# Sunway TaihuLight:
# Designing and Tuning Scientific Applications at the Scale of 10 Million Cores

Haohuan Fu

National Supercomputing Center in Wuxi

Department of Earth System Science, Tsinghua University

September 13th 2017 @ ICAS

# Outline

Sunway Machine: the Challenges and Opportunities

Scientific Computing with 10 Million Cores

Long Term Plan for Sunway TaihuLight

Sunway-I:

- CMA service, 1998

- commercial chip

- 0.384 Tflops

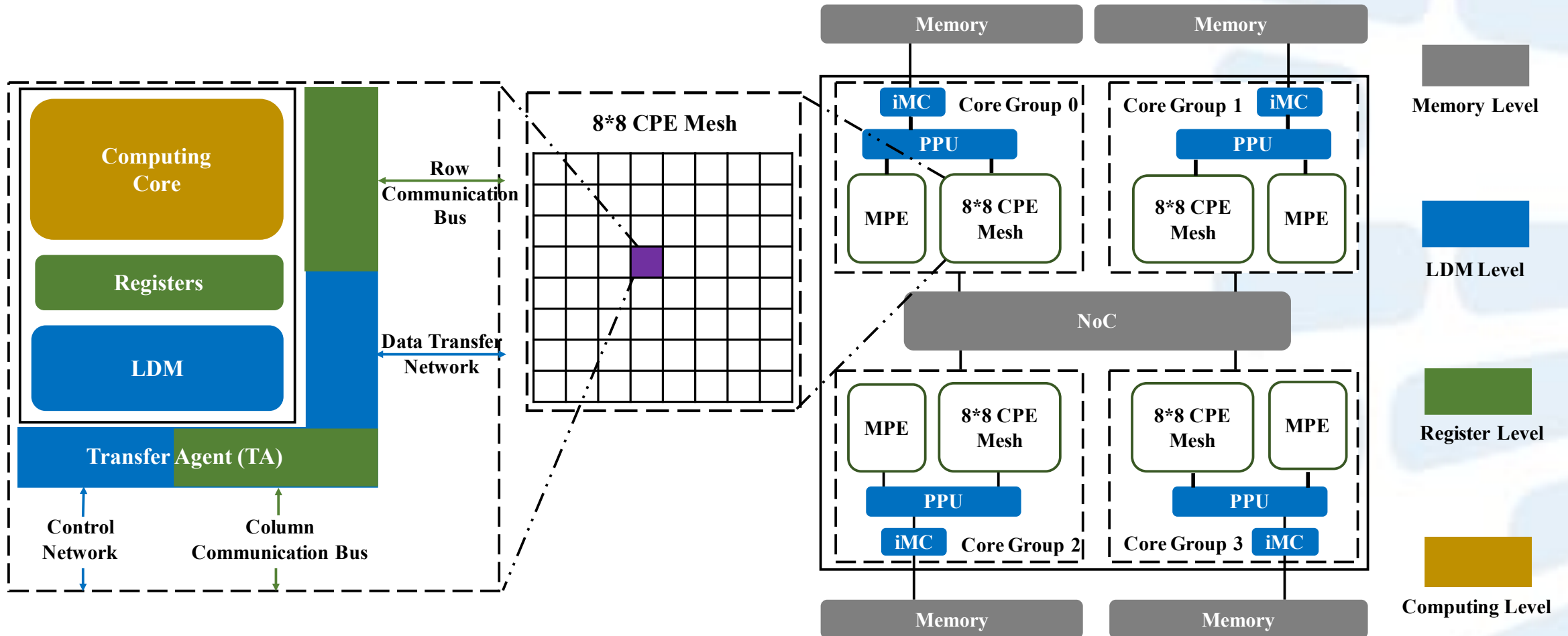- 48$^{th}$ of TOP500

Sunway BlueLight:

- NSCC-Jinan, 2011

- 16-core processor

- 1 Pflops

- 14$^{th}$ of TOP500

Sunway TaihuLight:

- NSCC-Wuxi, 2016

- 260-core processor

- 125 Pflops

- 1$^{st}$ of TOP500

# The Sunway Machine Family
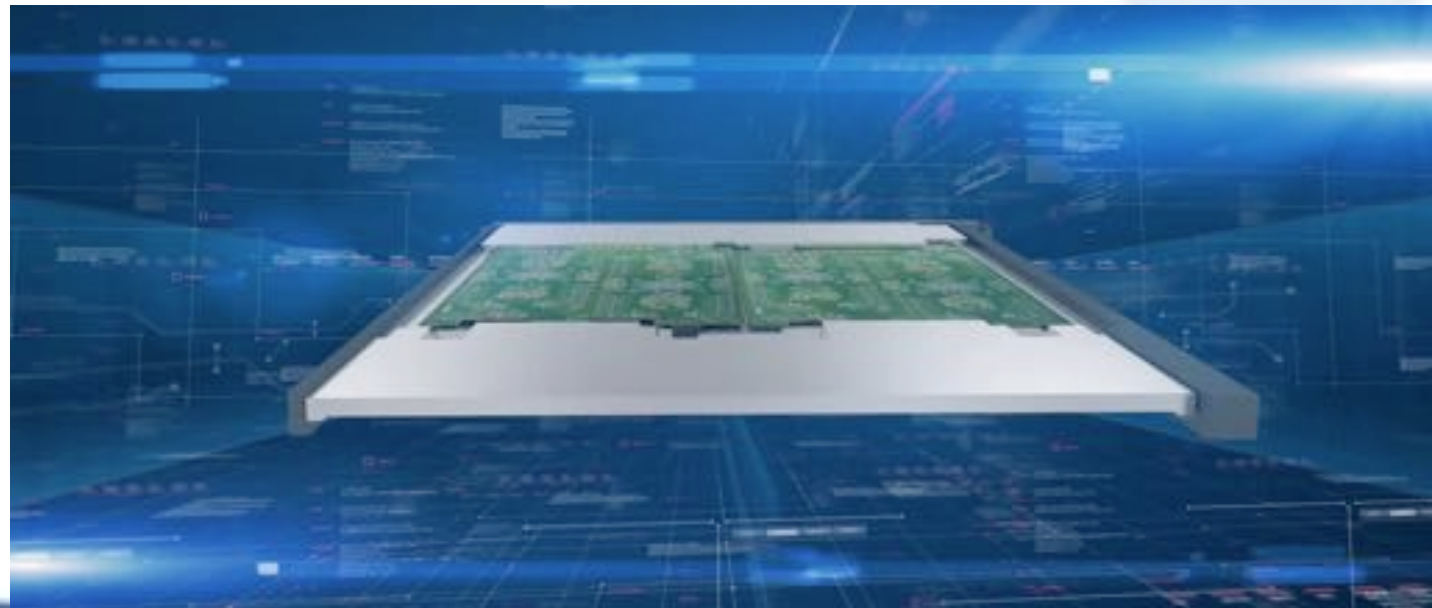
# SW26010: Sunway 260-Core Processor

# High-Density Integration of the Computing System

- A Five-Level Integration Hierarchy
  - ☐ computing node
  - ☐ computing board
  - ☐ super node
  - ☐ cabinet
  - ☐ entire computing system

# High-Density Integration of the Computing System

■ A Five-Level Integration Hierarchy

- ☐ <span style="color:red">computing node</span>
- ☐ <span style="color:red">computing board</span>
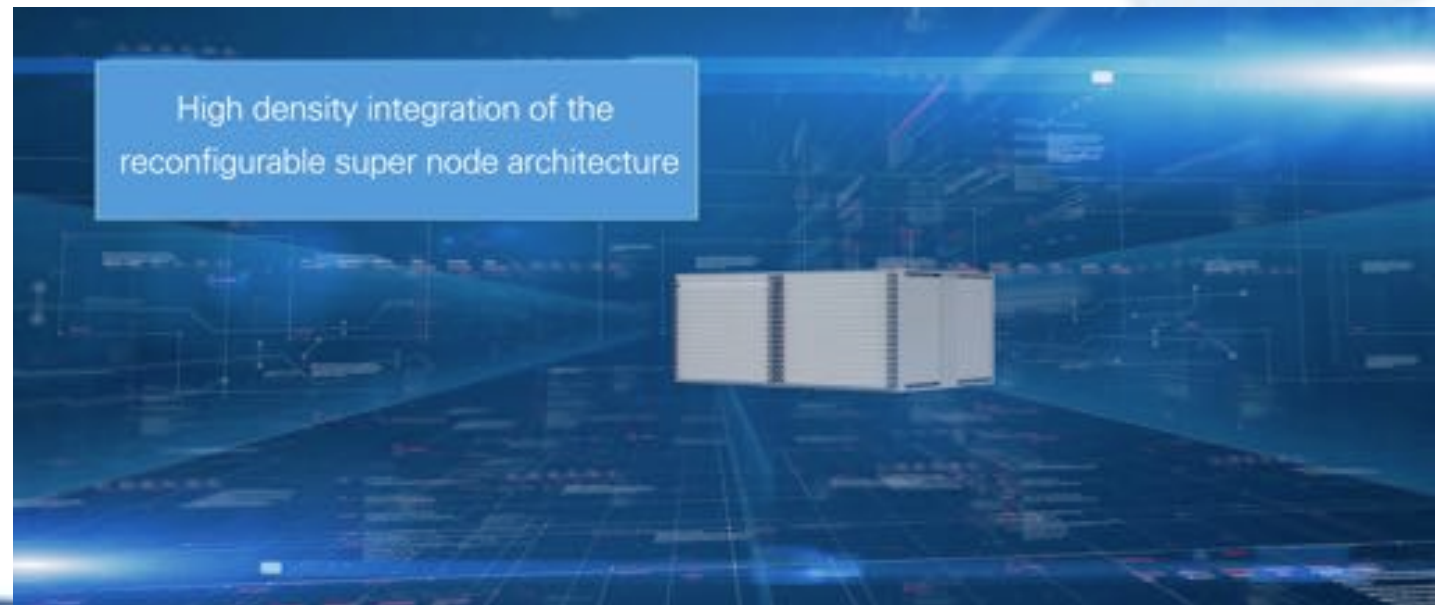- ☐ super node
- ☐ cabinet
- ☐ entire computing system

# High-Density Integration of the Computing System

- ## A Five-Level Integration Hierarchy
  - ☐ computing node
  - ☐ <span style="color:red">computing board</span>
  - ☐ <span style="color:red">super node</span>
  - ☐ cabinet
  - ☐ entire computing system

# High-Density Integration of the Computing System

- A Five-Level Integration Hierarchy
    - computing node
    - computing board
    - <span style="color:red">super node</span>
    - <span style="color:red">cabinet</span>
    - entire computing system



High density integration of the reconfigurable super node architecture

# High-Density Integration of the Computing System

■ A Five-Level Integration Hierarchy

  ❑ computing node

  ❑ computing board

  ❑ super node

  ❑ cabinet

  ❑ entire computing system

# How to Connect the 10 Million Cores?

$$40 \times 4 \times 256 \times 4 \times (1 + 8 \times 8) = 10{,}649{,}600$$

# How to Connect the 10 Million Cores?

$$40 \times 4 \times 256 \times 4 \times (1 + 8 \times 8) = 10,649,600$$

2D core array with row and column buses

# How to Connect the 10 Million Cores?

$$40 \times 4 \times 256 \times 4 \times (1 + 8 \times 8) = 10{,}649{,}600$$

2D core array with row and column buses

Network on Chip

# How to Connect the 10 Million Cores?

$$40×4×256×4×(1+8×8) = 10,649,600$$

2D core array with row and column buses

Network on Chip

Customized Network Board to Fully Connect 256 Nodes

# How to Connect the 10 Million Cores?

$$40 \times 4 \times 256 \times 4 \times (1 + 8 \times 8) = 10,649,600$$

2D core array with row
and column buses

Network on Chip

Customized Network Board to
Fully Connect 256 Nodes

Sunway Net

# Tweet Comments from Prof. Satoshi Matsuoka

**Satoshi Matsuoka**
@ProfMatsuoka

I was quite impressed with the engineering quality of TaihuLight, different from previous Chinese machines; now truly rivals US, Japan in SC twitter.com/profmatsuoka/s…

下午4:40 - 2016年11月3日 发自 東京 目黒区

# Tweet Comments from Prof. Satoshi Matsuoka

# Tweet Comments from Prof. Satoshi Matsuoka



**Satoshi Matsuoka**
@ProfMatsuoka

I was q
differen
Japan i

下午4:40

**Satoshi Matsuoka**
@ProfMatsuoka

TaihuLigh
dual-side
plate cool

下午5:57 - 2

**Satoshi Matsuoka**
@ProfMatsuoka

Also impressive was their software and application efforts. Contrary to my speculations OpenACC does work, used in many of their real apps.

下午5:59 - 2016年11月3日

# Tweet Comments from Prof. Satoshi Matsuoka

**Satoshi Matsuoka**
@ProfMatsuoka

Finally their design was cost&utility conscious. No expensive parts, quacky architecture, etc. Sunway apparently plans to sell the machine.

下午6:08 - 2016年11月3日

# Tweet Comments from Prof. Satoshi Matsuoka

**Satoshi Matsuoka**
@ProfMatsuoka

Finally their design was cost&utility conscious. No expensive parts, quacky architecture, etc. Sunway apparently plans to sell the machine.

下午6:08 - 2016年11月3日

Sunway Micro

# Outline

Sunway Machine: the Challenges and Opportunities

Scientific Computing with 10 Million Cores

Long Term Plan for Sunway TaihuLight

# Machine Capability Comparison

# Major Features to Consider

## Sunway TaihuLight

125 Pflops

10 million cores

user-controlled 64 KB LDM

32 GB and 136GB/s per node

22 flops/byte

MPE + CPE

register communication among CPEs

# Major Features to Consider

## Sunway TaihuLight

125 Pflops

10 million cores

user-controlled 64 KB LDM

32 GB and 136GB/s per node

22 flops/byte

MPE + CPE

register communication among CPEs

Intel KNL 7250 of Cori: 6.5 flops/byte
NVIDIA P100 of Piz Daint: 7.2 flops/byte

国家超级计算无锡中心

# Major Challenge #1: Scaling

**Sunway TaihuLight**

| 125 Pflops | 10 million cores | user-controlled 64 KB LDM |
| 32 GB and 136GB/s per node | 22 flops/byte | MPE + CPE | register communication among CPEs |

# Major Challenge #2: Memory Wall

## Sunway TaihuLight

| | | |
|---|---|---|
| 125 Pflops | 10 million cores | user-controlled 64 KB LDM |
| 32 GB and 136GB/s per node | 22 flops/byte | MPE + CPE | register communication among CPEs |

# Major Challenge #2: Memory Wall

## Sunway TaihuLight

125 Pflops

10 million cores

user-controlled 64 KB LDM

32 GB and 136GB/s per node

22 flops/byte

register communication among CPEs

Refactoring and Redesigning

# An (Incomplete) List of Full-Scale Applications

## 2016

Fully Implicit Solver for Atmospheric Dynamics

Surface Wave Modeling

Phase Field Simulations of Coarsening Dynamics

Atomistic Simulation of Silicon Nanowires

Run-away Electron Trajectory Simulation

Genome Functional Annotation and Homeotic Gene Building

Spacecraft CFD Numerical Simulation

## 2017

Extreme-scale Graph Processing Framework

Simulation of Planetary Rings

Simulations of Quantum Spin Liquid States via PEPS++

Molecular Dynamics Simulation of Condensed Covalent Materials

cryo-EM Macromolecule Structure Determination

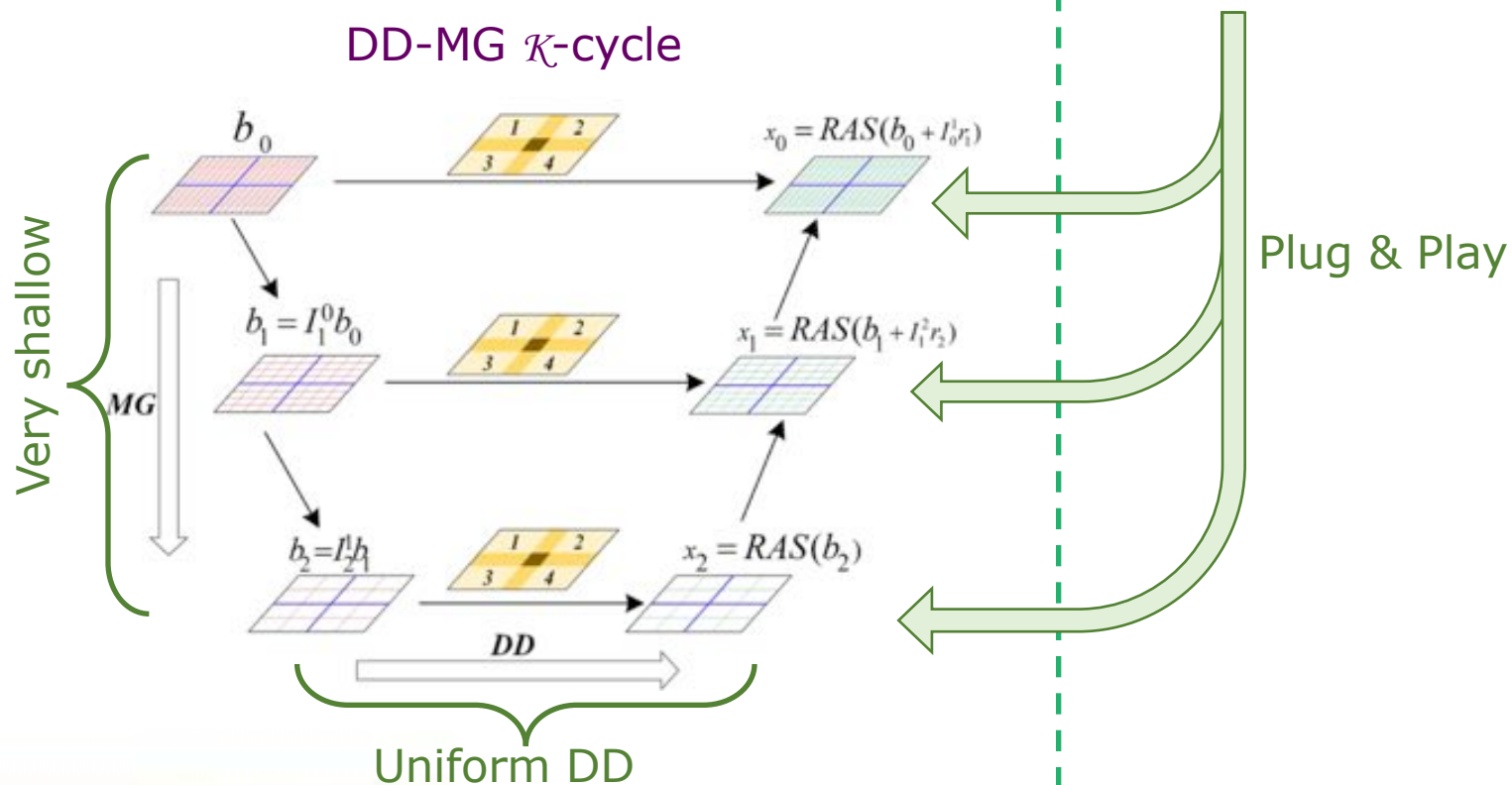Redesigning CAM-SE

Nonlinear Earthquake Simulation

# An (Incomplete) List of Full-Scale Applications

## 2016 Gordon Bell Finalists

- Fully Implicit Solver for Atmospheric Dynamics
- Surface Wave Modeling
- Phase Field Simulations of Coarsening Dynamics
- Atomistic Simulation of Silicon Nanowires
- Run-away Electron Trajectory Simulation
- Genome Functional Annotation and Homeotic Gene Building
- Spacecraft CFD Numerical Simulation

## 2017

- Extreme-scale Graph Processing Framework
- Simulation of Planetary Rings
- Simulations of Quantum Spin Liquid States via PEPS++
- Molecular Dynamics Simulation of Condensed Covalent Materials
- cryo-EM Macromolecule Structure Determination
- Redesigning CAM-SE
- Nonlinear Earthquake Simulation

# An (Incomplete) List of Full-Scale Applications

## 2016 Gordon Bell Prize

Fully Implicit Solver for Atmospheric Dynamics

Surface Wave Modeling

Phase Field Simulations of Coarsening Dynamics

Atomistic Simulation of Silicon Nanowires

Run-away Electron Trajectory Simulation

Genome Functional Annotation and Homeotic Gene Building

Spacecraft CFD Numerical Simulation

## 2017

Extreme-scale Graph Processing Framework

Simulation of Planetary Rings

Simulations of Quantum Spin Liquid States via PEPS++

Molecular Dynamics Simulation of Condensed Covalent Materials

cryo-EM Macromolecule Structure Determination

Redesigning CAM-SE

Nonlinear Earthquake Simulation

# An (Incomplete) List of Full-Scale Applications

## 2016 Gordon Bell Prize

Fully Implicit Solver for Atmospheric Dynamics

Surface Wave Modeling

Phase Field Simulations of Coarsening Dynamics

Atomistic Simulation of Silicon Nanowires

Run-away Electron Trajectory Simulation

Genome Functional Annotation and Homeotic Gene Building

Spacecraft CFD Numerical Simulation

## 2017 Gordon Bell Finalists

Extreme-scale Graph Processing Framework

Simulation of Planetary Rings

Simulations of Quantum Spin Liquid States via PEPS++

Molecular Dynamics Simulation of Condensed Covalent Materials

cryo-EM Macromolecule Structure Determination

Redesigning CAM-SE

Nonlinear Earthquake Simulation

163,840 processes  65 threads

racks  chips  core-groups  cores  total number of cores

$$40 \times 1{,}024 \times 4 \times 65 = 10{,}649{,}600$$

DD-MG $\mathcal{K}$-cycle

Very shallow

$b_0$

$x_0 = RAS(b_0 + I_0^1 r_1)$

$b_1 = I_1^0 b_0$

$x_1 = RAS(b_1 + I_1^2 r_2)$

MG

$b_2 = I_2^1 b_1$

$x_2 = RAS(b_2)$

DD

Uniform DD

Plug & Play

Now let's find a way to design a subdomain solver.

国家超级计算无锡中心

163,840 processes     65 threads

*racks*     *chips*     *core-groups*    *cores*     *total number of cores*

$$40 \times 1{,}024 \times 4 \times 65 = 10{,}649{,}600$$

DD-MG $\mathcal{K}$-cycle

Geometry-based pipelined ILU (GP-ILU)



$b_0$

$1$ $2$ $3$ $4$

$x_0 = RAS(b_0 + I_0^1 r_1)$

MG

$b_1 = I_1^0 b_0$

$1$ $2$ $3$ $4$

$x_1 = RAS(b_1 + I_1^2 r_2)$

$b_2 = I_2^1 b_1$

$1$ $2$ $3$ $4$

$x_2 = RAS(b_2)$

DD

Subdomain matrix of 1st-order with geometric index

Our goal of design:
1. Single sweep
2. Synchronization-free
3. Improved data-locality

$8 \times 8$

$\dfrac{reg\_size}{cell\_size}(num\_cores - 1) + blk\_height < dim\_z$

synchronization avoiding

# Strong-scaling results



The 3-km res run: 1.01 SYPD with 10.6M cores, dt=240s, I/O penalty <5%

# Weak-scaling results

Resolution (km)

DOFs=772B



The 488-m res run: 0.07 SYPD, 10.6M cores, dt=240s, 89.5X speedup over explicit

# Application (II): Porting CESM and Redesigning CAM-SE for Sunway TaihuLight



CAM5.0

POP2.0

CPL7

CLM4.0

CICE4.0

CESM1.2.0

Tsinghua + BNU **30+ Professors and Students**

- Four component models, millions lines of code
- Large-scale run on Sunway TaihuLight
    - **24,000** MPI processes
    - Over **one million** cores
- **10-20x** speedup for kernels
- **2-3x** speedup for the entire model

"Refactoring and Optimizing the Community Atmosphere Model (CAM) on the Sunway TaihuLight Supercomputer", in Proceedings of SC 2016.

# Application (II): Porting CESM and Redesigning CAM-SE for Sunway TaihuLight



CAM5.0

POP2.0

CPL7

CLM4.0

CICE4.0

CESM1.2.0

Tsinghua + BNU **30+ Professors and Students**

- Four component models, millions lines of code
- Large-scale run on Sunway TaihuLight
  - **24,000** MPI processes
  - Over **one million** cores
- **10-20x** speedup for kernels
- **2-3x** speedup for the entire model

"Refactoring and Optimizing the Community Atmosphere Model (CAM) on the Sunway TaihuLight Supercomputer", in Proceedings of SC 2016.

# Major Challenges

a high complexity in application, and a heavy legacy in the code base (millions lines of code)

an extremely complicated MPMD program with no hotspots (or hundreds of hotspots)

misfit between the in-place design philosophy and the new architecture

lack of people with interdisciplinary knowledge and experience

# OpenACC-based Refactoring of CAM

Pass tracers (u, v) to dynamics



Pass state variables and tracers

Pass state variables

- manual transformation of loops
- manual OpenACC parallelization and optimization on code and data structures

- tool based transformation of loops

# CAM model: scalability and speedup



- million core scale, 2.81 SYPD
- many-core refactoring for the entire model
- competitive simulation speed to the same model on NCAR Yellowstone

Simulation Speed (Described in Model Year Per Day(MYPD))

Number of CGs (each CG includes 1 MPE and 64 CPEs)

■ MPE only   ■ MPE+CPE for dynamic core   ■ MPE+CPE for both dynamic core and physics schemes

| Number of CGs | MPE only | MPE+CPE for dynamic core | MPE+CPE for both dynamic core and physics schemes |
|---|---|---|---|
| 1024 | 0.04 | | |
| 2400 | 0.15 | | |
| 4096 | 0.24 | | |
| 5120 | 0.25 | | |
| 7350 | 0.6 | | |
| 9600 | 0.78 | | |
| 12000 | 0.87 | 1.2 | 1.75 |
| 24000 | 1.54 | 1.62 | 2.81 |

# Athread-based Fine-grained Redesign

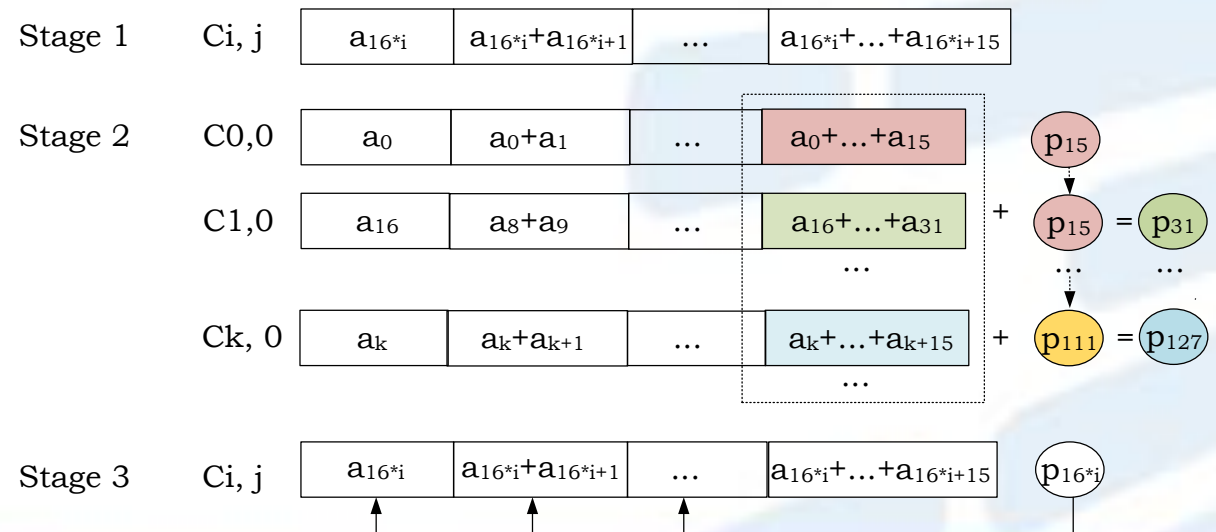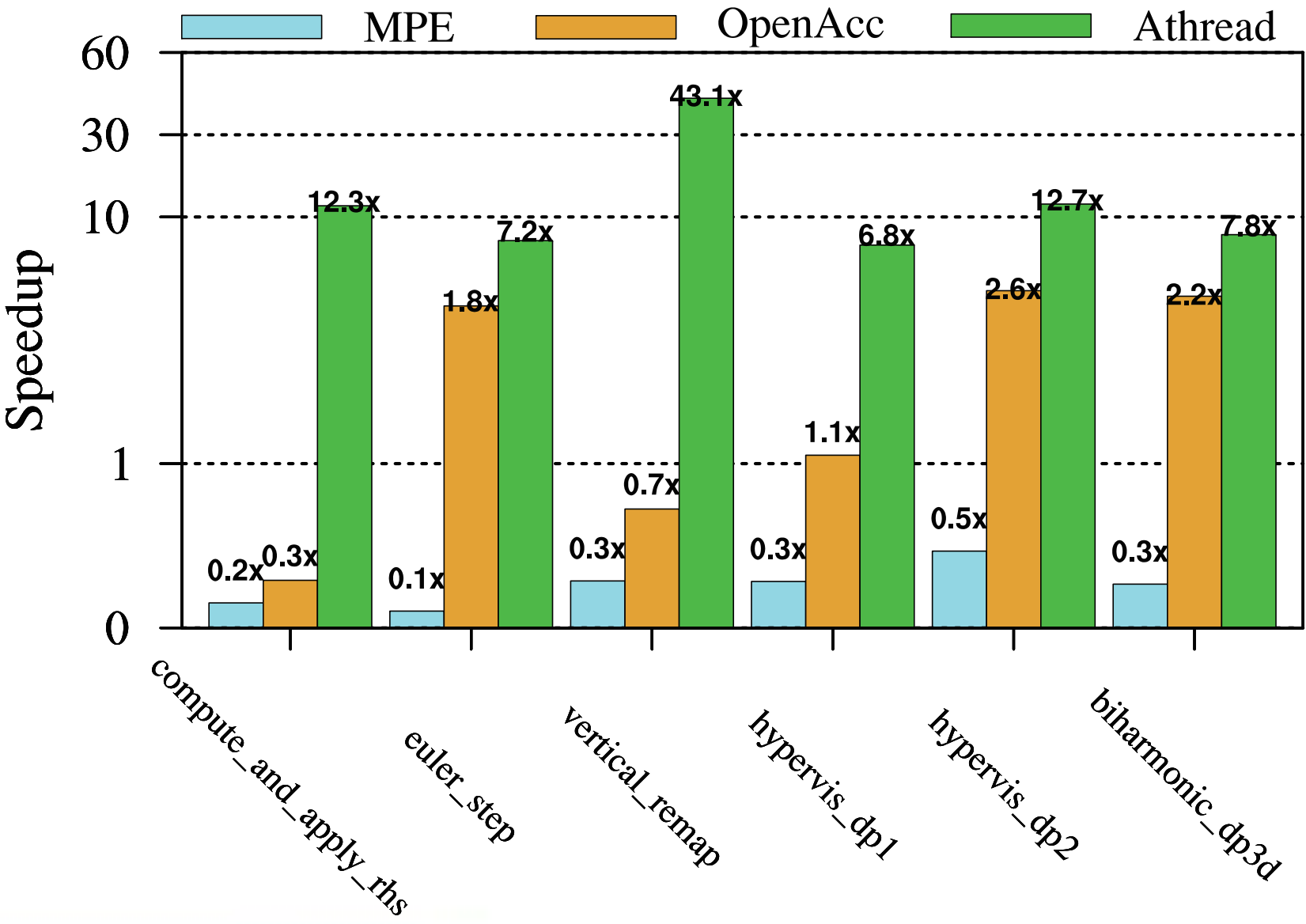- ■ Step 1: rewrite of Fortran OpenACC code to Athread C code
  - ☐ finer memory control through a specific DMA scheme
  - ☐ more efficient vectorization
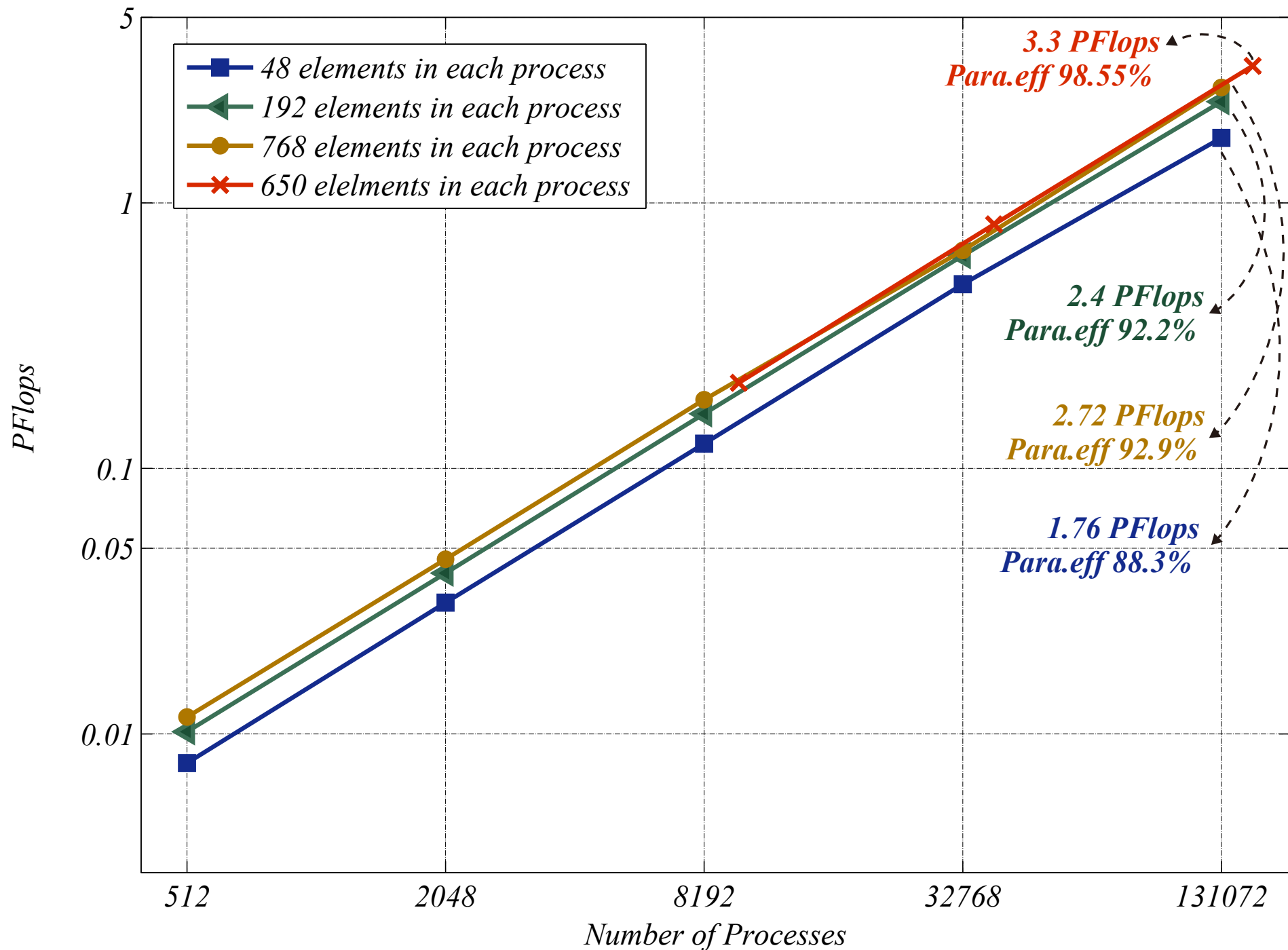
- ■ Step 2: register-communication based redesign
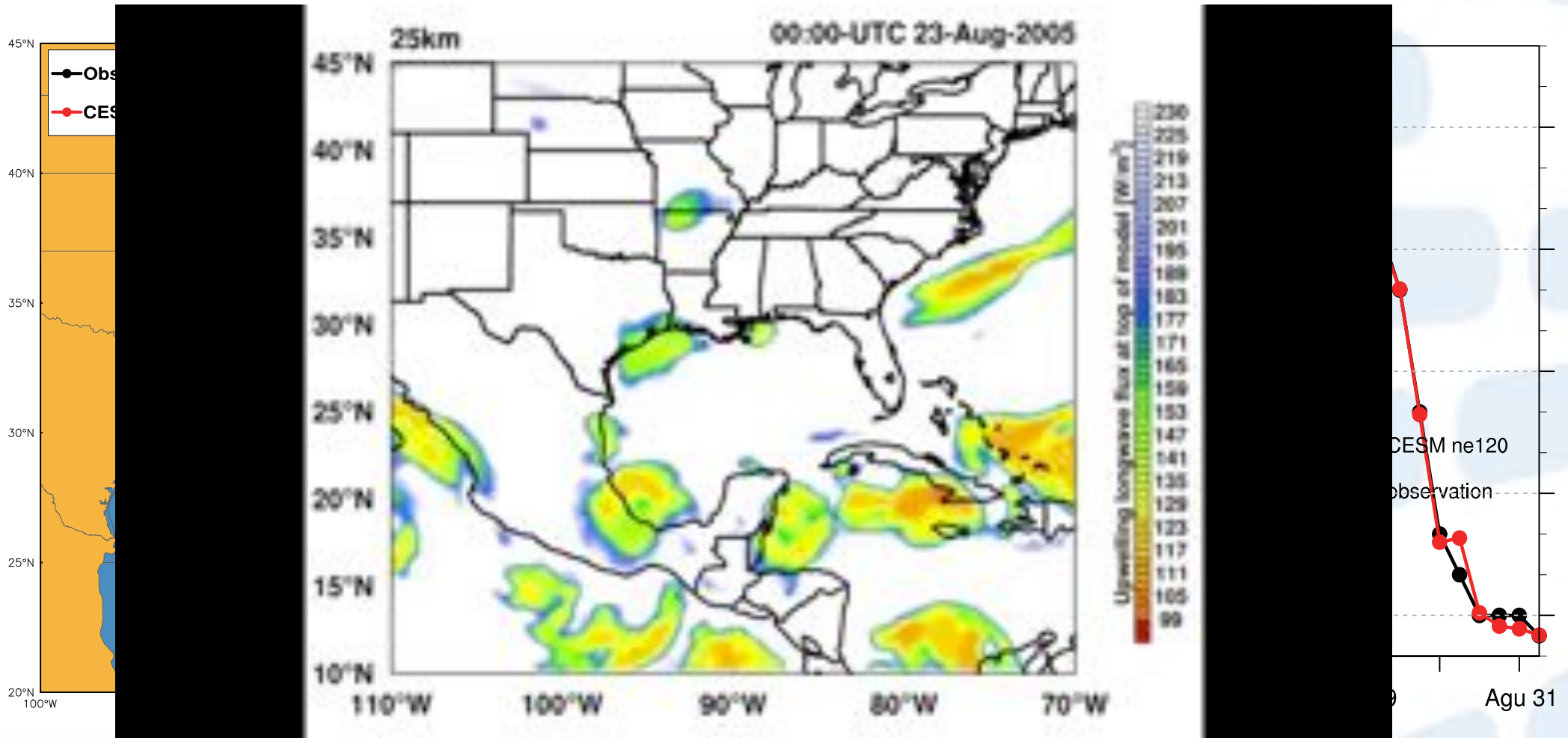  - ☐ remove data dependency
  - ☐ expose more parallelism

Performance Improvement through Redesign

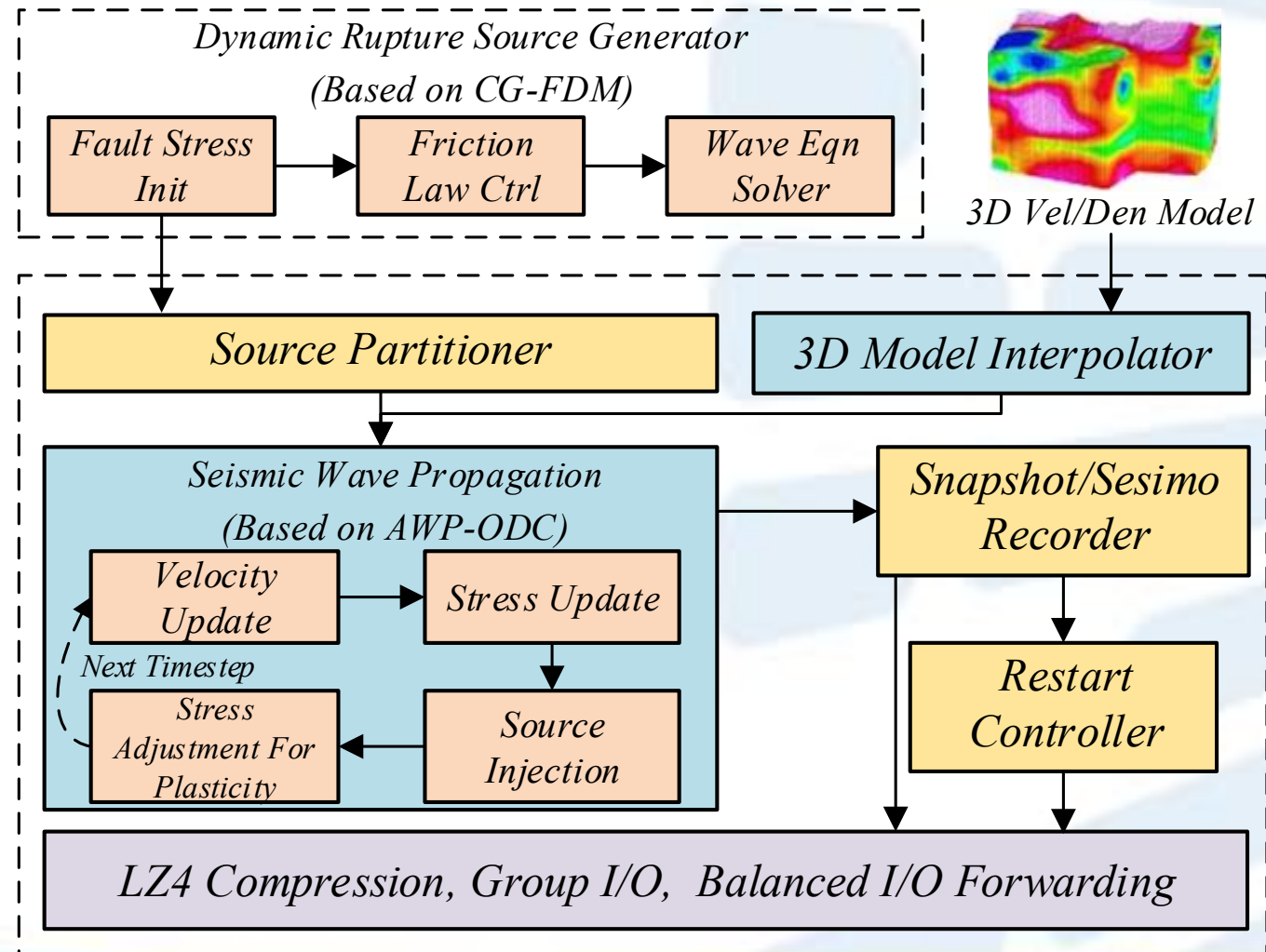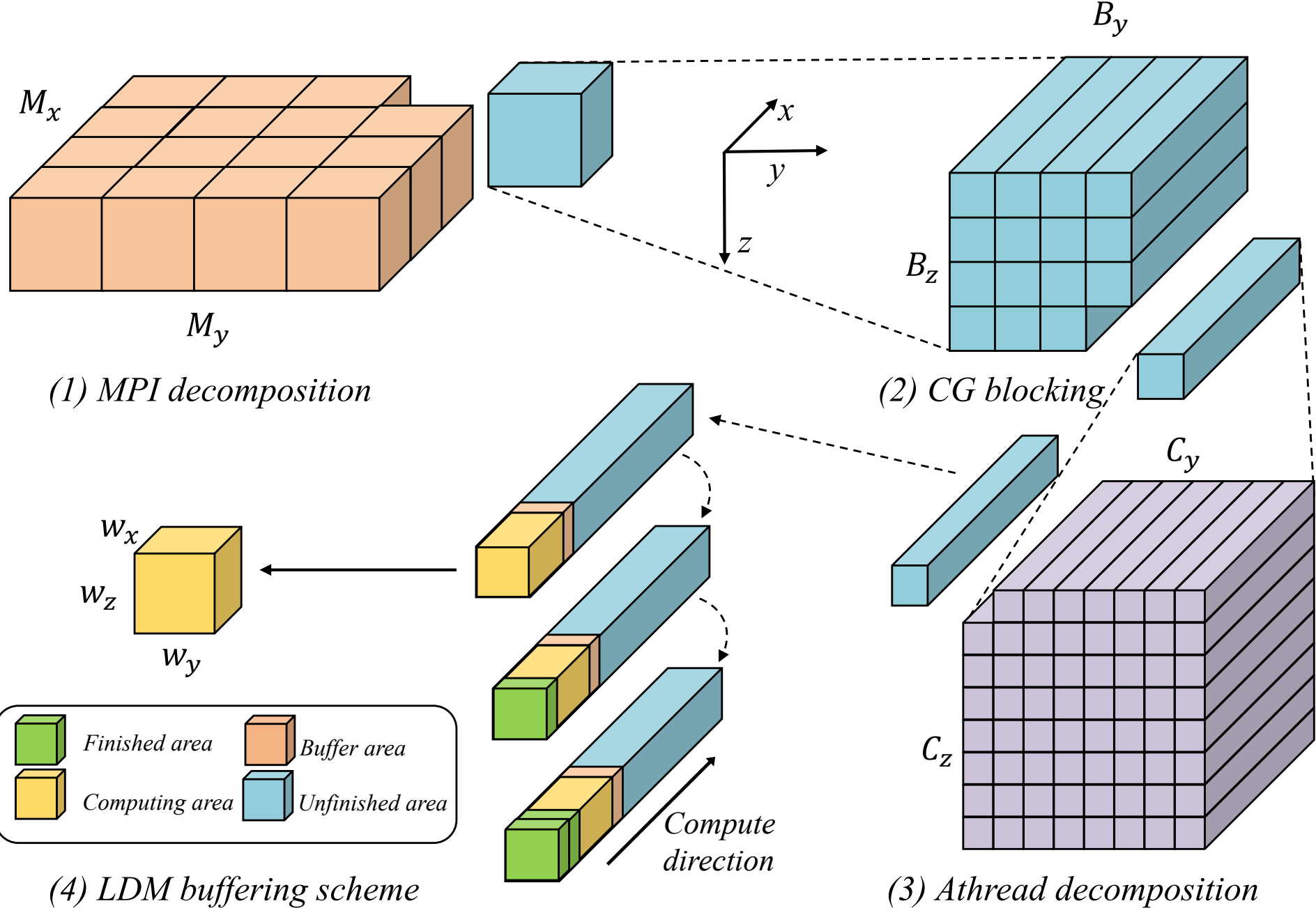Scaling the Dynamic Core to Millions of Cores
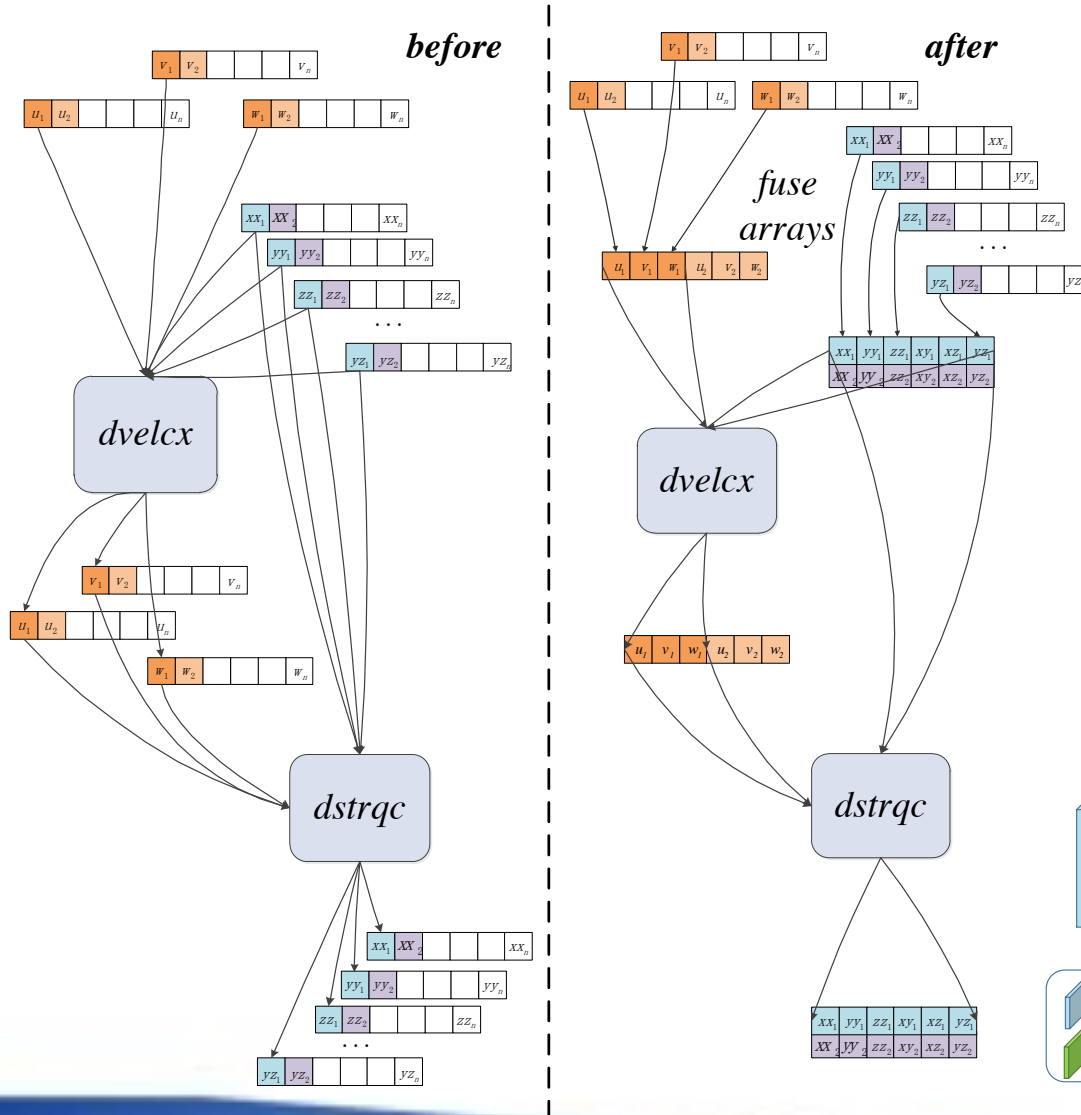
# Simulation of Hurricane Katrina

# Application (III): Nonlinear Earthquake Simulation on Sunway TaihuLight

- **Dynamic rupture source generator (originated from CG-FDM)**

- **Seismic wave propagation (originated from AWP-ODC)**

- **Other utilities:**
  - source partitioner
  - 3D Model Interpolator
  - Restart controller

Dynamic Rupture Source Generator
(Based on CG-FDM)

Fault Stress Init → Friction Law Ctrl → Wave Eqn Solver

3D Vel/Den Model

Source Partitioner

3D Model Interpolator

Seismic Wave Propagation
(Based on AWP-ODC)

Velocity Update → Stress Update

Next Timestep

Stress Adjustment For Plasticity ← Source Injection

Snapshot/Sesimo Recorder

Restart Controller

LZ4 Compression, Group I/O, Balanced I/O Forwarding

(1) MPI decomposition

(2) CG blocking

(3) Athread decomposition

(4) LDM buffering scheme

*Finished area*

*Buffer area*

*Computing area*

*Unfinished area*

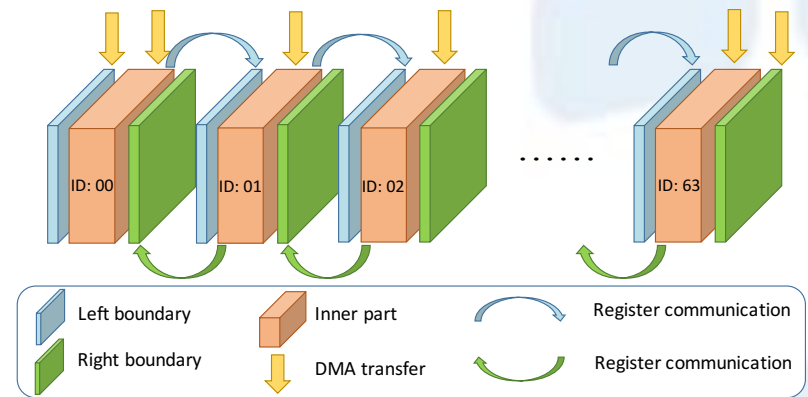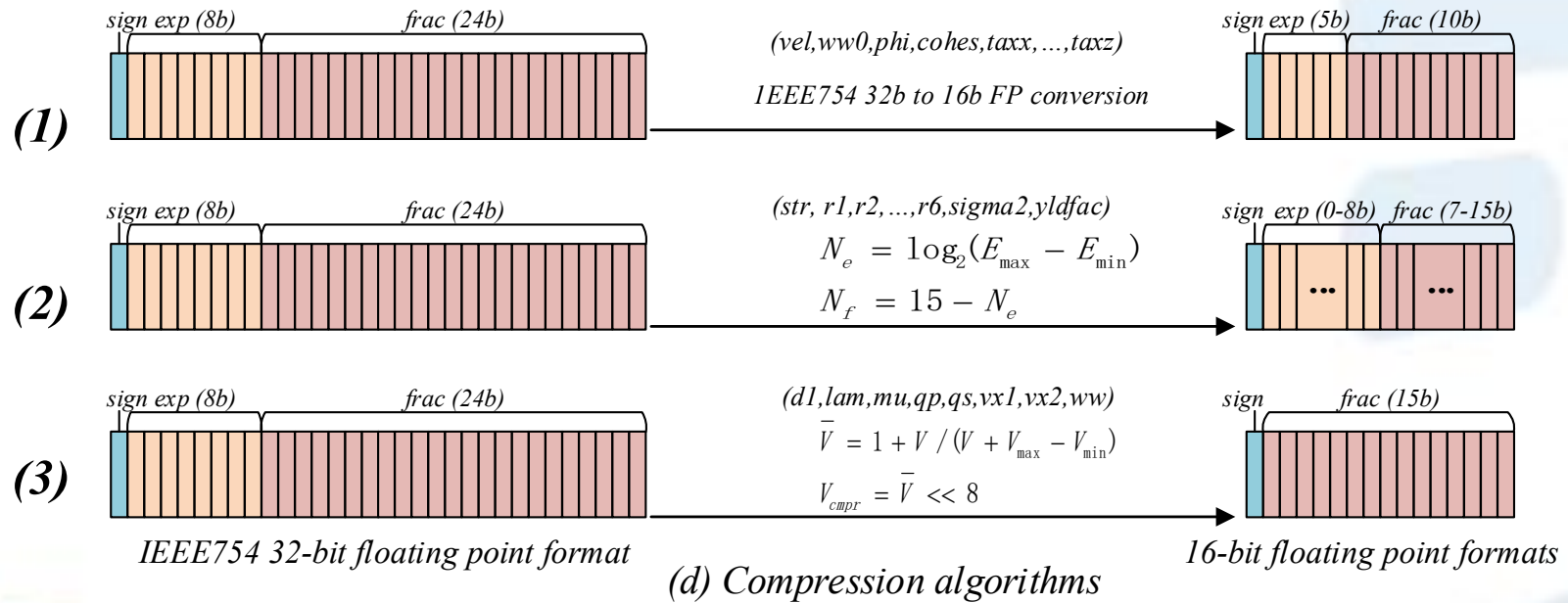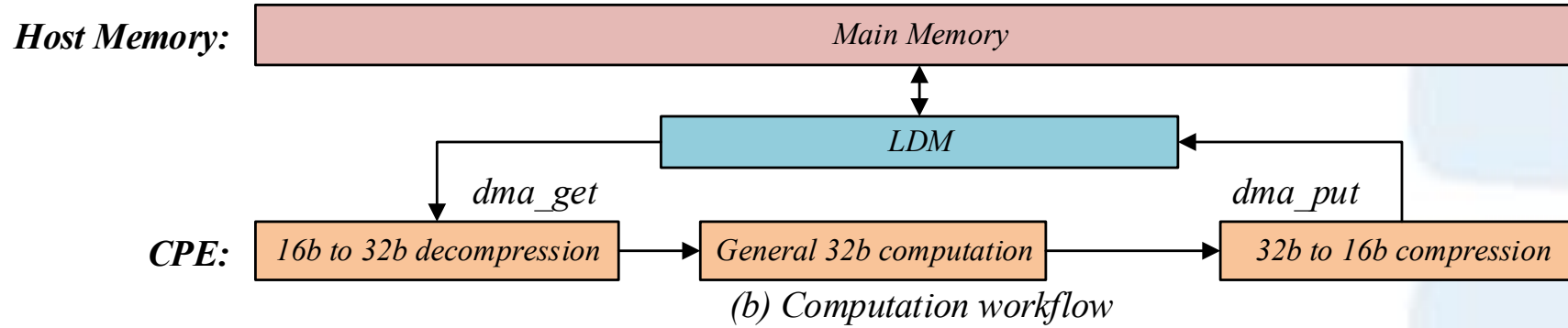*Compute direction*

**Multi-Level Domain Decomposition**

# A Balanced Memory Scheme



(1) array fusion,

(2) halo exchange through register communication,

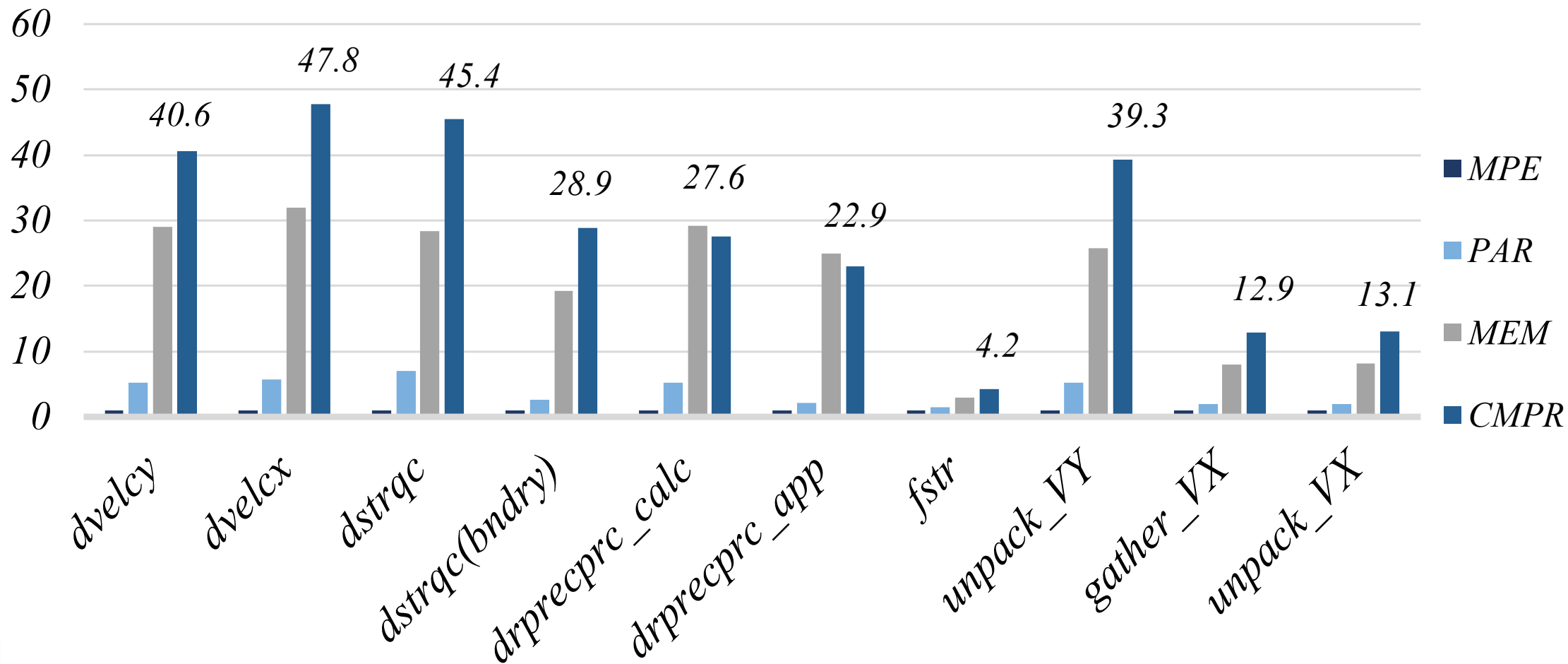(3) and optimized blocking configuration guided by an analytical model
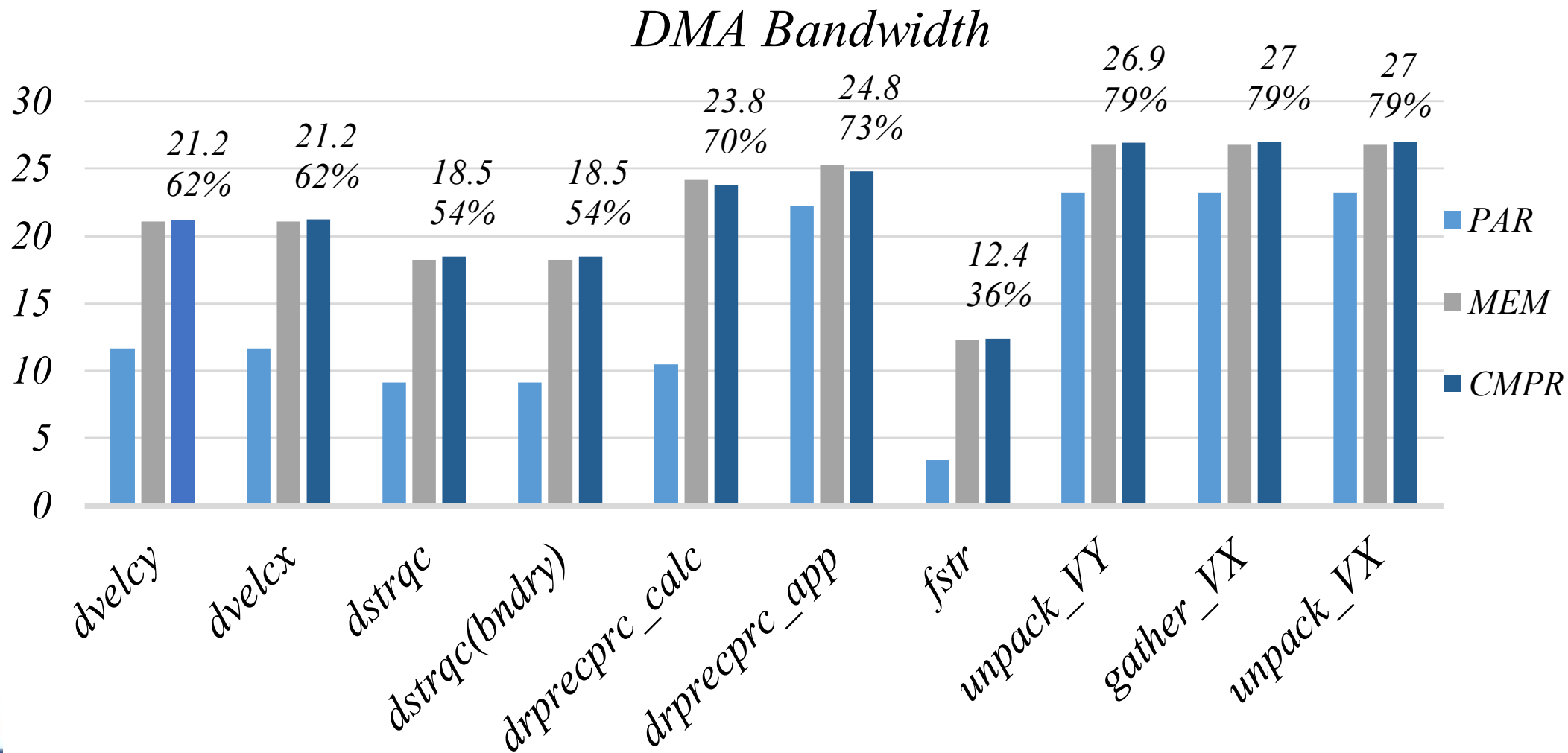
# On-the-fly Compression



(b) Computation workflow
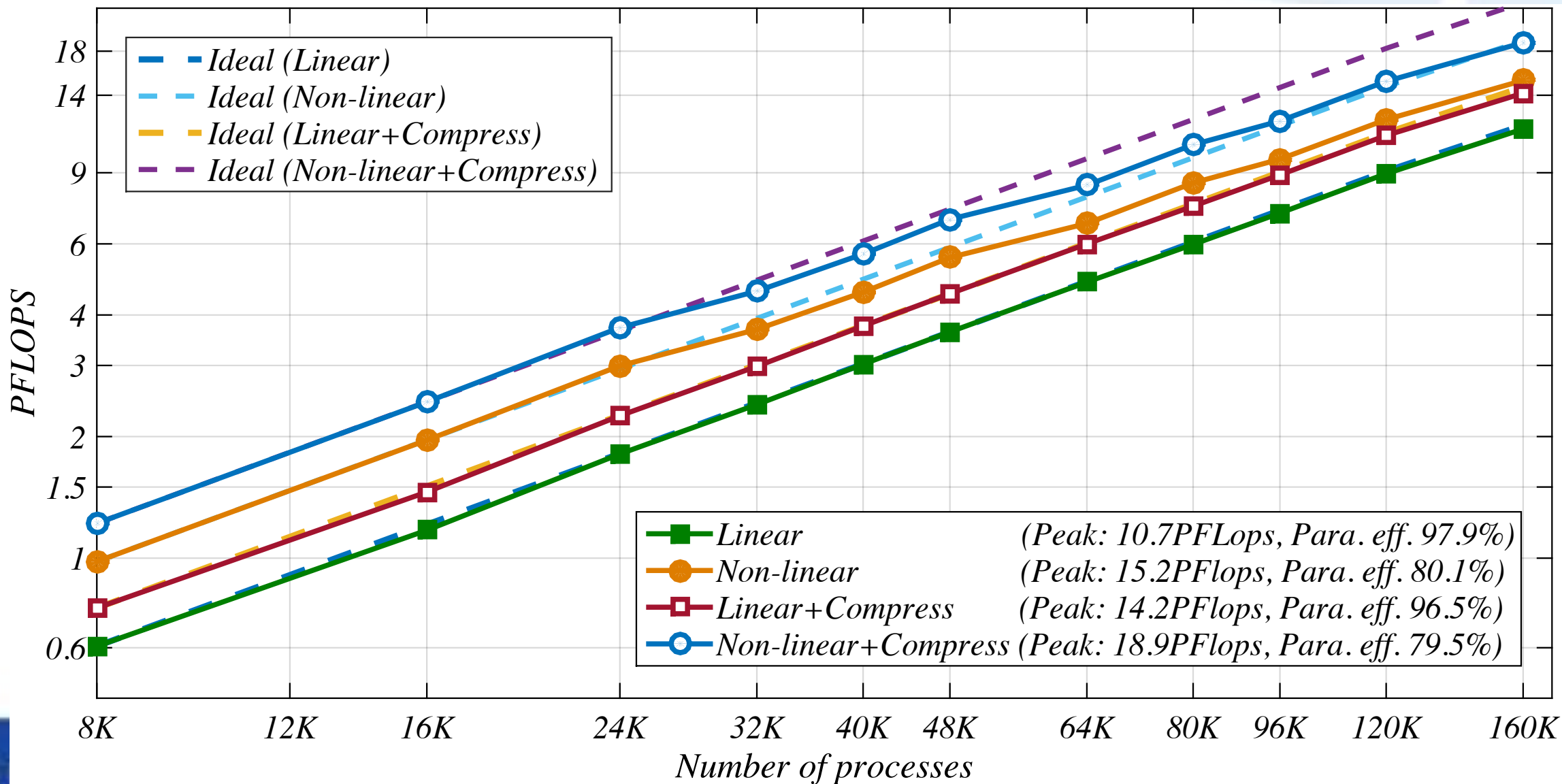
(1) $(vel,ww0,phi,cohes,taxx,...,taxz)$

$1EEE754$ 32b to 16b FP conversion

(2) $(str, r1,r2,...,r6,sigma2,yldfac)$

$$N_e = \log_2(E_{\max} - E_{\min})$$
$$N_f = 15 - N_e$$

(3) $(d1,lam,mu,qp,qs,vx1,vx2,ww)$

$$\overline{V} = 1 + V / (V + V_{\max} - V_{\min})$$
$$V_{cmpr} = \overline{V} << 8$$

IEEE754 32-bit floating point format

16-bit floating point formats

(d) Compression algorithms

# Speedup: 64 CPE vs 1 MPE

*Speedup*

# Memory Bandwidth Utilization



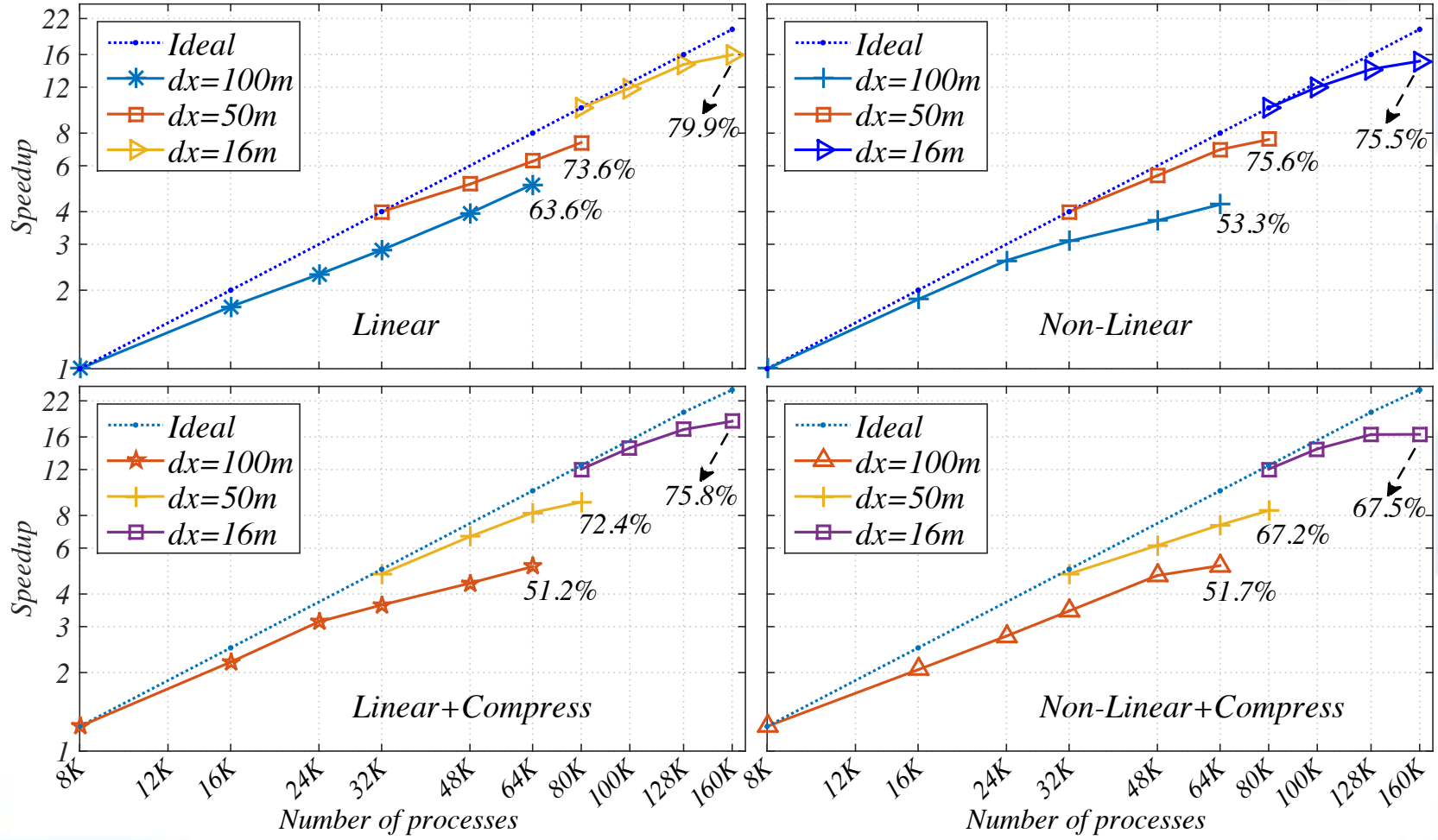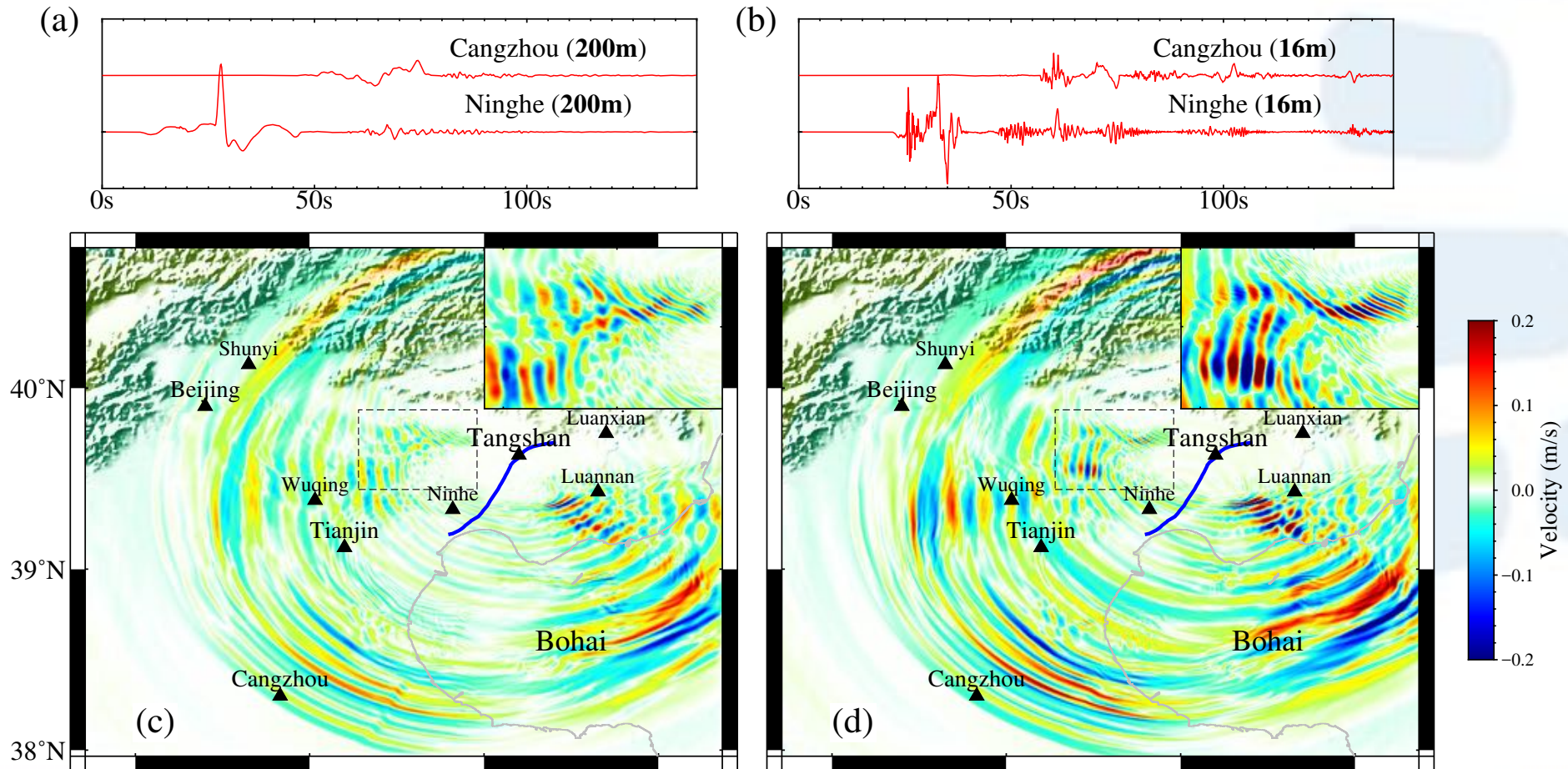DMA Bandwidth

# Weak Scaling

# Strong Scaling

# Simulation Results: 200m vs 16m
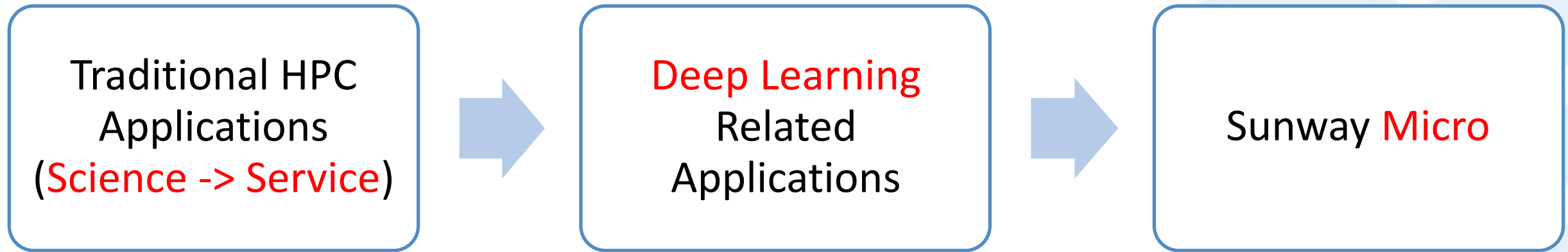
# Outline

Sunway Machine: the Challenges and Opportunities
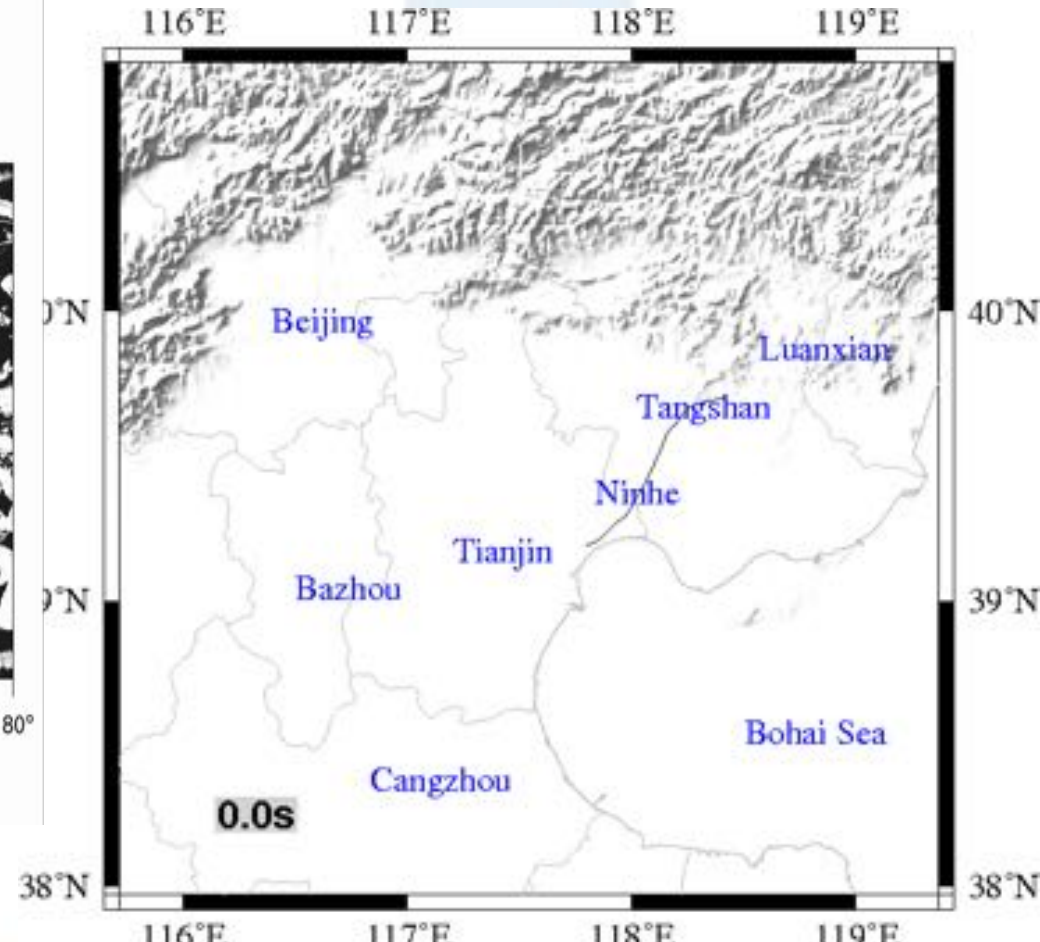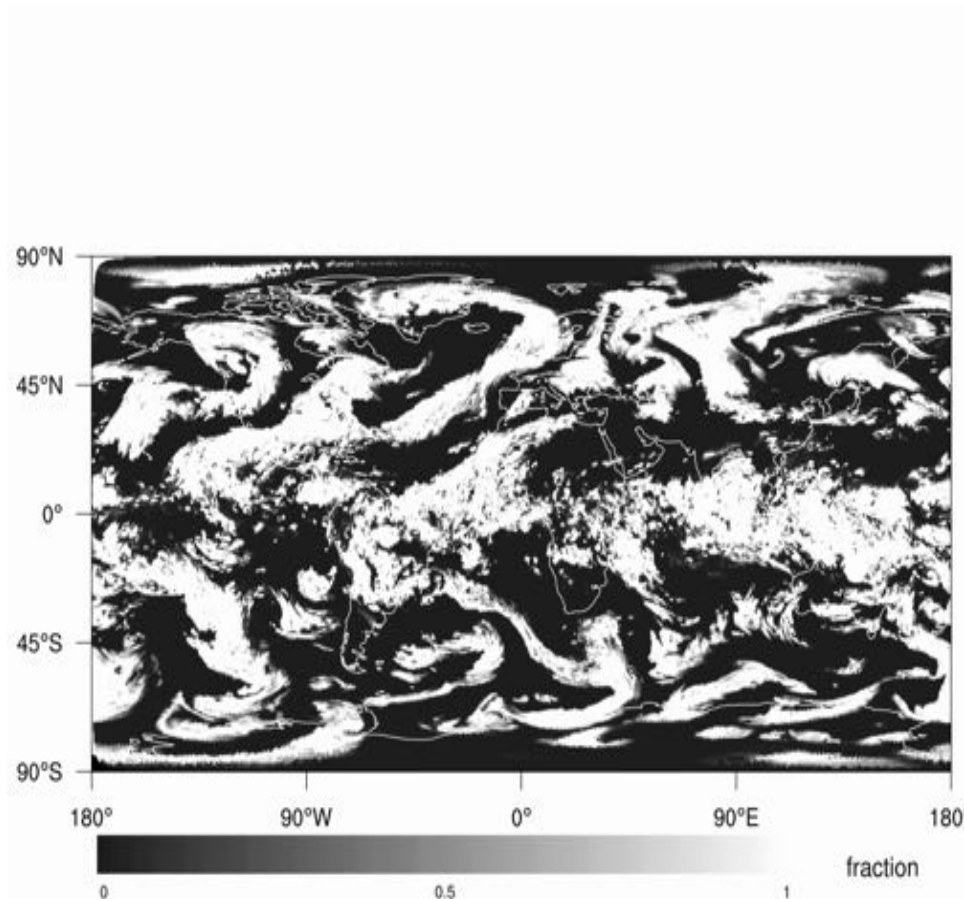
Scientific Computing with 10 Million Cores
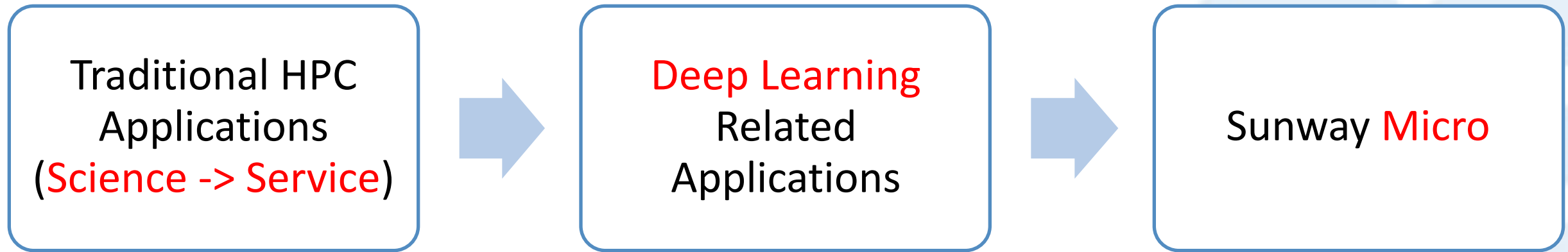
Long Term Plan for Sunway TaihuLight

# Long Term Plan

Traditional HPC Applications
(Science -> Service)

➤

Deep Learning Related Applications

➤

Sunway Micro

# Long Term Plan

Traditional HPC
Applications
(Science -> Service)

# Long Term Plan

Traditional HPC Applications
(Science -> Service)

➤

Deep Learning Related Applications

➤

Sunway Micro

# Long Term Plan



Training AlexNet with swCaffe

- total
- convolution
- fully connected

Traditional
Application
(Science -> S...)

Micro

# Long Term Plan

Traditional HPC Applications
(Science -> Service)

➡

Deep Learning Related Applications

➡

Sunway Micro

# Acknowledgements

- MOST, China: major sponsors of the HPC hardware and software development

- NRCPC: vendor of the machine

- NCAR: Rich Loft, John Dennis, Allison Baker, Haiying Xu (support and advice on the CAM-SE work)

- SCEC: Yifeng Cui, Steve Day, Daniel Roten, Kim Olsen, Josh Tobin, Alex Breuer, and Dawei Mu (discussion and advice on the earthquake simulation work)

国家超级计算无锡中心
National Supercomputing Center in Wuxi