

# Development of Element-by-Element Kernel Algorithms in Unstructured Implicit Low-Order Finite-Element Earthquake Simulation for Many-Core Wide-SIMD CPUs

Kohei Fujita<sup>1,2</sup>, Masashi Horikoshi<sup>3</sup>, Tsuyoshi Ichimura<sup>1,2,4</sup>, Larry Meadows<sup>5</sup>, Kengo Nakajima<sup>6,2</sup>, Muneo Hori<sup>7</sup>, and Lalith Maddeggedara<sup>1,2</sup>

<sup>1</sup> Earthquake Research Institute & Department of Civil Engineering,  
The University of Tokyo {fujita,ichimura,lalith}@eri.u-tokyo.ac.jp

<sup>2</sup> Center for Computational Science, RIKEN

<sup>3</sup> Core and Visual Computing Group, Intel K.K masashi.horikoshi@intel.com

<sup>4</sup> Center for Advanced Intelligence Project, RIKEN

<sup>5</sup> Data Center Group, Intel Corporation lawrence.f.meadows@intel.com

<sup>6</sup> Information Technology Center, The University of Tokyo  
nakajima@cc.u-tokyo.ac.jp

<sup>7</sup> Japan Agency for Marine-Earth Science and Technology horimune@jamstec.go.jp

**Abstract.** Acceleration of the Element-by-Element (EBE) kernel in matrix-vector products is essential for high-performance in unstructured implicit finite-element applications. However, the EBE kernel is not straightforward to attain high performance due to random data access with data recurrence. In this paper, we develop methods to circumvent these data races for high performance on many-core CPU architectures with wide SIMD units. The developed EBE kernel attains 16.3% and 20.9% of FP32 peak on Intel Xeon Phi Knights Landing based Oakforest-PACS and Intel Skylake Xeon Gold processor based system, respectively. This leads to 2.88-fold speedup over the baseline kernel and 2.03-fold speedup of the whole finite-element application on Oakforest-PACS. An example of urban earthquake simulation using the developed finite-element application is shown.

**Keywords:** finite-element method · random data access · many-core · SIMD.

## 1 Introduction

Energy efficiency is one of the requirements for exascale computing, and many-core CPUs with wide SIMD units are considered to be one of the architectures towards this goal. For example, the Post-K supercomputer that is scheduled for operation around 2021 is announced to equip many-core CPUs with 512 bit wide SIMD cores [1]. However, some applications are difficult to benefit from wide SIMD and many-core architectures due to random data access and

data recurrence, and thus algorithms enabling efficient use of these architectures are important towards exascale computing. The unstructured low-order implicit finite-element method is one of the applications that are not straightforward for efficient use of many-core wide SIMD CPUs. As the unstructured low-order implicit finite-element method is the de facto standard in structural simulations used in the manufacturing industry, the acceleration of this method is expected to lead to large ripple effects. In this paper, we focus on the unstructured low-order implicit finite-element method as a target for scalable algorithm development for many-core wide SIMD CPUs.

Many studies target implicit solvers on large scale systems. For example, the SC15 Gordon Bell Prize Paper [2] uses multi-grid solvers for simulation of extremely large mantle convection problems, and the SC16 Gordon Bell Prize Paper uses the Sunway Taihu Light Supercomputer to solve climate problems [3]. Also, an SC15 Gordon Bell Prize Finalist [4] solves an urban earthquake problem on the K computer. However, all of these simulations use structured or semi-structured mesh, which are both SIMD and multi-core friendly, for efficient computation of matrix-vector products for the implicit solvers. However, for solving general manufacturing or geoscience problems, use of pure unstructured mesh is required. The state-of-the-art of low-order unstructured finite-element method on CPU based systems is SC14 Gordon Bell Prize Finalist *GAMERA* [5]. Targeting architectures with high arithmetic capabilities and relatively low memory bandwidth, *GAMERA* uses a matrix-free kernel for computation of matrix vector products (i.e., the Element-by-Element kernel, EBE [6]). Here, instead of storing and reading the global matrix from memory, element matrices are generated and multiplied with the right hand side vector every time a matrix-vector product is called. This enables fast computation of matrix-vector products with good load balance on massively parallel systems with relatively low memory bandwidth, which lead to high application performance and high scalability. However, as *GAMERA* is targeted for the K computer which comprises 8 core CPUs with 2 wide SIMD cores [7], the computation algorithm used for the EBE kernel does not consider the wider SIMD and larger core counts of current and near future many-core CPUs. Thus, a new time-parallel computation algorithm *GHYDRA* was developed in a 2018 IPDPS paper [8], that enables use of uniformity of mesh in time domain to reduce random data access. This algorithm is suitable for recent wide SIMD CPUs; however, EBE kernel algorithms suitable for this type of architecture is required to fully exploit its performance. In this paper, we develop EBE kernel algorithms for efficient computation of the *GHYDRA* finite-element solver algorithm on many-core CPU systems, and compare with conventional EBE algorithms on the K computer, Intel Xeon Phi Knights Landing based Oakforest-PACS, and an Intel Skylake Xeon Gold CPU based system. We also show an urban earthquake simulation example using the unstructured finite-element solver accelerated by the developed EBE kernel algorithms.

## 2 Finite-element solver algorithm

### 2.1 Target problem

Earthquake simulations involve large domain nonlinear time-evolution problems with locally complex structures. For example, soft soil near surface have complex geometry and softens during strong ground motion under large earthquakes. Thus, we solve the target dynamic nonlinear continuum mechanics problem using a nonlinear dynamic 3D finite-element method with low-order solid elements because this method is suitable for modeling complex geometry and analytically satisfies the traction-free boundary condition at the surface. The target equation using Newmark- $\beta$  method ( $\beta = 1/4$ ,  $\delta = 1/2$ ) for time integration is:

$$\mathbf{A}^n \delta \mathbf{u}^n = \mathbf{b}^n, \quad (1)$$

where

$$\begin{cases} \mathbf{A}^n = \frac{4}{dt^2} \mathbf{M} + \frac{2}{dt} \mathbf{C}^n + \mathbf{K}^n, \\ \mathbf{b}^n = \mathbf{f}^n - \mathbf{q}^{n-1} + \mathbf{C}^n \mathbf{v}^{n-1} + \mathbf{M} \left( \mathbf{a}^{n-1} + \frac{4}{dt} \mathbf{v}^{n-1} \right), \end{cases}$$

and

$$\begin{cases} \mathbf{q}^n = \mathbf{q}^{n-1} + \mathbf{K}^n \delta \mathbf{u}^n, \\ \mathbf{u}^n = \mathbf{u}^{n-1} + \delta \mathbf{u}^n, \\ \mathbf{v}^n = -\mathbf{v}^{n-1} + \frac{2}{dt} \delta \mathbf{u}^n, \\ \mathbf{a}^n = -\mathbf{a}^{n-1} - \frac{4}{dt} \mathbf{v}^{n-1} + \frac{4}{dt^2} \delta \mathbf{u}^n. \end{cases} \quad (2)$$

Here,  $\delta \mathbf{u}$ ,  $\mathbf{u}$ ,  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{f}$  are incremental displacement, displacement, velocity, acceleration, and external force vectors, respectively.  $\mathbf{M}$ ,  $\mathbf{C}$ , and  $\mathbf{K}$  are the consistent mass, damping, and stiffness matrices, respectively.  $dt$  is the time increment and  $n$  is the time step. We use Rayleigh damping for  $\mathbf{C}$ , where the element damping matrix  $\mathbf{C}_e^n$  is calculated using the element consistent mass matrix  $\mathbf{M}_e$  and the element stiffness matrix  $\mathbf{K}_e^n$  as  $\mathbf{C}_e^n = \alpha \mathbf{M}_e + \beta \mathbf{K}_e^n$ . Here,  $\alpha$  and  $\beta$  are obtained by solving the following least squares equation:

$$\text{minimize} \left[ \int_{f_{\min}}^{f_{\max}} \left( h - \frac{1}{2} \left( \frac{\alpha}{2\pi f} + 2\pi f \beta \right) \right)^2 df \right],$$

where,  $f_{\max}$ ,  $f_{\min}$ , and  $h$  are the maximum and minimum target frequency and the damping ratio, respectively. Although arbitrary constitutive models can be used, we use the Ramberg Osgood model [9] and Masing rule [10] for the nonlinear constitutive modeling in urban earthquake simulations. In summary, time history response  $\mathbf{u}^n$  is computed by repeating the following steps.

1. Read boundary conditions.
2. Evaluate  $\mathbf{C}^n$  and  $\mathbf{K}^n$  based on the constitutive relations and strain at time-step  $n - 1$ .
3. Obtain  $\delta \mathbf{u}^n$  by solving Eq. 1.
4. Update Eq. 2 using  $\delta \mathbf{u}^n$ .

Since most of the computation cost is incurred in solving Eq. 1, we explain the details of the solver in the next subsection.

Standard solver algorithm	Time-parallel solver algorithm
<pre> 1: set <math>x_{-1} \leftarrow 0</math> 2: for( <math>i = 0; i &lt; n; i = i + 1</math> ){ 3:   guess <math>\bar{x}_i</math> using standard predictor 4:   set <math>b_i</math> using <math>x_{i-1}</math> 5:   solve <math>x_i \leftarrow A^{-1}b_i</math> with error tolerance <math>\frac{ Ax_i - b_i }{ b_i } &lt; \epsilon</math> using initial solution <math>\bar{x}_i</math> <b>(Computed using iterative solver with EBE kernel (1 vector))</b> 6: }</pre>	<pre> 1: set <math>x_{-1} \leftarrow 0</math> and <math>\bar{x}_i \leftarrow 0</math> (<math>i = 0, \dots, m - 2</math>) 2: for( <math>i = 0; i &lt; n; i = i + 1</math> ){ 3:   set <math>b_i</math> using <math>x_{i-1}</math> 4:   guess <math>\bar{x}_{i+m-1}</math> using standard predictor 5:   <math>b_j = \bar{b}_j</math> 6:   while (<math>\frac{ Ax_i - b_i }{ b_i } &gt; \epsilon</math>) do { 7:     guess <math>\bar{b}_j</math> using <math>\bar{x}_{j-1}</math> (<math>j = i + 1, \dots, i + m - 1</math>) 8:     refine solution <math>\{\bar{x}_j \leftarrow A^{-1}\bar{b}_j\}</math> with initial solution <math>\bar{x}_j</math> (<math>j = i, \dots, i + m - 1</math>) <b>(Computed using iterative solver with EBE kernel (m vectors))</b> 9:   } 10: } 11: <math>x_i = \bar{x}_i</math> 11: }</pre>

**Fig. 1.** Standard and time-parallel solver algorithm. The same solution is obtained within error tolerance  $\epsilon$  using both of the solver algorithms.

## 2.2 Overview of time-parallel finite-element solver algorithm

Solving Eq. 1 is a common target problem arising from low-ordered implicit finite-element methods used in solid continuum mechanics problems in engineering and science fields. In these low-ordered (e.g., second-ordered) implicit simulations, iterative solvers such as Conjugate Gradient (CG) methods are generally used. Since these solvers involve much random access and intensive memory access, it is difficult to attain performance on current systems.

The time-parallel algorithm in *GHYDRA* aims to accelerate this solver by reducing random access by using the uniformity of mesh in the time domain [8]. As shown in Fig. 1, instead of solving one time step at a time, several time steps are solved together in parallel. When solving  $m$  time steps together, the arithmetic count per iteration of the iterative solver becomes  $m$  times larger. However, as the solutions for future time steps can be used as accurate initial solutions, the total number of iterations becomes approximately  $1/m$  (see [8] for detailed comparison of convergence of this method). Thus, the total arithmetic count does not change by using this algorithm. However, since the mesh connectivity does not change in time, the random access involved in the kernel used to compute matrix-vector products can be reduced, which leads to shorter time-to-solution of the application with the same solution within the solver error threshold  $\epsilon$ . This time-parallel algorithm is different from parallel-in-time integration methods intended to improve the time-to-solution at the strong scaling limits, as the *GHYDRA* algorithm uses the uniformity of mesh in the time domain to reduce random data access and thus enables speedup of the application even within the strong scaling limits.

In *GHYDRA*, the time-parallel algorithm is combined with inexact preconditioned conjugate gradient method, multi-grid method, and mixed precision computation. Below we briefly summarize these components:

- Inexact preconditioned conjugate gradient method [11]: An inexact preconditioned conjugate gradient method is a preconditioned conjugate gradient method which uses another solver for solving the preconditioning matrix

equation instead of multiplying a fixed preconditioning matrix. As it is common to use another iterative solver for solving the preconditioning matrix equation, we call the original CG loop as the “outer loop” and the preconditioning solver as the “inner loop”. As the inner loop does not need to be solved exactly, this makes room for making improvements such as multi-grid and mixed precision arithmetic.

- Multi-grid method: We reduce the computation cost and communication size of the inexact preconditioned conjugate gradient method by using a multi-grid solver for the inner loop. Here we use a two step geometric multi-grid; we use the targeted second-ordered tetrahedral mesh FEmodel for the inner fine loop, and use the same mesh without intermediate nodes of each element for the inner coarse loop (i.e., a first-ordered tetrahedral mesh FEmodel<sub>c</sub>). We use a  $3 \times 3$  block Jacobi preconditioned conjugate gradient solver for solving the inner fine and inner coarse loops.
- Mixed precision arithmetic: Even for a target problem in FP64, we need only to solve the preconditioning matrix equation roughly. Thus, we compute the whole multi-grid preconditioner in FP32. This leads to halving memory access and communication size, and enable usage of high-performance FP32 arithmetic units equipped in recent CPUs.

When using the algorithm above, we can move most of the compute cost of the outer loop to the inner loops by using appropriate threshold values in the inner CG solvers. Thus, even though it may seem a complicated algorithm, the performance of the whole solver is dependent on kernels such as matrix-vector product kernels, inner product kernels, or saxpy kernels which are in common with standard CG solvers. In the following, we explain the most time consuming matrix-vector product kernel used in the inner fine loop, which is not straightforward to attain performance on many-core wide SIMD CPUs.

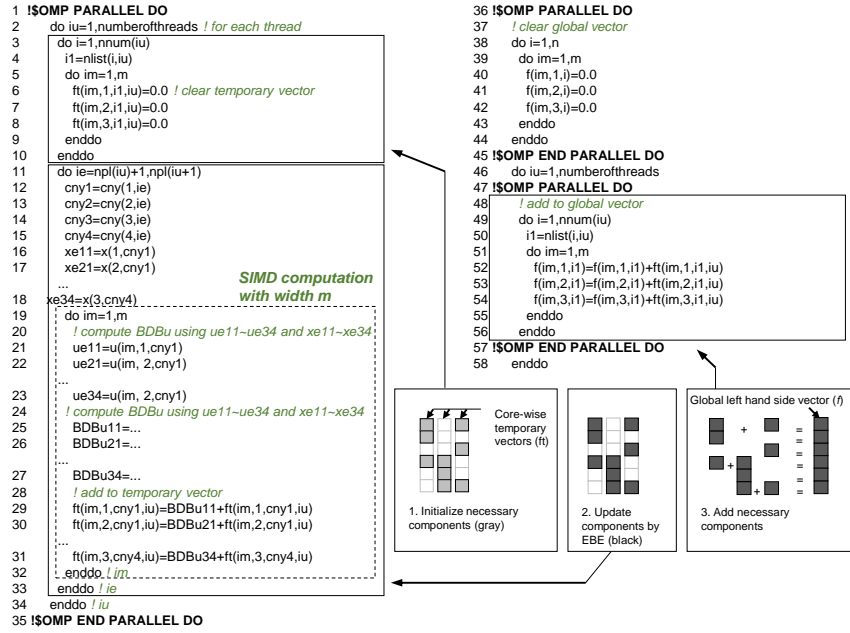
### 2.3 Baseline Element-by-Element kernel algorithm

As the relative memory transfer capability to floating point computation capability is becoming lower, using an algorithm that can reduce memory access is becoming important for fast time-to-solution. Thus, the Element-by-Element (EBE) method, which is a matrix-free matrix-vector multiplication method, is used for computation of  $\mathbf{f} = \mathbf{A}\mathbf{u}$ . Here,  $\mathbf{u}, \mathbf{f}$  are displacement and nodal force vectors, and  $\mathbf{A}$  is the global matrix that is generated by superimposing element matrices  $\mathbf{A}_i^e$ .  $\mathbf{A}_i^e$  can be computed by using coordinates of nodes ( $\mathbf{x}$ ) and element material properties. In standard matrix-vector product methods,  $\mathbf{f} = \mathbf{A}\mathbf{u}$  is computed by reading the global matrix  $\mathbf{A}$  from memory and multiplying it with  $\mathbf{u}$ . In the contrary, in the EBE method, matrix-vector products are computed by computing local matrix-vector products

$$\mathbf{f}_i^e = \mathbf{A}_i^e \mathbf{u}_i^e = \mathbf{A}_i^e \mathbf{Q}_i^{eT} \mathbf{u}, \quad (3)$$

and adding them up as

$$\mathbf{f} = \sum_i \mathbf{Q}_i^e \mathbf{f}_i^e. \quad (4)$$

Fig. 2. Baseline EBE kernel with  $m$  vectors

Here,  $\mathbf{Q}_i^e$  are matrices for mapping element-wise nodal values to global nodal values. Since the coordinates, displacement and force vectors  $\mathbf{x}$ ,  $\mathbf{u}$ ,  $\mathbf{f}$  can be kept on cache by reordering of nodes and elements, we can drastically reduce memory access when compared with methods reading the global matrix from memory. On the other hand, the EBE method involves more computation and consequently transfers memory access cost to computation cost.

Although the EBE method has low algorithmic Byte/FLOP and is potentially suitable for current computers, it is not straightforward for parallel computation due to the data recurrence for adding local force vectors of elements with shared nodes in Eq. 4. Figure 2 shows the multi-core EBE computation algorithm for *GHYDRA*. Here, temporary vectors  $\mathbf{ft}$  are allocated for each core, initialized on lines 3-10, and the core wise results are added to  $\mathbf{ft}$  in lines 11-33. Finally, the core-wise results are added to the global vector  $\mathbf{f}$  in lines 46-58. The innermost loop for time-parallelism ( $m$ ) can be computed using SIMD units. However, in practice, we use  $m = 4$  in earthquake problems as increasing  $m$  leads to deterioration of the accuracy of the predicted solutions and thus increase in the total number of arithmetic counts. Thus, only 1/4 of the SIMD lanes of 512 bit SIMD registers (i.e., 4 out of the 16 FP32 lanes) can be used. Although this is better than the non-time-parallel algorithm in which SIMD cannot be used completely due to the data recurrence involved in lines 29-31, this EBE imple-

```

1 ISOMP PARALLEL DO
2 do iu=1,numberofthreads ! for each thread
3 do i=1,nnum(iu)
4 i1=nilist(i,iu)
5 do im=1,m
6 ft(im,1,i1,iu)=0.0 ! clear temporary vector
7 ft(im,2,i1,iu)=0.0
8 ft(im,3,i1,iu)=0.0
9 enddo
10 enddo
11 ! block loop with blocksize NL/m
12 do ieo=np(iu)+1,np(iu+1),NL/m
13 ! load ue, xe
14 do ie=1,min(NL/m,np(iu+1)-ieo+1)
15 cny1=cny(1,ieo+ie-1)
16 cny2=cny(2,ieo+ie-1)
17 cny3=cny(3,ieo+ie-1)
18 cny4=cny(4,ieo+ie-1)
19 do im=1,m
20 ue11(im+(ie-1)*m)=u(im,1,cny1) SIMD
21 ue21(im+(ie-1)*m)=u(im,2,cny1) (width=m)
22 ... computation
23 ue34(im+(ie-1)*m)=u(im,3,cny4)
24 xe11(im+(ie-1)*m)=x(1,cny1)
25 ...
26 xe34(im+(ie-1)*m)=x(3,cny4)
27 enddo
28 enddo
29 ! compute BDBu
30 do ie=1,NL
31 BDBu11(ie)=... SIMD computation
32 BDBu21(ie)=...
33 ...
34 BDBu34(ie)=...
35 enddo
36 ! add to global vector
37 do ie=1,min(NL/m, np(iu+1)-ieo+1)
38 cny1=cny(1,ieo+ie-1)
39 cny2=cny(2,ieo+ie-1) SIMD
40 cny3=cny(3,ieo+ie-1) (width=m)
41 cny4=cny(4,ieo+ie-1) computation
42 do im=1,m
43 ft(im,1,cny1,iu)=BDBu11(im+(ie-1)*m)+f(im,1,cny1,iu)
44 ft(im,2,cny1,iu)=BDBu21(im+(ie-1)*m)+f(im,2,cny1,iu)
45 ...
46 ft(im,3,cny4,iu)=BDBu34(im+(ie-1)*m)+f(im,3,cny4,iu)
47 enddo
48 enddo ! ieo
49 ISOMP END PARALLEL DO
50 Add ft in to f (same as lines 36-58 of Fig. 2)

```

Fig. 3. EBE kernel with  $m$  vectors for wide-SIMD CPUs

mentation is expected to lead to poor performance on recent many-core CPUs with wider SIMD units. Development of EBE kernel algorithms that uses SIMD in full width is required for improving performance.

### 3 Developed EBE kernel algorithms for many-core wide SIMD CPUs

In this section, we show the developed algorithms for faster EBE computation on many-core wide SIMD CPUs.

We first develop an algorithm to utilize SIMD units in the main computation loop (Fig. 3). As the data recurrence in addition of results into the  $\mathbf{ft}$  vector was blocking the use of SIMD units, we split the innermost loop into two loops; i.e., the computation part of  $\mathbf{f}_i^e = \mathbf{A}_i^e \mathbf{u}_i^e$  (Fig. 3 lines 14-33) and the summation part of  $\mathbf{f}_i^e$  to  $\mathbf{f}$  (lines 35-45). This leads to use of SIMD operations to the computationally rich first loop. Here, we reduce the overhead of loop splitting by blocking the loop with small block size ( $NL$ , which is typically set to the SIMD vector length). This keeps the temporary buffers  $BDBu11-34$  on cache.

When using many cores, the size of the thread-wise temporary buffer  $\mathbf{ft}$  becomes large. In addition, the overhead of initializing  $\mathbf{ft}$  and adding the thread-wise results to the global vector  $\mathbf{f}$  is not negligible in the total kernel cost. Thus, we developed a thread partitioning method to eliminate the use of thread-wise temporary buffers (Fig. 4). The procedure in standard coloring methods are to color the whole mesh such that elements in each color does not have shared nodes (Fig. 5a). Instead, we partition the domain with prescribed thread numbers

```

1  !$OMP PARALLEL DO
2  !clear global vector
3  do i=1,n
4  do im=1,m
5  f(im,1,i)=0.0
6  f(im,2,i)=0.0
7  f(im,3,i)=0.0
8  enddo
9  enddo
10 !$OMP END PARALLEL DO
11 do icolor=1,ncolor !for each color or element set
12 !$OMP PARALLEL DO
13 do iu=1, numberofthreads
14 ! block loop with blocksize NL/m
15 do ieo=np1(icolor,iu)+1,np1(icolor,iu+1),NL/m
16 ! load ue, xe
17 do ie=1,min(NL/m,np1(icolor,iu+1)-ieo+1)
18 cny1=cny(1,ieo+ie-1)
19 cny2=cny(2,ieo+ie-1)
20 cny3=cny(3,ieo+ie-1)
21 cny4=cny(4,ieo+ie-1)
22 do im=1,m
23 ue11(im+(ie-1)*m)=u(im,1,cny1)
24 ue21(im+(ie-1)*m)=u(im,2,cny1)
25 ...
26 ue34(im+(ie-1)*m)=u(im,3,cny4)
27 xe11(im+(ie-1)*m)=x(1,cny1)
28 xe21(im+(ie-1)*m)=x(2,cny1)
29 ...
30 xe34(im+(ie-1)*m)=x(3,cny4)
31 enddo
32 enddo
33 enddo
34 enddo
35 enddo
36 enddo
37 enddo
38 !add to global vector
39 do ie=1,min(NL/m, np1(icolor,iu+1)-ieo+1)
40 cny1=cny(1,ieo+ie-1)
41 cny2=cny(2,ieo+ie-1)
42 cny3=cny(3,ieo+ie-1)
43 cny4=cny(4,ieo+ie-1)
44 do im=1,m
45 f(im,1,cny1)=BDBu11(im+(ie-1)*m)+f(im,1,cny1)
46 f(im,2,cny1)=BDBu21(im+(ie-1)*m)+f(im,2,cny1)
47 ...
48 f(im,3,cny4)=BDBu34(im+(ie-1)*m)+f(im,3,cny4)
49 enddo
50 enddo ! ieo
51 enddo ! iu
52 !$OMP END PARALLEL DO
53 enddo ! icolor

```

**Fig. 4.** Coloring/thread partitioning of EBE kernel with  $m$  vectors for wide-SIMD CPUs. The same code can be used for both coloring/thread partitioning methods shown in Fig. 5.

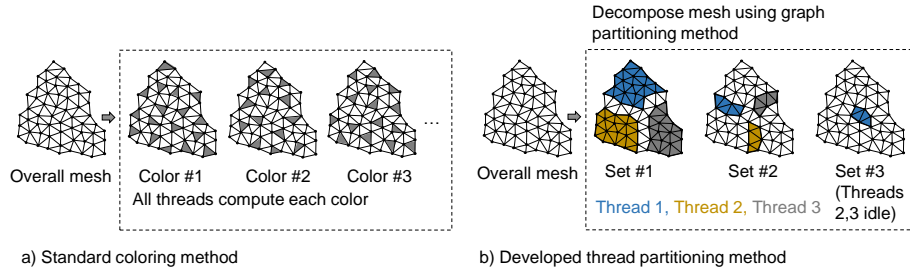
(Fig. 5b). Here, we partition the overall mesh into the number of threads using a graph partitioning method (METIS [12]). Then, we remove elements that share nodes between the thread partitions. The remaining elements becomes the first set of elements to be computed (Set #1 of Fig. 5b). The removed elements are partitioned again to recursively decompose the mesh into sets (Set #2 and #3 of Fig. 5b). As will be shown in Section 4, the proposed thread partitioning method enables better cache reuse of nodal values  $\mathbf{u}$ ,  $\mathbf{x}$  and  $\mathbf{f}$  when compared with the case of standard coloring methods.

In practice, we use  $m = 4$  in earthquake problems as increasing  $m$  leads to deterioration of the accuracy of the predicted solutions and thus increase in the total number of arithmetic counts. On the other hand, 16 FP32 floats can be packed in a single 512 bit SIMD register; thus, we can only use the first 1/4 of the SIMD capability in the vector load/store computation. In order to accelerate this part, we used 16 wide SIMD for loading from  $\mathbf{u}$  while using 4 wide SIMD for storing to local vectors  $\mathbf{ueXX}$  (Fig. 6a). Similarly, we used 4 wide SIMD for loading local vectors  $\mathbf{BDBuXX}$  while using 16 wide SIMD for loading and storing  $\mathbf{f}$  (Fig. 6b). This enables reduction of inefficient 4 wide SIMD accesses and thus improvement in overall kernel performance.

## 4 Performance Measurements

We measure performance of the developed kernels and the accelerated finite-element application on the K computer, Intel Xeon Phi Knights Landing based





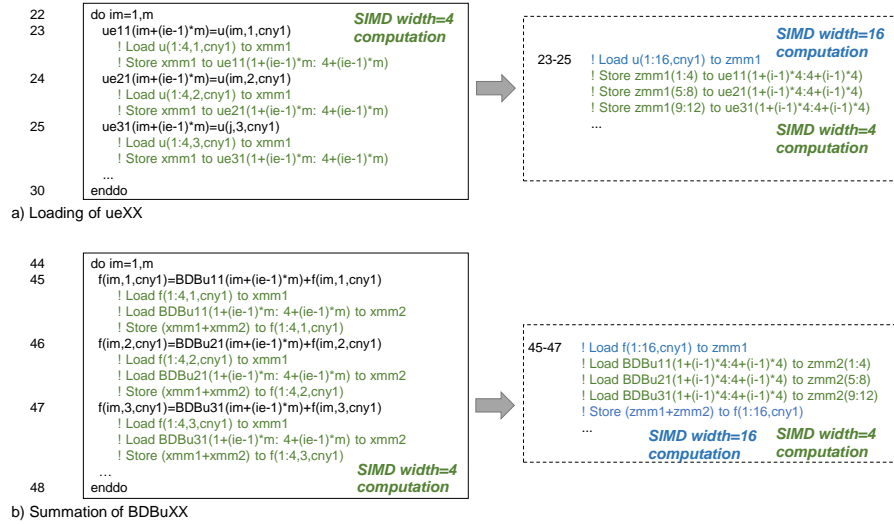
**Fig. 5.** Coloring/thread partitioning methods for the EBE kernel

**Table 1.** Performance measurement environment. Memory bandwidth on OFP is measured values of STREAM benchmark, while others are hardware peak values.

	K computer	Oakforest-PACS	Intel Skylake Xeon Gold based server
Nodes	8	1	1
Sockets/node	1	1	2
Cores/socket	8	68	20
FP32 SIMD width	2	16	16
Clock frequency	2.0 GHz	1.4 GHz	2.4 GHz
Total peak FP32 FLOPS	1024 GFLOPS	6092 GFLOPS	6144 GFLOPS
Total DDR bandwidth	512 GB/s	80.1GB/s	255.9 GB/s
Total MCDRAM bandwidth	-	490 GB/s	-

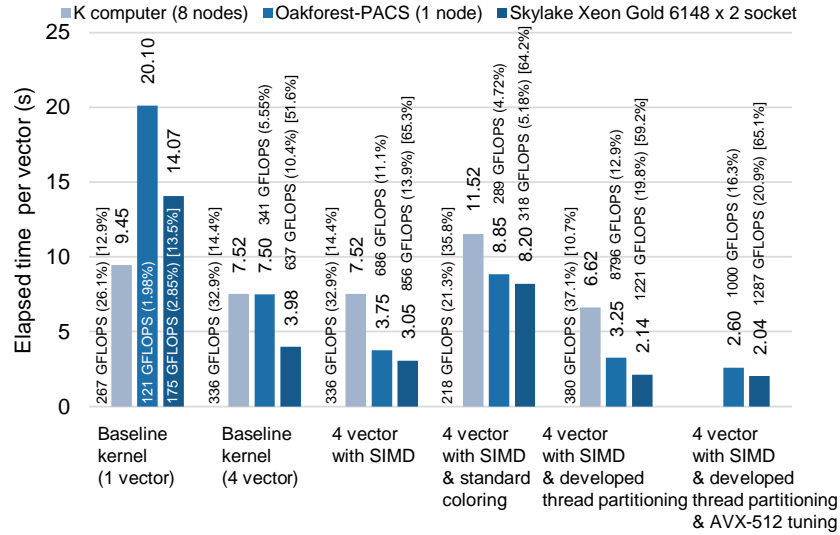
Oakforest-PACS, and an Intel Skylake Xeon Gold CPU based system. The latter two systems are examples of many-core systems with 512-bit SIMD. Table 1 summarizes the configuration of the systems. K computer (K) [7] consists of 82,944 compute nodes, each with a single eight-core SPARC64 VIIIfx CPU. Each core has two sets of SIMD FMA arithmetic units of width 2. Oakforest-PACS (OFP) [13] is a supercomputer system introduced by the Joint Center for Advanced HPC, which was established by the University of Tokyo and the University of Tsukuba. The system comprises 8,208 nodes with a 68-core Intel Xeon Phi 7250 (Knights Landing) CPU [14], 96 GB of DDR4 RAM, and 16 GB of stacked 3D MCDRAM. Intel Skylake Xeon Gold CPU based system comprises two 20-core Intel Skylake Xeon Gold 6148 CPU [16] and 192 GB of DDR4 RAM. Both the Xeon Phi 7250 CPU and Skylake Xeon Gold 6148 CPU support AVX-512 SIMD instructions.

We first measure performance of the EBE kernel. For the K computer, we use 8 nodes each with 1 MPI process with 8 OpenMP threads per process. For OFP, we use quadrant/flat mode, and run 8 MPI processes with 8 OpenMP threads per process on a single node. Each thread is bound to one CPU core (hyperthreading is not used), and we use `numactl --preferred=1` option such that memory is preferentially allocated to MCDRAM. For the Skylake system,



**Fig. 6.** Tuning of  $m$  vector EBE kernel for 512 bit SIMD architecture. Here, `xmm` indicate 128 bit FP32 registers and `zmm` indicate 512 bit FP32 registers.

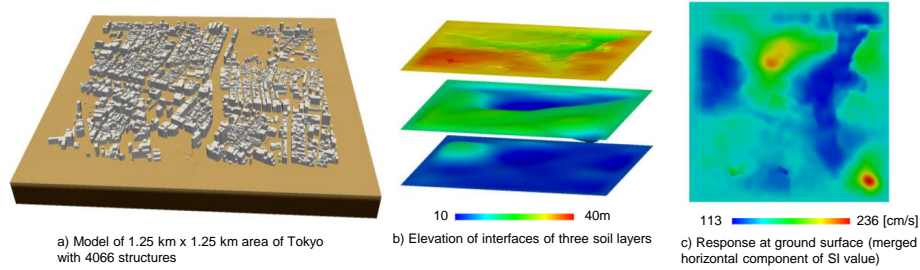
we use 8 MPI processes with 5 OpenMP threads per process on a single node. Each MPI process on K/OFP/Skylake runs on a mesh with 6,534,144 second-order tetrahedral elements and 9,044,560 nodes and elapsed time for computing 529 times of matrix-vector products is measured. Figure 7 shows the EBE kernel performance. Comparing the K computer and OFP for the baseline kernel with one vector, we can see that high performance of 26.1% FP32 peak efficiency is attained for the K computer; however, only 1.98% is attained on OFP. The kernel performance is improved by using the time-parallel algorithm with  $m = 4$ ; a 9.45 s / 7.52 s = 1.25-fold speedup was obtained on K computer and 20.1 s / 7.50 s = 2.68-fold speedup was obtained on OFP, respectively, for elapsed time per vector. This performance is expected to be further improved by using the developed EBE algorithms. First we see the effectiveness of the kernel algorithm enabling use of SIMD for the main computation part. By using SIMD by loop splitting and blocking for the main computation part of the EBE kernel, we can shorten elapsed time from 7.50 s to 3.75 s on OFP. Note that loop splitting and blocking is disabled for the K computer as the 2 wide FP32 SIMD of the K computer can be directly applied to the innermost time-parallel  $m$  loop. Next we see the effectiveness of using coloring and thread aware partitioning. When using standard coloring generated by serial greedy algorithm, the mesh was decomposed into 41 colors, resulting in total of 263.9 GB memory transfer per node incurred for computation of the EBE kernel on the K computer. This is more memory transfer than the baseline algorithm with 69.3 GB memory transfer, as the cache reuse for `u`, `x` and `f/ft` is disabled; which lead to longer



**Fig. 7.** EBE kernel performance. Elapsed time per vector is shown. Numbers in brackets indicate efficiency to FP32 peak, and numbers in square brackets indicate memory bandwidth efficiency to hardware peak of K computer and Skylake system. 8 nodes of the K computer (1 MPI process  $\times$  8 OpenMP threads per node), 1 node of OFP (8 MPI processes  $\times$  8 OpenMP threads), and a Skylake system (8 MPI processes  $\times$  5 OpenMP threads) is used for computation.

elapsed time. When using the developed thread partitioning method, the mesh was decomposed into 5 colors. With more cache reuse, the total memory transfer was reduced to 45.3 GB, which resulted to 7.52 s / 6.62 s = 1.13-fold speedup on the K computer and 3.75 s / 3.25 s = 1.15-fold speedup on OFP. Furthermore, by using the AVX-512 tuning for reducing random data access, the elapsed time was decreased to 2.60 s on OFP. This lead to high performance of 16.3% of peak FP32 performance on 64 cores of OFP’s Intel Xeon Phi Knights Landing CPU. We can see that the effective use of SIMD and circumventing random access to temporary vectors lead to significant speedup of 7.50 s / 2.60 s = 2.88-fold from the baseline algorithm with  $m = 4$  vectors on OFP. Compared with the baseline algorithm without time-parallelism ( $m = 1$ ), this is 20.1 s / 2.60 s = 7.73-fold speedup on OFP. The developed EBE algorithm is also effective for the Skylake system, and lead to 3.98 s / 2.04 s = 1.95-fold from the baseline algorithm with  $m = 4$  vectors. Compared with the baseline algorithm without time-parallelism ( $m = 1$ ), this is 14.1 s / 2.04 s = 6.89-fold speedup on the Skylake system.

Next we see the performance of developed application on an urban earthquake problem targeting a 1.25 km  $\times$  1.25 km area of Tokyo. Using digital elevation map of [17] and soil information of [18], we constructed a soil model consisting of three soil layers. We discretized this problem with minimum element size of 1 m, which lead to a problem size of 1,022,620,536 degrees-of-freedom, 252,738,195

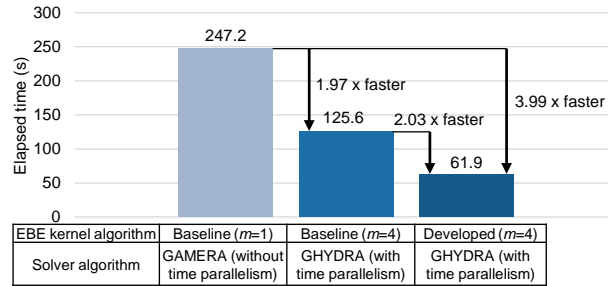


**Fig. 8.** Urban earthquake problem settings and computation results. Earthquake wave propagation in the three-layered soil structure is computed in this application example.

second-order tetrahedral elements and 340,873,512 nodes (Fig.8a,b). We input the wave observed during 1995 Kobe Earthquake [19] with time stepping  $dt = 0.01$  s. The mesh is partitioned with METIS into 1,152 partitions and computed using 144 nodes of OFP (8 MPI processes per node). Figure 9 shows the elapsed time for solving the first 25 time steps using *GHYDRA* with the baseline EBE kernel and *GHYDRA* with the developed EBE kernel (both with  $m = 4$ ). We also show performance of *GAMERA*, which can be considered as the non-time parallel version of *GHYDRA* with the baseline EBE kernel ( $m=1$ ). We can see that the application was accelerated by  $125.6 \text{ s} / 61.9 \text{ s} = 2.03$ -fold by using the developed EBE kernel algorithms leading to high peak performance of 11.6% of FP64 peak. Together with the  $247.2 \text{ s} / 125.6 \text{ s} = 1.97$ -fold speedup using time-parallelism, the application was accelerated by total of 3.99-fold when compared to the SC14 Gordon Bell Prize Finalist solver *GAMERA*. We can see that the acceleration of the EBE kernel with algorithms suitable for wide SIMD and multi-cores are effective for attaining high performance and shorter time-to-solution of the overall application. With less energy-to-solution, we can afford to conduct more detailed simulations of larger areas of interest, which is expected to contribute to earthquake disaster mitigation.

## 5 Closing remarks

In this paper we developed algorithms to accelerate the Element-by-Element kernel in unstructured low-order implicit finite-element methods on systems with many-core wide SIMD CPUs. By using the developed algorithm on the Intel Xeon Phi Knights Landing based Oakforest-PACS system, the elapsed time of the EBE kernel and total unstructured finite-element application was accelerated by 2.88-fold and 2.03-fold, respectively. The developed EBE kernel algorithms were also effective for the K computer and an Intel Skylake Xeon Gold CPU based system. These insights are expected to enable high performance on other large-scale many-core wide CPU based supercomputer systems, enabling energy efficient finite-element computation towards exascale computing.



**Fig. 9.** Elapsed time for solving the first 25 time steps of application problem. Measured on 144 nodes of OFP.

## Acknowledgments

Our results were obtained using the Oakforest-PACS at the Joint Center for Advanced HPC and the K computer at the Center for Computational Science, RIKEN (hp170249, hp180217). We acknowledge support from the Japan Society for the Promotion of Science (17K14719, 26249066, 25220908, 18H05239).

## References

1. Outline of the Development of the Post-K computer [Online]. <https://www.r-ccs.riken.jp/en/postk/project/outline>
2. J. Rudi, A. C. I. Malossi, T. Isaac, G. Stadler, M. Gurnis, P. W. J. Staar, Y. Ineichen, C. Bekas, A. Curioni, and O. Ghattas. “An Extreme-Scale Implicit Solver for Complex PDEs: Highly Heterogeneous Flow in Earth’s Mantle.” Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC’15). ACM, 2015, p. 5.
3. C. Yang, W. Xue, H. Fu, H. You, X. Wang, Y. Ao, F. Liu, L. Gan, P. Xu, L. Wang, G. Yang, and W. Zheng. “10M-Core Scalable Fully-Implicit Solver for Nonhydrostatic Atmospheric Dynamics.” Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC’16). IEEE Press, 2016, p. 6.
4. T. Ichimura, K. Fujita, P. E. B. Quinay, L. Madgededara, M. Hori, S. Tanaka, Y. Shizawa, H. Kobayashi, and K. Minami. “Implicit Nonlinear Wave Simulation with 1.08T DOF and 0.270T Unstructured Finite Elements to Enhance Comprehensive Earthquake Simulation.” Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC’15). ACM, 2015, p. 4.
5. T. Ichimura, K. Fujita, S. Tanaka, M. Hori, M. Lalith, Y. Shizawa, and H. Kobayashi. “Physics-based urban earthquake simulation enhanced by 10.7 BlnDOF x 30 K time-step unstructured FE non-linear seismic wave simulation.” Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis (SC’14). IEEE Press, 2014, pp. 15–26.

6. J. M. Winget and T. J. Hughes. "Solution algorithms for nonlinear transient heat conduction analysis employing element-by-element iterative strategies." *Computer Methods in Applied Mechanics and Engineering*, **52**(1-3), 1985, pp. 711-815.
7. H. Miyazaki, Y. Kusano, N. Shinjou, F. Shoji, M. Yokokawa, and T. Watanabe. "Overview of the K computer system." *FUJITSU Sci. Tech. J.*, **48**(3), 2012, pp. 302-309.
8. T. Ichimura, K. Fujita, M. Horikoshi, L. Meadows, K. Nakajima, T. Yamaguchi, K. Koyama, H. Inoue, A. Naruse, K. Katsushima, M. Hori and L. Maddegedara. "A Fast Scalable Implicit Solver with Concentrated Computation for Nonlinear Time-Evolution Problems on Low-Order Unstructured Finite Elements," 2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS), Vancouver, BC, 2018, pp. 620-629. doi: 10.1109/IPDPS.2018.00071
9. I. M. Idriss, R. Dobry, and R. D. Sing. "Nonlinear Behavior of Soft Clays during Cyclic Loading." *Journal of the Geotechnical Engineering Division*, **104**(ASCE 14265), 1978.
10. G. Masing. "Eigenspannungen und Verfestigung beim Messing." *Proceedings of the 2nd International Congress of Applied Mechanics*, 1926, pp. 332-335.
11. G. H. Golub and Q. Ye. "Inexact conjugate gradient method with inner-outer iteration." *SIAM Journal on Scientific Computing*, **21**(4), 1997, pp. 1305-1320.
12. G. Karypis and V. Kumar, "A fast and high quality multilevel scheme for partitioning irregular graphs," *SIAM Journal on scientific Computing*, vol. 20(1), 1998, pp. 359-392.
13. Oakforest-PACS, [Online].  
<http://www.cc.u-tokyo.ac.jp/system/ofp/index-e.html>
14. A. Sodani, "Knights landing (KNL): 2nd Generation IntelR Xeon Phi processor," 2015 IEEE Hot Chips 27 Symposium (HCS), Cupertino, CA, 2015, pp. 1-24. doi: 10.1109/HOTCHIPS.2015.7477467
15. STREAM triad benchmark, [Online]. <http://www.cs.virginia.edu/stream/>
16. A. Kumar, "The New Intel Xeon Processor Scalable Family (Formerly Skylake-SP)," 2017 IEEE Hot Chips 29 Symposium (HCS), Cupertino, CA, 2017.
17. 5m mesh digital elevation map, Tokyo ward area, Geospatial Information Authority of Japan, [Online].  
<https://www.gsi.go.jp/MAP/CD-ROM/dem5m/index.htm>
18. National Digital Soil Map, The Japanese Geotechnical Society, [Online].  
[https://www.jiban.or.jp/?page\\_id=432](https://www.jiban.or.jp/?page_id=432)
19. Strong ground motion of The Southern Hyogo prefecture earthquake in 1995 observed at Kobe JMA observatory, Japan Meteorological Agency, [Online].  
[https://www.data.jma.go.jp/svd/eqev/data/kyoshin/jishin/hyogo\\_nanbu/dat/H1171931.csv](https://www.data.jma.go.jp/svd/eqev/data/kyoshin/jishin/hyogo_nanbu/dat/H1171931.csv)