



Scaling Geoscience Applications on Sunway Supercomputer

*Lin Gan Assistant Professor, Tsinghua University, Beijing
Assistant Director, NSCC-Wuxi, Jiangsu*

大寒

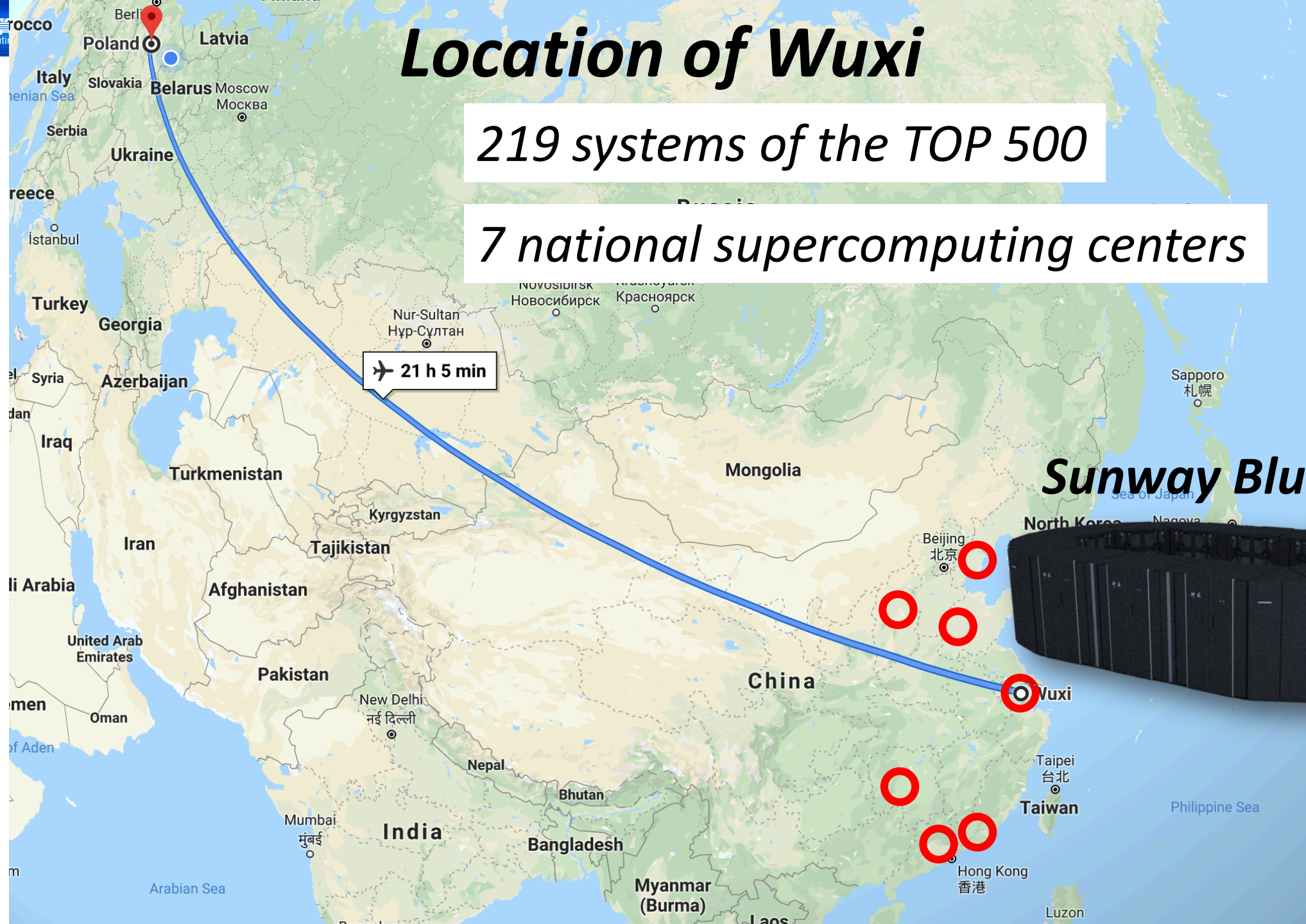




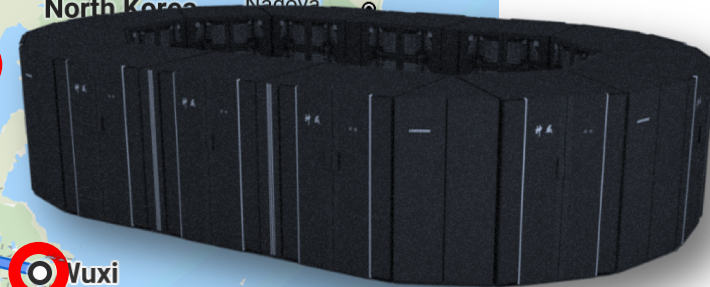
Location of Wuxi

219 systems of the TOP 500

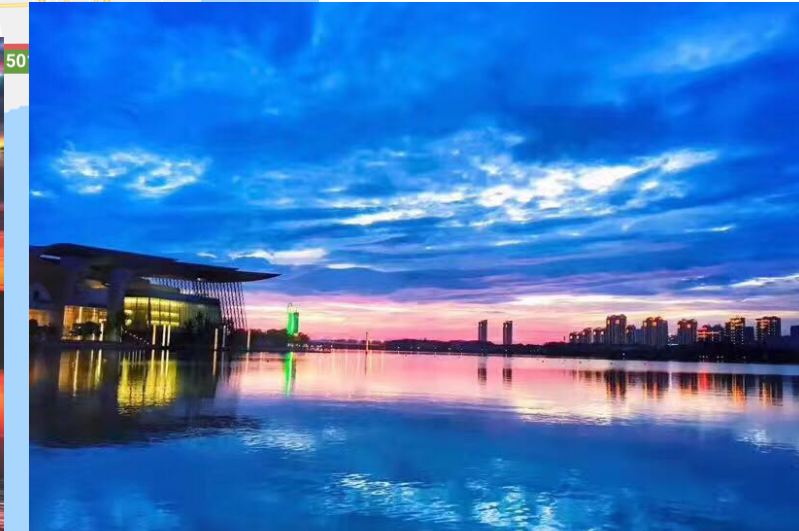
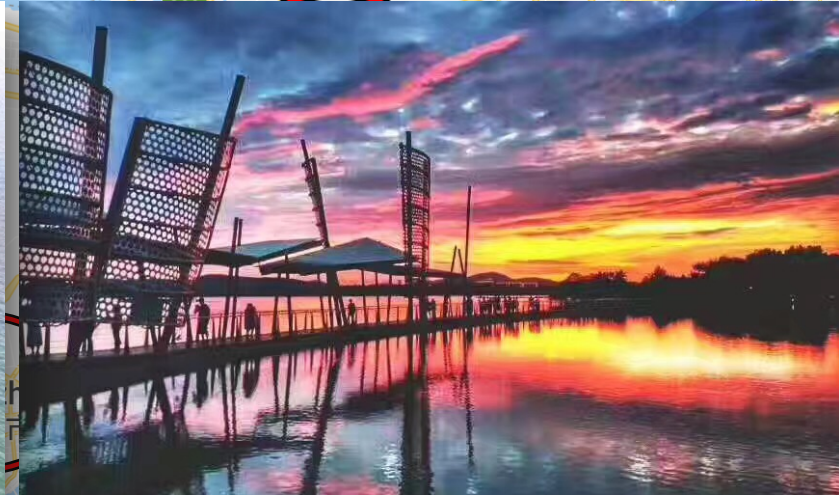
7 national supercomputing centers



Sunway BlueLight



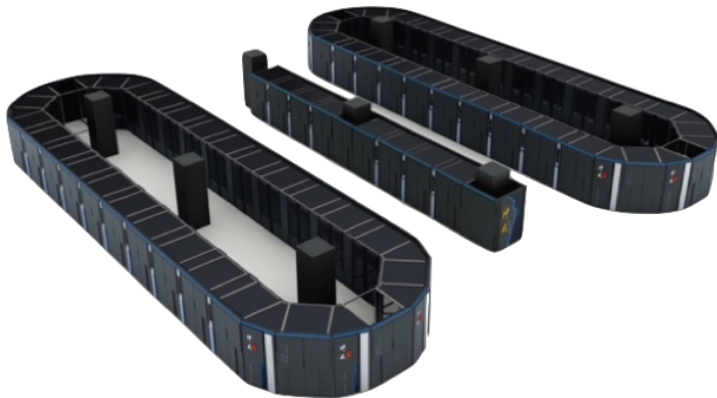






国家超级计算无锡中心
National Supercomputing Center in Wuxi

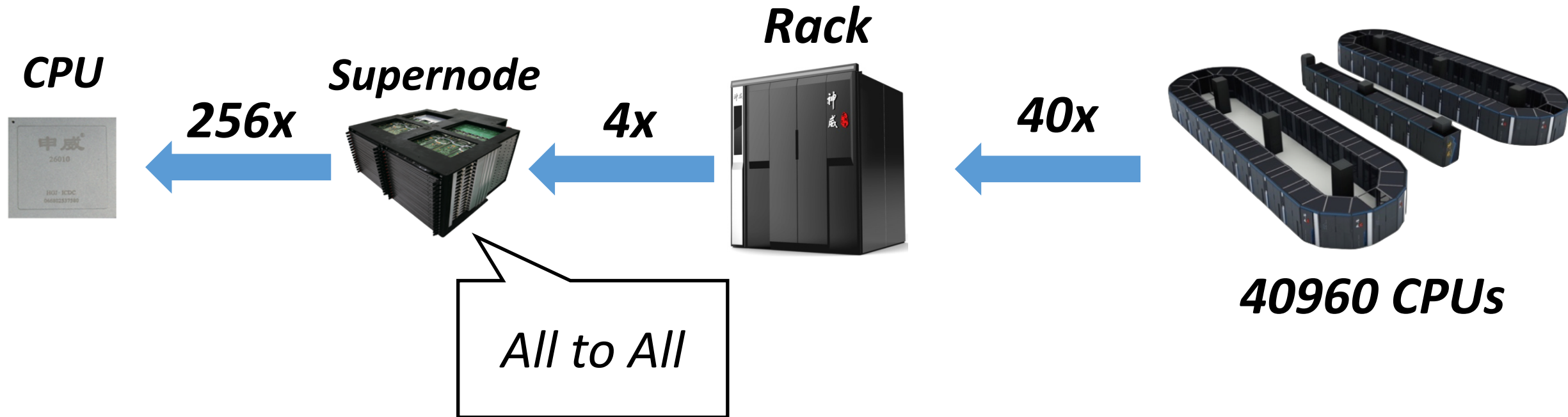
- ***National Supercomputing Center in Wuxi***
 - *Sponsored by MOST and state Gov.*
 - *Operated by Tsinghua University*
 - *Goal: world-leading SC center for important academic and industrial applications*



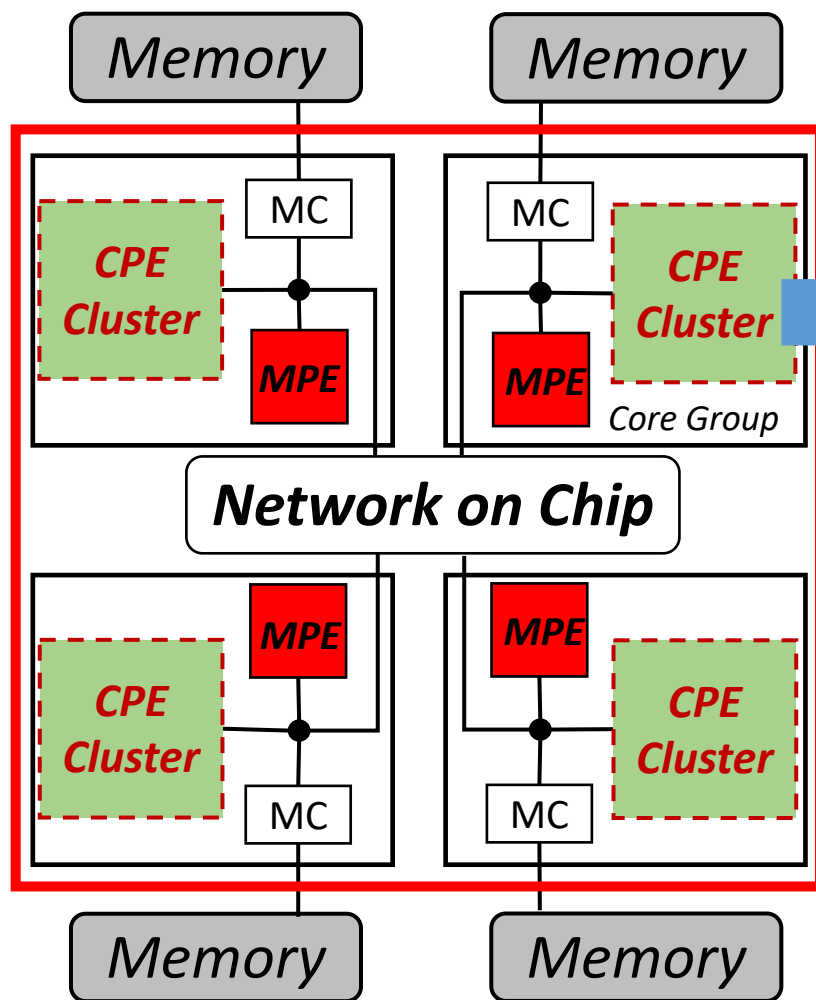
- ***Sunway TaihuLight Supercomputer***
 - *125 PFlops peak*
 - *93 PFlops LINPACK (No.3 in TOP 500)*
 - *6.05 GFlops/Watt*
 - *40,960 homegrown 260-core CPUs*



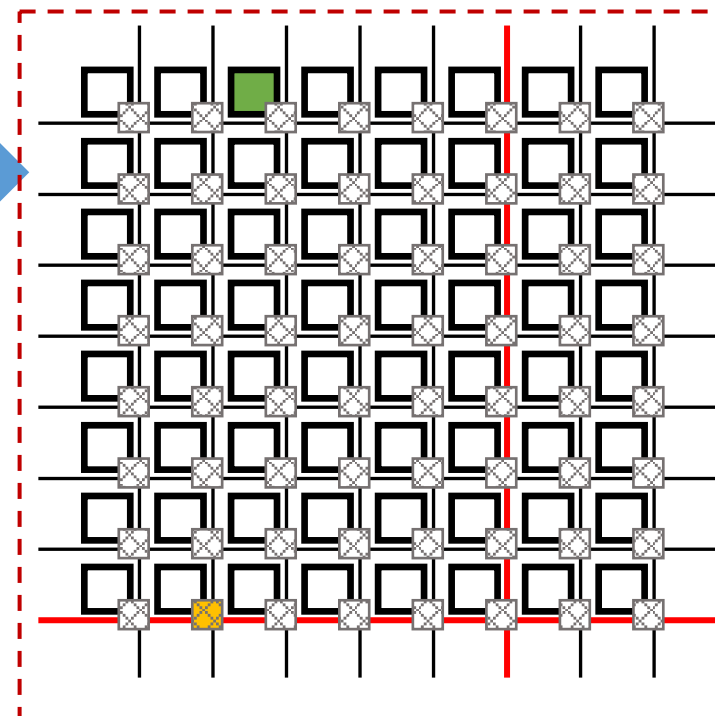
System Hierarchy



SW26010 CPU



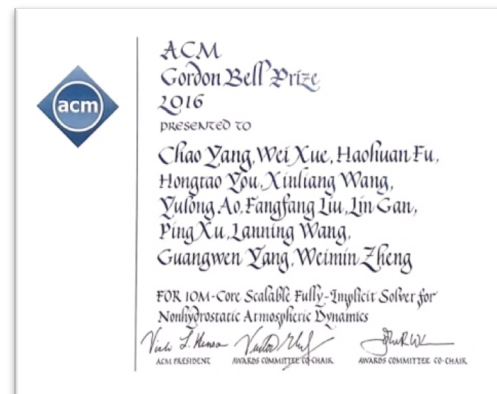
64 CPEs (Slave Cores)

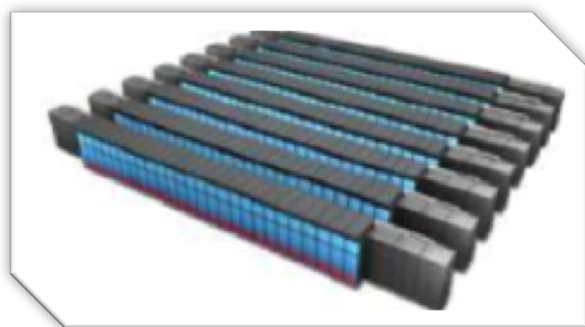
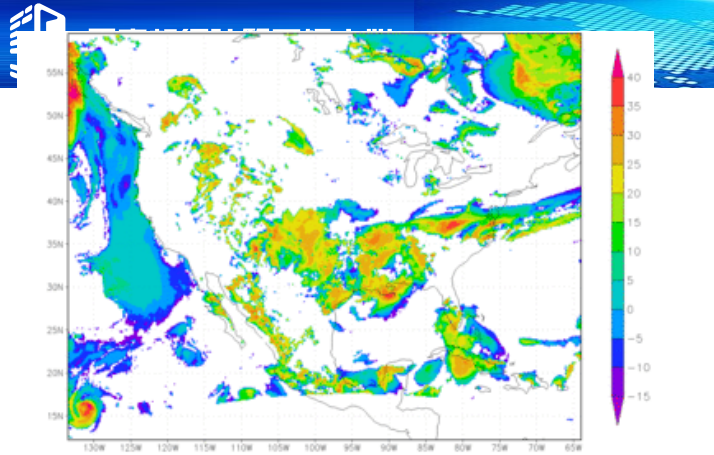




Application Overview

Since 2016, over **200** large-scale applications from over **300** research institutes covering **19** application domains, **22** full-scale applications, **40** half-scale, **100** million-core-scale, **6** Gordon Bell Finalists, and **2** Gordon Bell Prizes.





Climate Modeling

surface Exploration Geophysics

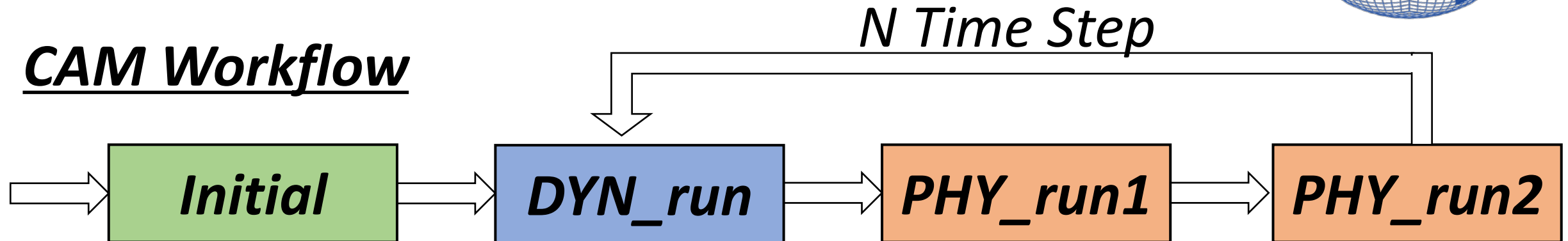


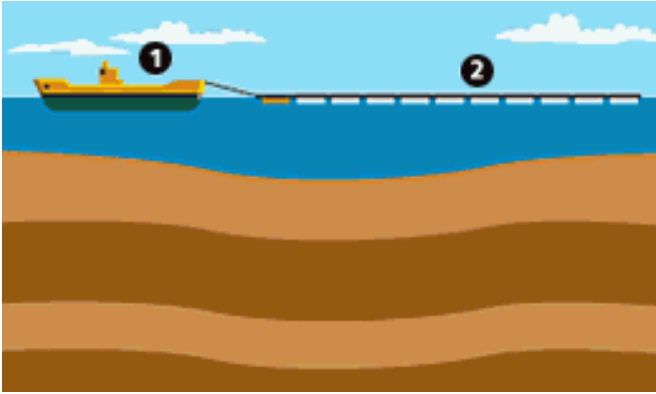
Community Atmosphere Model (CAM)

- Atmospheric model for CESM (climate model) developed by **the National Center for Atmosphere Research (NCAR)**
- One of the most computationally consuming part of CESM
- A popular & killer HPC application
- **Dynamics core (DNY)** and **Physics scheme (PHY)**



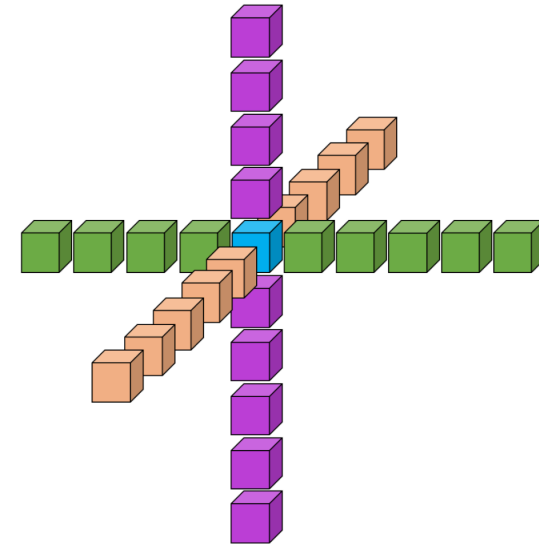
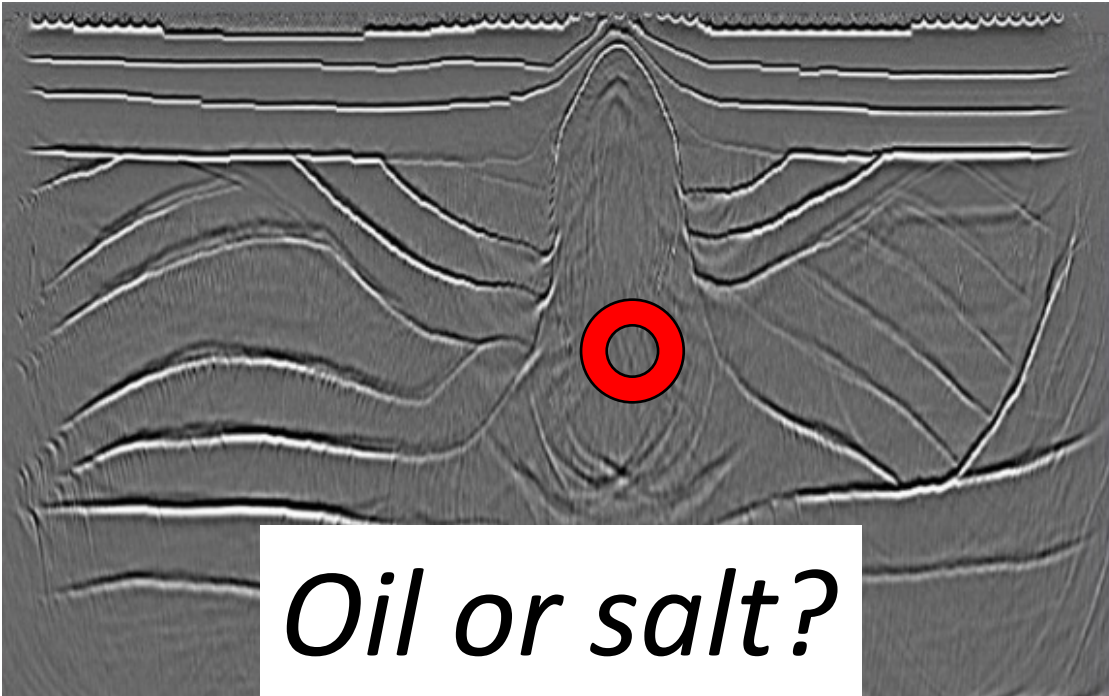
CAM Workflow





• *Elastic Reverse Time Migration (RTM)*

- *Major hot spot for migration algorithm*
- *Evolving to be more complex*



RTM stencil



Geoscience Applications to Sunway TaihuLight



OpenACC Refactoring



Athread Redesigning



OpenACC

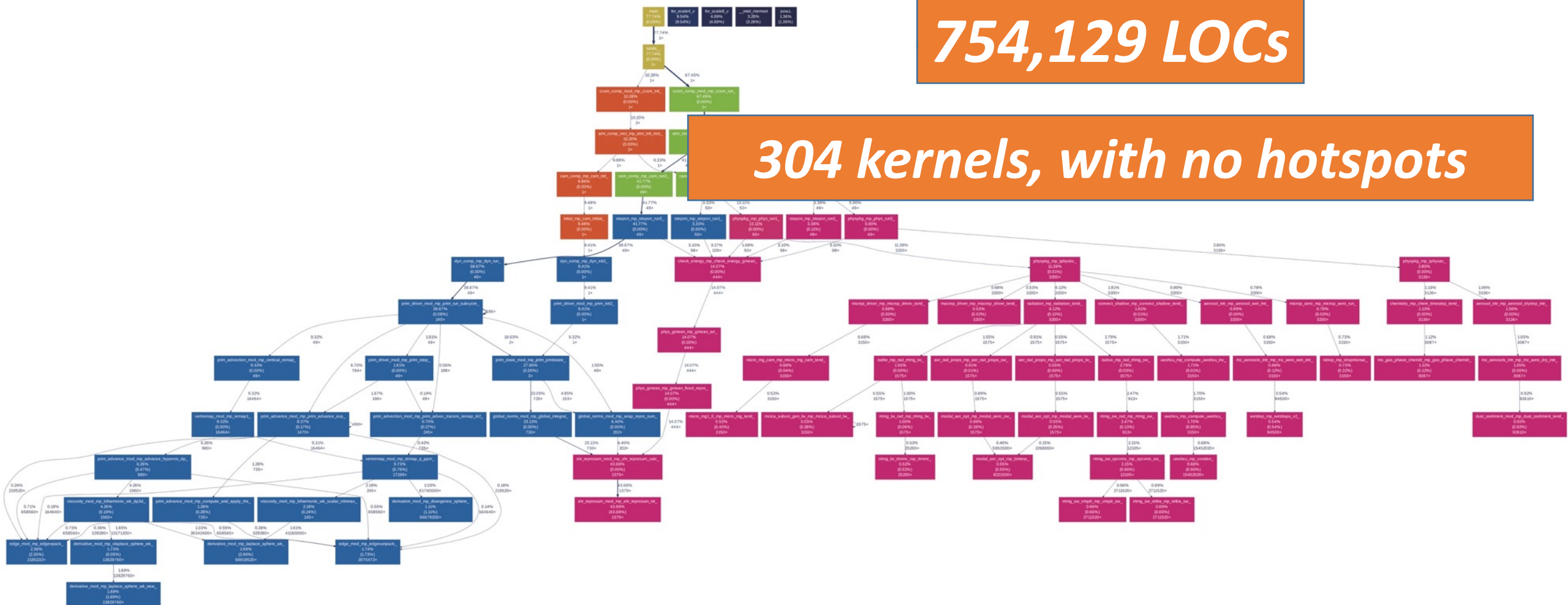
- *Sunway OpenACC Compiler*
 - *A customized version based on OpenACC 2.0*
 - *Directive-based and source-to-source compiler*



Workflow of CAM

754,129 LOCs

304 kernels, with no hotspots



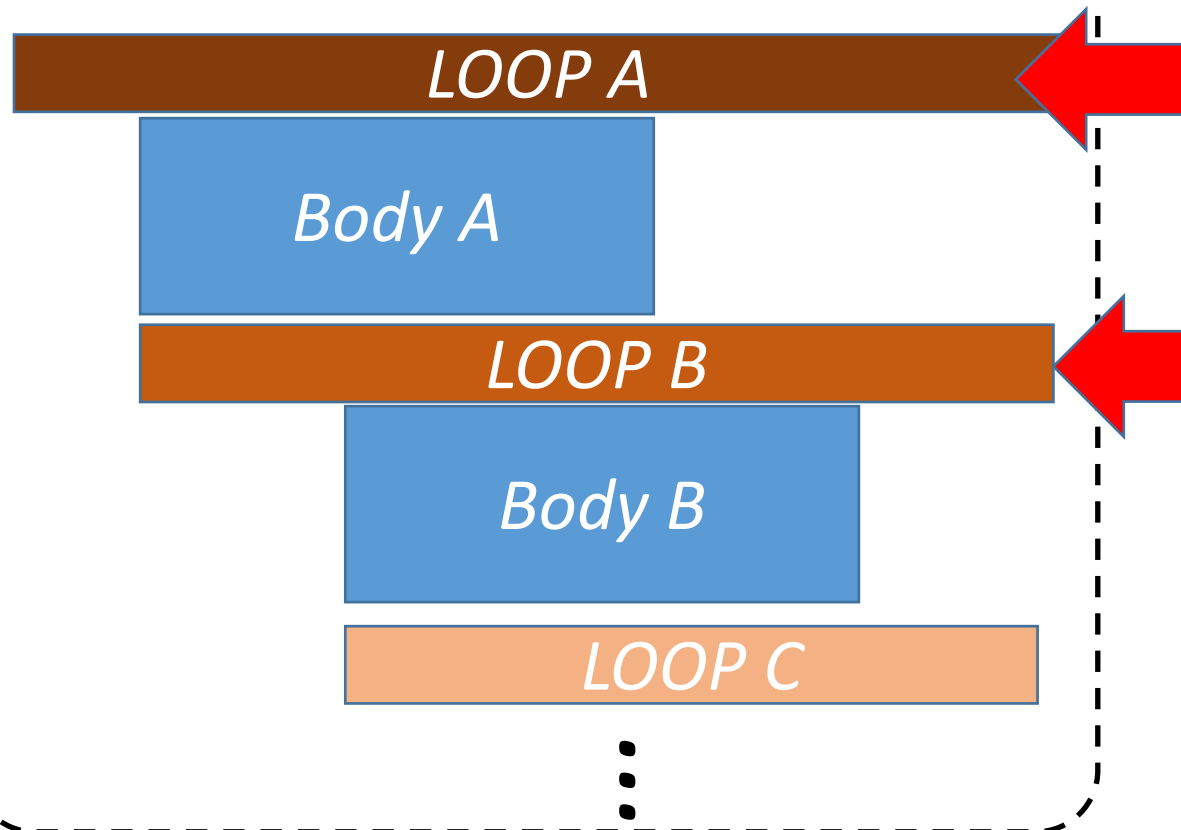


OpenACC Refactoring

- *Sunway OpenACC Compiler*
 - *A customized version based on OpenACC 2.0*
 - *Directive-based and source-to-source compiler*
- *To efficiently apply OpenACC on numerous Loops*
 - *Loop transformation tool to expose the right level of parallelism and data size*

OpenACC for Refactoring (exp. 1)

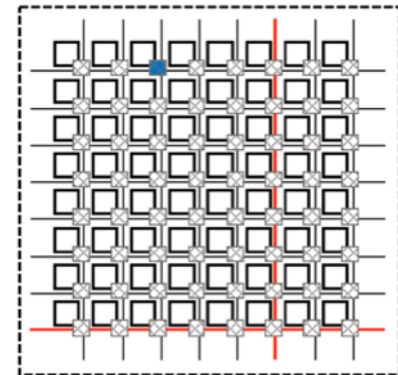
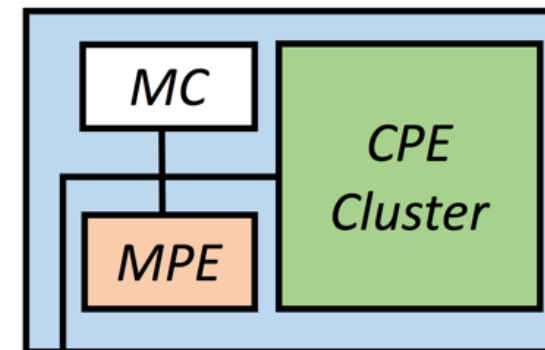
Loop Abstraction



Originally designed for multi-core system with cache

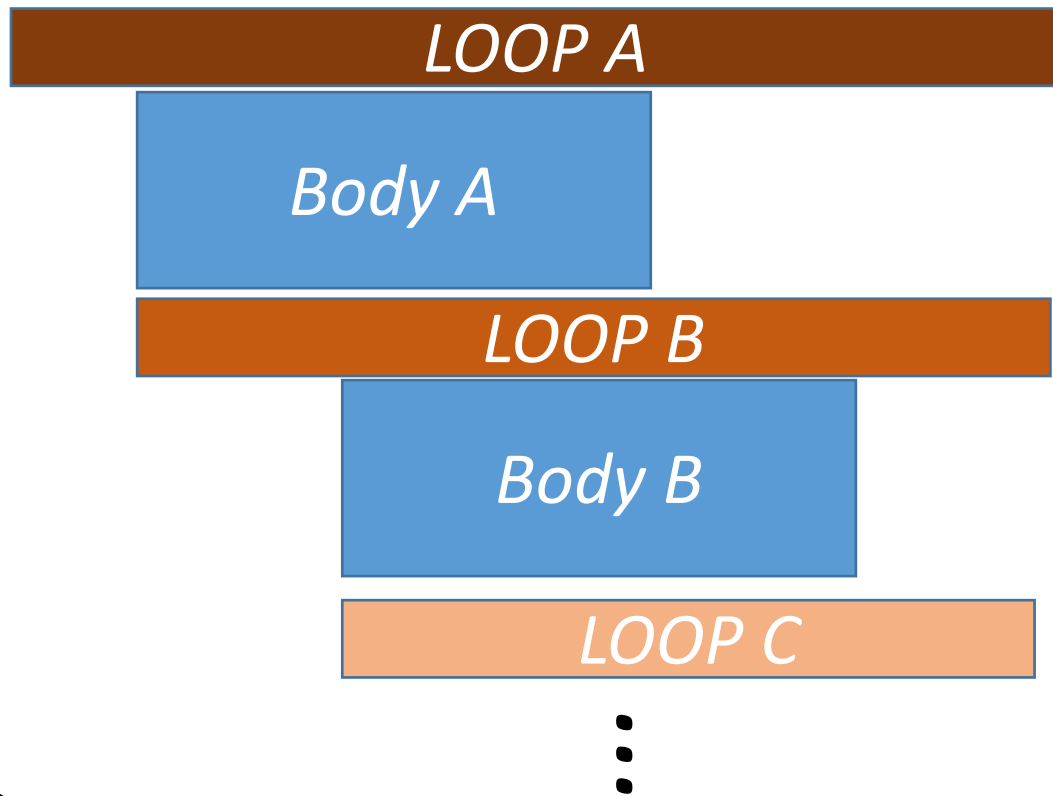
Small loop size

Fastest direction

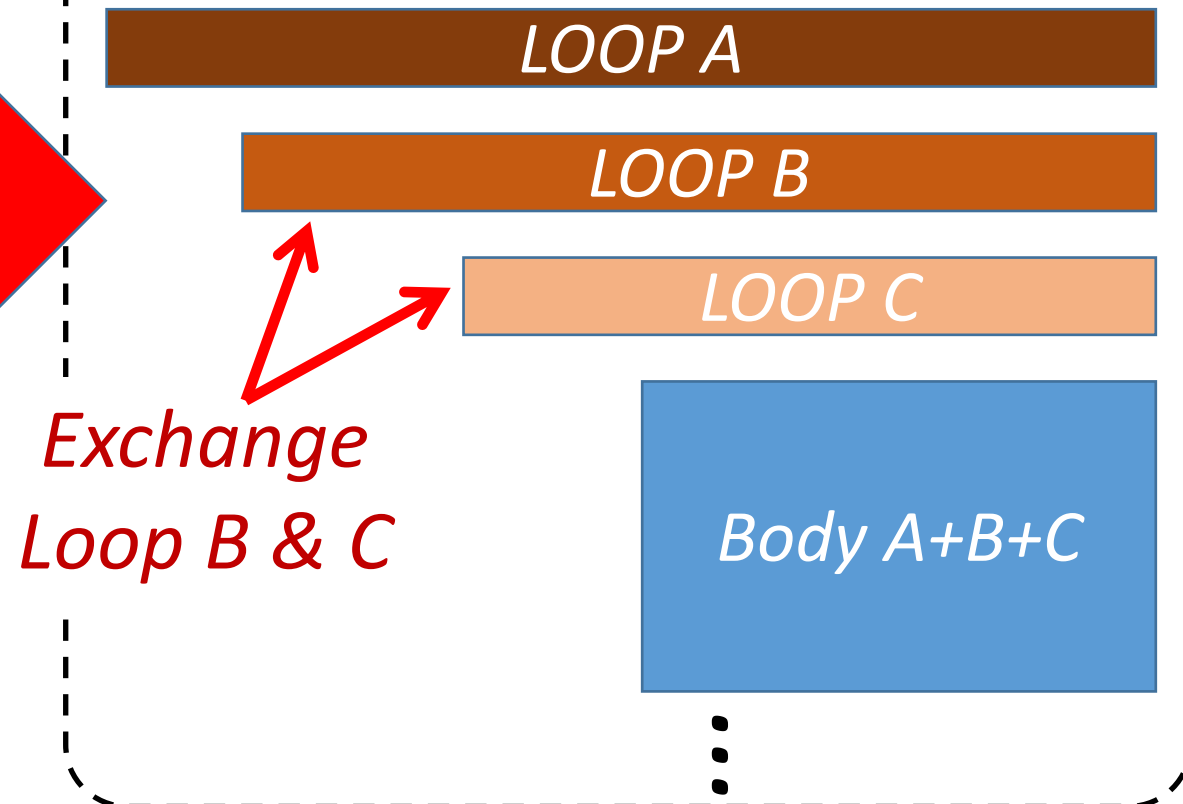


OpenACC for Refactoring (exp. 1)

Loop Abstraction

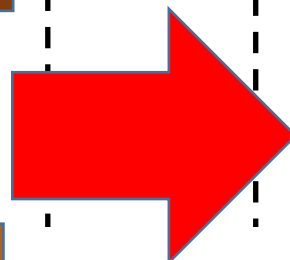
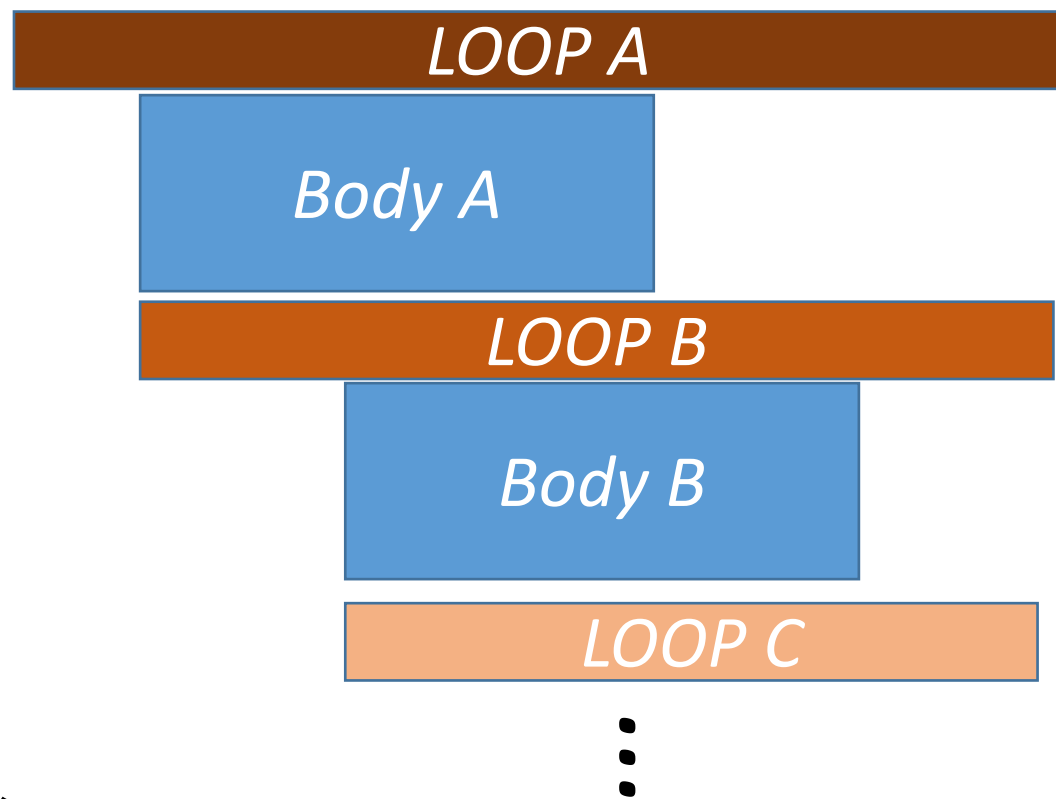


Loop Abstraction



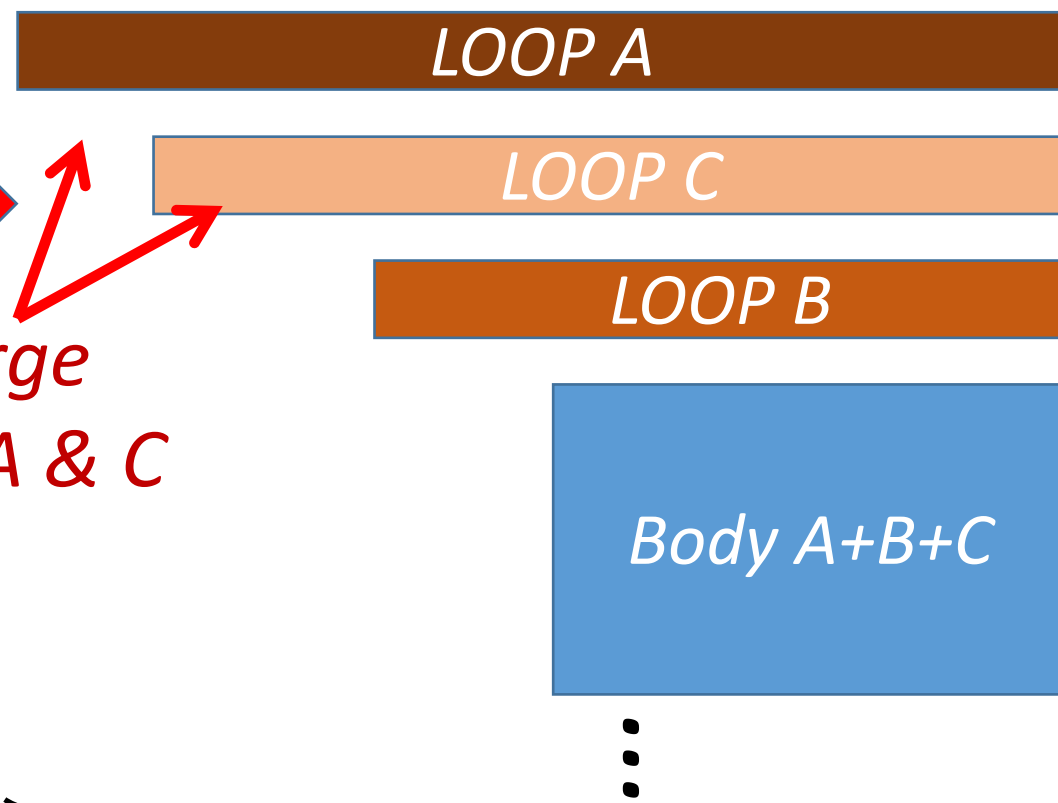
OpenACC for Refactoring (exp. 1)

Loop Abstraction



*Merge
Loop A & C*

Loop Abstraction



OpenACC for Refactoring (exp. 1)

Loop Abstraction

#OpenACC Directive

LOOP A

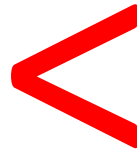
Body A

LOOP B

Body B

LOOP C

⋮



Loop Abstraction

#OpenACC Directive

LOOP (A + C)

LOOP B

Body A+B+C

⋮



OpenACC for Refactoring (exp. 2)

Loop Abstraction

LOOP A

Func 1

Func 2

Func 3

128 KB data



64 KB LDM of each CPE core

OpenACC for Refactoring (exp. 2)

Loop Abstraction

#OpenACC Directive

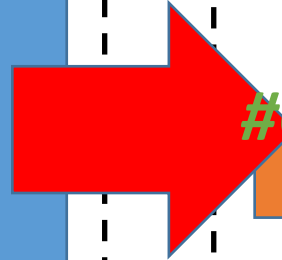
LOOP A

Func 1

Func 2

Func 3

128 KB data



Loop Abstraction

#OpenACC Directive

LOOP A1

Func 1

#OpenACC Directive

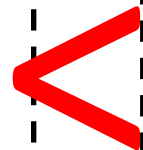
LOOP A2

Func 2

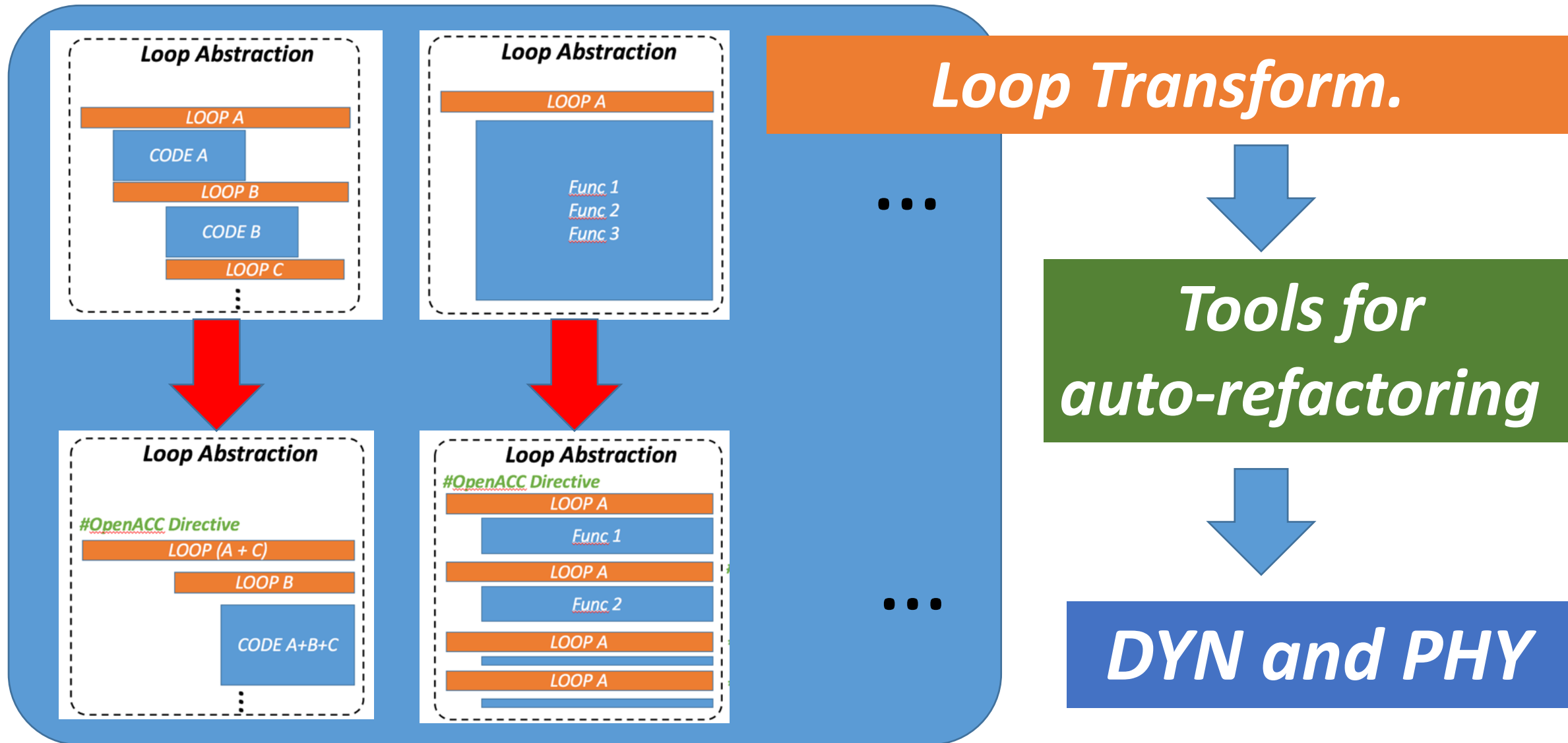
#OpenACC Directive

LOOP A3

Func 3



Tools for OpenACC Implementation





Restrictions

CAM Model

- *Difficulty in data locality to fit with the 64 KB LDM*
- *Computation patterns that are unfriendly to vectorization*
- *Absence of overlapping schemes to hide communications*

OpenACC

- *OpenACC removes options to achieve finer control of computing or memory operations or register communication*
- *Threading overhead becomes a huge issue for no-hotspot module*

Redesign is necessary



Geoscience Applications to Sunway TaihuLight



OpenACC Refactoring

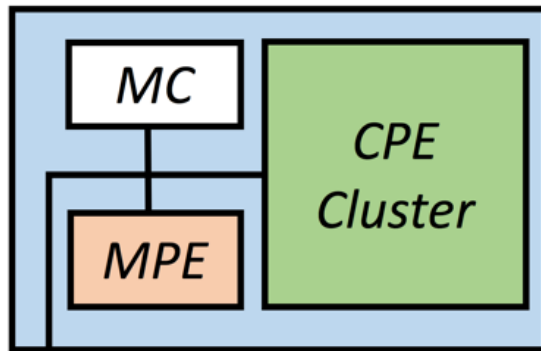


Athread Redesigning

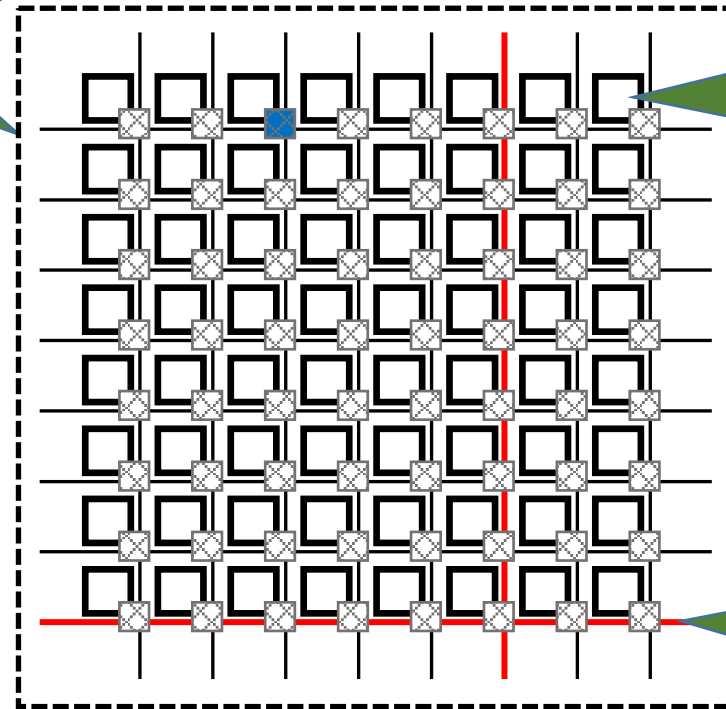
Ahtread: Fine-Grained Parallel Approach

- *Take more aggressive operations to redesign*

Parallel scheme using 64 CPES



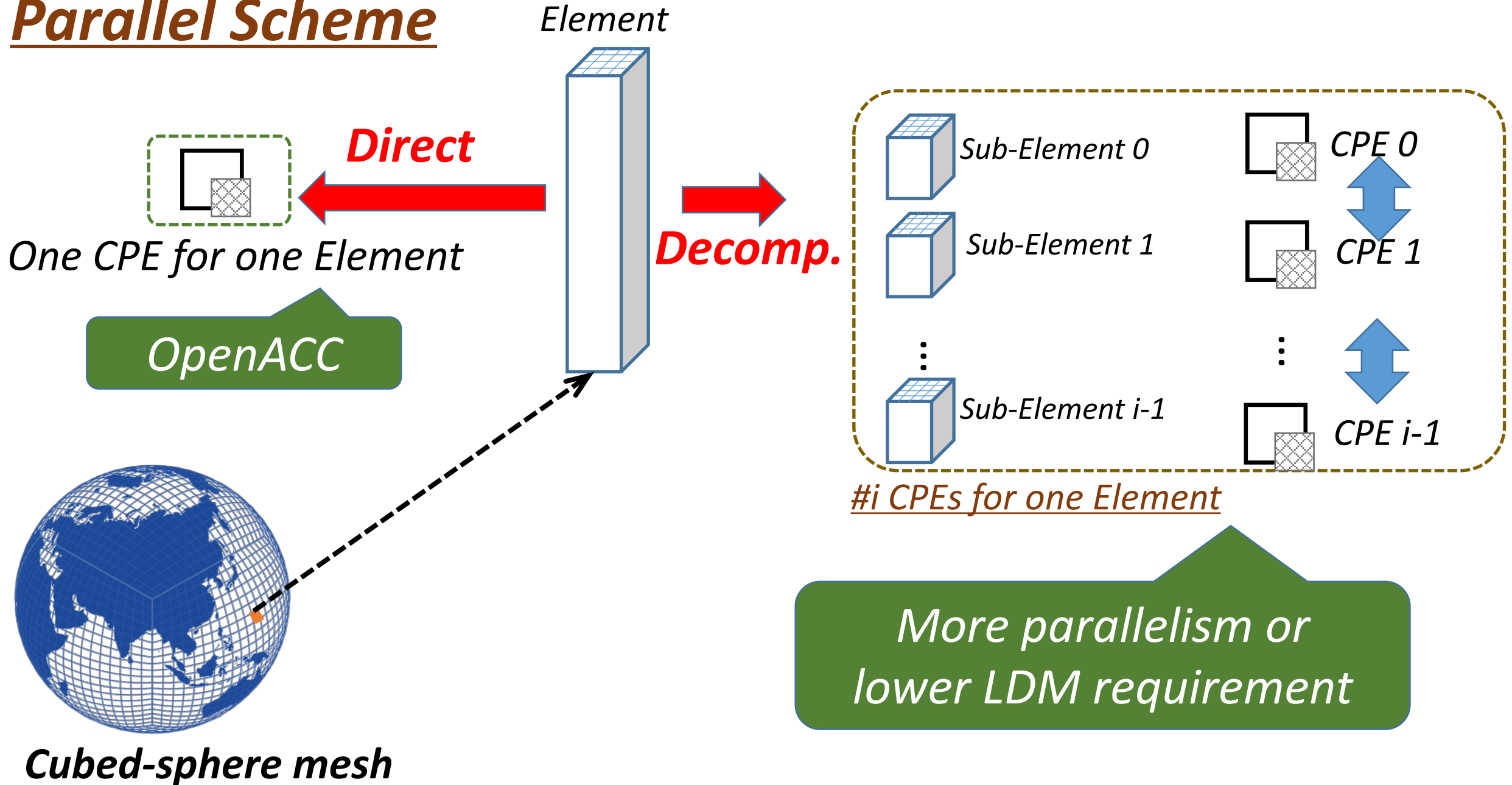
Heterogeneous arch.



One CPE
64 KB Local data
memory (LDM)

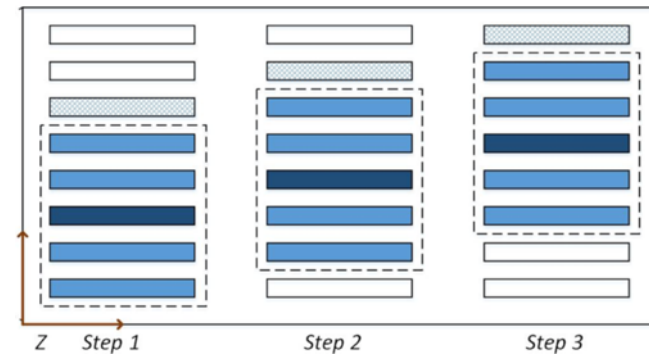
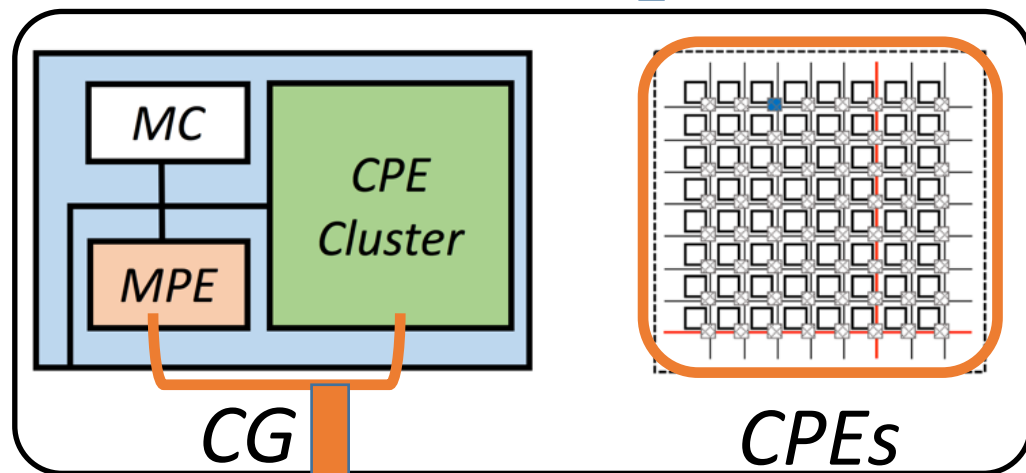
Between CPE
Fast register
communication

Parallel Scheme



Overlapping Strategies

Heterogeneous
MPE for inter-Comm.
CPEs for Comp.

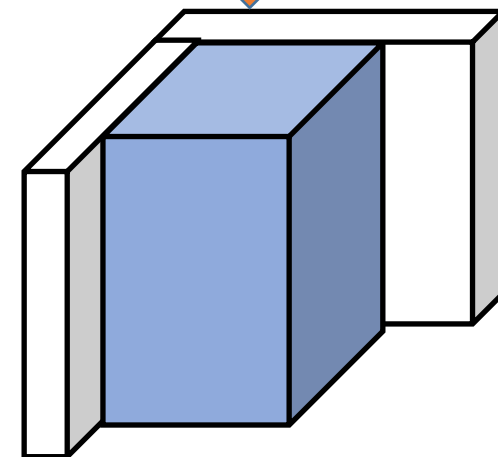


Pipeline Scheme

CPEs Computation

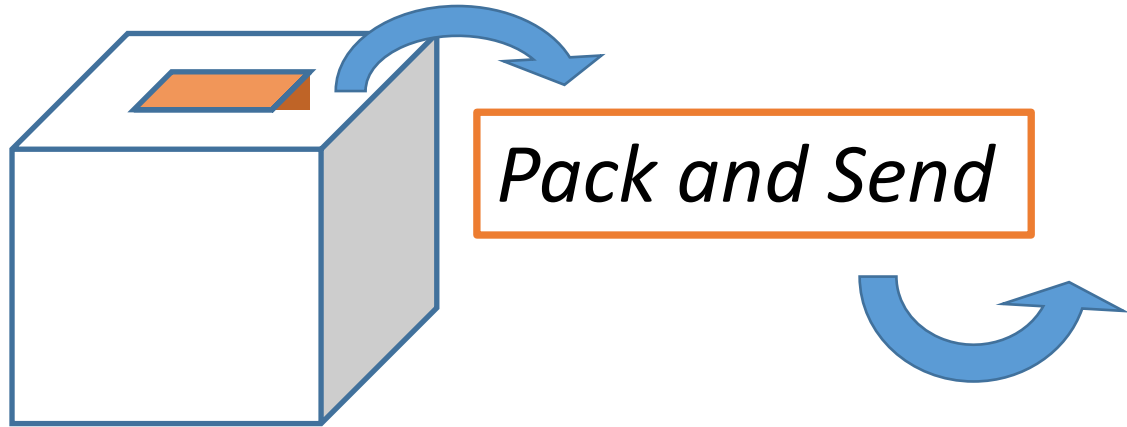
MPE Boundary Exchange

Stencil-like
Hybrid decomp.
into Inner and
Outer parts for
halo overlapping

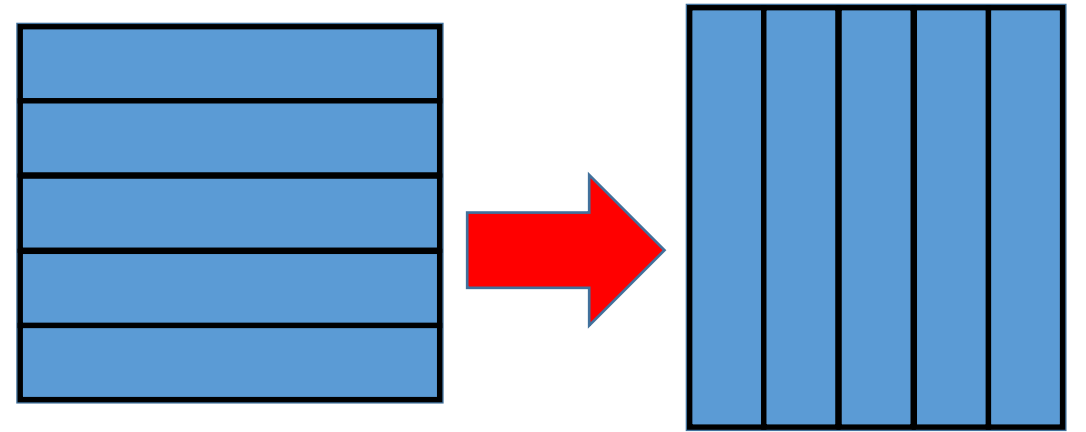


Array Transposition by Register Communication

Boundary Exchange

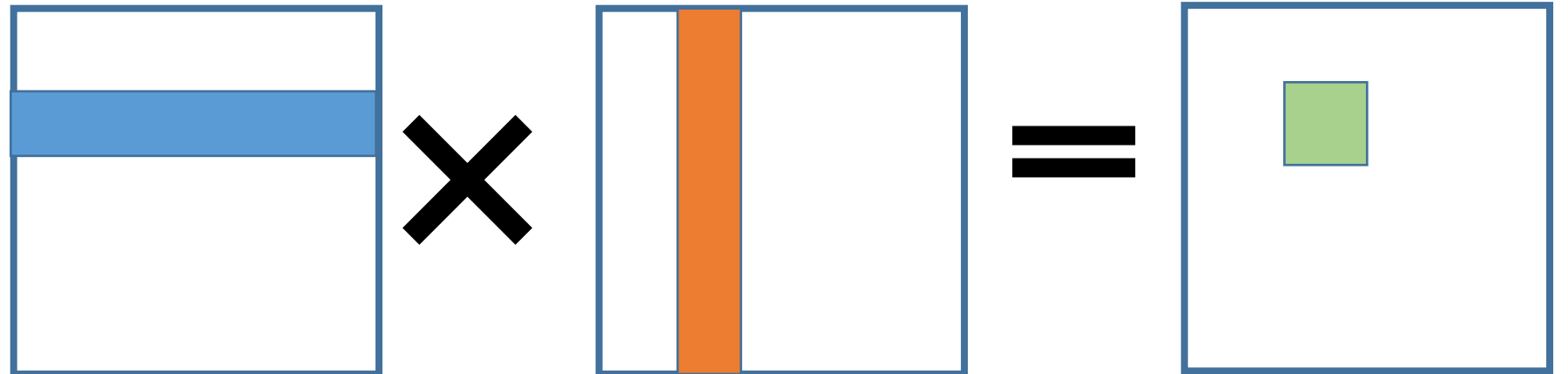


On-the-fly Matrix Transposition

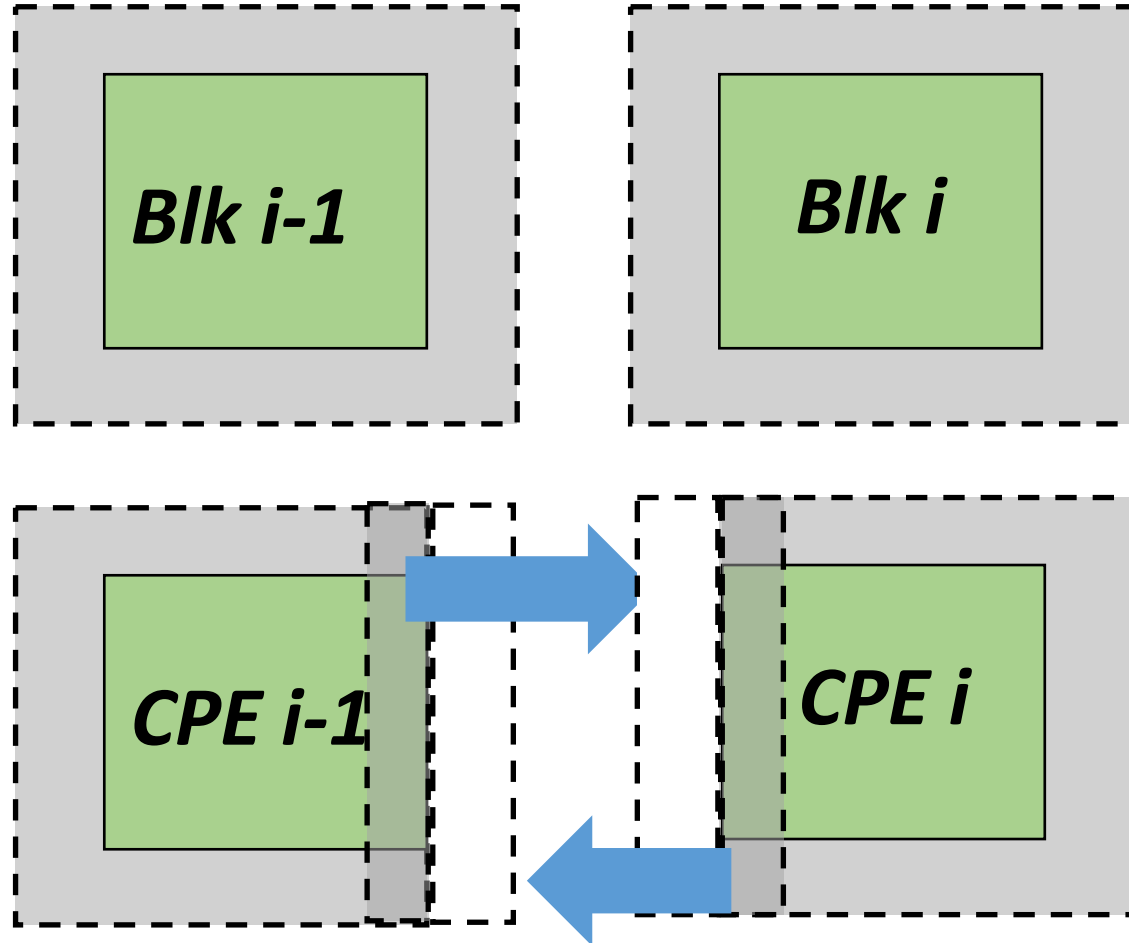


Matrix-Matrix Multiplication

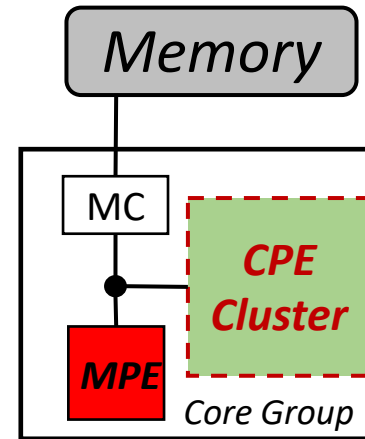
Register
communication
(~10cycles)



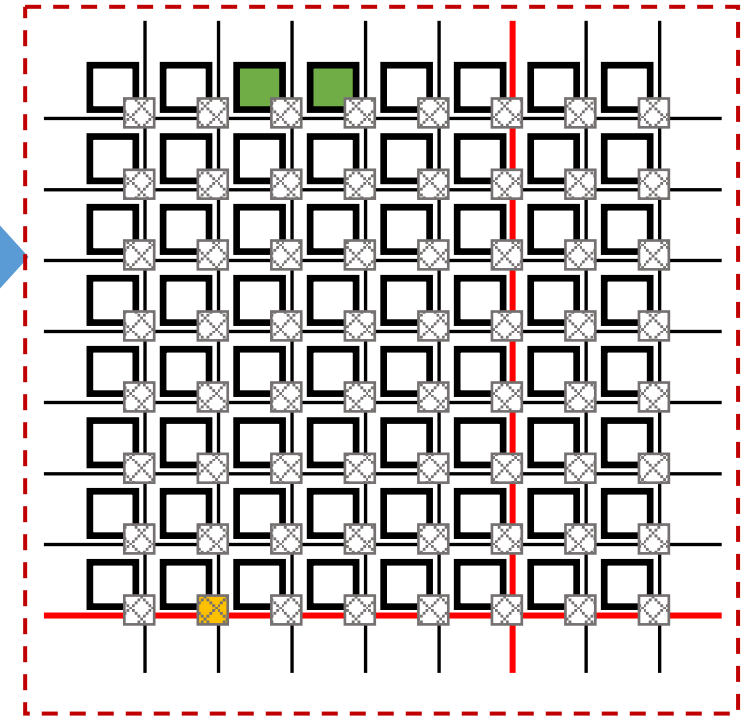
Data Sharing via Register Communication



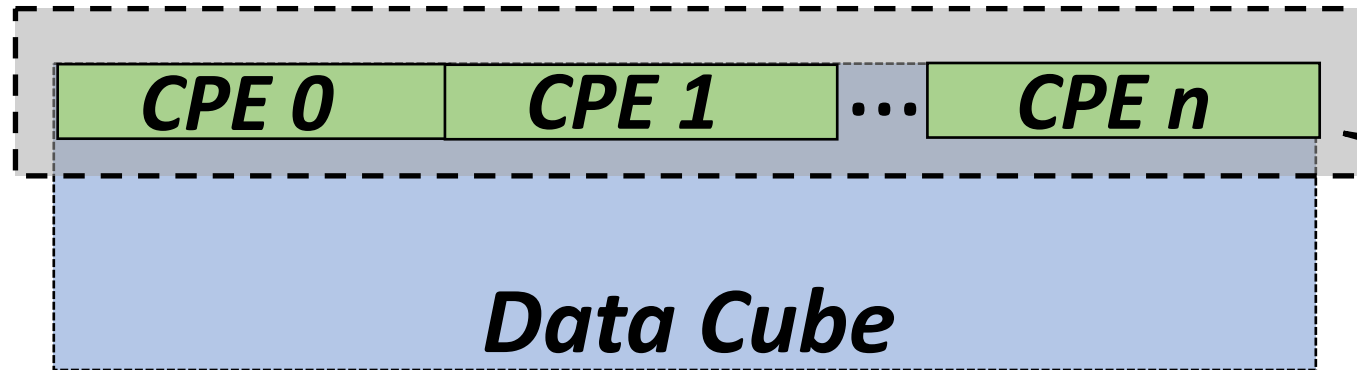
Register Comm.



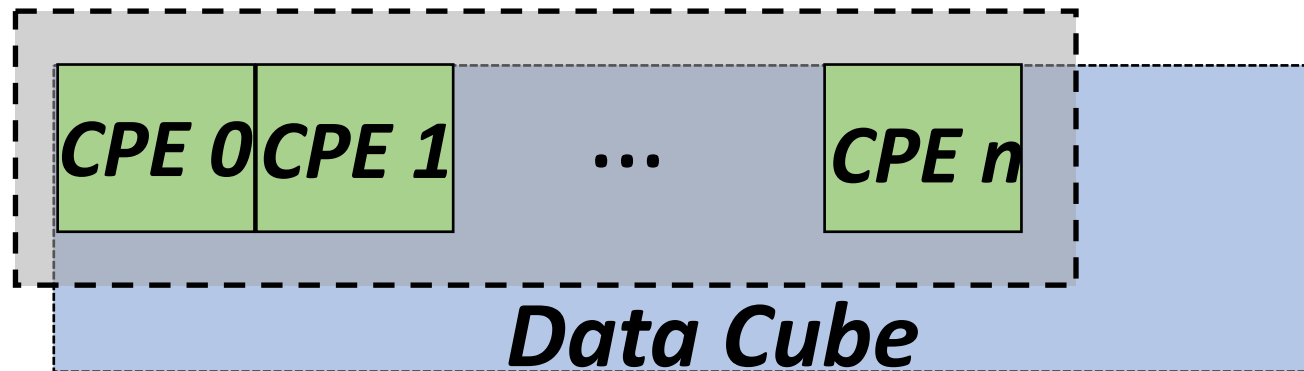
64 CPEs (Slave Cores)



Locality-aware Memory Design



One CPE to control a **stride/block** data

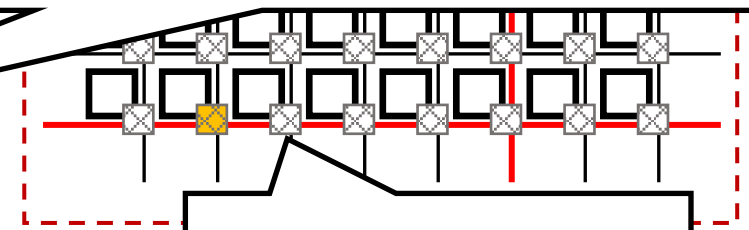


64 CPEs (Slave Cores)

Better memory

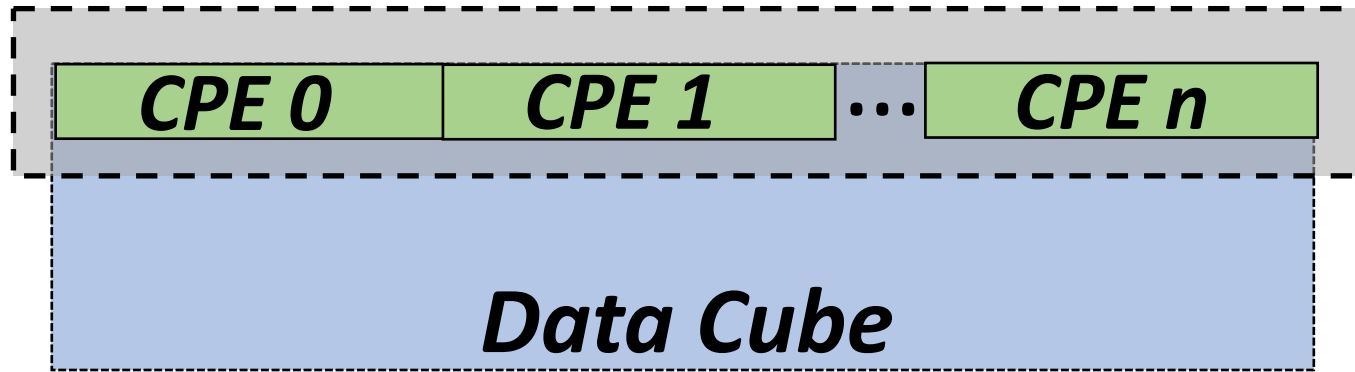
Trade-off

Less redundancy

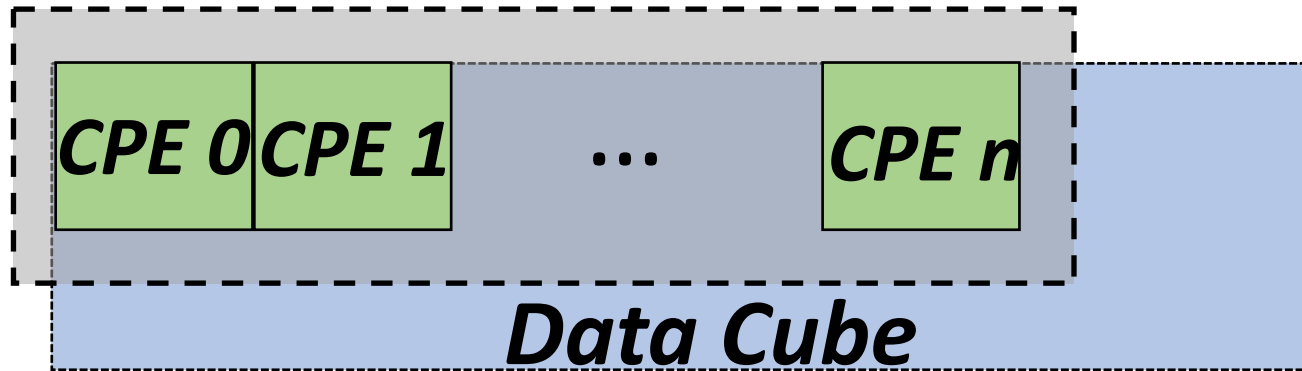


64 KB LDM

Locality-aware Memory Design



One CPE to control a **stride/block** data

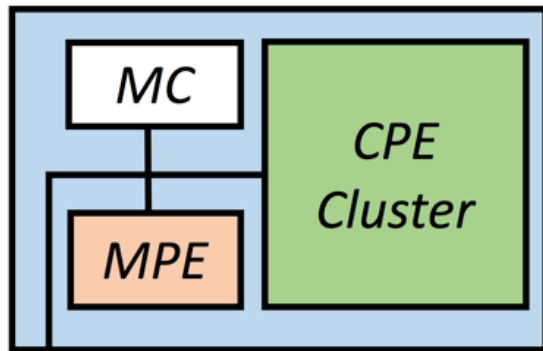


**Cache-like Mechanism
for better locality**

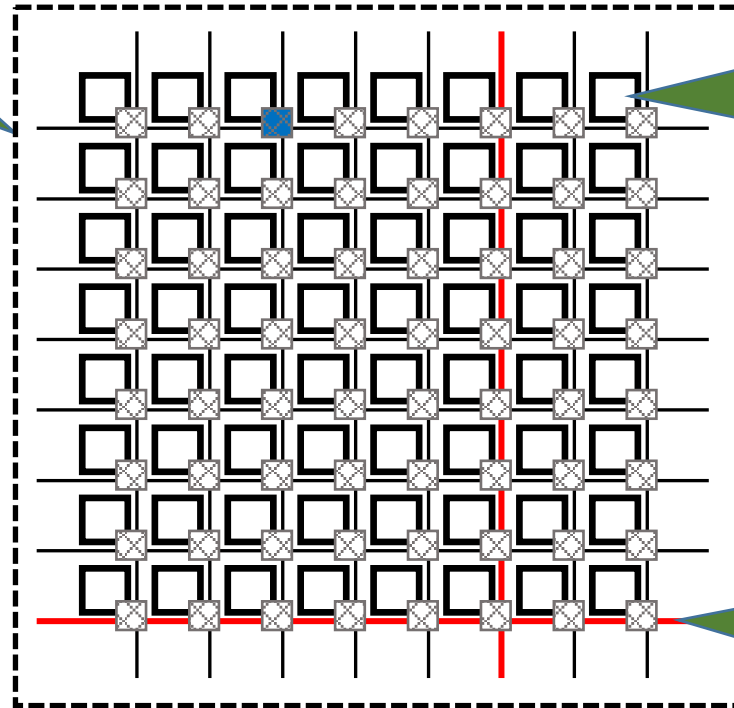
Systematic Redesign Solution by Athread

- *Take more aggressive operations to redesign*

Parallel scheme using 64 CPES



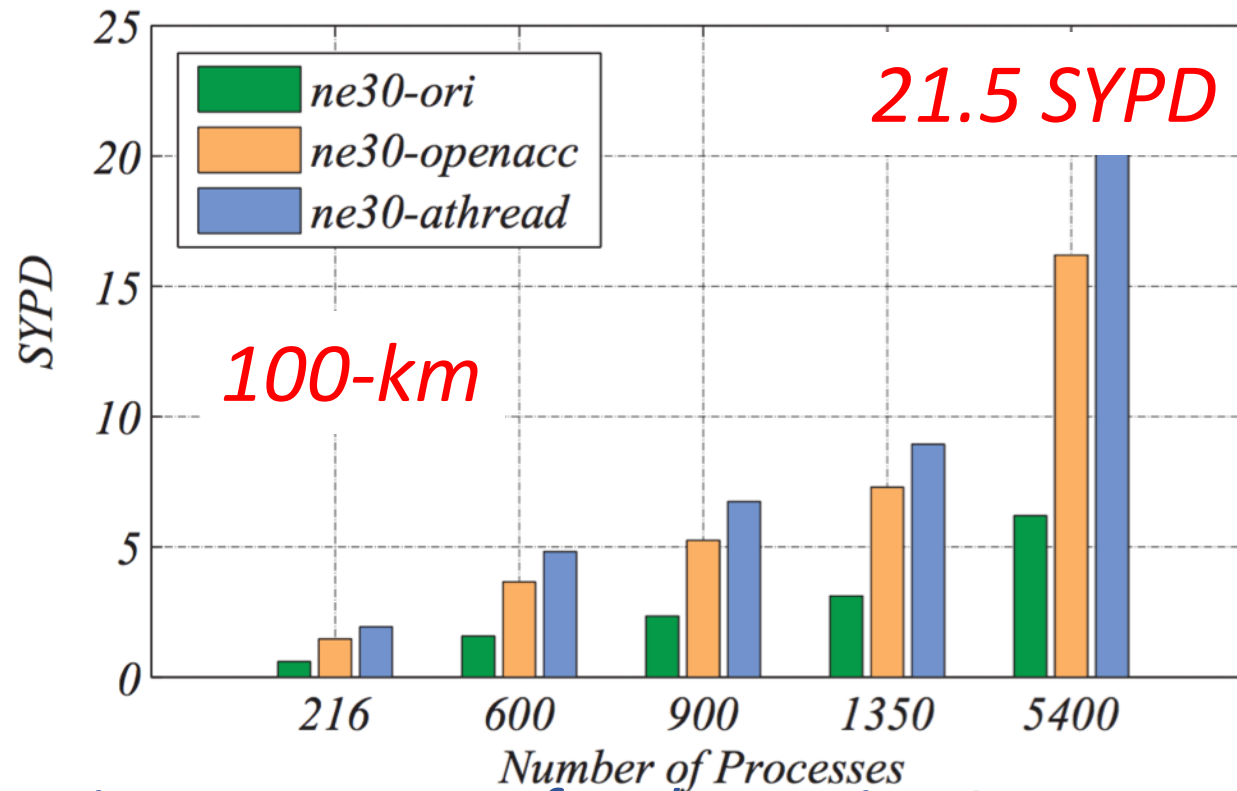
Heterogeneous arch.



One CPE
64 KB Local data
memory (LDM)

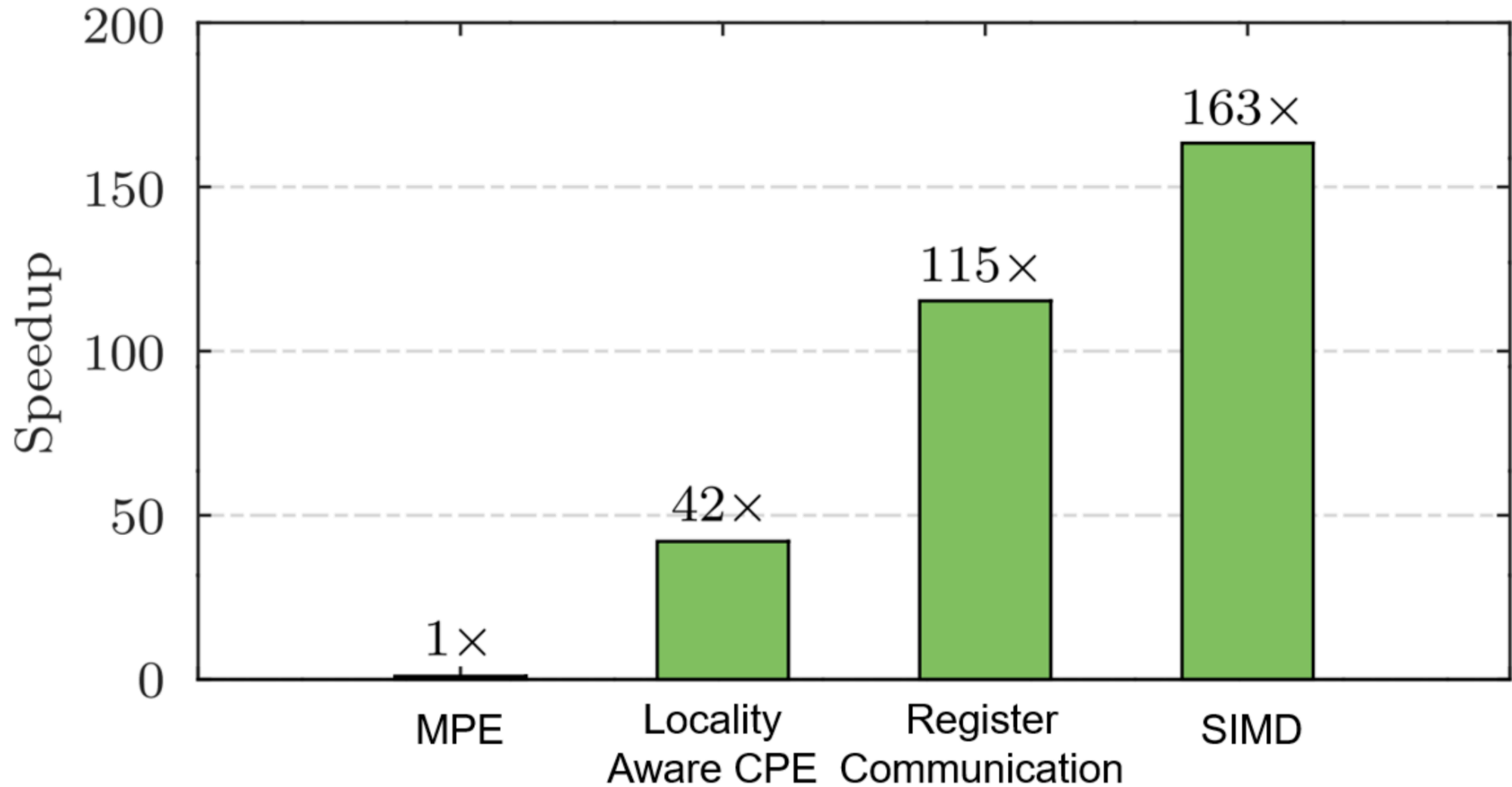
Between CPE
Fast register
communication

Performance Speedups for CAM on TaihuLight

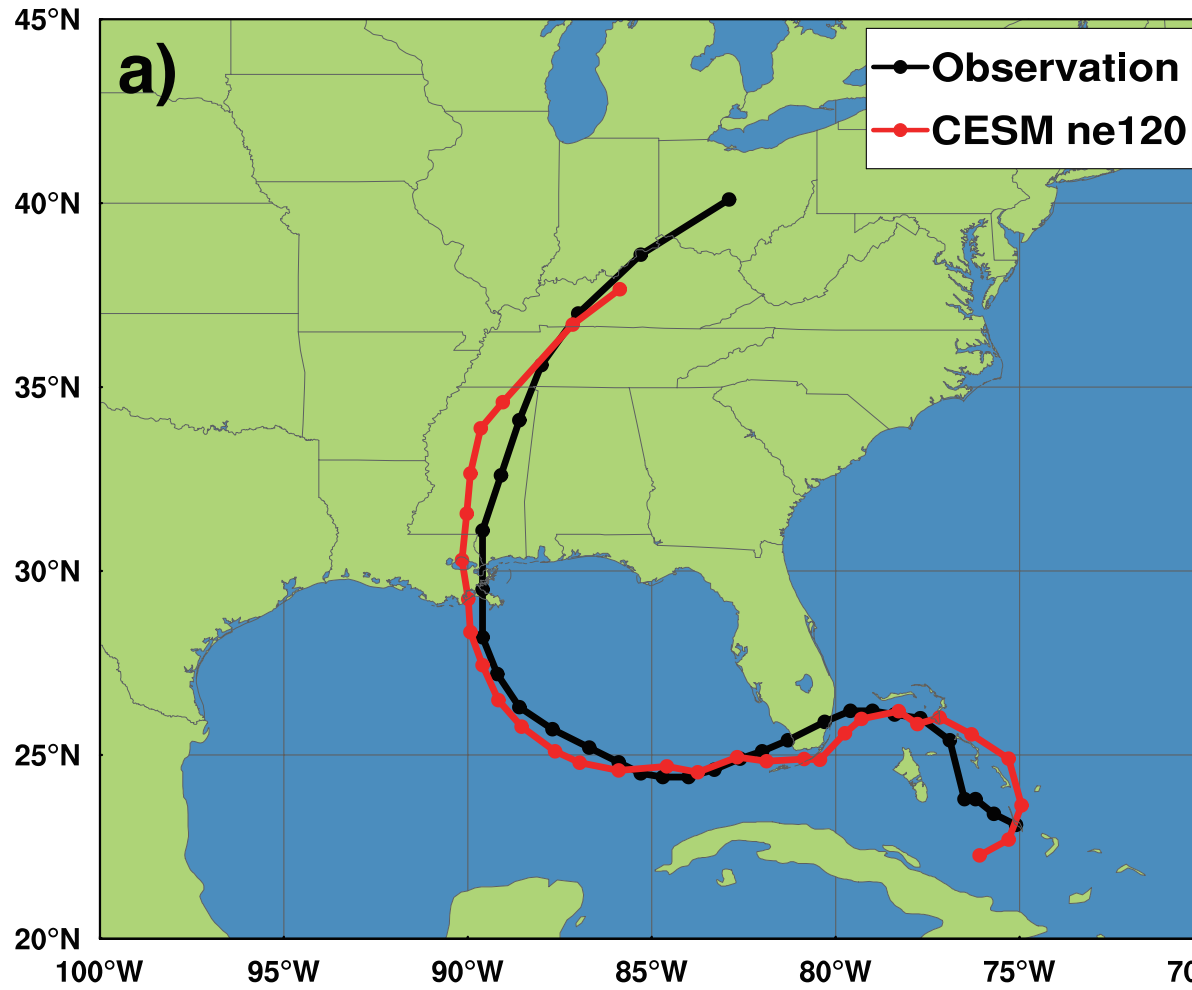


The performance improvements for the entire CAM model in ne30 and ne120. ori refers to the original version based on MPE, openacc refers to the usage of OpenACC directive, and athread refers to the further usage of Athread.

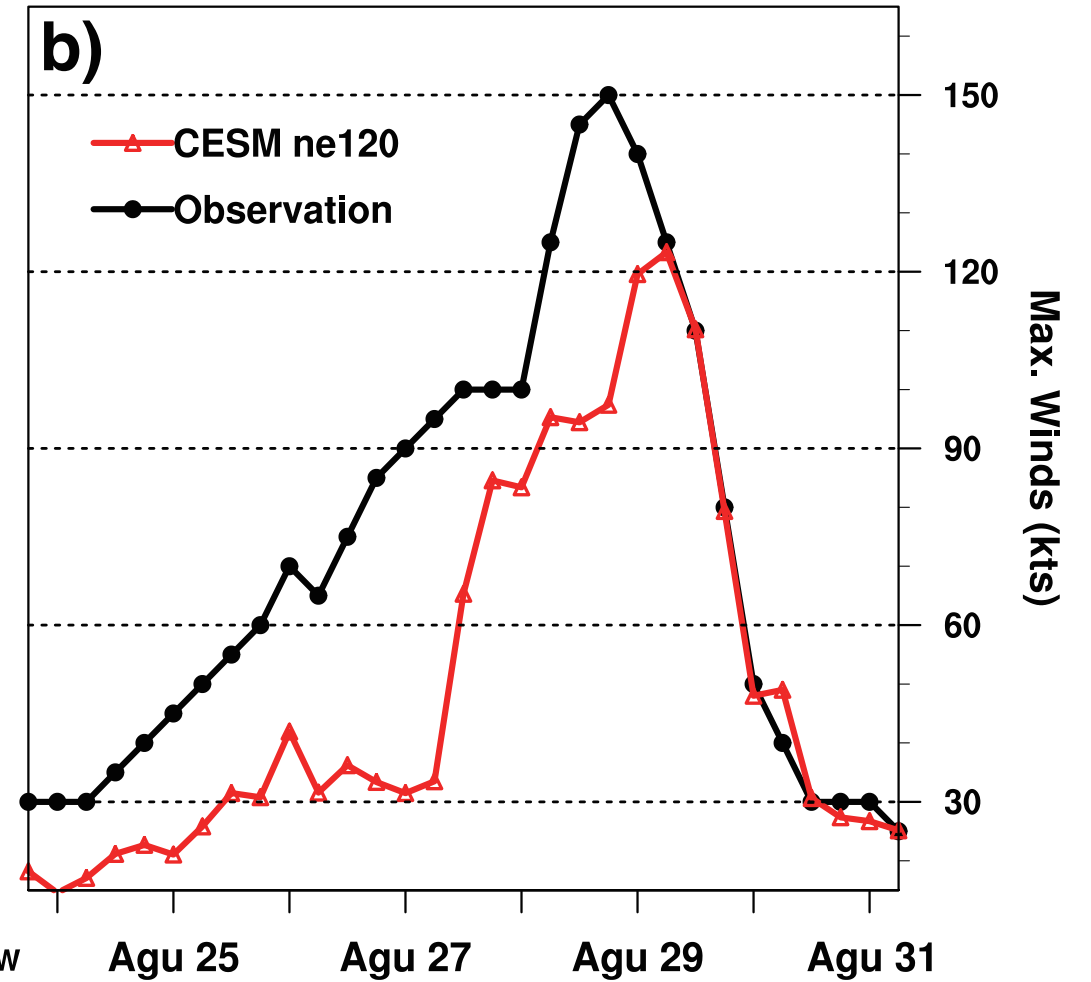
Performance Speedup of Elastic RTM



Hurricane Simulation – Katrina (2005)

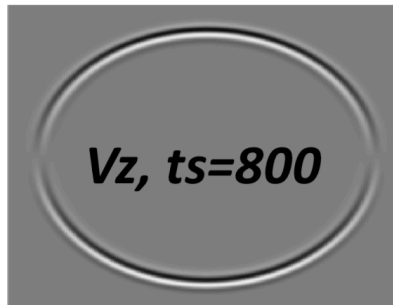
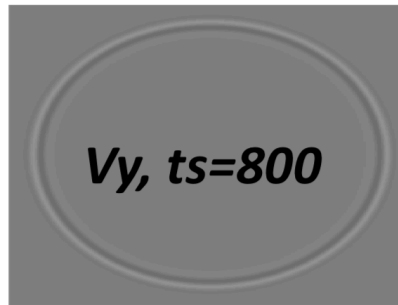
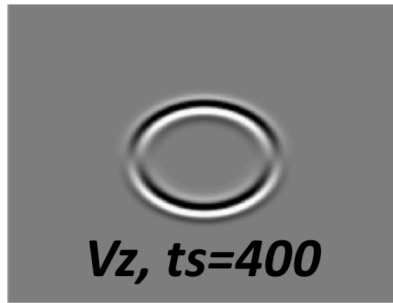
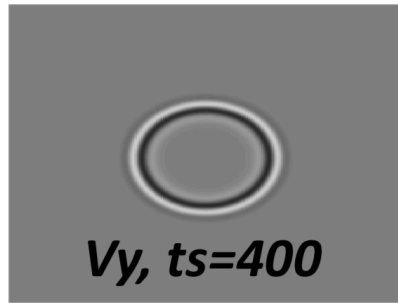
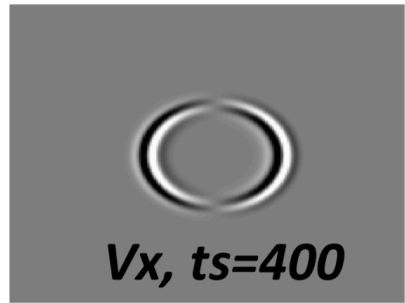


Track



Intensity

Validation of Elastic RTM



Wave propagation of three particle velocities at step 400 and 800, views from xoy, yoz, and xoz

Drho-Layer Benchmark

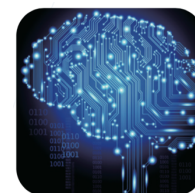
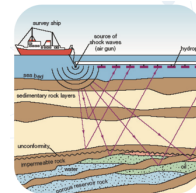
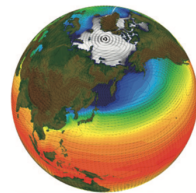
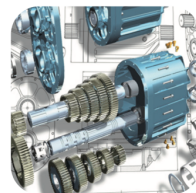


Summary

- **Solutions for dealing with *real-world* applications**
 - *OpenACC Refactoring, and the automatic tools for loop transformation*
 - *Athread Redesign, with different solutions towards architectural features*
- **Proved thread-level comm. a good option**
- **Virtual cache: make up the con of the Sunway architecture**
- **Provide better support for applications & current and upcoming (Exascale) many-core supercomputers**



Big Computing & Big Data



Software Ecosystem

xMath

Thunder

VASP

swLBM

Petsc

OpenFoam

LAMMPS

swAllies

bowtie

Star

Gromacs

swRTM

Deep Learning



C

C++

Fortran

OpenACC

Athread

gprof

Jperf

swLU





Scaling Geoscience Applications on Sunway Supercomputer

Lin Gan

Assistant Professor, Tsinghua University, Beijing
Assistant Director, NSCC-Wuxi, Jiangsu

大寒

