# Implicit Low-Order Unstructured Finite-Element Multiple Simulation Enhanced by Dense Computation using OpenACC
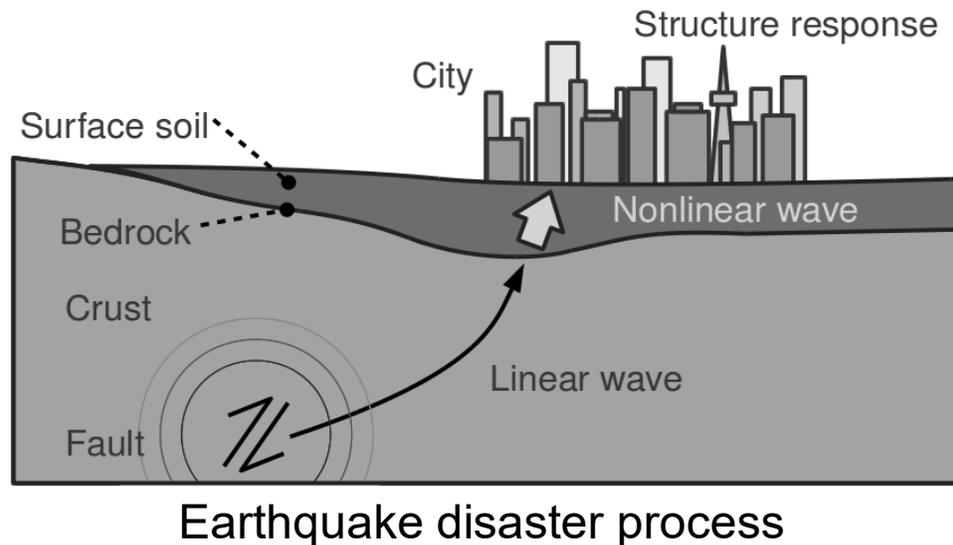
Takuma Yamaguchi, Kohei Fujita, Tsuyoshi Ichimura,
Muneo Hori, Lalith Maddegedara, Kengo Nakajima
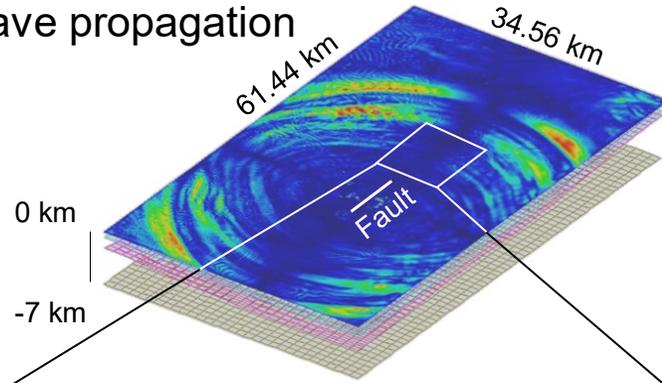
THE UNIVERSITY OF TOKYO

# Introduction

- Contribution of HPC to earthquake mitigation highly anticipated from society
- We are developing comprehensive earthquake simulation that simulate all phases of earthquake disaster by use of full K computer system
  - Simulate all phases of earthquake by speeding up core solver
  - **SC14 Gordon Bell Prize Finalist, SC15 Gordon Bell Prize Finalist & SC16 Best Poster & SC17 Best Poster Finalist**
- Ported this solver to GPU environment using OpenACC in WACCPD 2016 (**Best Paper**)
- Today's topic is enhancement of this GPU solver, and report performance on Pascal and Volta GPUs



Earthquake disaster process



K computer: 8 core CPU x 82944 node system with peak performance of 10.6 PFLOPS

# Comprehensive earthquake simulation



a) Earthquake wave propagation
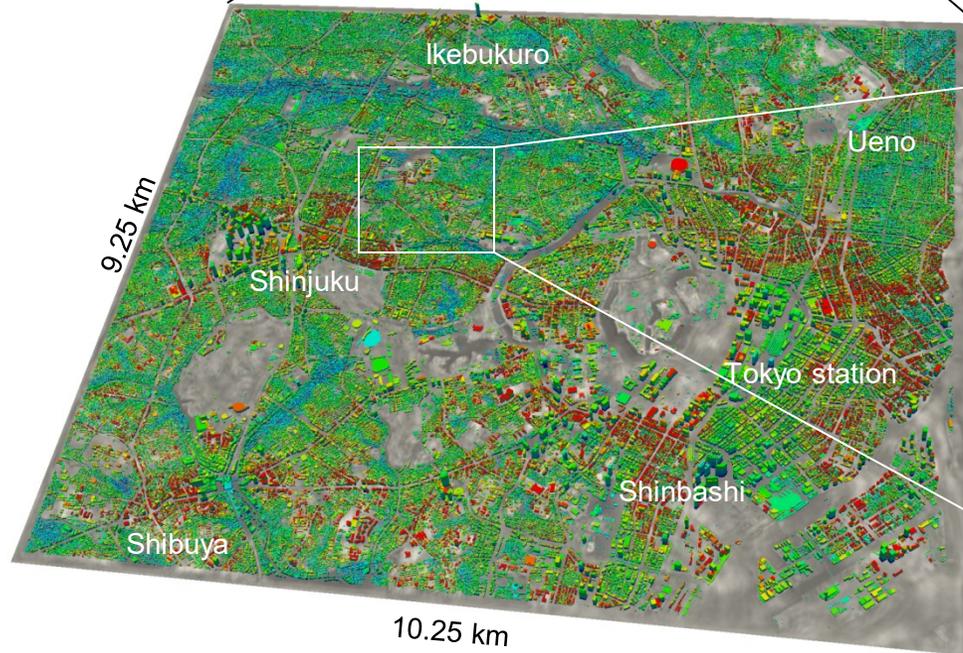
34.56 km

61.44 km

0 km

Fault

-7 km

9.25 km

Ikebukuro

Ueno

Shinjuku

Tokyo station

Shinbashi

Shibuya

10.25 km
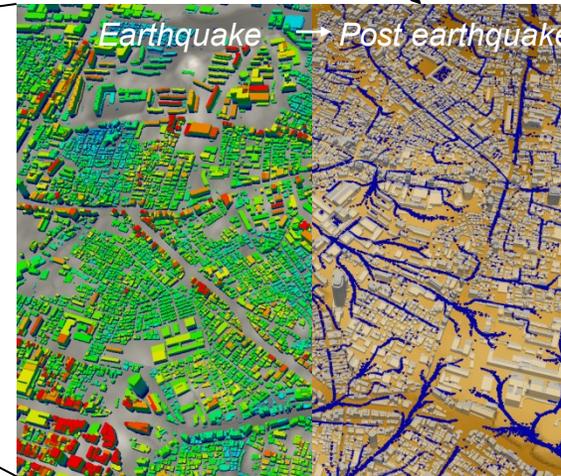
b) City response simulation

c) Resident evacuation

Two million agents evacuating to nearest safe site

Earthquake → Post earthquake

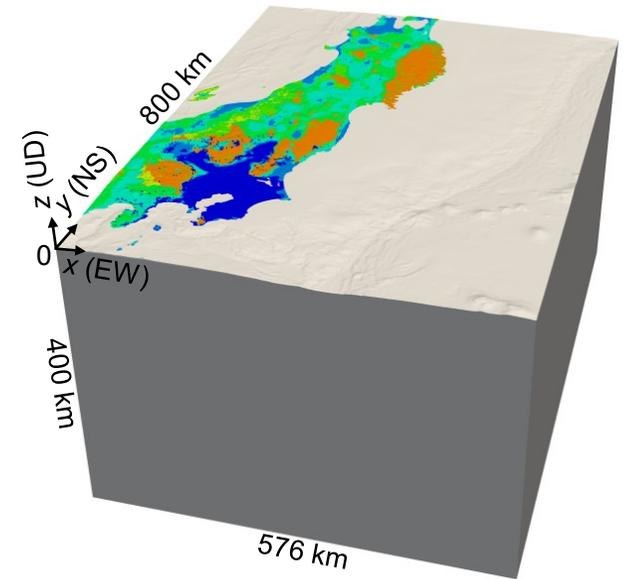Large finite-element simulation enabled by developed solver

3

City simulation



Image Landsat
Data SIO, NOAA, U.S. Navy, NGA, GEBCO
Image IBCAO

Google earth

Visualized by CYBERNET SYSTEMS CO., LTD

# Target problem: Earth's crust deformation problem

- Compute elastic response to given fault slip
  - Many case analysis required for inverse analyses and Monte Carlo simulations
- Compute using finite-element method: solve large matrix equation many times
  - Involves many random data access & communication
- Difficulty of problem
  - Attaining load balance & peak-performance & convergency of iterative solver & short time-to-solution at same time
  - Smart use of compute precision space, constraints in solver search space according to physical solution space required



Earth's crust deformation problem

$$Ku = f$$ ← Outer force vector

Unknown vector with up to 1 trillion degrees of freedom

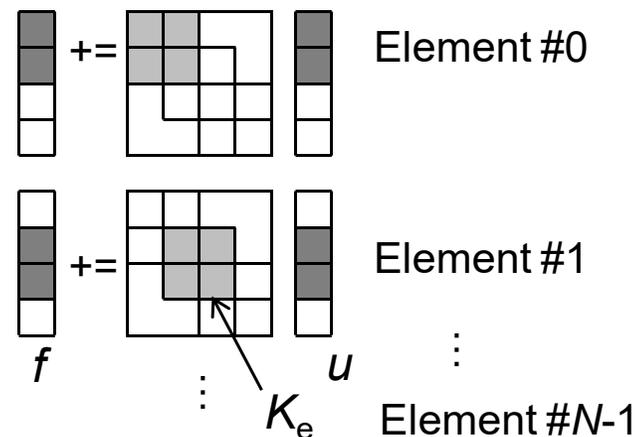Sparse, symmetric positive definite matrix

# Designing scalable & fast finite-element solver

- Design algorithm that can obtain equal granularity at O(million) cores
  - Matrix-free matrix-vector product (Element-by-Element method) is promising: Good load balance when elements per core is equal
    - Also high-peak performance as it is on-cache computation
- Combine Element-by-Element method with multi-grid, mixed precision arithmetic, and adaptive conjugate gradient method
  - Scalability & peak-performance good (core computation kernels are Element-by-Element), convergency good, time-to-solution good
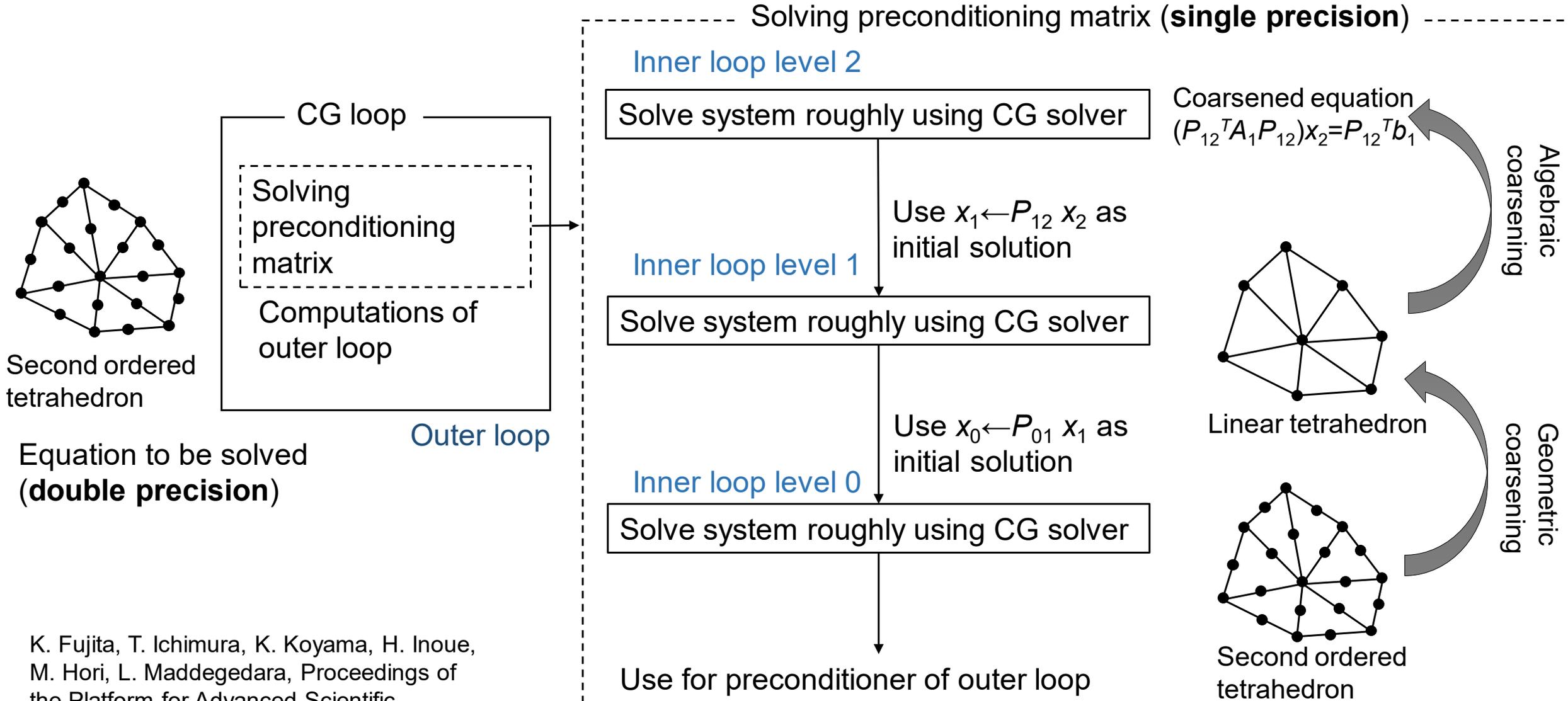
Element-by-Element method

$f = \Sigma_e\ P_e\ K_e\ P_e^\top\ u$

[$K_e$ is generated on-the-fly]

+=    Element #0

+=    Element #1

$f$    $u$    ⋮

⋮    $K_e$    Element #$N$-1

6

# Solver algorithm



Second ordered tetrahedron

Equation to be solved (**double precision**)

CG loop
- Solving preconditioning matrix
- Computations of outer loop

Outer loop

Solving preconditioning matrix (**single precision**)

Inner loop level 2

Solve system roughly using CG solver

Use $x_1 \leftarrow P_{12} x_2$ as initial solution

Inner loop level 1

Solve system roughly using CG solver

Use $x_0 \leftarrow P_{01} x_1$ as initial solution

Inner loop level 0

Solve system roughly using CG solver

Use for preconditioner of outer loop

Coarsened equation
$(P_{12}{}^T A_1 P_{12})x_2 = P_{12}{}^T b_1$

Algebraic coarsening

Linear tetrahedron

Geometric coarsening

Second ordered tetrahedron

# Performance on K computer

- Developed solver significantly faster than
  - PCGE (standard CG solver algorithm; preconditioning with 3x3 block diagonal matrix)
  - SC14 Gordon Bell Prize finalist solver (base solver for WACCPD 2016 GPU solver)
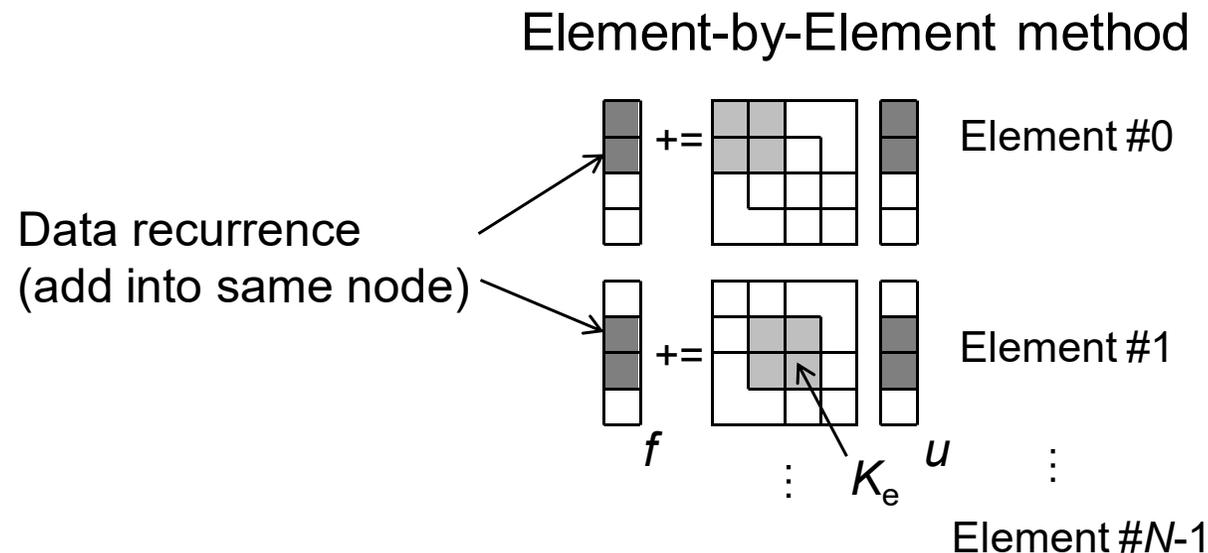- Use this as a base for GPU solver



- 94.8% scalability from 9216 to 294912 cores
- 4 times peak performance of HPCG benchmark (HPCG on K computer: 5.3% in double precision)

# Introduction of GPU computations

- Further speedup of the simulation by introducing GPUs
  - Good load balance, Reduced computation cost & data transfer size is also beneficial for GPUs
  - High performance can be obtained using OpenACC with low development cost

- GPU architecture is different from CPU architecture
  - Latency bound especially when we conduct random memory access
  - Relatively smaller cache size

- Finite-element applications tend to be memory bandwidth bound
- ➜Simple porting of the CPU code is not sufficient

# Key kernel: Element-by-Element kernel

- Most costly kernel; involves data recurrence

- Algorithm for avoiding data recurrence on **CPUs**
  - Use temporary buffers per core & per SIMD lane
  - Suitable for small core counts with large cache capacity

- Algorithm for avoiding data recurrence on **GPUs**
  - Last year, we developed algorithm using atomics and achieved high performance
  - However, random access becomes bottleneck…
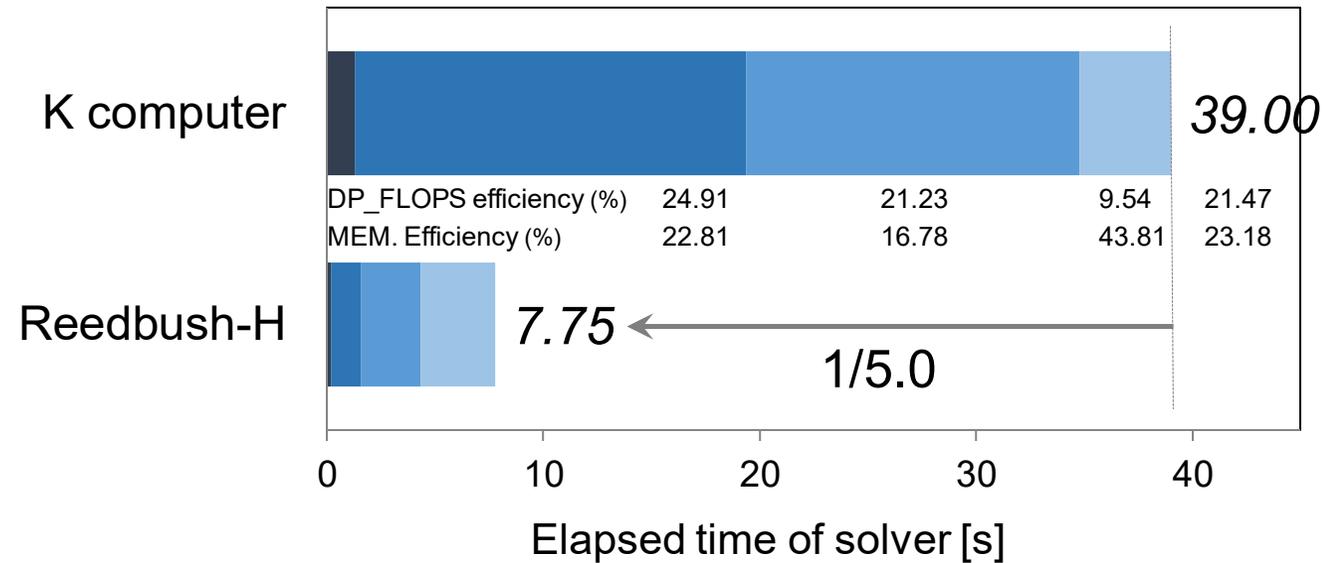
Element-by-Element method

Data recurrence
(add into same node)

Element #0

Element #1

$f$    $K_e$   $u$

Element #$N$-1

# Performance in simple porting

DOF: 125,177,217, # of elements: 30,720,000

## Computational Environment

|  | K computer | Reedbush-H |
|---|---|---|
| # of nodes | 20 | 10 |
| CPU/node | 1 x SPARC64 VIIIfx | 2 x Intel Xeon E5-2695 v4 |
| GPU/node | -- | 2 x NVIDIA P100 |
| Hardware peak FLOPS /process | 128 GFLOPS | 5.30 TFLOPS |
| Memory bandwidth /process | 64 GB/s | 732 GB/s |

■ Outer   ■ Inner level 0   ■ Inner level 1   ■ Inner level 2

**K computer**

*39.00*

DP_FLOPS efficiency (%)   24.91   21.23   9.54   21.47
MEM. Efficiency (%)   22.81   16.78   43.81   23.18

**Reedbush-H**   *7.75* ← 1/5.0

0   10   20   30   40

Elapsed time of solver [s]

- Simple porting achieved 5.0 times speedup
- However, there is some room for improvement
  - Memory bandwidth is 11 times larger

# Strategy for Introduction of OpenACC

- To attain optimal performance, algorithm/implementation suitable for GPUs should differ from that for CPUs

Thereby, we

1. Modify the solver algorithm to suit the GPU architectures
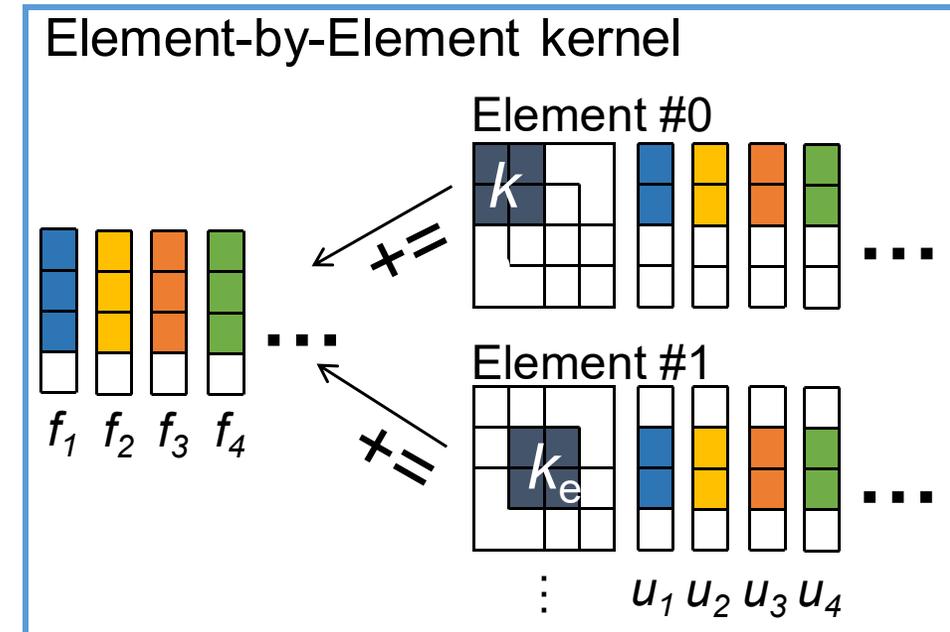2. Port the solver to GPUs using OpenACC

# Modification of Algorithm for GPUs

- Reduce random memory accesses

- Target applications (Inverse analyses, Monte Carlo method etc.) solve many systems of equations
  - Same stiffness matrix
  - Different right-hand side input vectors

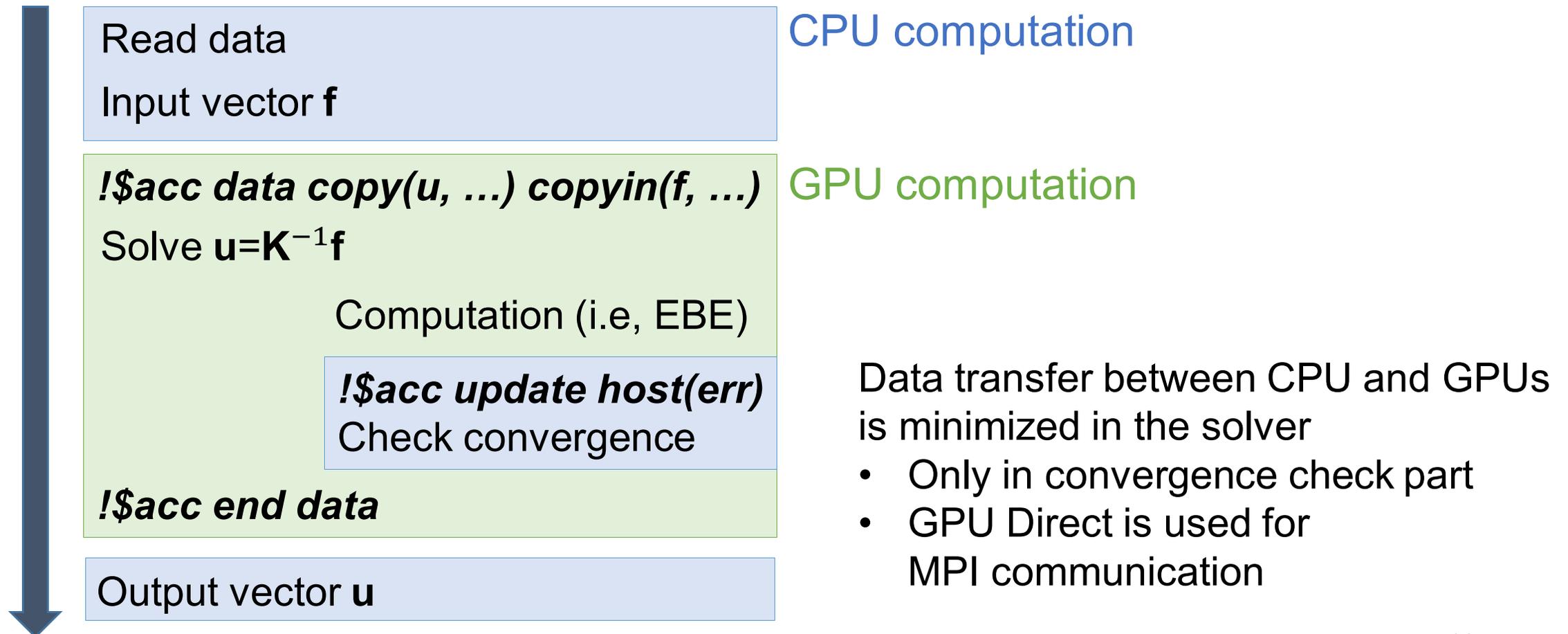- Multiple equations at the same time

$$K[u_1, u_2, u_3, ..., u_{16}]^T = [f_1, f_2, f_3, ..., f_{16}]^T$$

Instead of $Ku_1 = f_1, Ku_2 = f_2, Ku_3 = f_3, ...$



Element-by-Element kernel

# Introduction of OpenACC – 1/3

Control of data transfer

| | |
|---|---|
| Read data<br>Input vector **f** | CPU computation |
| **!$acc data copy(u, …) copyin(f, …)**<br>Solve $\mathbf{u}=\mathbf{K}^{-1}\mathbf{f}$ | GPU computation |
| Computation (i.e, EBE) | |
| **!$acc update host(err)**<br>Check convergence | |
| **!$acc end data** | |
| Output vector **u** | |

Data transfer between CPU and GPUs is minimized in the solver
- Only in convergence check part
- GPU Direct is used for MPI communication

# Introduction of OpenACC – 2/3

Insertion of some directives for parallel computation

Example for Element-by-Element multiplication

- Assign 16 threads for one element

- Introduce atomic functions to avoid data race

```
1   !$acc parallel loop collapse(2) present(…
2   do i_ele = 1, n_element
3   do i_vec = 1, 16
4   cny1 = connect(1, i_ele)
5        ⋮
6   cny10 = connect(10, i_ele)
7
8   u0101 = u(i_vec, 1, cny1)
9        ⋮
10  u1003 = u(i_vec, 3, cny10)
11
12  Ku01 = ...
13       ⋮
14  Ku30 = …
15
16  !$acc atomic
17  r(i_vec, 1, cny1) = r(i_vec, 1, cny1) + Ku01
18       ⋮
19  !$acc atomic
20  r(i_vec, 3, cny10) = r(i_vec, 3, cny10) + Ku30
21  enddo
22  enddo
23  !$acc end parallel
```

# Introduction of OpenACC – 3/3

Minor tuning for OpenACC parameters

• The allocation of gang, worker and vector

• The length of vector

Optimize fine-grain control of parallelism

(Not large effect on performance)

# Performance of the proposed solver

DOF: 125,177,217, # of elements: 30,720,000

## Computational Environment

|  | K computer | Reedbush-H |
|---|---|---|
| # of nodes | 20 | 10 |
| CPU/node | 1 x SPARC64 VIIIfx | 2 x Intel Xeon E5-2695 v4 |
| GPU/node | -- | 2 x NVIDIA P100 |
| Hardware peak FLOPS /process | 128 GFLOPS | 5.30 TFLOPS |
| Memory bandwidth /process | 64 GB/s | 732 GB/s |

**Legend:** ■ Outer ■ Inner level 0 ■ Inner level 1 ■ Inner level 2

**K computer 1 vector:** 39.00

DP_FLOPS efficiency (%)   24.91   21.23   9.54   21.47
MEM. Efficiency (%)   22.81   16.78   43.81   23.18

**Reedbush-H 1 vector:** 7.75    1/5.0

**K computer 16 vectors:** 26.43

**Reedbush-H 16 vectors:** 2.75    1/9.6    1/14.2

Elapsed time of solver, per vector [s]

# The speedup of each kernel

| Kernel | | Elapsed time per vector(s) | | Speedup |
|---|---|---|---|---|
| | | 1 vector | 16 vectors | |
| Element-by-Element computation | 1st order (FP32) | 0.948 | 0.584 | *1.62* |
| | 2nd order (FP32) | 0.687 | 0.401 | *1.71* |
| | 2nd order (FP64) | 0.044 | 0.025 | *1.78* |
| SpMV | | 1.465 | 0.091 | *16.10* |
| Dot product | | 0.213 | 0.522 | *0.41* |
| Total | | 7.75 | 2.75 | *2.82* |

- Reduction in random memory access in EBE kernels
- Total computation time for SpMV is constant
  - Bound by reading global matrix for memory
- Dot product is not efficiently computed
  - OpenACC cannot use arrays for reduction option
  - Using scalars (tmp1,tmp2,…,tmp16) causes stride memory access
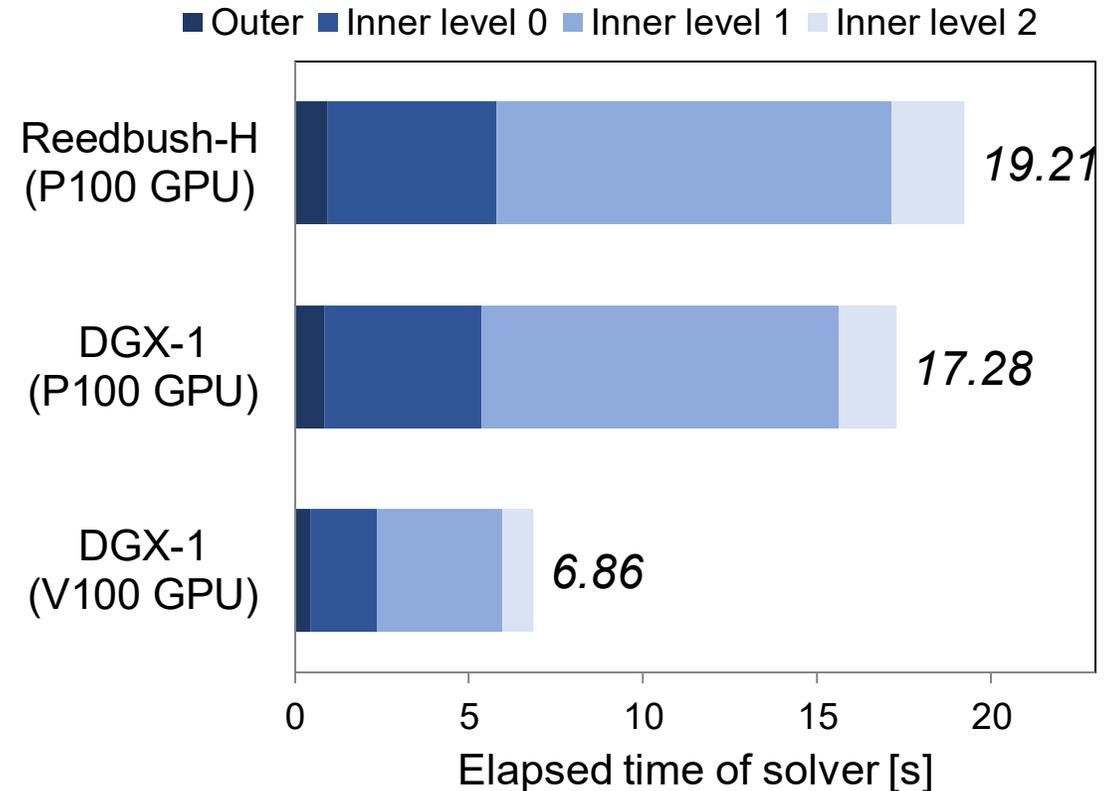
# Weak Scaling

Reedbush-H: P100 GPU x 240

| Model | DOF | # of elements | PCGE iterations | # of GPUs |
|-------|-----|---------------|-----------------|-----------|
| No. 1 | 125,177,217 | 30,720,000 | 4,928 | 20 |
| No. 2 | 249,640,977 | 61,440,000 | 4,943 | 40 |
| No. 3 | 496,736,817 | 122,880,000 | 4,901 | 80 |
| No. 4 | 992,038,737 | 245,760,000 | 4,905 | 160 |
| No. 5 | 1,484,953,857 | 368,640,000 | 4,877 | 240 |



■ Outer  ■ Inner level 0  ■ Inner level 1  ■ Inner level 2

| | | | | total |
|---|---|---|---|---|
| No.1 | 12 + 311 + 1300 + 2161 | 1.57 / 16.76 / 21.21 / 7.93 | | 47.47 |
| No.2 | 12 + 311 + 1346 + 2477 | 1.60 / 16.95 / 22.05 / 9.03 | | 49.63 |
| No.3 | 11 + 280 + 1312 + 1924 | 1.50 / 15.54 / 21.61 / 6.99 | | 45.64 |
| No.4 | 11 + 281 + 1430 + 2367 | 1.51 / 15.83 / 23.68 / 8.97 | | 49.99 |
| No.5 | 11 + 280 + 1341 + 2201 | 1.54 / 16.13 / 23.27 / 8.95 | | 49.89 |

Elapsed time of solver [s]

# Performance in using V100 GPUs

DOF: 38,617,017, # of elements: 9,440,240

Computational Environment

| | Reedbush-H | DGX-1 (P100/V100) | |
|---|---|---|---|
| # of nodes | 4 | 1 | |
| CPU/node | 2 x Intel Xeon E5-2695 v4 | 2 x Intel Xeon E5-2698 v4 | |
| GPU/node | 2 x NVIDIA P100 | 8 x NVIDIA P100 | 8 x NVIDIA V100 |
| Hardware peak FLOPS / process | 5.30 TFLOPS | 5.30 TFLOPS | 7.5 TFLOPS |
| Memory bandwidth /process | 732 GB/s | 732 GB/s | 900 GB/s |



Legend: ■ Outer ■ Inner level 0 ■ Inner level 1 ■ Inner level 2

Reedbush-H (P100 GPU): 19.21
DGX-1 (P100 GPU): 17.28
DGX-1 (V100 GPU): 6.86

Elapsed time of solver [s]

Higher performance than expected from hardware capability
- improved L1/L2 caches

# Application Example

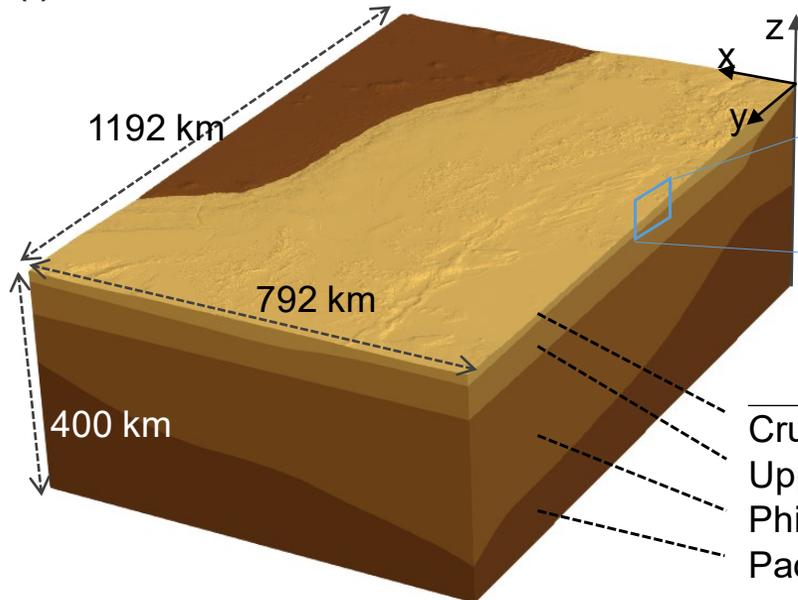Estimation of coseismic slip distribution in 2011 Tohoku Earthquake
DOF: 409,649,580
# of input vectors: 368 = 23 sets × 16 vectors
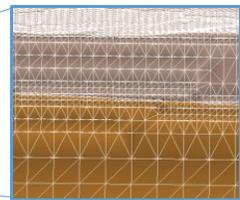Computation Environment: 64 x P100 GPUs (32 nodes of Reedbush-H)
Computation time: 828s for 23 sets of analyses
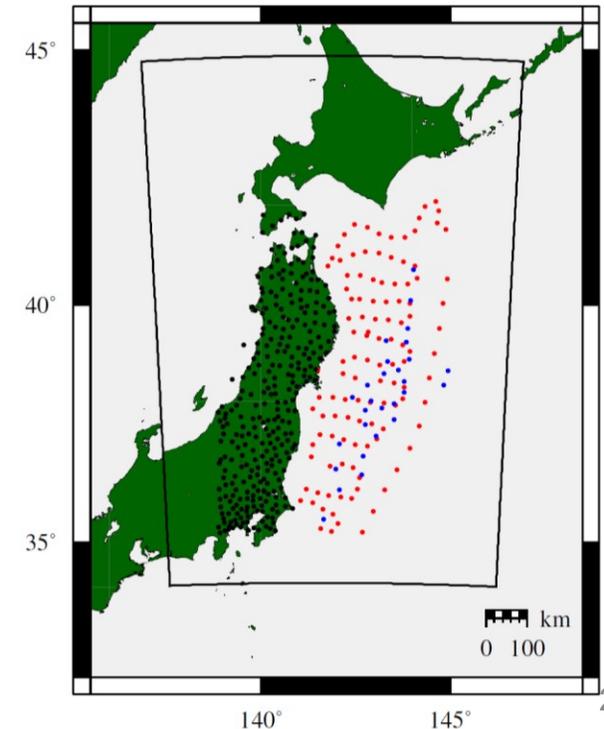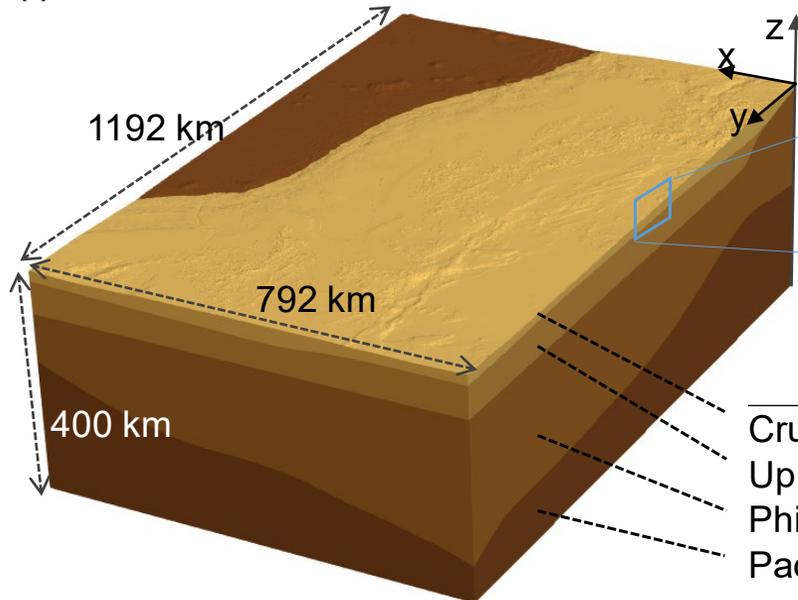 (29 times better in performance than previous studies)

Target Domain

(i) Whole FE model

1192 km

792 km

400 km

z

x

y

(ii) Close-up view

(iii) Material properties

|  | Vp (m/s) | Vs (m/s) | $\rho$ (kg/m$^3$) |
|---|---|---|---|
| Crust layer | 5664 | 3300 | 2670 |
| Upper-mantle layer | 8270 | 4535 | 3320 |
| Philippine Plate | 6686 | 3818 | 2600 |
| Pacific Plate | 6686 | 3818 | 2600 |

# Application Example

Estimation of coseismic slip distribution in 2011 Tohoku Earthquake
DOF: 409,649,580
# of input vectors: 368 = 23 sets × 16 vectors
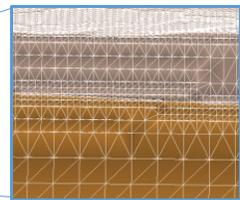Computation Environment: 64 x P100 GPUs (32 nodes of Reedbush-H)
Computation time: 828s for 23 sets of analyses
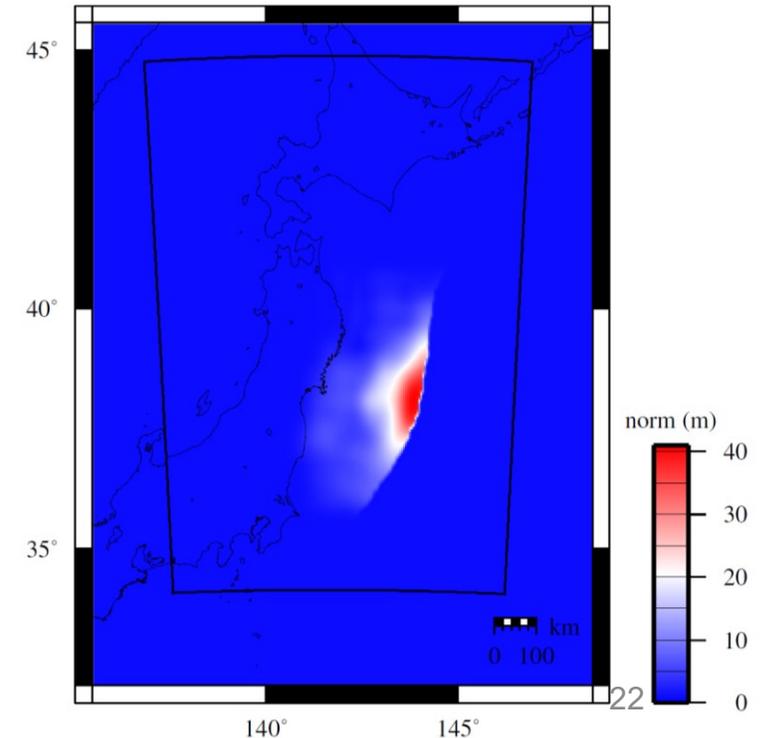 (29 times better in performance than previous studies)

Estimated Distribution

(i) Whole FE model



1192 km

792 km

400 km

(ii) Close-up view



(iii) Material properties

| | Vp (m/s) | Vs (m/s) | $\rho$ (kg/m$^3$) |
|---|---|---|---|
| Crust layer | 5664 | 3300 | 2670 |
| Upper-mantle layer | 8270 | 4535 | 3320 |
| Philippine Plate | 6686 | 3818 | 2600 |
| Pacific Plate | 6686 | 3818 | 2600 |

# Conclusion

- Accelerate the unstructured low-order finite element solvers by OpenACC
  - Design the solver appropriate for GPU computations
  - Port the codes to GPUs

- Obtain high performance with low development costs
  - 14.2 times speedup on P100 GPUs from the original solver on CPU-based K computer
  - Computation in low power consumption

- Improvement in reliability of earthquake simulation
  - Many-case simulation within short time