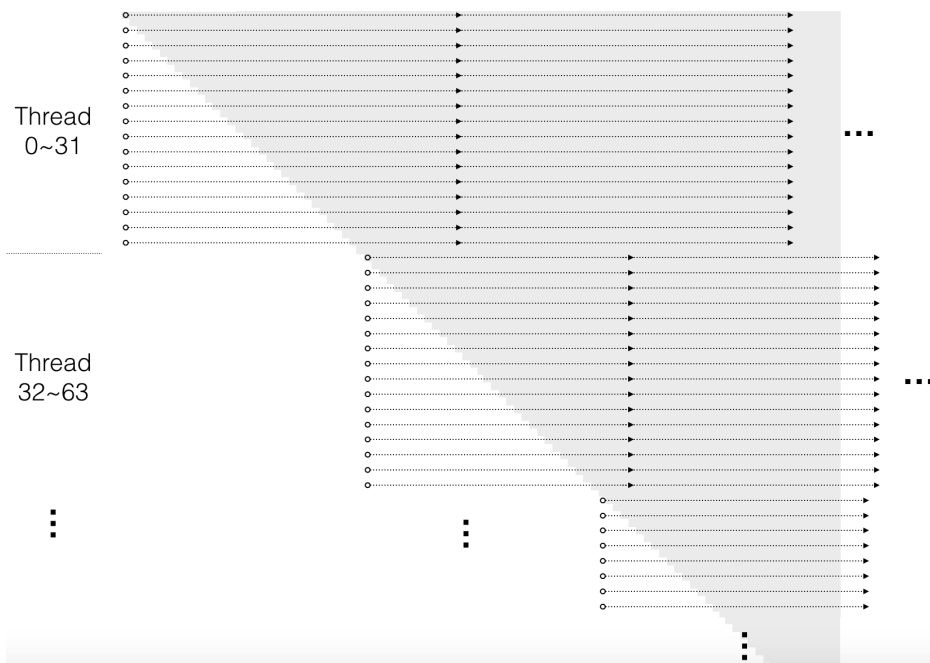


Figure 4: Upper triangle part of self calculation.

3.5.2. Decision and lower triangle part

The computation of the lower triangle part of the self-interaction table, together with decision of the trial in the host-group, should be performed at last. As shown in Fig. 5, when thread 0 finishes all the summations for ΔE_0 , it can make a decision by comparing with the random number $X'.w$ under acceptance ratio $\min\{1, \exp(-\beta\Delta E_0)\}$. If $X'.w < \min\{1, \exp(-\beta\Delta E_0)\}$, thread 0 replaces \mathbf{x}_0 by \mathbf{x}'_0 on the shared memory. Otherwise no coordinate update will be made. After thread 0 updates the coordinate of 0'th particle, the remaining threads $tid > 0$ in this block calculate the interactions with the 0'th particle together, and add the results to their own ΔE_{tid} . Then thread 1 updates the coordinate of particle 1 and the updated energy of particle 1 is propagated to the other threads. This procedure is repeated for other threads until all the host-particles in this block are updated.

The last step is to copy the information array of the host-particles from the shared memory to the global memory and set the corresponding flag “BlockState” to be one. The “BlockState” is an array with B elements, each of which indicates the state of corresponding block. Initially all the elements in this array is 0. When

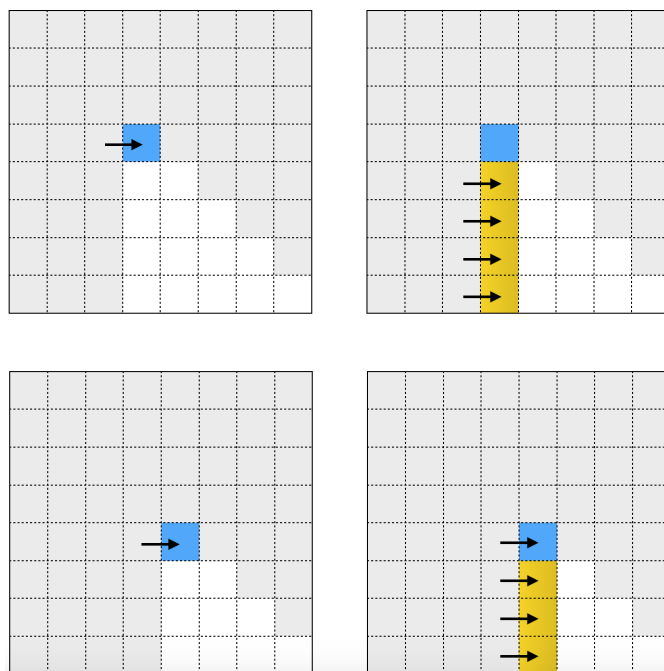


Figure 5: Lower triangle part of self calculation.

a thread block has written the updated coordinates of its host-particles to the global memory, the corresponding flag $\text{BlockState}[\text{bid}]$ is set to 1. Since CUDA is a weak order programming language, the writing instruction for coordinates and the setting instruction for “BlockState” should be separated by a “syncthreads” or a “threadfence” instruction to ensure the required executing order viewed in other blocks. The pseudo-code for this part is shown in algorithm 4.

3.6. Global procedure

Fig. 6 shows the global procedure of a thread block. After initialization, threads in this block firstly compute the inter-group interactions with groups $\text{bid} + 1, \text{bid} + 2, \dots, B - 1$. Then they compute the internal upper-triangle interactions.

The second global loop is to evaluate the interactions with groups $0, 1, \dots, \text{bid} - 1$. In serial implementation, particles in these groups are updated before the host-particles. To realize this dependency in the GPU implementation, before loading information of a guest-group J , $J < \text{bid}$, all threads in this block should wait until the coordinates of particles in guest-group J are updated. Therefore, in each step J of this loop, one of the threads in this block firstly runs a subloop to check and wait until $\text{BlockState}[J]$ is 1. Finally, threads make the decisions and evaluate the lower-triangle part of self-interaction table.

A remaining issue is to reinitialize the flag-array to 0s for the next grid. We

Algorithm 4 Decision and lower triangle part

```
Y[tid] ← X;
__syncthreads();
for l := 0 to S - 1 do
  if (tid = l) and (X'.w < exp(-λΔE)) then
    X'.w ← X.w;
    Y[l] ← X';
  end if
  __syncthreads();
  if (tid > l) then
    ΔE ← ΔE + ΔU(X, X', Y[l]);
  end if
end for
Xdev[bid*S+tid] ← Y[tid];
__syncthreads();
Thread 0 do: Set BlockState[bid] to be 1
```

assign this additional work to block $B - 1$, which is the last one to terminate. Threads in it set all the elements of “BlockState” to 0s so that the next grid can use these indicators directly. The pseudocode for the global procedure is shown in algorithm 5.

The above inter-block communications via flags are **blocked communications** since receivers will not carry out follow-up tasks until receiving the required messages. In a parallel program with blocked communication, we always need to avoid a deadlock state where receivers are waiting for each other and cannot terminate without external force. A case of deadlock in GPU is that all the resident blocks in SMXs are waiting for a message from a block which is not executed yet. In this case, the inactive sender block cannot move forward until one of the resident blocks terminates so that there is enough space to launch this sender block. Therefore, the resident blocks and the inactive sender blocks are waiting for each others. In our implementation, deadlock can be fortunately avoided because each thread block receives messages only from its previous blocks.

3.7. Multi-Thread-Per-Particle Treatment

In above algorithm, one thread is responsible for one particle’s trial. This one-thread-per-particle assignment cannot take full use of device if the number of particles is small because insufficient number of threads are carried out simultaneously to fully take advantage of the capability of GPU. To deal with small system, we take a multi-thread-per-particle assignment, in which the change of energy due to the trial of one particle is evaluated by two or more adjacent threads.

Algorithm 5 Global procedure

Initialization
for $J := \text{bid}+1$ **to** $B - 1$ **do**
 Calculate the energy changes with group J
end for
 Calculate the energy changes of upper triangular part
for $J := 0$ **to** $\text{bid}-1$ **do**
 Thread 0 do: Wait and check until $\text{BlockState}[J]=1$
 $_ \text{syncthread}()$;
 Calculate the energy changes with group J
end for
 Decision and lower triangular part
if $\text{bid}=B-1$ **then**
 Set all the elements of BlockState to 0
end if

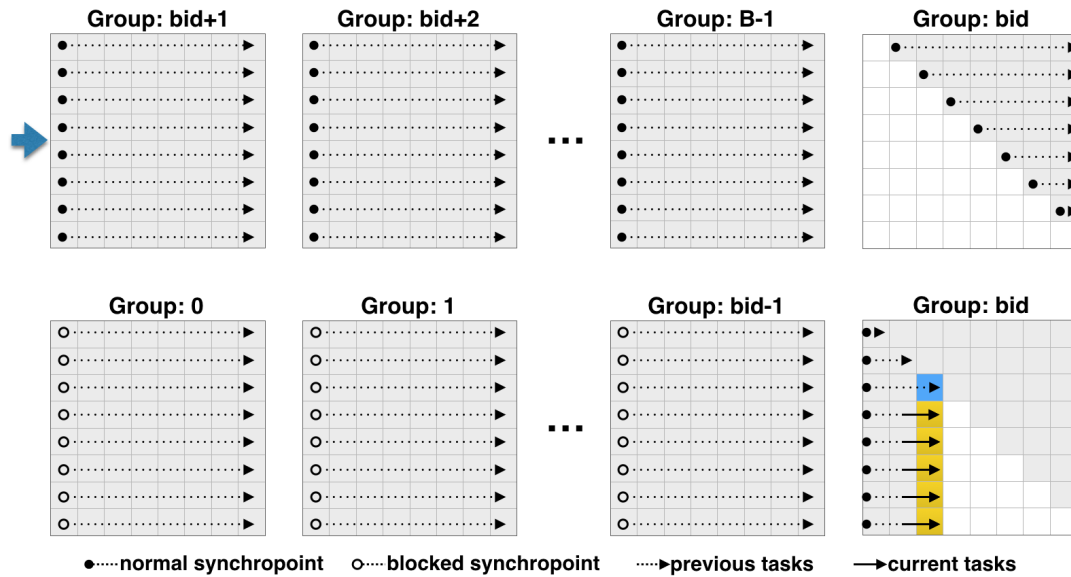


Figure 6: Global procedure in Monte Carlo kernel.

In a program with η -thread-per-particle scheme, the size of thread block is η times of the group size. The energy change of the host particle k (k is the relative index within its group) will be computed by a bundle of threads $tid = \eta k$, $tid = \eta k + 1, \dots, tid = \eta(k + 1) - 1$, with interleaved loops as show in Fig. 7. This interleaved approach is performed in the inter-group evaluation and upper triangle part of self-table. Before the decisions and evaluation of the lower triangle part of self-table, we use warp-reductions to sum up results in all the slave threads of a particle to one. For convenience, η is selected to be 2, 4, 8 and 16, leading 1-step, 2-step, 3-step and 4-step warp shuffle instructions in the reductions.

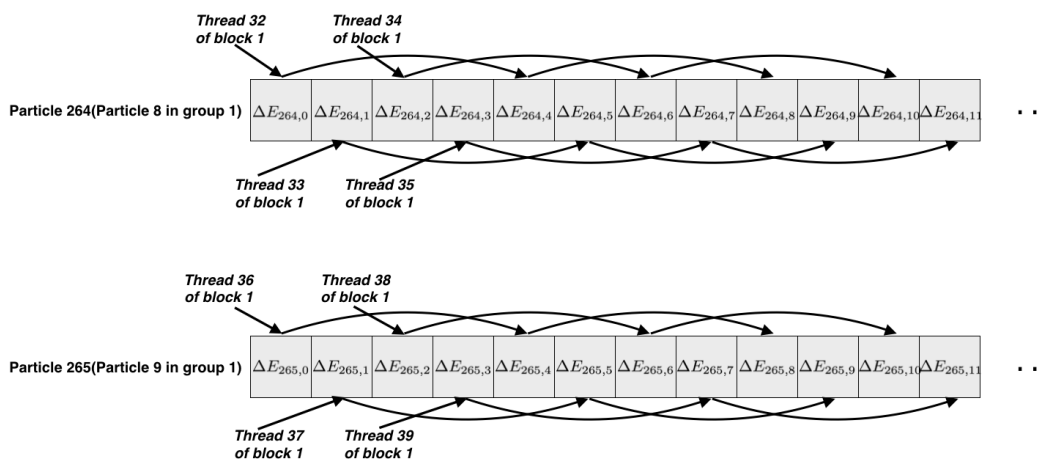


Figure 7: Multi-thread-per-particle assignment for the intergroup evaluation. Each thread block contains 1024 threads. The size of particle group is 256.

4. Benchmark

We benchmark the performance of our Brush Metropolis Algorithm on GPU by simulating the primitive model of electrolyte with various numbers of particles, and make a comparison with two other MC implementations: a sequential CPU code and a GPU code with **parallel reduction scheme**.

In the CPU code, the evaluation and summation of the pairwise interactions are performed sequentially without using any acceleration methods. In the GPU code with parallel reduction scheme, each kernel tries to move one particle. There are $N - 1$ pairwise interactions per trial to be calculated and summed, between the selected particle and the rest of particles in the system. We distribute these $N - 1$ evaluation to N threads while the thread for the selected particle remains idle. In this reduction kernel, the size of thread block is 1,024, which is the maximum

number of threads per block in Tesla K20. The reduction within one block is done by warp shuffle instructions [46], which is an optimized version of parallel reduction [47] for Tesla K20. The intermediate results of warp sums are stored in the shared memory temporarily. We perform warp reduction again with these intermediate results so that the partial sums in this block are integrated. The sub-total of each block is to be added into the global memory atomically. When a block finishes adding the sub-total to the global memory, it should increase a global counter by 1 atomically. This global counter indicates the number of completed thread blocks. Finally, the last block decides acceptance of the proposed particle move. Before the last block updates the coordinate of the selected particle, one of threads in the last block checks this global counter ceaselessly until it indicates completion of summation. This method is straightforward and easy to implement, but it has some drawbacks that limit its performance. Firstly, in the parallel reduction there is a great thread divergency. For a reduction within warp(32 threads), the mean number of active threads is only $(32 + 16 + 8 + 4 + 2 + 1)/6 \approx 10.5$, which means about 2/3 of computational resources are idle. Secondly, each time when the kernel is invoked, the system needs to take some time to initialize the environment. In this naive method program invokes kernel for N times in a cycle (N is the number of particles) whereas our Brush Metropolis method just invokes kernel once per cycle. Thus this naive method accumulates a considerable time for initialization. Finally, in each time when the kernel is invoked, the information of particles needs to be reloaded, leading to many global memory transactions and thus low cache efficiency.

For the Brush Metropolis GPU code, we mainly measure the performance of one-thread-per-particle approach, with number of particles from 8,192 up to 1,048,576. In addition, we carry out benchmarks of multi-thread-per-particle approach with number of particles from 8,192 to 32,768. We find that the size of block 1,024 exhibits best performance. However, for the multi-thread-per-particle implementation, the number of slave threads per particle η with best performance changes with the size of system. We measure the execution time for 4 small systems, and find the best η for each system, as shown in Table. 1

| N | Size of thread block | η |
|-------|----------------------|--------|
| 8192 | 1024 | 16 |
| 16384 | 1024 | 8 |
| 24576 | 1024 | 4 |
| 32768 | 1024 | 4 |

Table 1: Parameters

We carry out all benchmarks on a Linux machine with 2×8 Intel Xeon E5-2630

(2.3GHz) CPU cores and two NVIDIA Tesla K20 GPUs. The operating system is Ubuntu 14.04, with the host C-code compiler GCC 4.8.4 and the GPU code compiler CUDA 7.0. The GPU driver's version is 346.96. The Monte Carlo simulation program on CPU is compiled with optimization option `-O3`. Both GPU codes are compiled with option `-arch=sm_35`. Since the GPU implementation of the Brush Metropolis method costs lots of registers, we take option `-maxrregcount=65` as the register usage per thread to yield the best performance.

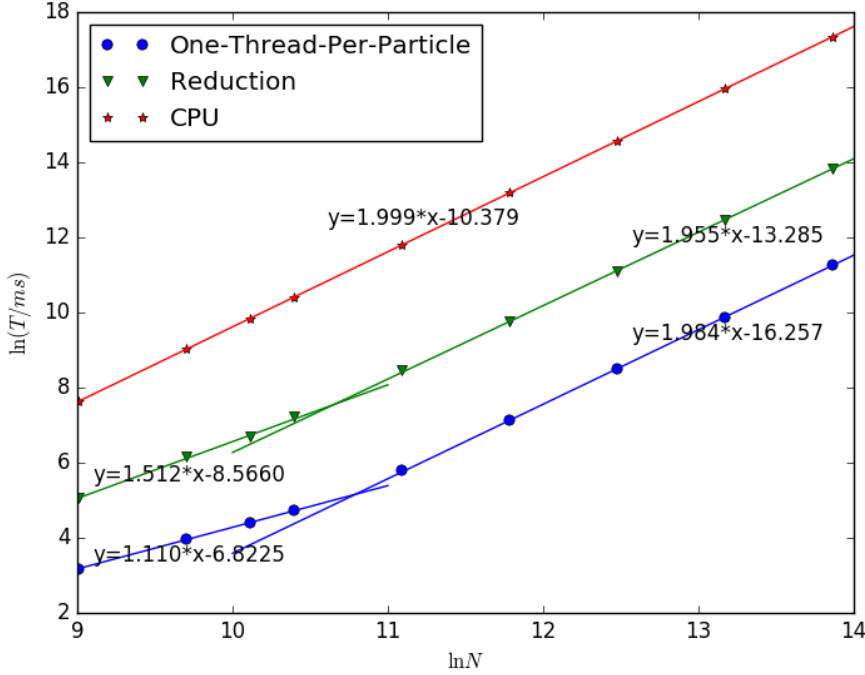


Figure 8: Execution time in log-scale.

The execution time in log-scale for all three algorithms is showed in Fig.8. It clearly shows that the time complexity of the Brush Metropolis Algorithm with one-thread-per-particle assignment is between the step complexity $O(N)$ and the work complexity $O(N^2)$. (The step complexity measures the computational complexity assigned to each thread, including the cost of synchronization.) For particle number less than 50,000, the time complexity scales approximately linearly with the step complexity. For even larger particle numbers, the workload of device is nearly saturated so that the time complexity scales with the work complexity. The GPU implementation with parallel reduction has a time complexity between $O(N \log N)$ and $O(N^2)$.

The speedup of the Brush Metropolis GPU code (comparing with the CPU

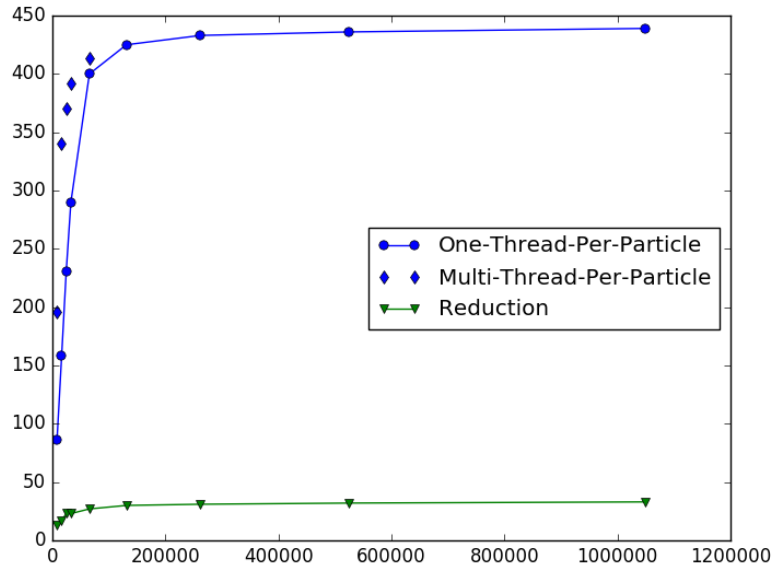


Figure 9: Speedup of the Brush Metropolis GPU code and the convention GPU code using parallel reduction.

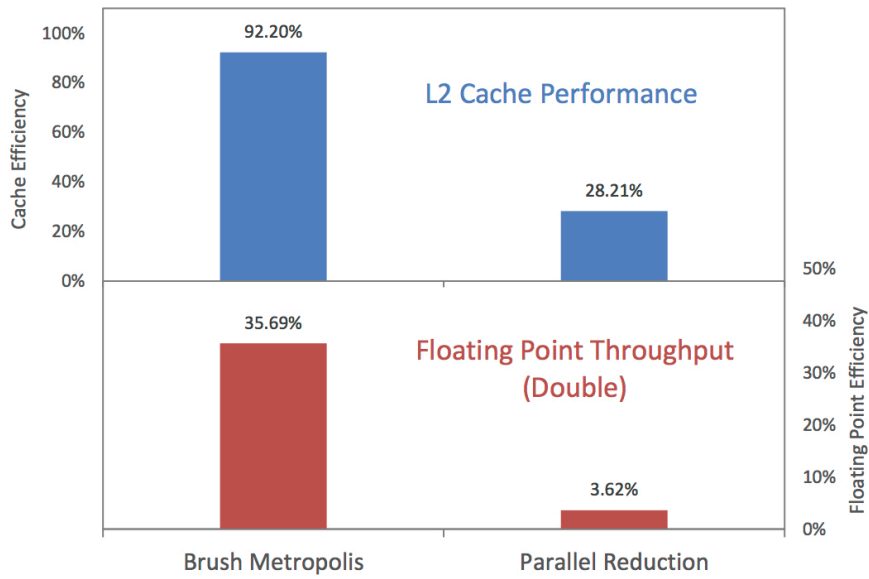


Figure 10: Performance comparison between GPU implementations based on parallel reduction and Brush Metropolis on a system with 131,072 particles. Performance data obtained from nvprof. Due to significantly higher L2 cache efficiency, our GPU implementation has better floating point throughput than that of parallel reduction implementation.

code) as a function of system size is showed in Fig.9. For large systems the implementation with one-thread-per-particle saturates to approximately 440. In general the Brush Metropolis method exhibits significantly better speedup than the implementation with parallel reduction method, because of less global memory transactions and higher L2 cache efficiency, as well as higher floating point throughput in comparison with parallel reduction method, as shown in Fig.10. The main factor that prevents us from achieving even higher floating point efficiency is the Metropolis-style acceptance and rejection, which generate a divergence path for each proposed move.

5. Simulation of Linear Response Properties of Dense Electrolytes

In this section, we use our Brush Metropolis GPU code to carry out a large scale MC simulation of the *primitive model* electrolytes, which are modeled as charged hard spheres. In the experiment the system contains 1,048,576 ions, and more than a hundred of Debye length. Because of the huge amount of available data, we are able to measure all pair correlation functions up to *extremely precision and very long scale*. Using these correlation functions, we determine various renormalized parameters that characterize the linear response properties of the electrolyte, including renormalized valences of the constituent ions, the renormalized Debye length, and the renormalized dielectric constant. These results demonstrate unequivocally that the properties of system are beyond the classical Poisson-Boltzmann theory (PB). We emphasize that in order to compute precisely the long scale properties of electrolytes, it is essentially important to simulate very large system sizes *without making any approximation in long scales*. This is the main advantage of our GPU code, compared with methods using multipole expansions.

5.1. Charge renormalization in concentrated electrolytes

The classical Poisson-Boltzmann theory (PB) [50] predicts that the mean potential around a fixed ion with charge Q is given by

$$\phi = \frac{Q e^{-\kappa r}}{4\pi\epsilon r}, \quad (3)$$

whereas the effective interaction, i.e., two-ion potential of mean force (PMF), between two ions Q_1, Q_2 is

$$U_{12} = \frac{Q_1 Q_2 e^{-\kappa r}}{4\pi\epsilon r}. \quad (4)$$

In the above equations, κ is the *bare* inverse Debye length, which, according to PB, is related to bulk ion densities \bar{n}_\pm and charges q_\pm via:

$$\kappa^2 = \frac{\beta}{\epsilon} (\bar{n}_+ q_+^2 + \bar{n}_- q_-^2). \quad (5)$$

Because the classical PB theory ignores correlation effects, it is not applicable in the concentrated regime. In this regime, the correct far field behaviors of mean potential and two-ion PMF are [48, 49]:

$$\phi = \frac{Q_R e^{-\kappa_R r}}{4\pi\epsilon_R r}, \quad (6)$$

$$U_{12} = \frac{Q_1^R Q_2^R e^{-\kappa_R r}}{4\pi\epsilon_R r}, \quad (7)$$

where $Q_R, Q_1^R, Q_2^R, \kappa_R, \epsilon_R$ are, respectively, the renormalized charges, renormalized inverse Debye length, and renormalized dielectric constant, which are different from their bare values. There is an exact relation between the renormalized Debye length and renormalized charges of constituent ions:

$$\left(\frac{\kappa_R}{\kappa}\right)^2 = \frac{q_+^R - q_-^R}{q_+ - q_-}. \quad (8)$$

Eqs. (6)-(8) are the main results of the *dressed-ion theory* [48]. We shall compute renormalized valences of constituent ions, renormalized Debye length and renormalized dielectric constant, and finally verify the relation Eq. (8) using large-scale MC simulations.

5.2. Simulation Methodology

To compute all renormalized parameters of dense electrolytes, we perform large scale simulations of electrolytes and measure all pair correlation functions. All our simulations are carried out in the Center for High Performance Computing (HPC) of Shanghai Jiaotong University. To fully take advantage of computation resources, we perform one individual simulation on each GPU card. The number of cards employed in simulation varies according to the ion densities. Typically one data point needs 2 GPU cards, each of which carries out about 2,000 iterations with about 5 days. In particular, as the system approaches the charge oscillation regime (where ion densities are high), more Monte Carlo cycles are needed and thus more cards are used. For each simulation(MPI Process), the memory cost on host is 263MB, and the cost of global memory on GPU is about 32MB. The system contains $1,024 \times 1,024 = 1,048,576$ particles with room temperature $T = 300K$, and relative dielectric constant of the solvent is chosen to be $\epsilon = 78.3$ and the Bjerrum length $b = 7.117\text{\AA}$. We use a spherical simulation domain with hard wall boundary conditions. To eliminate influences from the boundary, ions that are less than ten Debye lengths away from the boundary are *not* used for data collection. The initial state is generated by setting particles uniformly inside the simulation domain with hard core repulsions. The proposed position of the selected particle is uniformly generated inside a cubic center around its original position

with size about $L/3$ from experience, where L is the size of simulation domain. In the warming iterations we output the total energy to monitor whether the system equilibrates or not. We find that almost all the systems we study can equilibrate in about 10 cycles. Therefore we are sure that for all the systems the correlation between successive configurations is weak. In Fig. 11 we present a typical autocorrelation function of total energy for an equilibrated 3:-1 electrolyte with classical debye length 28.53\AA , whose autocorrelation time is about 2.3. Here the autocorrelation function of total energy is defined as

$$C_n = \frac{\overline{(E_k - \bar{E})(E_{k+n} - \bar{E})}}{\overline{(E_k - \bar{E})^2}}. \quad (9)$$

Here E_k denotes the total energy at step k and \bar{E} stands for the mean energy at equilibrium. The over-line can be performed by averaging variable over both ensemble (MPI procedures) and time (steps). The integrated autocorrelation time is defined as

$$T_c = \frac{1}{2} + \sum_{k=1}^{\infty} C_k. \quad (10)$$

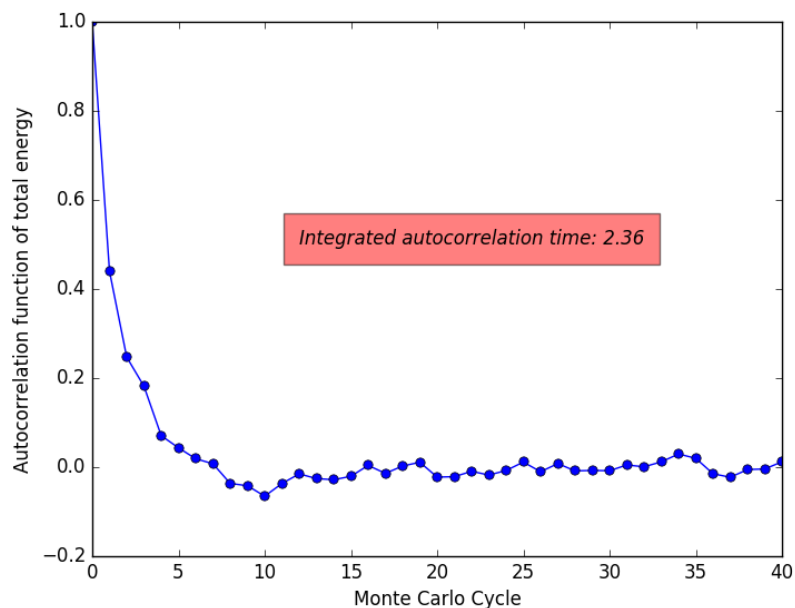


Figure 11: Autocorrelation function of a 3:-1 electrolyte.

A typical pair correlation function is showed in Fig. 12.

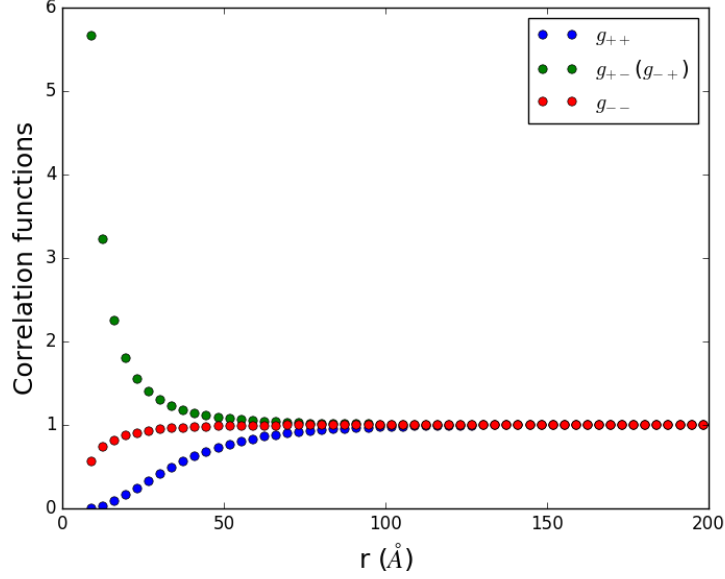


Figure 12: A typical radial correlation function of 3:-1 electrolyte. The radius of ions is 3.75\AA and the classical debye length is 28.53\AA .

The renormalized parameters $\kappa_R, \epsilon_R, q_{\pm}^R$ can be obtained from tails of pair correlation functions $g_{\pm\pm}(\mathbf{r})$ as follows. Firstly the two-ion PMF are obtained via:

$$U_{\pm\pm}(\mathbf{r}) = -k_B T \log g_{\pm\pm}(\mathbf{r}). \quad (11)$$

Now, let us fix a positive/negative ion q_+/q_- at the origin in the bulk electrolyte. The average charge density (excluding the charge q_{\pm} fixed at the origin) can also be computed by:

$$\rho_{\pm}^q(\mathbf{r}) = n_+ q_+ g_{+\pm}(r) + n_- q_- g_{-\pm}(r). \quad (12)$$

Now the mean potential $\phi_{\pm}(\mathbf{r})$ around the fixed ion is given by Eq. (6) with Q_R replaced by q_{\pm}^R in the far field. The mean charge density $\rho_{\pm}^q(\mathbf{r})$ can be obtained using the exact Poisson equation. In the far field, they decay in the form of screened Coulomb potential:

$$\rho_{\pm}^q(\mathbf{r}) = -\epsilon \nabla^2 \phi_{\pm}(\mathbf{r}) \sim \frac{\epsilon \kappa_R^2 q_{\pm}^R e^{-\kappa_R r}}{4\pi \epsilon_R r}. \quad (13)$$

We can take the logarithm of Eq. (13) and obtain:

$$\log [r \rho_{\pm}^q(\mathbf{r})] \sim \log \left[\frac{\epsilon \kappa_R^2 q_{\pm}^R}{4\pi \epsilon_R} \right] - \kappa_R r. \quad (14)$$

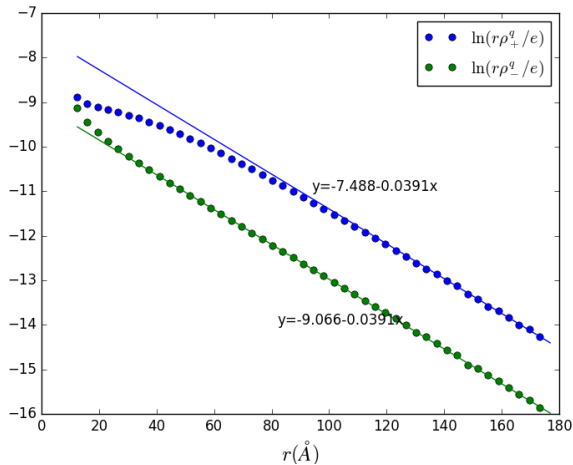


Figure 13: Plotting $\ln(rp_{\pm}^q/e)$ v.s. r for 3:-1 electrolyte. The radius of ions is 3.75\AA and the bare debye length (as predicted by PB) is 28.53\AA . The renormalized Debye length is 25.57\AA .

We can therefore plot the l.h.s. (measured by simulations) as a function of radius r , fit the data to straight-lines *in the far field*, and extract the renormalized inverse Debye length κ_R from their slopes. This is illustrated in Fig. (13). Note that two straight-lines have the same slope and give a renormalized Debye length $\kappa_R^{-1} = 25.57\text{\AA}$, manifestly different from the bare Debye length $\kappa^{-1} = 28.53\text{\AA}$.

According to Eq. (7), the far field asymptotics of the two-ion PMF is

$$U_{\pm\pm}(\mathbf{r}) \sim \frac{q_{\pm}^R q_{\pm}^R e^{-\kappa_R r}}{4\pi\epsilon_R r}. \quad (15)$$

Taking the ratio of Eqs. (13) and (15), we find the following relation valid in the far field:

$$q_{\pm}^R = -\epsilon\kappa_R^2 \frac{U_{\pm\pm}(\mathbf{r})}{\rho_{\pm}^q(\mathbf{r})}. \quad (16)$$

Since all quantities in the r.h.s. are known, we can use this relation to determine the renormalized charges q_{\pm}^R of positive and negative ions. In fact we have two independent ways to compute the renormalized charges. Let us write them out explicitly:

$$q_+^R = -\epsilon\kappa_R^2 \frac{U_{++}(\mathbf{r})}{\rho_+^q(\mathbf{r})} = -\epsilon\kappa_R^2 \frac{U_{-+}(\mathbf{r})}{\rho_-^q(\mathbf{r})}, \quad (17a)$$

$$q_-^R = -\epsilon\kappa_R^2 \frac{U_{+-}(\mathbf{r})}{\rho_+^q(\mathbf{r})} = -\epsilon\kappa_R^2 \frac{U_{--}(\mathbf{r})}{\rho_-^q(\mathbf{r})}. \quad (17b)$$

The data for one particular simulation are shown in Fig. 14, from which we extract the renormalized charges of positive and negative ions to be 4.125 and -0.884 re-

spectively. Note that these are substantially different from the bare charges, which are 3 and -1 respectively. This charge renormalization arises as a consequence of ionic correlations and signify the failure of the classical PB theory.

Using Eq. (13) we can compute the renormalized dielectric constant ϵ_R in terms of κ_R , q_{\pm}^R , and $\rho_{\pm}^q(\mathbf{r})$ (again in two independent ways):

$$\epsilon_R = -\epsilon \frac{\kappa_R^2 q_{\pm}^R e^{-\kappa_R r}}{4\pi r \rho_{\pm}^q(\mathbf{r})}. \quad (18)$$

Finally, we also use the computed κ_R and q_{\pm}^R to test the validity of the exact relation Eq. (8). The results are displayed in Fig. 15(a). Additionally, we plot q_{\pm}^R and κ_R in Figs. 15(b) and 15(c) respectively measured in two independent ways and show that they are consistent with each other, within computational errors. All these numerical tests unambiguously demonstrate the validity and internal consistency of the dressed-ion theory, Eqs. (6)-(8).

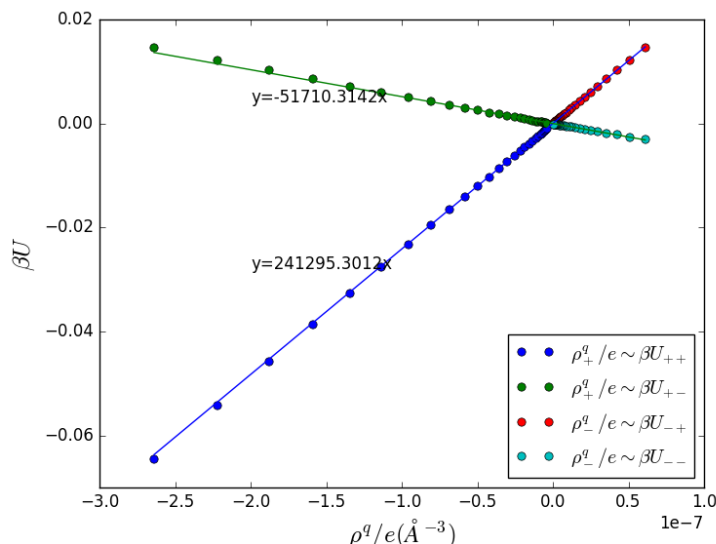


Figure 14: Plotting βU v.s. ρ^q/e for 3:-1 electrolyte. The radius of ions is 3.75\AA and the classical debye length is 28.53\AA . Using Eq. (17) and fitting the long range data to straight-lines, we obtain the renormalized valences of the positive ion and of the negative ions to be 4.125 and -0.884 respectively, which are different from the bare valences 3 and -1 respectively. This again demonstrates the failure of classical PB theory.

6. Conclusion

In this work, we have developed an efficient GPU code for large-scale Monte Carlo simulation of Coulomb many-body systems, which parallelizes the sequential

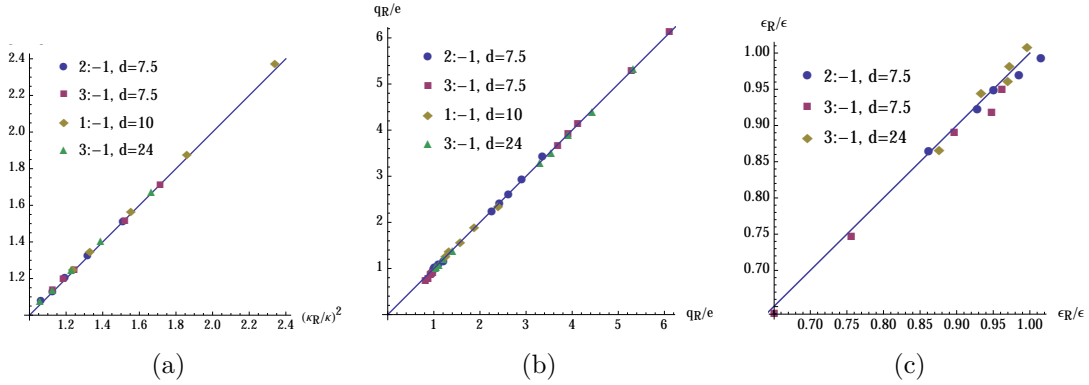


Figure 15: (a): Verification of the exact relation Eq. (8) using MC simulation. d is the ion diameter (in the unit of \AA). Vertical axis: $(q_+^R - q_-^R) / (q_+ - q_-)$. (b): The renormalized valences of ions, computed using in two independent ways, see Eqs. (17). The purpose of this panel is to demonstrate that two independent computations give the same result, within computational errors. (c): Renormalized dielectric constant can also be computed in two ways, see (18). The fact that these renormalized parameters are different from their bare values demonstrates that the properties of concentrated electrolytes are beyond the classical Poisson-Boltzmann theory.

updating scheme and achieves an acceleration of 440 over the sequential CPU code, without sacrificing accuracy. We have further applied this method to precisely measure the long scale linear response properties of dense asymmetric electrolytes and have demonstrated that they are beyond the classical Poisson-Boltzmann theory. Further applications of this method will be reported in future publications.

Y.H. Liang and X.J. Xing acknowledge financial support from NSFC (grant No. 11174196 and 91130012). Y.H. Li acknowledges support from NSF (grant No. 1066471). The authors also thank Beijing Computational Science Research Center (BCSRC) for hospitality, where part of this work is done. This work is supported by Center for HPC, Shanghai Jiao Tong University.

References

- [1] Andrew W. Appel. An efficient program for many-body simulation. **SIAM J. Sci. Stat. Comput.**, 6(1):85, Jan 1985.
- [2] Josh Barnes and Piet Hut. A hierarchical $o(n \log n)$ force-calculation algorithm. **Nature**, 324(4):446, Dec 1986.
- [3] L. Greengard and V. Rokhlin. A fast algorithm for particle simulations. **Journal of Computational Physics**, 73(2):325–348, 1987.
- [4] Gan Zecheng and Xu Zhenli. Efficient implementation of the barnes-hut octree algorithm for monte carlo simulations of charged systems. **Science China, Mathematics**, 57(7):1331–1340, July 2014.

- [5] Yihao Liang, Zhenli Xu, and Xiangjun Xing. A multi-scale Monte Carlo method for electrolytes **New J. Phys.**, 17(2015), 083082.
- [6] Zhenli Xu, Yihao Liang, Xiangjun Xing. Mellin Transform and Image Charge Method for Dielectric Sphere in an Electrolyte. **SIAM J. Appl. Math.**, **73**(4), 1396-1415.(2013).
- [7] Y. Lin, A. Baumketner, S. Deng, Z. Xu, D. Jacobs, and W. Cai. An image-based reaction field method for electrostatic interactions in molecular dynamics simulations of aqueous solutions. **J. Chem. Phys.**, 131:154103, 2009.
- [8] NVIDIA. CUDA C Programming Guide. NVIDIA, 7.0 edition, Mar 2015.
- [9] Lars Nyland, Mark Harris, and Jan Prins. Fast n-body simulation with cuda. http://http.developer.nvidia.com/GPUGems3/gpugems3_ch31.html.
- [10] Carolyn L. Phillips, Joshua A. Anderson, and Sharon C. Glotzer. Pseudo-random number generation for brownian dynamics and dissipative particle dynamics simulations on gpu devices. **Journal of Computational Physics**, 230:7191–7201, 2011.
- [11] Tsuyoshi Hamada and Toshiaki Iitaka. The chamomile scheme: An optimized algorithm for n-body simulations on programmable graphics processing units. **arXiv:astro-ph/0703100v1**, 2007.
- [12] Evghenii Gaburov, Jeroen Bédorf, and Simon Portegies Zwart. Gravitational tree-code on graphics processing units: implementation in cuda. In **Procedia Computer Science**, volume 1, pages 1119–1127. International Conference on Computational Science, ICCS 2010, 2012.
- [13] Juekuan Yang, Yujuan Wang, and Yunfei Chen. Gpu accelerated molecular dynamics simulation of thermal conductivities. **Journal of Computational Physics**, 221:799–804, 2007.
- [14] Trung Dac Nguyen, Carolyn L. Phillips, Joshua A. Anderson, and Sharon C. Glotzer. Rigid body constraints realized in massively-parallel molecular dynamics on graphics processing units. **Computer Physics Communications**, 182:2307–2313, 2011.
- [15] Weiguo Liu, Bertil Schmidt, Gerrit Voss, and Wolfgang Müller-Wittig. Accelerating molecular dynamics simulations using graphics processing units with cuda. **Computer Physics Communications**, 179:634–641, 2008.

- [16] D.C. Rapaport. Enhanced molecular dynamics performance with a programmable graphics processor. **Computer Physics Communications**, 182:926–934, 2011.
- [17] Andreas W. Götz, Mark J. Williamson, Dong Xu, Duncan Poole, Scott Le Grand, and Ross C. Walker. Routine microsecond molecular dynamics simulations with amber on gpus. 1. generalized born. **Journal of Chemical Theory and Computation**, 8:1542–1555, 2012.
- [18] Romelia Salomon-Ferrer, Andreas W. Götz, Duncan Poole, Scott Le Grand, and Ross C. Walker. Routine microsecond molecular dynamics simulations with amber on gpus. 2. explicit solvent particle mesh ewald. **Journal of Chemical Theory and Computation**, 9:3878–3888, 2013.
- [19] Romelia Salomon-Ferrer, David A. Case, and Ross C. Walker. An overview of the amber biomolecular simulation package. **WIREs Computational Molecular Science**, 3:198–210, Mar 2013.
- [20] Scott Le Grand, Andreas W. Götz, and Ross C. Walker. Spfp: Speed without compromise—a mixed precision model for gpu accelerated molecular dynamics simulations. **Computer Physics Communications**, 184:374–380, 2013.
- [21] Michela Taufer, Narayan Ganesan, and Sandeep Patel. Gpu-enabled macromolecular simulation: Challenges and opportunities. **Computing in Science & Engineering**, pages 56–65, 2013.
- [22] Narayan Ganesan Brad A. Bauer Timothy R. Lucas Sandeep Patel and Michela Taufer. Structural, dynamic, and electrostatic properties of fully hydrated dmpe bilayers from molecular dynamics simulations accelerated with graphical processing units (gpus). **Journal of Computational Chemistry**, 32:2958–2973, 2011.
- [23] William Michael Brown, Steven James Plimpton, Peng Wang, Pratul K. Agarwal, Scott Hampton, and Paul Stewart Crozier. Porting lammmps to gpus. In **SciTech Connect**, Savannah, GA., 2010. Proposed for presentation at the SOS 14 Conference.
- [24] W. Michael Brown, Peng Wang, Steven J. Plimpton, and Arnold N. Tharrington. Implementing molecular dynamics on hybrid high performance computers – short range forces. **Computer Physics Communications**, 182:898–911, 2011.

- [25] John E. Stone, James C. Phillips, Peter L. Freddolino, David J. Hardy, Leonardo G. Trabuco, and Klaus Schulten. Accelerating molecular modeling applications with graphics processors. **Wiley InterScience**, 2007.
- [26] James C. Phillips, John E. Stone, and Klaus Schulten. Adapting a message-driven parallel application to gpu-accelerated clusters. In **Proceedings of the 2008 ACM/IEEE Conference on Supercomputing**, page 1, Austin, Texas, 2008. IEEE/ACM, IEEE Press.
- [27] Peter H. Colberg and Felix Höfling. Highly accelerated simulations of glassy dynamics using gpus: Caveats on limited floating-point precision. **Computer Physics Communications**, 182:1120–1129, 2011.
- [28] Peter Eastman and Vijay S. Pande. Efficient nonbonded interactions for molecular dynamics on a graphics processing unit. **Journal of Computational Chemistry**, 31:1268–1272, 2010.
- [29] Joshua A. Anderson, Chris D. Lorenz, and A. Travesset. General purpose molecular dynamics simulations fully implemented on graphics processing units. **Journal of Computational Physics**, 227:5342–5359, 2008.
- [30] Hoomd-blue. <http://codeblue.umich.edu/hoomd-blue/>.
- [31] Gromacs. <http://www.gromacs.org>.
- [32] Acemd. <https://www.acellera.com>.
- [33] Benjamin Block, Peter Virnau, and Tobias Preis. Multi-gpu accelerated multi-spin monte carlo simulations of the 2d ising model. **Computer Physics Communications**, 181:1549–1556, 2010.
- [34] Tobias Preis, Peter Virnau, Wolfgang Paul, and Johannes J. Schneider. Gpu accelerated monte carlo simulation of the 2d and 3d ising model. **Journal of Computational Physics**, 228:4468–4477, 2009.
- [35] Tal Levy, Guy Cohen, and Eran Rabani. Simulating lattice spin models on graphics processing units. **Journal of Chemical Theory and Computation**, 6:3293–3301, 2010.
- [36] Joshua A. Anderson, Eric Jankowski, Thomas L. Grubb, Michael Engel, and Sharon C. Glotzer. Massively parallel monte carlo for many-particle simulations on gpus. **Journal of Computational Physics**, 254:27–38, 2013.

- [37] Ashraf Yaseen and Yaohang Li. Accelerating knowledge-based energy evaluation in protein structure modeling with graphics processing units. **Journal of Parallel and Distributed Computing**, 72:297–307, 2012.
- [38] Ashraf Yaseen and Yaohang Li. A Load-Balancing Workload Distribution Scheme for Three-Body Interaction Computation on Graphics Processing Units (GPU). **Journal of Parallel and Distributed Computing**, 87:91–101, 2016.
- [39] Jason Mick, Eyad Hailat, Vincent Russo, Kamel Rushaidat, Loren Schwiebert, and Jeffrey Potoff. Gpu-accelerated gibbs ensemble monte carlo simulations of lennard-jonesium. **Computer Physics Communications**, 184:2662–2669, 2013.
- [40] Eyad Hailat, Vincent Russo, Kamel Rushaidat, Jason Mick, Loren Schwiebert, Kamel Rushaidat, Jason Mick, Loren Schwiebert, and Jeffrey Potoff. Parallel monte carlo simulation in the canonical ensemble on the graphics processing unit. **International Journal of Parallel, Emergent and Distributed Systems**, 29(4):379–400, 2014.
- [41] Jihan Kim and Jocelyn M. Rodgers and Manuel Ath enes and Berend Smit Molecular Monte Carlo Simulations Using Graphics Processing Units: To Waste Recycle or Not? **Journal of Chemical Theory and Computation**, 7:3208–3222, 2011
- [42] Jihan Kim and Berend Smit Efficient Monte Carlo Simulations of Gas Molecules Inside Porous Materials **Journal of Chemical Theory and Computation**, 8:2336–2343, 2012
- [43] Jihan Kim and Richard L. Martin and Oliver RÃijbel and Maciej Haranczyk and Berend Smit High-Throughput Characterization of Porous Materials Using Graphics Processing Units **Journal of Chemical Theory and Computation**, 8:1684–1693, 2012
- [44] William H. Press and Saul A. Teukolsky. **Numerical Recipes**. 9780521431088. Cambridge University Press, 3rd edition, August 2007.
- [45] Wenjian Yu, Kuangya Zhai, Hao Zhuang, Junqing Chen/ Accelerated floating random walk algorithm for the electrostatic computation with 3-D rectilinear-shaped conductors **Simulation Modelling Practice and Theory**, 34:20–36, 2013

- [46] Justin Luitjens. Faster parallel reductions on kepler. <http://devblogs.nvidia.com/paralleforall/faster-parallel-reductions-kepler/>, Feb 2014.
- [47] Mark Harris. Optimizing parallel reduction in cuda. <http://developer.download.nvidia.com/assets/cuda/files/reduction.pdf>, 2008.
- [48] Roland Kjellander and D John Mitchell. An exact but linear and poisson—boltzmann-like theory for electrolytes and colloid dispersions in the primitive model. **Chemical physics letters**, 200(1):76–82, 1992.
Roland Kjellander and D. John Mitchell. Dressed ion theory for electrolyte solutions: A Debye-Hückel-like reformulation of the exact theory for the primitive model. **The Journal of Chemical Physics**, 101(1):603–626, 1994.
- [49] Mingnan Ding, Yihao Liang, Bing-sui Lu, and Xiangjun Xing. Charge Renormalization and Charge Oscillation in Asymmetric Primitive Model. Submitted to **Journal of Statistical Physics**.
- [50] D. Andelman. *Electrostatic Properties of Membranes: The Poisson-Boltzmann Theory.*, chapter 12. Structure and Dynamics of Membranes Generic and Specific Interactions. ELSEVIER, Amsterdam, 1995.