

SIMD Vectorization for the Lennard-Jones Potential with AVX2 and AVX-512 instructions

Hiroshi Watanabe* and Koh M. Nakagawa

*The Institute for Solid State Physics, The University of Tokyo,
Kashiwanoha 5-1-5, Kashiwa, Chiba 277-8581, Japan*

(Dated: October 23, 2018)

This work describes the SIMD vectorization of the force calculation of the Lennard-Jones potential with Intel AVX2 and AVX-512 instruction sets. Since the force-calculation kernel of the molecular dynamics method involves indirect access to memory, the data layout is one of the most important factors in vectorization. We find that the Array of Structures (AoS) with padding exhibits better performance than Structure of Arrays (SoA) with appropriate vectorization and optimizations. In particular, AoS with 512-bit width exhibits the best performance among the architectures. While the difference in performance between AoS and SoA is significant for the vectorization with AVX2, that with AVX-512 is minor. The effect of other optimization techniques, such as software pipelining together with vectorization, is also discussed. We present results for benchmarks on three CPU architectures: Intel Haswell (HSW), Knights Landing (KNL), and Skylake (SKL). The performance gains by vectorization are about 42% on HSW compared with the code optimized without vectorization. On KNL, the hand-vectorized codes exhibit 34% better performance than the codes vectorized automatically by the Intel compiler. On SKL, the code vectorized with AVX2 exhibits slightly better performance than that with vectorized AVX-512.

I. INTRODUCTION

Since Alder and Wainwright performed molecular dynamics (MD) simulations for the first time [1], MD has been an important tool for exploring the wide variety of fields in science. Starting from the first MD simulation with 32 atoms, the number of atoms involved in MD has continued to increase owing to the development of computational power and has reached hundreds of billions to trillions of atoms [2, 3]. It is expected that the size and complexity of systems will continue to increase with the development of computers. Since the increase in the CPU frequency stopped in the early 2000s [4], the performance development of CPUs has mainly depended on increasing the number of CPU cores and the SIMD width. SIMD, which is short for single instruction multiple data, enables us to process multiple data with a single instruction. The width of SIMD is the number of data that can be processed simultaneously and it is determined by the bit length of registers. Beginning with SSE, which supports 128-bit registers, 256-bit registers for AVX and 512-bit registers for AVX-512 are available, where SSE and AVX are shorts for Streaming SIMD extensions and advanced vector extensions, respectively. Since the theoretical peak performance of a CPU is the value when the application uses the vector width to the full, the utilization of SIMD is crucial for the performance of applications on modern CPU architecture. However, it is not trivial to transform a scalar kernel to a SIMD-vectorized one. Additionally, the optimal method of vectorization is different for each instruction set architecture

(ISA). Portability is one of the important issues for vectorization. The SIMD capabilities are usually provided as an extension of the existing ISA. SIMD and related instructions are different for each architecture. To address this issue, several approaches have been proposed. One of them is a directive-based approach [5]. By using directives, one can utilize GPGPU or many-core processors effectively while retaining portability. Another is a framework approach. Karpiński and McDonald provided a framework called Unified Multi/Many-Core Environment (UME), which allows a programmer to utilize the SIMD capabilities without detailed knowledge of the specific architecture [6]. With these approaches, one can develop an application that exhibits high performance without increasing software complexity.

Despite the above efforts, SIMD vectorization remains a difficult and cumbersome task, since the performance would be unsatisfactory simply by vectorizing the existing scalar code. To utilize the SIMD capability of modern CPUs, it is necessary to combine SIMD vectorization with an optimal data layout and other optimization techniques. In this paper, we describe the SIMD vectorization of the force calculation for the Lennard-Jones (LJ) potential with AVX2 and AVX-512 on several types of CPU. The force calculation is the most time-consuming part of MD, and therefore, the efficiency of the kernel directly determines the performance of MD. Since the force-calculation kernel is a typical example of the vectorization of a loop with indirect access, there have been many reports on the vectorization of MD codes, including that of the many-body potential [7]. Additionally, the important kernels of widely used MD packages, such as Gromacs [8], LAMMPS [9], and NAMD [10], have already been effectively vectorized and are available. However, the previous works mainly focused on the performance of the whole application instead of that of the

*Corresponding author; Electronic address:
hwatanabe@issp.u-tokyo.ac.jp

vectorized kernel. While the performance of the application is ultimately most important, the total performance depends on various factors such as the communication and the cost of thread synchronization. Such factors make it difficult to evaluate the impact of vectorization. Therefore, we focus on force-calculation kernel in this paper. In particular, we investigate in detail the impact of optimization techniques when used with vectorization. We also discuss the difference between the vectorization with AVX2 and AVX-512 and between CPU architectures. Since the computational intensity of LJ is low, *i.e.*, a number of floating point operations is small relative to the amount of memory access, the memory access is the main bottleneck. Therefore, most of our efforts are devoted to optimizing memory access. The optimization techniques described in this paper can be applied to other cases with low computational intensity. Although mixed-precision calculations are known to improve the performance of MD, we use double-precision numbers for all calculations throughout this work for simplicity.

The rest of the article is organized as follows. The benchmark conditions are described in the next section. Then we give a brief introduction to the optimization techniques employed before vectorization in Sec. III. Vectorization with AVX2 is described in Sec. IV and that with AVX-512 and related optimizations are described in Sec. V. The benchmark results on SKL are shown in Sec. VI. Section VII is devoted to a summary and discussion. The associated code is available at [11].

II. BENCHMARK CONDITIONS

We consider three types of CPU architecture; Intel Haswell (HSW), Knights Landing (KNL), and Skylake (SKL), which we refer to as HSW, KNL, and SKL, respectively. We performed benchmark simulations on the systems of the Institute for Solid State Physics of the University of Tokyo for HSW and SKL and of the Information Technology Center of the University of Tokyo for KNL. The conditions for the benchmark simulations are as follows. The simulation box is a cube with a linear size of 100σ , where σ is the diameter of an LJ atom. Atoms are placed at the face-centered-cubic lattice. The cutoff length is 3.0σ . The number of atoms is 119164, and consequently, the number density is 1.0. Only force calculations are performed and the time to calculate the force 100 times is measured. The program is executed as a single-threaded process on a single CPU core. The software is compiled using Intel C++ compiler. The details of the CPU, the versions of the compiler, and the compilation options are listed in Table I.

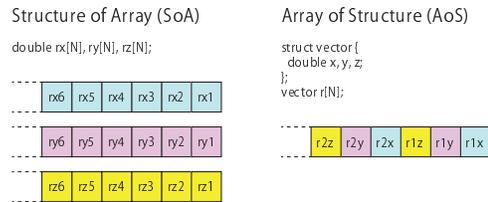


FIG. 1: (Color online) Data layout. The numbers represent the indices of atoms. Left: Structure of Array (SoA). The x coordinates of the atoms are contiguously arranged in the memory. The same is true for the y and z coordinates. Right: Array of Structure (AoS). Coordinates are stored for each atom. Examples of data definition in C language are also shown for both cases.

III. OPTIMIZATION BEFORE VECTORIZATION

In this section, we describe several optimization techniques employed before SIMD vectorization since these techniques affect the SIMD vectorization. While some of the techniques have been described in a past paper [12], we discuss them here to make this paper self-contained. Since we consider a three dimensional system, a position vector and velocity vector each contain three elements. When a system contains N atoms, we have to store N position and velocity vectors in memory. There are two ways to arrange such data in memory, Structure of Array (SoA) and Array of Structure (AoS). See Fig. 1 for the difference between the two layouts. Here, we adopt SoA as the data layout, *i.e.*, the data are arranged so that the x -coordinates of the i -atom and $(i + 1)$ -atom are adjacent. The results of optimizations are shown in Fig. 2. We describe each optimization below.

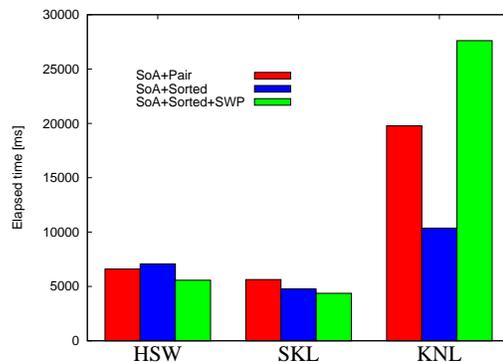


FIG. 2: (Color online) Effects of optimizations before vectorization.

Name	Processor	Vector ISA	Compiler	Options
HSW	Intel Xeon E5-2680 v3	AVX2	icpc (ICC) 16.0.4	-xHOST
KNL	Intel Xeon Phi 7250	AVX2, AVX-512	icpc (ICC) 18.0.1	-axMIC-AVX512
SKL	Intel Xeon Gold 6148	AVX2, AVX-512	icpc (ICC) 18.0.1	-xCORE-AVX512 -qopt-zmm-usage=high

TABLE I: Details of CPUs, versions of the compiler, and compilation options used for benchmark simulations. The common compiler option is `-O3 -std=c++11 -w2 -w3 -diag-disable:remark -restrict` and only CPU-dependent options are shown.

Name	Description
Pair	Naive implementation using a pair list.
Sorted	Reduce the memory access by sorting a list.
AoS-4	Adopt the AoS data structure with padding.

TABLE II: List of optimizations.

A. Lennard-Jones Potential

We consider a classical MD simulation of the LJ potential with truncation. The force calculation of the two-body potential consists of double loops. The loop counter of the outer loop is denoted by i and that of the inner loop is denoted by j . The loop counters i and j correspond to the indices of atoms. Consider the atom pair (i, j) , for which we calculate the force between them. We call the atom with index i the i -atom and the other the j -atom. The coordinates of the i - and j -atoms are denoted by \vec{q}_i and \vec{q}_j , respectively, and the distance between them is given by $r = |\vec{q}_j - \vec{q}_i|$. Then the LJ potential is described by

$$V(r) = 4\varepsilon \left[\left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right], \quad (1)$$

where ε is the well depth and σ is the atomic diameter. Hereafter, we set σ and ε to unity. Since the potential decays rapidly as the distance increases, it is wasteful to calculate the force between pairs at long distances. Therefore, a cutoff distance is introduced and only the interactions of pairs within the distance are considered. In this paper, we adopt the simple truncation, *i.e.*, we introduce a cutoff distance r_c and we ignore force between pairs of atoms at a distance greater than r_c . This truncation modifies the potential function by adding a constant so that $V(r_c) = 0$. Note that this simple method may exhibit problems in terms of energy conservation. For practice use, it is better to adopt a truncation method with not only the potential but also the force becoming zero at the truncation distance [13–15]. While we here adopt the simple truncation, the optimization and vectorization techniques described in this paper can be applied to other truncation methods with minor modifications.

The force calculation for a single pair is shown in Algorithm 1, for which we count the number of arithmetic operations. For example, the operation $\vec{q}_j - \vec{q}_i$ involves three subtractions. The total number of arithmetic operations required to calculate the force between a single pair is 27 addition/multiplications and one division. To

Algorithm 1 Force calculation of LJ potential

```

1:  $\vec{r} \leftarrow \vec{q}_j - \vec{q}_i$ 
2:  $r^2 \leftarrow \vec{r} \cdot \vec{r}$ 
3:  $r^6 \leftarrow r^2 \times r^2 \times r^2$ 
4:  $r^{14} \leftarrow r^6 \times r^6 \times r^2$ 
5:  $df \leftarrow (48 - 24 \times r^6) \times dt/r^{14}$ 
6:  $\vec{p}_i \leftarrow \vec{p}_i - df \times \vec{r}$ 
7:  $\vec{p}_j \leftarrow \vec{p}_j - df \times \vec{r}$ 

```

calculate the force between a single pair, 48 bytes must be read and 24 bytes must be written since it is necessary to load the coordinates and momenta of the atoms and write back the momenta. Considering the above, the force calculation of the LJ system becomes memory-bound.

B. Verlet List and Bookkeeping Method

Since the potential function is truncated, we must construct a pair list including only pairs for which the distance is shorter than the cutoff distance. While the trivial implementation to construct the list has a complexity of $O(N^2)$ for the total number of atoms N , the complexity of the computation can be reduced to $O(N)$ by adopting the linked-list method [16, 17]. In the construction of the pair list, we register each pair within a search length r_s , which is set at longer than r_c . Then the list has a margin of $r_s - r_c$ and we can reuse the list for several time steps. This technique is called the bookkeeping method, which greatly reduces the time required to construct the list [18]. By monitoring the fastest atom in the system, we can safely determine whether we can continue to use the list or whether we must rebuild it [19]. By adopting the bookkeeping method, pairs beyond the cutoff distance can be registered in the pair list. Therefore, we must exclude such pairs from the calculation of forces using the pair list. This fact affects the subsequent SIMD vectorization.

C. Sorting by Indices

An interacting pair can be denoted by a pair of indices. Therefore, the pair list is represented by a list of index pairs (see Fig. 3 (a)). The pseudocode for calculating forces using the pair list is shown in Algorithm 2. The force-calculation kernel contains a single loop and it re-

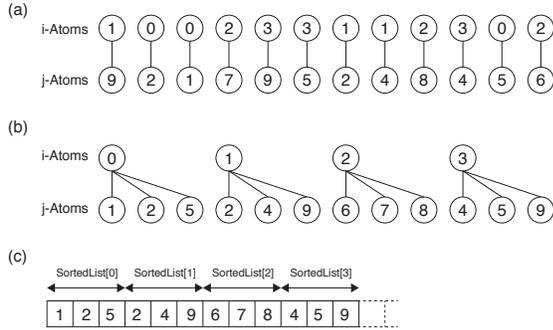


FIG. 3: Data layout of a pair list. (a) List of pairs. (b) Sorted list. The j -atoms interacting with the same i -atom are grouped together. We expect that the data of i -atoms are on registers. (c) Array representation of (b). The array `SortedList[i]` denotes the list of j -atoms that interact with i -atom.

quires memory access both for i - and j -atoms. Since the number of arithmetic operations is small as described before, the memory access becomes the bottleneck, resulting poor performance. To reduce the number of memory accesses, we sort the pair list into group j -atoms interacting with the same i -atom (see Fig. 3 (b)). The array representation of the sorted list is shown in Fig. 3 (c). The array `SortedList[i]` denotes the array that stores the indices of j -atoms interacting with i -atoms. This list can be constructed by a counting sort and the computational complexity of the construction is $O(N)$. Using the sorted list, the force-calculation kernel contains double loops. Then the information of the i -atom, which is denoted by the loop counter of the outer loop, is stored in registers. Then only the memory access of j -atoms is required in the inner loop. The pseudocode of the force calculation using the sorted list is shown in Algorithm 3. Hereafter, we refer to this optimization to “Sorted”. This optimization works effectively on KNL as shown in Fig. 2.

Algorithm 2 Calculating the force in a simple manner

```

1: for all pairs  $(i, j)$  in PairList do
2:   Load  $\vec{q}_i$ 
3:   Load  $\vec{q}_j$ 
4:   if  $|\vec{q}_j - \vec{q}_i| < r_c$  then
5:     Load  $\vec{p}_i$ 
6:     Load  $\vec{p}_j$ 
7:     Calculate force between  $i$ - and  $j$ -particles.
8:     Update  $\vec{p}_i$  and  $\vec{p}_j$ 
9:     Store  $\vec{p}_i$ 
10:    Store  $\vec{p}_j$ 
11:   end if
12: end for

```

Algorithm 3 Calculating the force with a sorted list

```

1: for  $i = 0$  to  $N - 1$  do
2:   Load  $\vec{q}_i$ 
3:   Load  $\vec{p}_i$ 
4:   for all  $j$  in SortedList[i] do
5:     Load  $\vec{q}_j$ 
6:     if  $|\vec{q}_j - \vec{q}_i| < r_c$  then
7:       Load  $\vec{p}_j$ 
8:       Calculate force between  $i$ - and  $j$ -particles.
9:       Update  $\vec{p}_i$  and  $\vec{p}_j$ 
10:      Store  $\vec{p}_j$ 
11:     end if
12:   end for
13:   Store  $\vec{p}_i$ 
14: end for

```

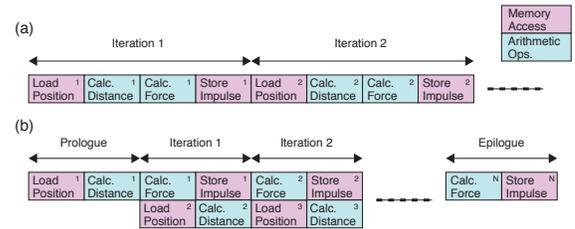


FIG. 4: (Color online) Software pipelining.

D. Software Pipelining

Software pipelining (SWP) is one of the loop optimization techniques that reforms the loop to reduce the execution time [20]. While SWP is often used to hide the latency of instructions so that pipelines are kept busy without stalling, we adopt this technique to increase the number of instructions per cycle (IPC) rather than to avoid pipeline hazards. The structure of the force calculation has a double-loop form where the outer loop is for i -atoms and the inner loop is for j -atoms. In the body of the inner loop, there are four different operations, A, B, C , and D , which are A : load the position of a j -atom, B : calculate the distance between i - and j -atoms, C : calculate the force and update the momenta, D : store the momentum of the j -atom. Suppose n is the number of j -atoms interacting with the i -atom, then the inner loop has the form $\{ABCD\}^n$. Since the four operations are mutually dependent, the operations should be performed sequentially (see Fig. 4 (a)). To increase parallelism, we change the inner loop from $\{ABCD\}^n$ to $AB\{CDAB\}^{n-1}CD$ (see Fig. 4 (b)) so that the memory accesses and the arithmetic operations in the body of the inner loop are no longer mutually dependent. Since the recent CPUs are superscalar, independent memory accesses and arithmetic operations can be executed simultaneously, increasing IPC. Since the modification of the loop causes an increase in the number of instructions, the overall performance is improved when the increase in IPC is larger than the increase in the number of in-

structions. The impact of this optimization technique, therefore, strongly depends on the environment, such as the CPU architecture, the version of the compiler, and so forth.

With the Intel 16.0.4 C++ compiler and on HSW, we find that SWP improves IPC by 55% while the number of instructions increases by 22%, and therefore, the total performance gain is about 21%. We confirm that SWP significantly improves the performance on SKL with the Intel 17.0.6 C++ compiler. However, the gain in the performance vanishes when we use the Intel 18.0.1 C++ compiler, since the performance without SWP is improved by using a later version of the compiler. The difference between the 17.0.6 and the 18.0.1 compilers is the vector optimization capability. While the 17.0.6 compiler used xmm registers only, the 18.0.1 compilers vectorized the code with using zmm registers. As the result, the code compiled by the 18.0.1 compilers is faster than that by the 17.0.6 compiler. Note that, we used the identical compiler options for both cases. If we apply SWP by hand, then the loop was not vectorized since the loop body becomes too complicated. However, the performance is improved since IPC increases. Whether the performance improves by SWP depends on the balance between performance improvement by improving IPC and performance degradation due to inhibition of vectorization. In SKL, the performances with and without SWP happened to be nearly the same. However, the performance on KNL becomes worse with SWP. This is because the auto-vectorization by the Intel compiler works efficiently for the code “SoA+Sorted” on KNL. The auto-vectorization by the compiler on KNL is discussed in Sec. V A.

IV. SIMD VECTORIZATION WITH AVX2

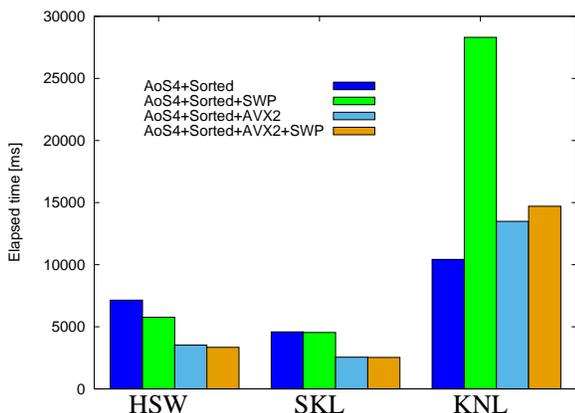


FIG. 5: (Color online) Effects of optimizations with AVX2.

In this section, we describe the vectorization with AVX2 instructions. We can use a 256-bit-width register with AVX2. Therefore, four pairs of forces can ideally

be calculated simultaneously. Since the innermost loop of the MD code contains indirect accesses, it is not straightforward to pack data from the memory to the register and to unpack data from the register to the memory. Since the force calculations of the LJ potential are relatively light, the packing and unpacking processes may be the bottleneck. While AVX2 includes a gather instruction, it is very slow. Additionally, AVX2 does not have a scatter instruction. Therefore, some ingenuities are required to pack and unpack the data. The results for the codes vectorized with AVX2 are shown in Fig. 5. In the following, we describe how to vectorize the force calculation with AVX2.

A. AoS with Padding

The key idea in the vectorization of the force calculation is to make the data structure the Array of Structures (AoS) with padding [5]. The atom data are aligned so that the data structure has 256-bit boundaries. Since the coordinates and momenta have three elements, one double precision floating number is inserted as a padding for every three elements (see Fig. 6). We refer to this layout as “AoS4” since each structure contains four elements including a padding. This data structure has two advantages over SoA. One is the cache efficiency. This data structure ensures that the data for a single atom will be on a single cache line with 256-bit width. The other is the memory-access efficiency. Since the data are aligned on 256-bit boundaries, the coordinates or momenta of a single atom can be moved to the YMM register by a single vector-load instruction (`vmoveupd`). We find that the latter advantage is more effective for increasing the speed. After loading the data of the coordinates into a register, we calculate the relative coordinate vector between i - and j -atoms (Fig. 7 (a)). Performing the above process four times, we obtain four registers containing relative coordinate vectors. Then we transpose the four vectors to obtain three vectors (Fig. 7 (b)). By calculating the sum of the squares of the vectors, we obtain a single register that contains the relative distances of the four pairs (Fig. 7 (c)).

Since we adopt the bookkeeping method, the pair list contains the pairs with distances greater than the truncation length. This list can be treated by masking operations. We first calculate impulses of all pairs. Next, we set the impulses between pairs that are outside the interaction range to zero (see Fig. 8). Then we load the momenta of the j -atoms, update them using the calculated impulses, and write them back.

The results of the above vectorization are denoted by “AoS4+Sorted+AVX2” in Fig. 5. The increase in speed compared with the scalar code with SWP, *i.e.*, the improvements from “AoS4+Sorted” to “AoS4+Sorted+AVX2”, are 51% on HSW and 44% on SKL. While the vectorization with AVX2 works effectively on HSW and SKL, the vectorized code becomes

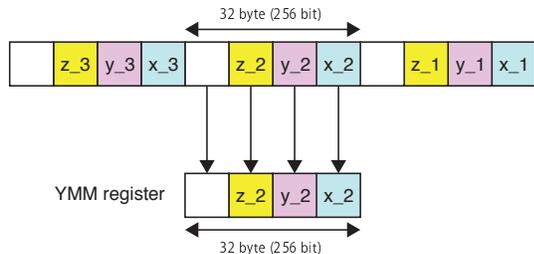


FIG. 6: (Color online) Array of Structures (AoS) with padding. This data layout allows us to load simultaneously three elements of coordinates to a YMM register.

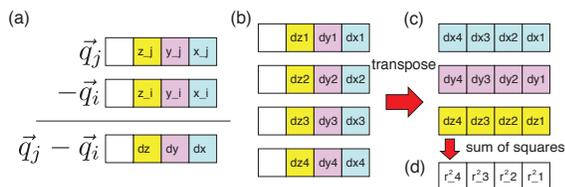


FIG. 7: (Color online) Calculation of relative distances. (a) By calculating the difference between the coordinates of the i - and j -atoms stored in each YMM register, the relative position vector is stored in the YMM register. (b) The calculation of the relative position vector is performed four times for four j -atoms. (c) The four YMM registers are transposed to obtain three registers. (d) The sums of the squares of the three registers are calculated to obtain four squared relative distances of four pairs.

significantly slower than the original code on KNL. This is because automatic vectorization by the compiler works effectively with AVX-512 on KNL. This issue will be discussed later.

B. Software Pipelining

We can apply the SWP technique to the vectorized code with AVX2 in a similar manner to the scalar code. However, the performance improvements are moderate. While the performance is improved by 6% on HSW, the effect on SKL is not significant. Compared with the fastest scalar codes which is “AoS4+Sorted+SWP”, the performance gains by vectorization are about 42% on HSW and 44% on KNL. The performance on KNL deteriorates upon applying SWP.

V. SIMD VECTORIZATION WITH AVX-512

In this section, we describe the vectorization with AVX-512 instructions. The AVX-512 includes gather, scatter, and mask operations, which are useful for the vectorization of the loop involving indirect access. Here

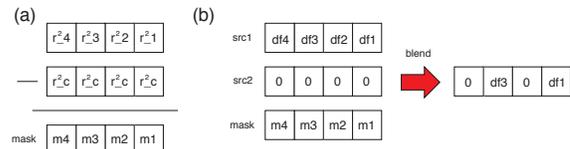


FIG. 8: Conditional execution using masks. The impulses between the pairs with distances greater than the truncation length are set to zero.

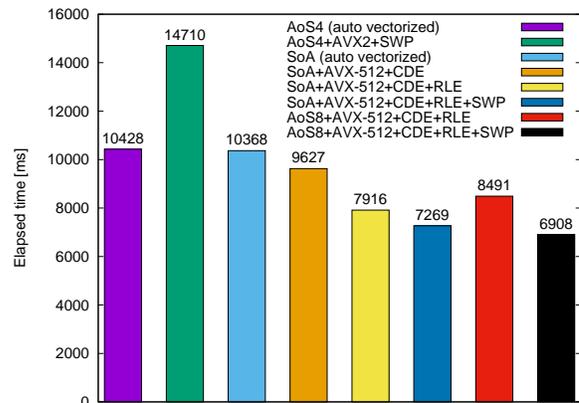


FIG. 9: (Color online) Results on KNL. The meanings of the abbreviations in the figure are AoS4 (Array of Structures with 4 members), AoS8 (Array of Structures with 8 members), SoA (Structure of Arrays), CDE (collision detection elimination), RLE (remainder loop elimination), and SWP (software pipelining).

we describe the efficiency of the optimization on KNL. The results are shown in Fig. 9. The results on SKL are discussed in the next section.

A. Auto-Vectorization by Compiler

While the Intel compiler cannot vectorize the benchmark program automatically with AVX2 without directives, it can vectorize the program with AVX-512. Actually, the compiler vectorized the code in Algorithm 3 as follows.

1. Unroll the inner loop eight times.
2. Gather the positions and momenta of j -atoms with `vgatherdpd`.
3. Calculate the forces and update the momenta of j -atoms.
4. Detect conflicts between the indices of j -atoms with `vpconflictd`.
5. Scatter the updated momenta with `vscatterdpd`.

Here, the compiler inserted the codes for conflict detection since it did not know that there would be no conflict between the indices of j -atoms interacting with i -atoms. However, the penalty for inserting conflict detection is not expensive as shown later. Inversely, the compiler cannot vectorize the code without conflict detection instructions. In that sense, the conflict detection support in AVX-512 is crucial for auto-vectorization by the compiler as suggested by Höhnerbach *et al.* [7].

The code vectorized by the compiler with AVX-512 exhibits better performance than the code vectorized by hand with AVX2. SWP reduces the performance since it interferes with the optimization by the compiler.

B. AoS to SoA

Suppose x_j is the x -coordinate of a j -atom. The distance between x_j and x_{j+1} in memory is 8 bytes for SoA, while it is 32 bytes for AoS (see Fig. 10). The gather (`vgatherdpd`) and scatter (`vscatterdpd`) instructions afford a scale factor that takes a values of 1, 2, 4, or 8. When the data layout is SoA, then the indices of j -atoms can be directly used for gather/scatter instructions with a scale factor of 8. However, bit shifts are required for AoS, and therefore, the bit shifts impose some penalty on performance. The efficiency of this optimization strongly depends on the code. Actually, the improvement of the performance by changing the data layout is 1% for the codes vectorized by the compiler. However, it improves the performance by 7% for the codes vectorized by hand with some additional optimizations, such as the remainder loop elimination and software pipelining, which are described later.

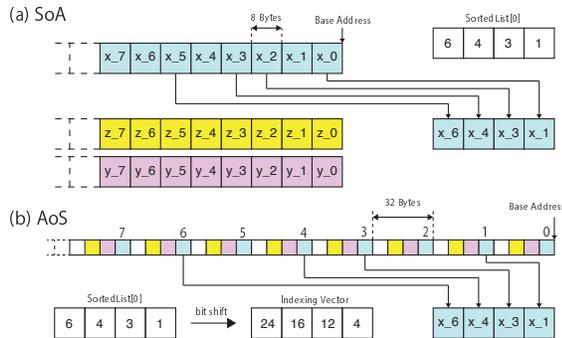


FIG. 10: (Color online) Gather procedures. While eight double-precision floating elements can be gathered with AVX-512, only four are shown for visibility. (a) Gather process for the SoA layout. By using 8 as the scale factor of the gather instruction, the data can be gathered by using the sorted list as is. (b) Gather process for the AoS layout. Since each element is separated by 32 bytes and the maximum scale factor is 8, it is necessary to carry out the gather instruction using bit-shifted indices of the sorted list.

C. Collision Detection Elimination

As described above, the compiler inserted unnecessary codes for the conflict detection. For the case of collision, the compiler prepares a code storing the momenta sequentially which will be never called. Since there are no conflicts between the indices, the conflict detection can be eliminated if we call gather and scatter instructions directly by using intrinsic functions. By this conflict detection elimination (CDE), the performance is improved by 7% (from "SoA+AVX-512+CDE" to "SoA+AVX-512+CDE+RLE"). Note that, it may be possible to tell the compiler that there are no conflicts using directives. However, as far as we tried, the collision detection codes could not be eliminated by adding directives.

D. Remainder Loop Elimination

The number of iterations of the inner loop of the force calculation is equals the number of atoms existing within the interaction distance, which depends on the density of atoms. For the benchmark condition adopted in the manuscript, the number of j -atoms interacting with one i -atom is approximately 55 to 80. If we unroll the inner loop eight times, the number of iterations of the vectorized loop kernel becomes 6 to 10. Since the number of iterations of the remaining loop is up to 7, the computational time to process the remainder loop is comparable to the time required for the vectorized loop kernel. Therefore, we eliminate the remainder loop by masking. The pseudocode for the remainder loop elimination (RLE) is shown in Algorithm 4. The key idea is to prepare a vector of the loop counter (\hat{d} in the pseudocode). Then we can make a mask \hat{m}_{loop} for the remainder loop. Note that we must also to make a mask \hat{m}_{cutoff} for the truncation of the interaction. The forces are zeroed by a mask \hat{m}_{total} which is the logical conjunction of \hat{m}_{loop} and \hat{m}_{cutoff} . The scatter of the momenta of j -atoms should be masked by \hat{m}_{loop} since there is a possibility of conflict due to the RLE. RLE improves the performance by 18% (from "SoA+AVX-512+CDE" to "SoA+AVX-512+CDE+RLE").

E. Software Pipelining

The performance of the vectorized kernel loop can be improved by the SWP technique described in Sec. III D. The performance improvement by SWP strongly depends on the optimization methods used together. While SWP does not improve the efficiency of the code vectorized with AVX2, it improves the performance of the code vectorized with AVX-512 by 9% .

Algorithm 4 Remainder loop elimination

```

1:  $\hat{t} \leftarrow \{8, 8, 8, 8, 8, 8, 8, 8\}$ 
2:  $\hat{v}_c \leftarrow \{r_c, r_c, r_c, r_c, r_c, r_c, r_c, r_c\}$ 
3: for  $i = 0$  to  $N - 1$  do
4:    $\hat{q}_i \leftarrow$  Position of  $i$ -Atom ▷ Broadcast
5:    $\hat{d} \leftarrow \{7, 6, 5, 4, 3, 2, 1, 0\}$ 
6:    $n \leftarrow |\text{SortedList}[i]|$ 
7:    $k \leftarrow 0$ 
8:    $\hat{n} \leftarrow \{n, n, n, n, n, n, n, n\}$ 
9:   while  $k < n$  do
10:     $\hat{q}_j \leftarrow$  Positions of  $j$ -Atoms ▷ Gather
11:     $\hat{v}_d \leftarrow$  Distance between  $\hat{q}_j$  and  $\hat{q}_i$ 
12:    Calculate forces between  $i$ - and  $j$ -atoms
13:     $\hat{d} \leftarrow \hat{d} + \hat{t}$ 
14:     $\hat{m}_{\text{cutoff}} \leftarrow \_mm512\_cmp\_pd\_mask(\hat{v}_c, \hat{v}_r, \_CMP\_LE\_OS)$ 
15:     $\hat{m}_{\text{loop}} \leftarrow \_mm512\_cmp\_epi64\_mask(\hat{d})$ 
16:     $\hat{m}_{\text{total}} \leftarrow \_mm512\_kand(\hat{m}_{\text{cutoff}}, \hat{m}_{\text{loop}})$ 
17:    Zero the forces with the mask  $\hat{m}$ 
18:    Store  $\vec{p}_j$  with the mask  $\hat{m}_{\text{loop}}$ 
19:     $k \leftarrow k + 8$ 
20:  end while
21:  Store  $\vec{p}_j$ 
22: end for

```

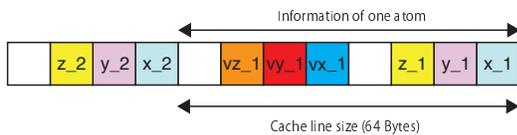
F. AoS with Eight Elements

FIG. 11: (Color online) AoS data layout with eight elements (AoS8). Each structure has eight members, three of them are coordinates, three of them are momenta, and two of them are padding, respectively. The information of one atom is placed on one cache line.

From the viewpoint of the number of instructions, SoA is more advantageous than AoS since AoS requires additional instructions to prepare the vector of indices for gather/scatter instructions. However, the AoS data layout is more advantageous than that of SoA from the viewpoint of cache efficiency. In the previous section, we introduced AoS with four elements: three of them were the coordinates or momenta of atoms and the other one was padding. Here, we adopt AoS with eight elements so that all the information of one atom is placed within one cache line (see Fig. 11). The optimal data layout depends on the optimization methods used together. For the code applying CDE and RLE, SoA is superior to AoS8. However, AoS8 exhibits better performance than SoA with SWP since SWP significantly improves the performance of the code with AoS8.

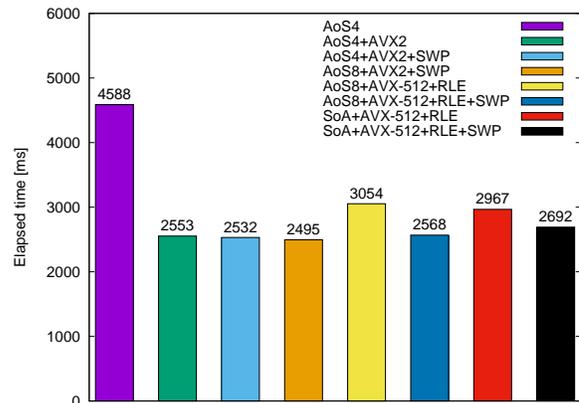
VI. RESULTS ON SKYLAKE

FIG. 12: (Color online) Results on SKL. The meanings of the abbreviations in the figure are AoS4 (Array of Structures with 4 members), AoS8 (Array of Structures with 8 members), SoA (Structure of Arrays), AVX2 (hand-vectorization with AVX2), AVX-512 (hand-vectorization with AVX-512), CDE (collision detection elimination), RLE (remainder loop elimination), and SWP (software pipelining).

The results on SKL are shown in Fig. 12. The code with the AoS having four members and sorting by indices is denoted by “AoS4”. This code is automatically vectorized by the Intel compiler using ZMM registers, but gather/scatter instructions are not used and the performance is not satisfactory. The results with “AVX2” or “AVX-512” are vectorized by hand using AVX2 or AVX-512, respectively. We find that the code with AoS that is vectorized with AVX2 exhibits the best performance, which is denoted by “AoS8+AVX2+SWP” in the figure. The performance improvement compared with the code automatically vectorized by the compiler, *i.e.*, from “AoS4” to “AoS8+AVX2+SWP”, is about 46%.

VII. SUMMARY AND DISCUSSION

In this paper, we investigated the effect of combining different data layouts and optimization methods used with vectorization for the truncated-LJ-force kernel on different architectures, HSW, KNL, and SKL. It was found that the choice of data layout is crucial for vectorization using AVX2 on HSW. While AoS8 was found to exhibit the best performance on all the architectures studied in this work, the differences in performance between AoS and SoA are less significant when the code is vectorized with AVX-512. We found that the data layout AoS8 may interfere with optimization by the compiler. In the force-calculation kernel, the coordinates are only read, but the momenta involve writing back. Since each structure in AoS8 includes both coordinates and momenta, the compiler cannot determine which members can be modified. Therefore, AoS4 is the optimal

choice from the viewpoint of portability. KNL and SKL support both AVX2 and AVX-512. While the code vectorized with AVX-512 is much faster than that vectorized with AVX2 on KNL, the code with AVX2 is slightly faster than that with AVX-512 on SKL. Although SWP is an optimization technique that can be used together with vectorization with any instruction set, the gain in performance by SWP is not significant in this work. However, we found that SWP works efficiently with another compilers such as GCC or with other architectures [12]. Although its improvement in performance strongly depends on the optimization ability of the compiler, the SWP is worth considering as an optimization method to be used with vectorization. In this work, we did not consider the compiler directives in detail since it was easier to write codes with builtin-functions than searching for effective combinations of compiler directives. This situation may change for force calculations of more complex interactions.

While we performed vectorization faithfully to the original scalar code, the performance can be improved by using reciprocal approximations. Although reciprocal approximation instructions also exist in SSE and AVX2, more accurate instructions are available in AVX-512. For example, `vrcp28pd` computes an approximate reciprocal with 28-bit accuracy. Therefore, full precision can be obtained with the first-order correction. However, `vrcp28pd` is included in the AVX-512ER instruction set, which is supported by KNL but not SKL. Instead, `vrcp14pd` is available with SKL but it requires the second-order correction to obtain full precision. Us-

ing these instructions, the performance of vectorization on SKL and KNL may be changed, which should be considered as a future issue.

While we were writing this paper, Intel announced discontinuance of Xeon Phi family [21]. Therefore, some part of this study became obsolete. However, we think that the knowledge obtained in this study is still useful. It is almost certain that the number of CPU cores and SIMD width will increase without increasing the operating frequency. Then the code vectorized with masked gather and scatter instructions are likely to be more efficient as it was on KNL, while the code vectorized with AVX2 was faster than that with AVX-512 on SKL. In any case, it is difficult to predict the optimal combination of optimization techniques, and trial and error are necessary. We hope that our work will assist vectorizations and optimizations on the future architectures.

Acknowledgements

The authors would like to thank S. Mitsunari and H. Noguchi for fruitful discussions. This work was supported by JSPS KAKENHI Grant Number 15K05201 and by the MEXT project as “Exploratory Challenge on Post-K Computer” (Frontiers of Basic Science: Challenging the Limits). The computations were carried out using the facilities of the Information Technology Center of the University of Tokyo and the Institute for Solid State Physics of the University of Tokyo.

-
- [1] B. Alder and T. Wainwright, *J. Chem. Phys.* **27**, 1208 (1957).
- [2] T. C. Germann and K. Kadau, *Int. J. Mod. Phys. C* **19**, 1315 (2008).
- [3] H. Watanabe, M. Suzuki, and N. Ito, *Comput. Phys. Commun.* **184**, 2775 (2013).
- [4] K. Rupp, <https://github.com/karlrupp/microprocessor-tide4>.
- [5] W. M. Brown, J.-M. Y. Carrillo, N. Gavhane, F. M. Thakkar, and S. J. Plimpton, *Comput. Phys. Commun.* **195**, 95 (2015).
- [6] P. Karpiński and J. McDonald, in *Proceedings of the 8th International Workshop on Programming Models and Applications for Multicores and Manycores, PMAM17* (ACM Press, 2017), pp. 21–28.
- [7] M. Höhnerbach, A. E. Ismail, and P. Bientinesi, in *SC16: International Conference for High Performance Computing, Networking, Storage and Analysis* (IEEE, 2016).
- [8] M. J. Abraham, T. Murtola, R. Schulz, S. Páll, J. C. Smith, B. Hess, and E. Lindahl, *SoftwareX* **1-2**, 19 (2015).
- [9] H. C. Edwards, C. R. Trott, and D. Sunderland, *J. Parallel Distrib. Comput.* **74**, 3202 (2014).
- [10] J. C. Phillips, R. Braun, W. Wang, J. Gumbart, E. Tajkhorshid, E. Villa, C. Chipot, R. D. Skeel, L. Kalé, and K. Schulten, *J. Comput. Chem.* **26**, 1781 (2005).
- [11] https://github.com/kaityo256/lj_simd.
- [12] H. Watanabe, M. Suzuki, and N. Ito, *Prog. Theor. Phys.* **126**, 203 (2011).
- [13] S. D. Stoddard and J. Ford, *Phys. Rev. A* **8**, 1504 (1973).
- [14] J. D. Broughton and G. H. Gilmer, *J. Chem. Phys.* **79**, 5095 (1983).
- [15] B. L. Holian, A. F. Voter, N. J. Wagner, R. J. Ravelo, S. P. Chen, W. G. Hoover, C. G. Hoover, J. E. Hammerberg, and T. D. Dontje, *Phys. Rev. A* **43**, 2655 (1991).
- [16] B. Quentrec and C. Brot, *J. Comput. Phys.* **13**, 430 (1973).
- [17] R. W. Hockney and J. W. Eastwood, *Computer Simulation Using Particles* (Taylor & Francis Ltd, 1988), ISBN 0852743920.
- [18] L. Verlet, *Phys. Rev.* **159**, 98 (1967).
- [19] M. Isobe, *Int. J. Mod. Phys. C* **10**, 1281 (1999).
- [20] V. H. Allan, R. B. Jones, R. M. Lee, and S. J. Allan, *ACM Comput. Surv.* **27**, 367 (1995).
- [21] <http://qdma.intel.com/dm/i.aspx/9C54A9A7-BF37-4496-B268-BD>