

# Preliminary Implementation of GRAPES global model on Sunway Taihu light

Zhiyan JIN

Numerical Weather Prediction Center, China Meteorological Administration

Wei Xue, Ping Xu, Hongsong Meng, Yongbin Jiang, Zhao Liu  
Tsinghua University and National Supercomputing Center in Wuxi

18<sup>th</sup> Workshop on HPC in Meteorology, ECMWF, 24-28 September 2018



# Outline

1. Introduction to GRAPES global model
2. Introduction to Sunway TaihuLight
3. Refactorization of critical computing kernels of GRAPES
  - ✓ Semi-Lagrange interpolation
  - ✓ Helmholtz solver
  - ✓ Halo communication optimization
  - ✓ Stencil kernel optimization
4. Optimizations of GRAPES with OpenACC\*
5. Evaluation
6. Summary

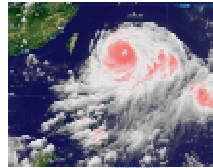


# Introduction of GRAPES global model



# Meteorological Centers and Research Institutions in CMA Campus

National Satellite Met. Center



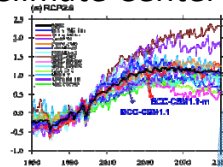
National Met. Information Center



National Meteorological Center



National Climate Center



CMA Met. Observation Center



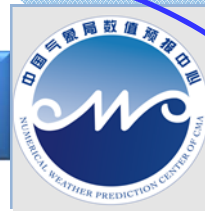
CMA Public Met. Service Center



CMA Weather Modification Center



CMA Numerical Weather Prediction Center



Chinese Academy of Met. Sciences  
8 CMA Specialized Research Institutes

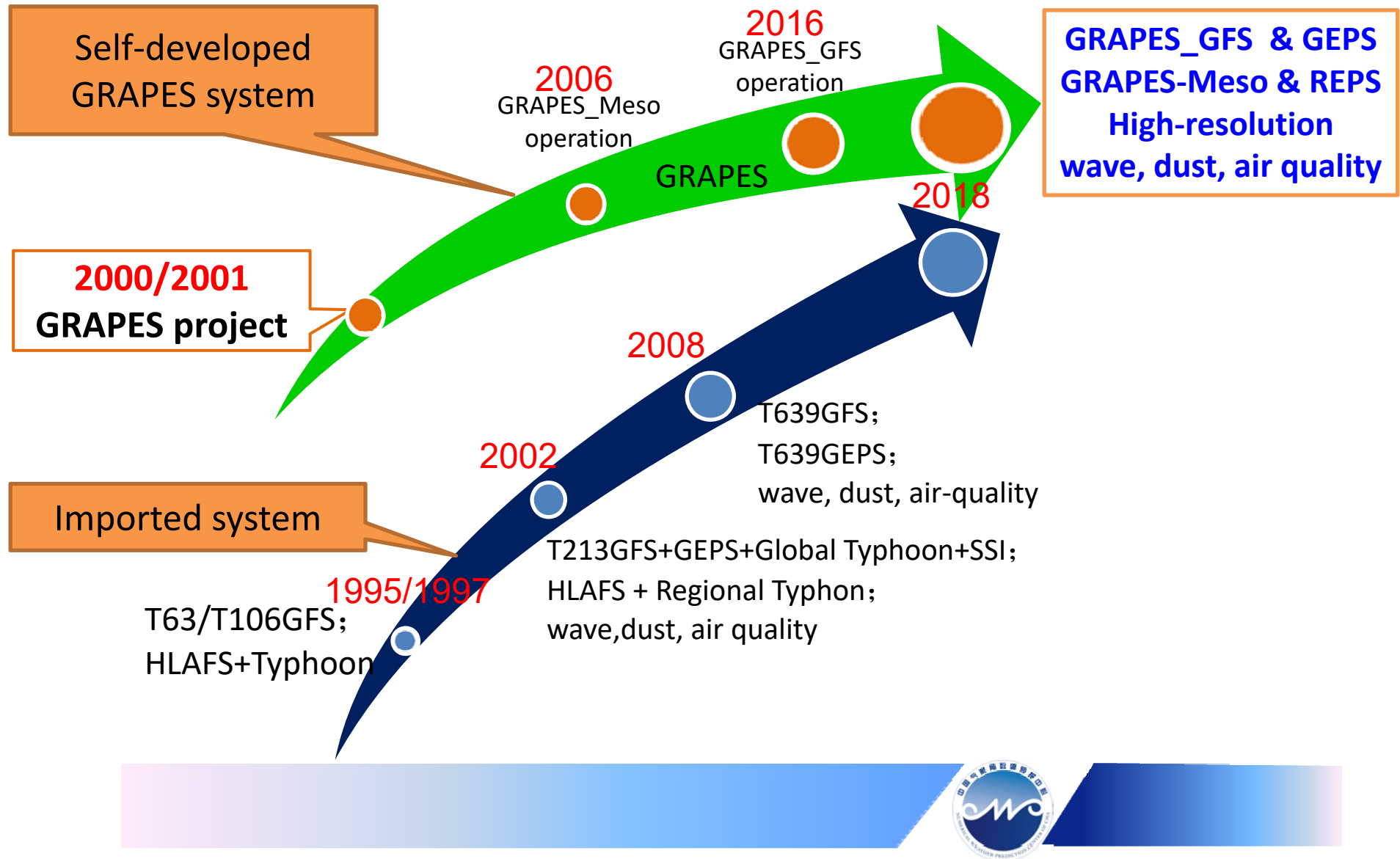


**Mission: 0-10day weather fcst.**



# History of CMA NWP Operational systems

—from imported to self-developed core techniques/systems



# About GRAPES

## GRAPES

(Global/Regional Assimilation PrEdiction System)

- **Model**

- Fully compressible equations with shallow atmosphere approximation
- Regular Lat/lon, Arakawa-C with V at poles
- Terrain-following Z, Charney-Phillips staggering
- 2TL-SISL time integration
- PRM scalar advection (conservation & monotonicity)

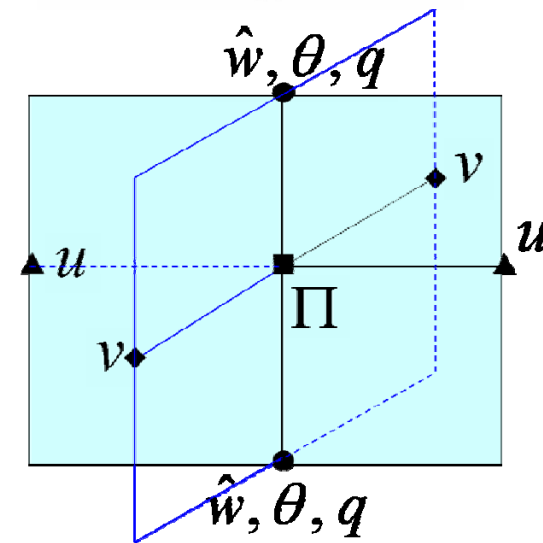
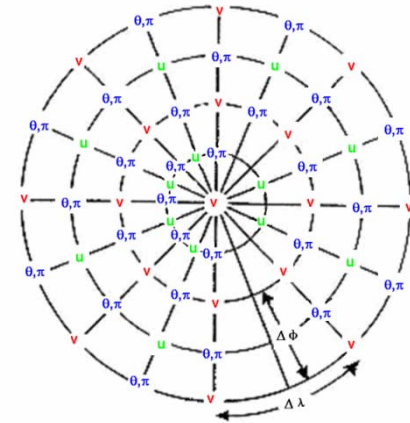
- **Assimilation**

- Unified 3/4DVAR framework
- Incremental analysis
- Digital filter, initialization in 3DVAR, weak constraint in 4DVAR



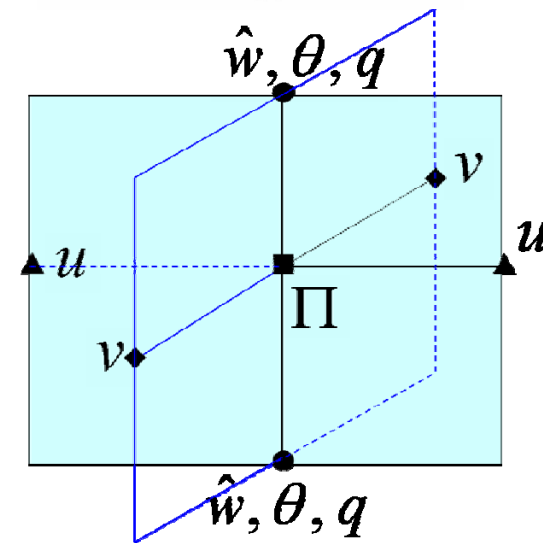
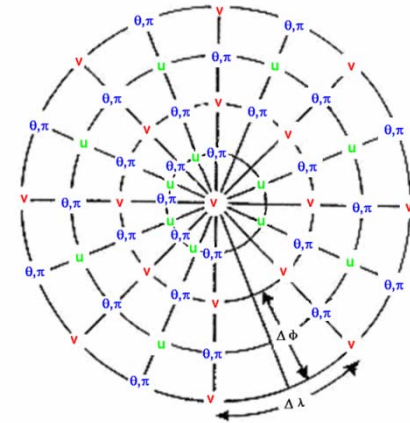
# GRAPES Dynamic Core

- ◆ **Fully compressible equations**
  - ◆ Height-based terrain-following coordinate
  - ◆ Option of Hydrostatic and Non-hydrostatic
  - ◆ Modified Arakawa C lat-lon horizontal grid
  - ◆ Charney-Phillips vertical grid
- ◆ **Off-centered 2-time-level semi-implicit semi-Lagrangian (SISL) time-stepping**
- ◆ **3D vector form of SISL formulation**
- ◆ **PRM for scalar advection**
- ◆ **Preconditioned GCR for Helmholtz Eq.**
- ◆ **Spherical & polar effects of trajectory calculation**
- ◆ **Polar filter**



# GRAPES Dynamic Core

- ◆ **Fully compressible equations**
  - ◆ Height-based terrain-following coordinate
  - ◆ Option of Hydrostatic and Non-hydrostatic
  - ◆ Modified Arakawa C lat-lon horizontal grid
  - ◆ Charney-Phillips vertical grid
- ◆ **Off-centered 2-time-level semi-implicit semi-Lagrangian (SISL) time-stepping**
- ◆ **3D vector form of SISL formulation**
- ◆ **PRM for scalar advection**
- ◆ **Preconditioned GCR for Helmholtz Eq.**
- ◆ **Spherical & polar effects of trajectory calculation**
- ◆ **Polar filter**





# Physics package

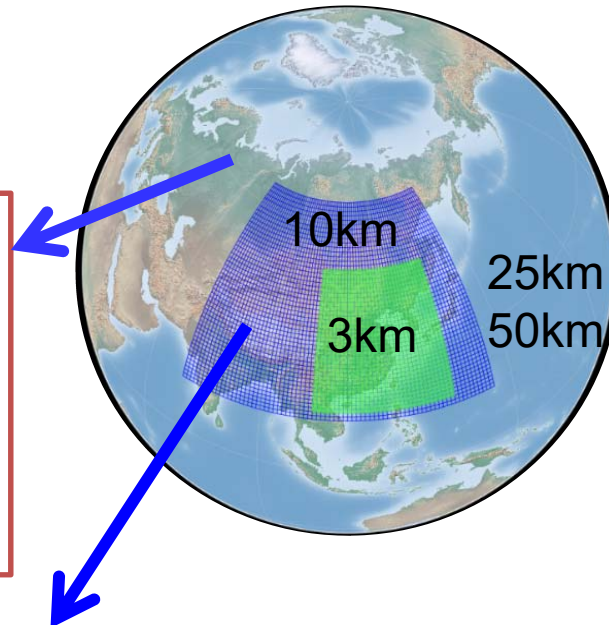
- WRF physics for meso-scale application
- Physics for global forecast
  - Radiation:
    - RRTMG LW(√4.71)/SW(√3.61)
  - Cumulus:
    - Simplified Arakawa Schubert
  - Microphysics: CMA two-moment microphysics
  - Cloud: Prognostic
  - Land surface: CoLM
  - Gravity wave drag:
    - Kim & Arakawa 1995; Lott & Miller 1997; Alpert, 2004
  - Small scale orographic form drag : Beljaars, Brown & Wood(2004)



# Operational NWP Configs at CMA

## Global: GRAPES\_GFS/GEPS

- 25/50km deterministic/ensemble(M30)
- 60 vertical levels ( $\sim 3\text{hPa}$  top)
- 10-day forecast twice daily
- 4DVAR-100km inner loop



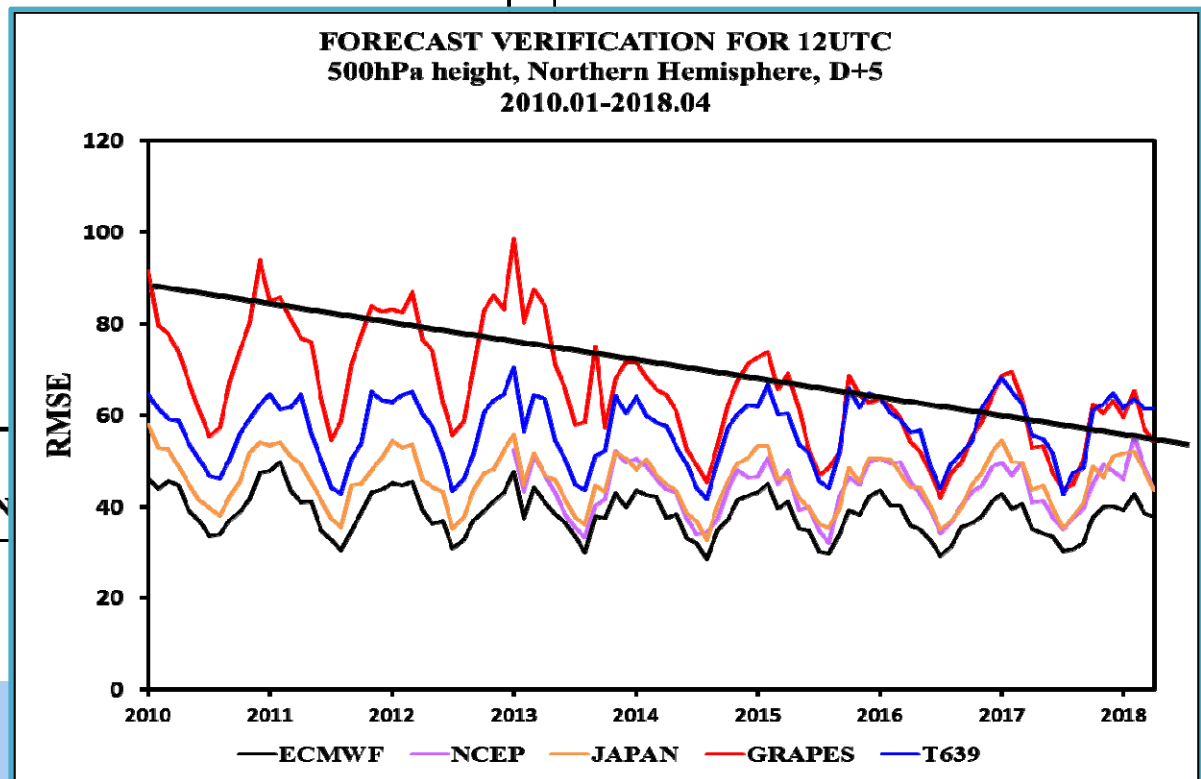
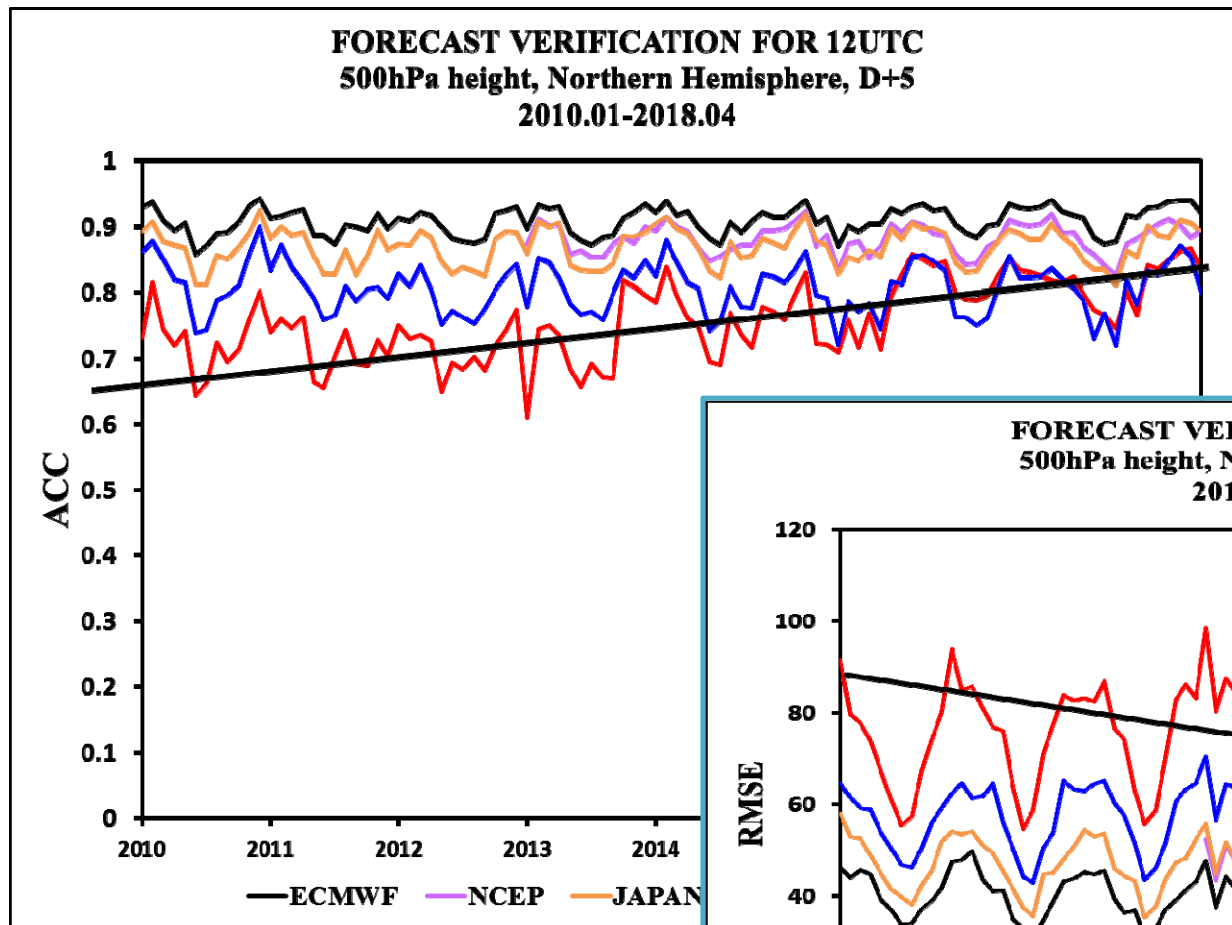
## East Asia: GRAPES\_Meso/REPS

- 3+10/10km deterministic/ensemble(M15)
- 50 vertical levels ( $\sim 50\text{hPa}$  top)
- 24-hour(eight times/day)/120-hour forecast (two times/day)
- 3DVAR



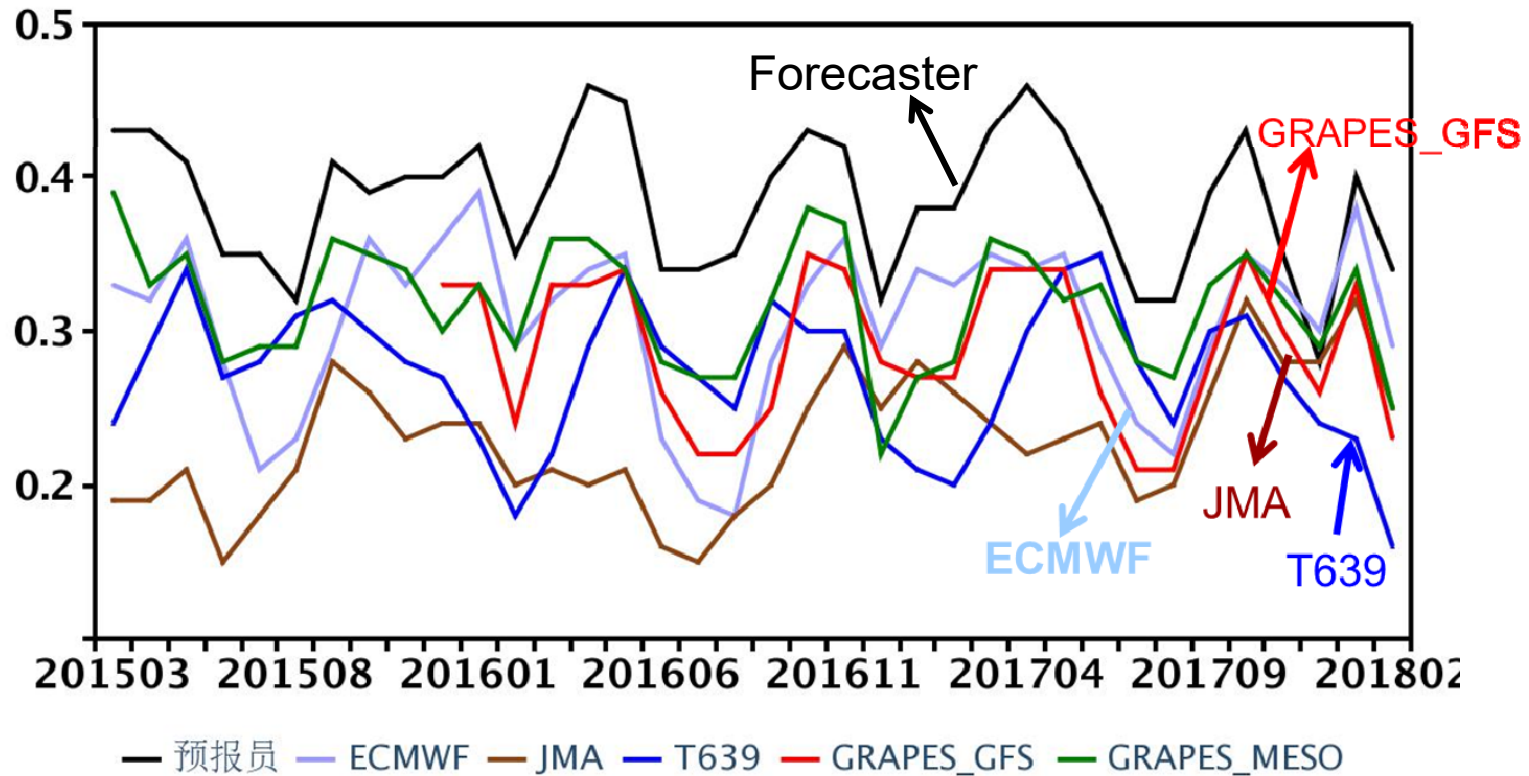
# ACC & RMSE of 500hPa height 5-day forecast

Northern Hemisphere



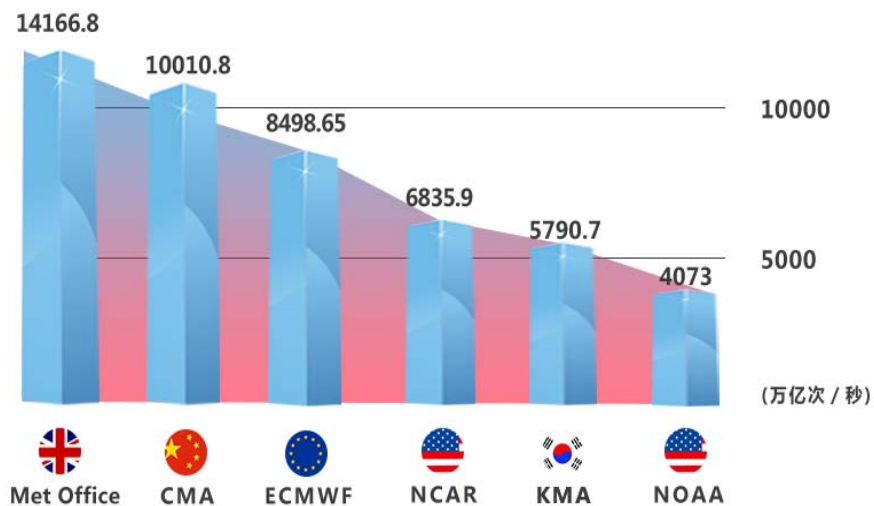
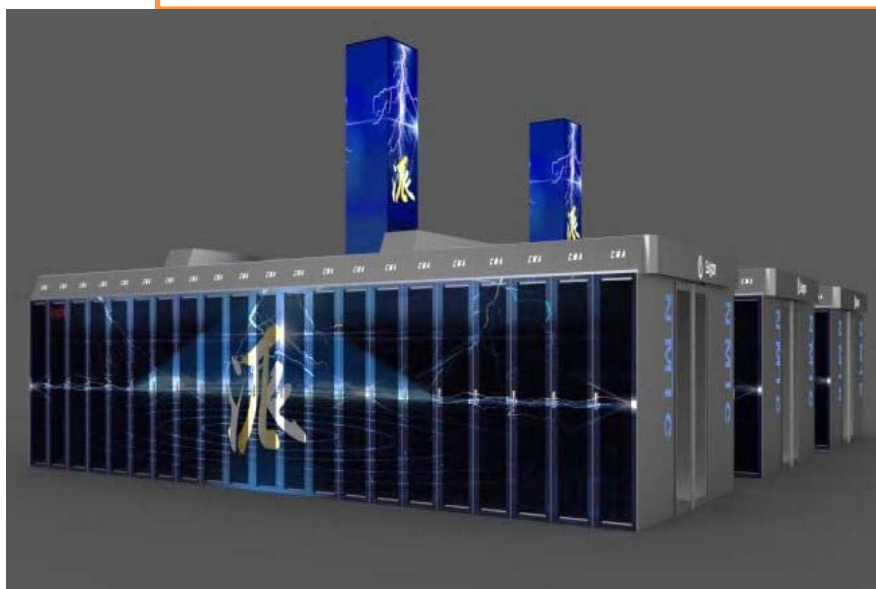
# Comparison of precipitation forecast over China among ECMWF, JMA & CMA GRAPES\_GFS, Meso

## ETS score of 48h forecast of rain belt



# New CMA HPC

- 2 computers, peak performance: 8189.5 TFLOPS
  - Parastor 300 storage: 23,088TB
  - Node/CPU: 3076 nodes, 98432 cores
    - Based on Intel Xeon Gold 6124 (2.66GHz 16 cores) processor
    - 2 CPUs/node (16 cores/CPU)
  - 100Gb/s EDR InfiniBand inter connection
  - RedHat Enterprise Linux Server V7.4
- and
- 4xIntel KNL 7250(68c 1.4GHz) X 6, 73.1TFlops
  - 2xTesla P100 GPU X 24, 289.5TFlops

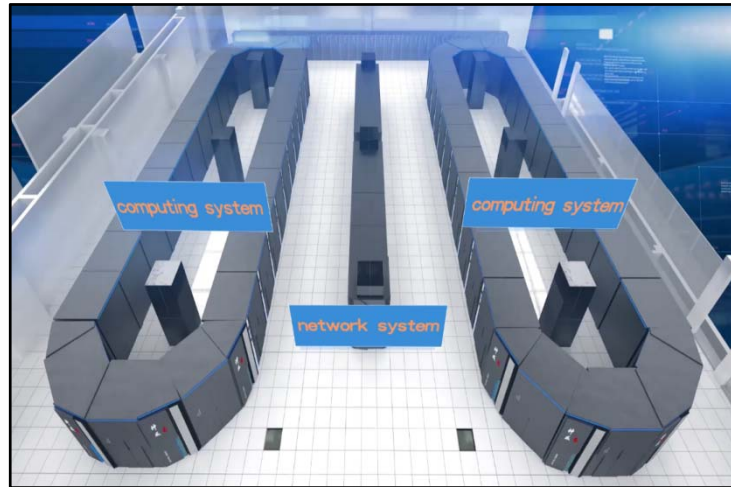


# Heterogeneous and Many-core architecture becomes mainstream

System/ Launch date	Rpeak (PFLOPS)	Rmax (PFLOPS)	Power Efficiency (MFLOPS/W)	Cores	Architecture
<b>Summit/ 201806</b>	187.66	122.30	13889.05	2,282,544	<b>GPU Acceleration</b>
<b>TaihuLight 201606</b>	125.44	93.02	6051.13	10,649,600	<b>Heterogeneous many-core</b>
<b>Sierra/ 201806</b>	119.19	71.61	/	1,572,480	<b>GPU Acceleration</b>
<b>Tianhe-2A 201806</b>	100.68	61.44	3324.56	4,981,760	<b>Matrix2000 Acceleration</b>
<b>ABCI/ 201806</b>	32.58	19.88	/	391,680	<b>GPU Acceleration</b>
<b>Piz Daint 201706</b>	25.33	19.59	8622.36	361,760	<b>GPU Acceleration</b>
<b>Titan 201211</b>	27.11	17.59	2142.77	560,640	<b>GPU Acceleration</b>
<b>Seq 201</b>	Evaluation of the potential of current operational model refactorization has to be done as early as possible				



# Target Platform: Sunway TaihuLight

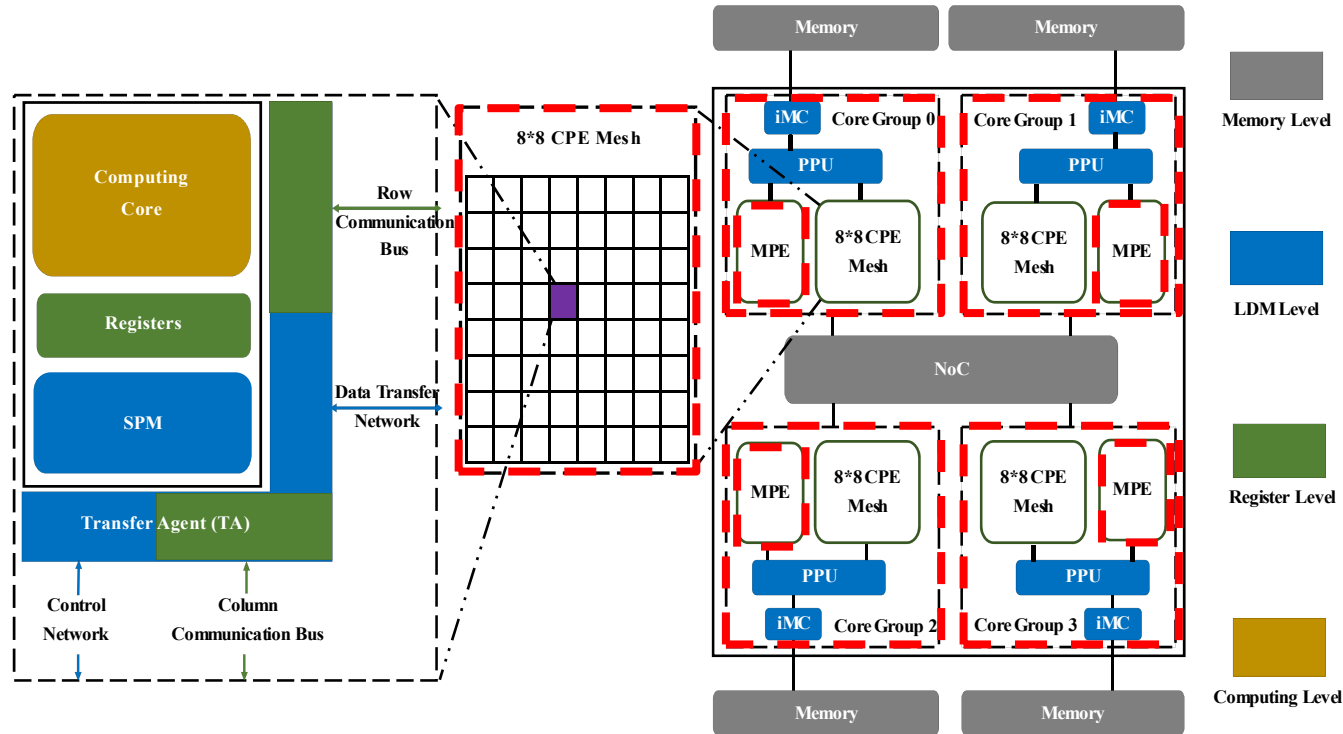


## Entire System

Peak Performance	125 PFlops
Linpack Performance	93 Pflops / 74.4%
Total Memory	1310.72 TB
Total Memory Bandwidth	5591.45 TB/s
# nodes	40,960
# cores	10,649,600



# SW26010 Processor

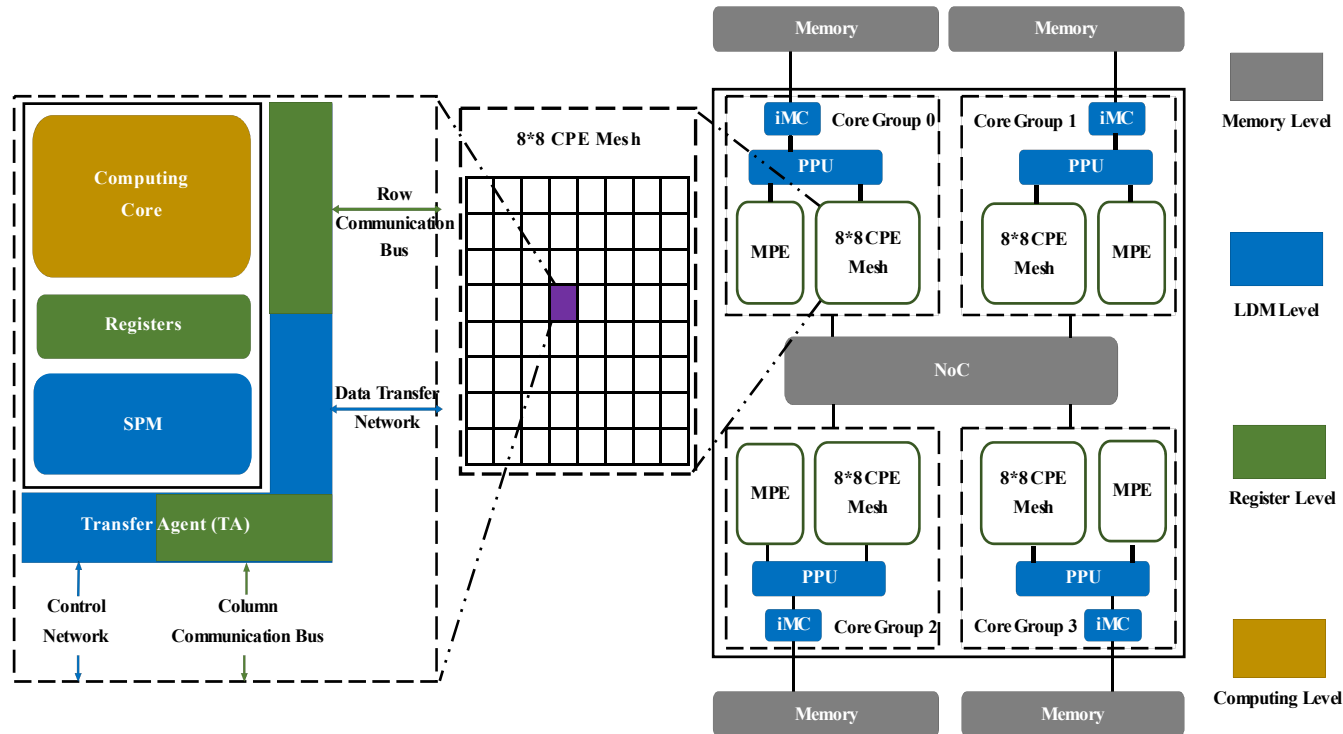


Name	MPE	CPE
Frequency (GHz)	1.45	1.45
Cores	1	64
SIMD width	4 double / 8 integer	4 double / 8 integer
Peak Gflop/s in D.P.	23.2	742.4
Data L1 / L2 cache	32KB / 256KB	64KB (SPM) / -
Instruction L1 / L2 cache	32KB / +	16KB / 64KB
Memory capacity (GB)	32 (shared)	





# SW26010 Processor

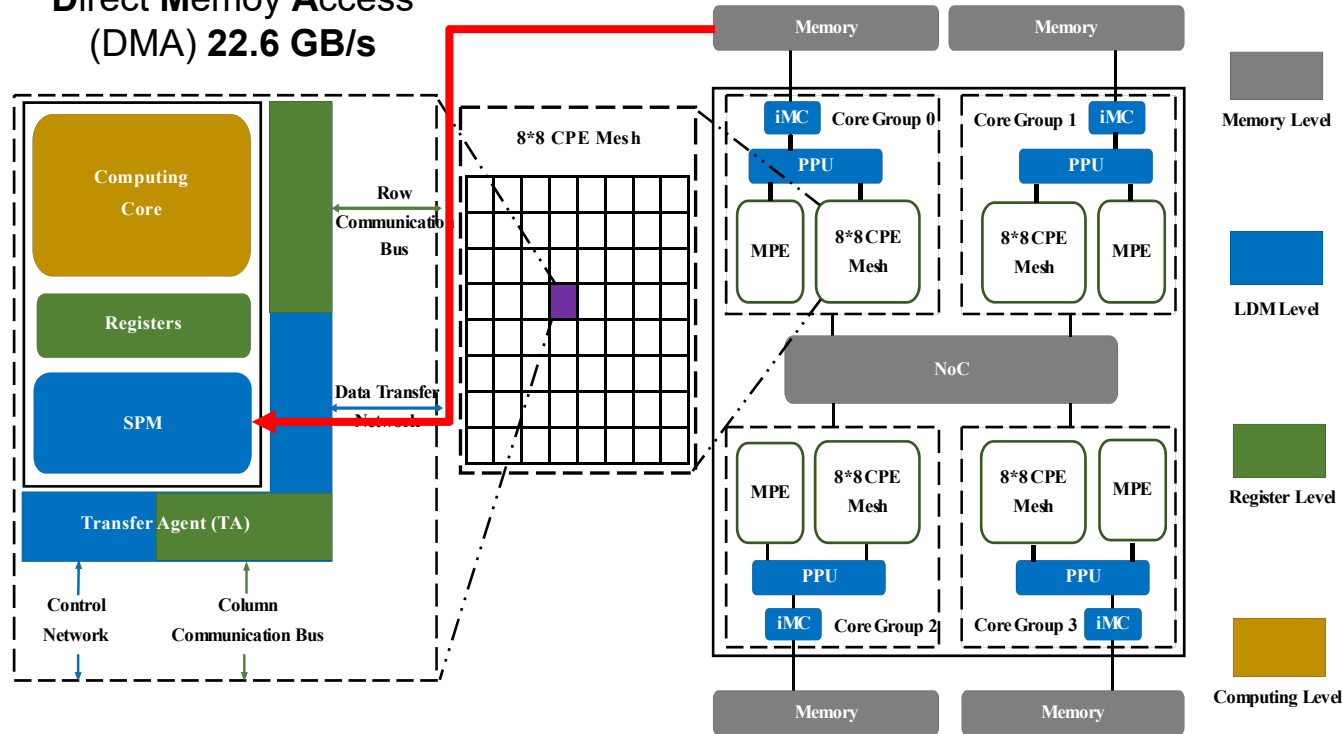


Name	MPE	CPE
Frequency (GHz)	1.45	1.45
Cores	1	64
SIMD width	4 double / 8 integer	4 double / 8 integer
Peak Gflop/s in D.P.	23.2	742.4
Data L1 / L2 cache	32KB / 256KB	64KB (SPM) / -
Instruction L1 / L2 cache	32KB / +	16KB / 64KB
Memory capacity (GB)	32 (shared)	



# SW26010 Processor

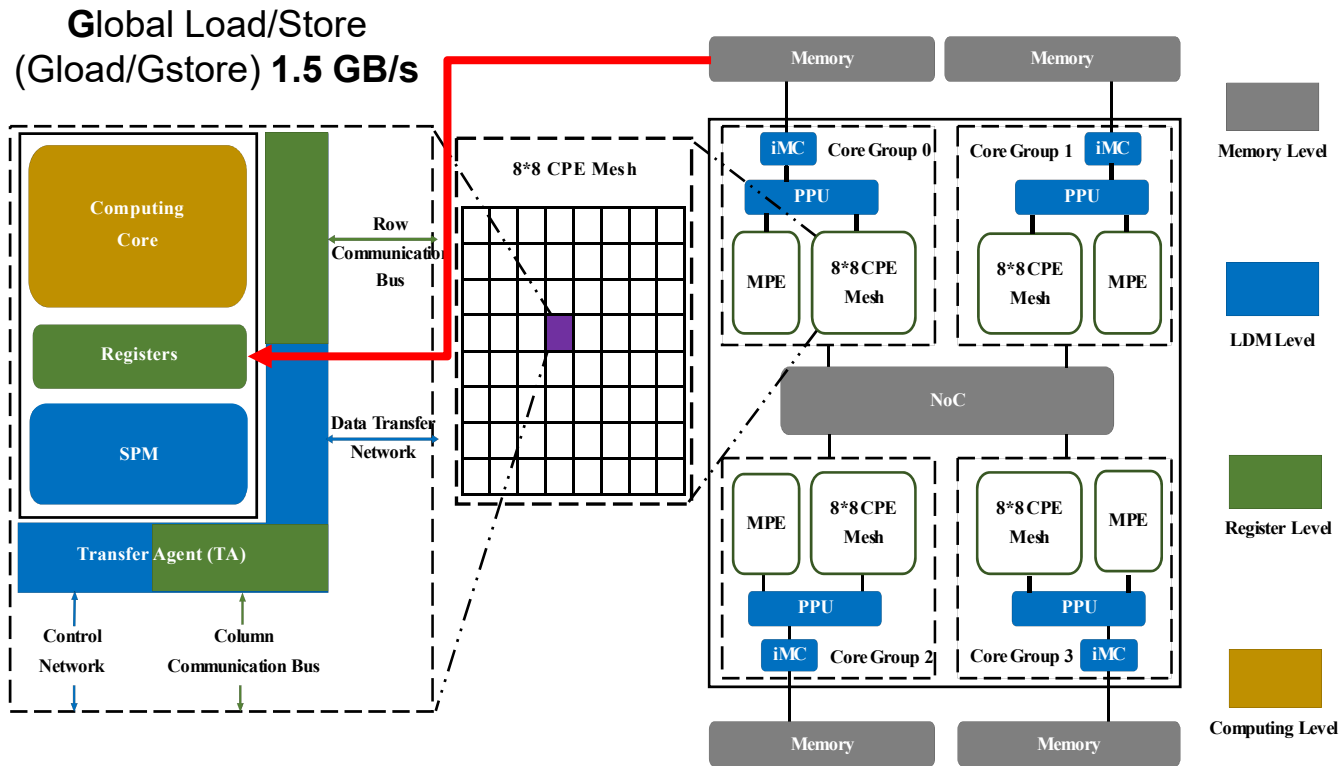
Direct Memory Access  
(DMA) 22.6 GB/s



Name	MPE	CPE
Frequency (GHz)	1.45	1.45
Cores	1	64
SIMD width	4 double / 8 integer	4 double / 8 integer
Peak Gflop/s in D.P.	23.2	742.4
Data L1 / L2 cache	32KB / 256KB	64KB (SPM) / -
Instruction L1 / L2 cache	32KB / +	16KB / 64KB
Memory capacity (GB)	32 (shared)	



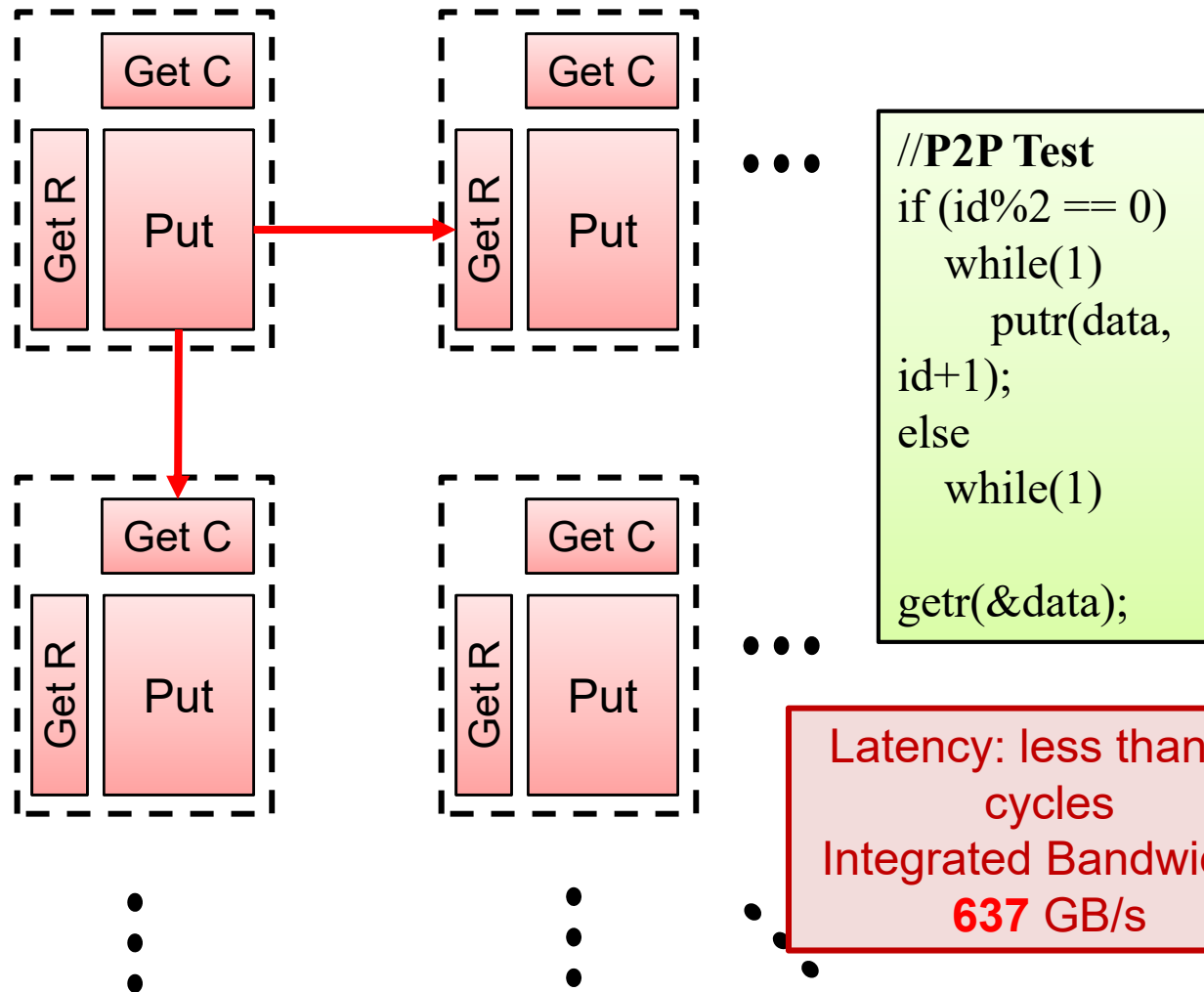
# SW26010 Processor



Name	MPE	CPE
Frequency (GHz)	1.45	1.45
Cores	1	64
SIMD width	4 double / 8 integer	4 double / 8 integer
Peak Gflop/s in D.P.	23.2	742.4
Data L1 / L2 cache	32KB / 256KB	64KB (SPM) / -
Instruction L1 / L2 cache	32KB / +	16KB / 64KB
Memory capacity (GB)	32 (shared)	



# Register Communication of SW26010



Xu, Zhigeng, James Lin, and Satoshi Matsuoka. "Benchmarking SW26010 Many-Core Processor." *Parallel and Distributed Processing Symposium Workshops (IPDPSW), 2017 IEEE International*. IEEE, 2017.

# Summary of SW26010 Processor

- Heterogeneous architecture
- Manual cache system (SPM)
- Direct memory access (DMA)
- Limited register communication

Different computing resources have been fully controlled by developers along with high development efforts



# Principal Programming Model on TaihuLight

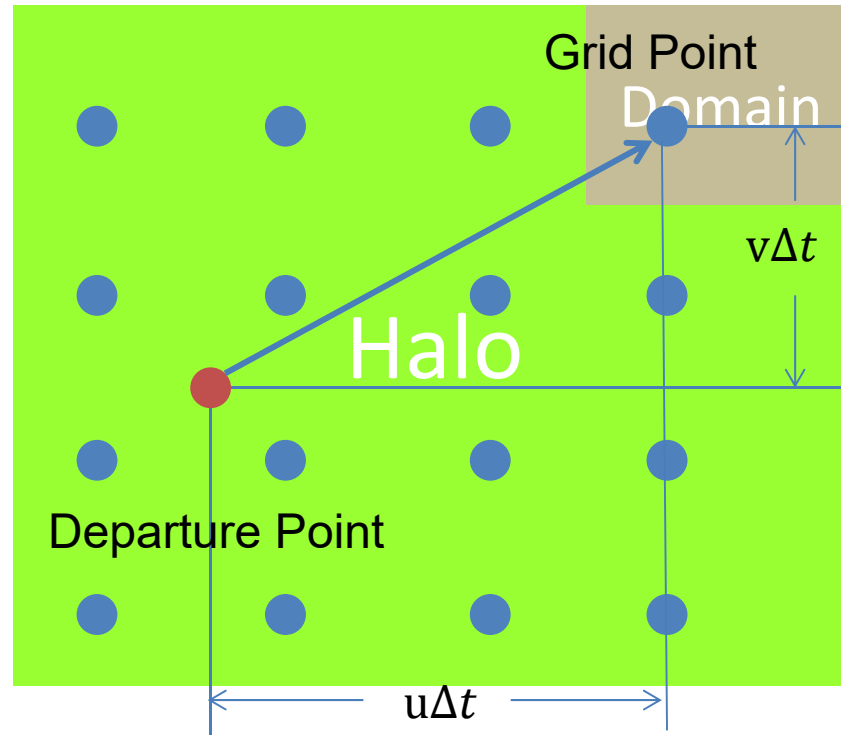
- MPI+X
  - X : OpenACC\* / Athread
  - One MPI process manages to run on one management core (MPE)
  - OpenACC\* is directive-based programming tool for SW26010
    - OpenACC2.0 based
    - Extensions for the architecture of SW26010
    - Supported by SWACC/SWAFORT compiler
    - OpenACC\* conducts data transfer between main memory and on-chip memory (SPM), and distributes the kernel workload across compute cores (CPEs)
  - Athread is the threading library to manage thread on compute core (CPE), which is used in OpenACC\* implementation



# Semi-Lagrange Interpolation



# Semi-Lagrange Interpolation

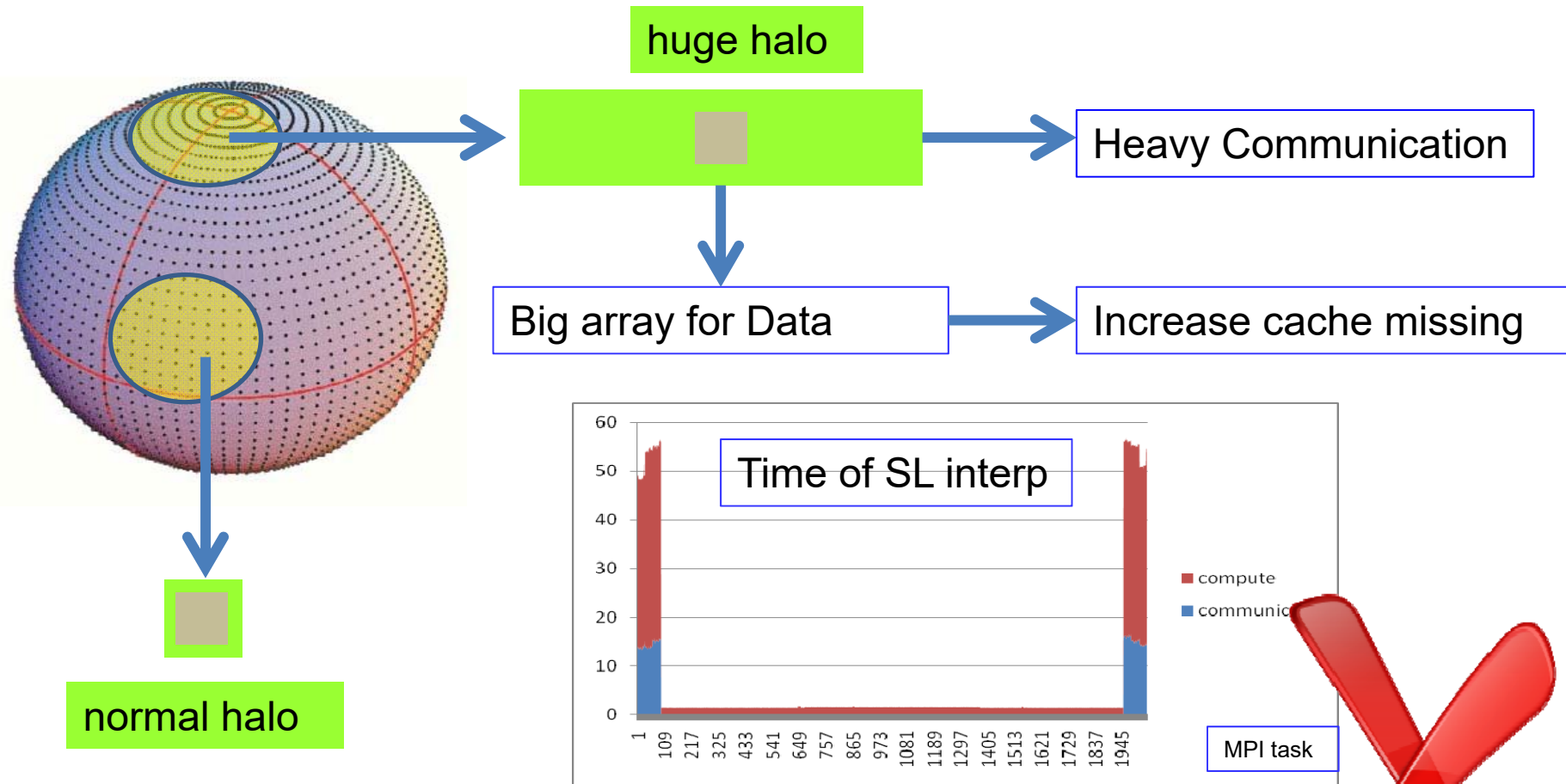


If the departure point is outside the domain,  
halo is needed for interpolation





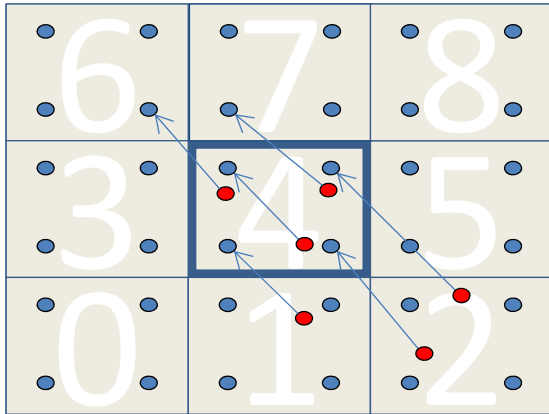
# Semi-Lagrange interpolation



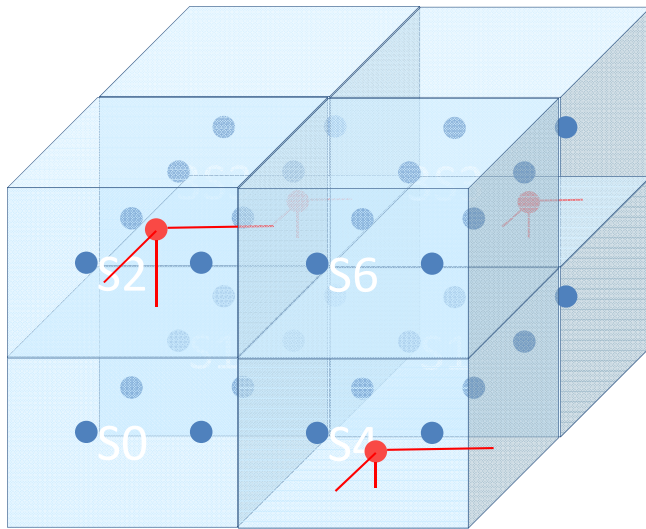
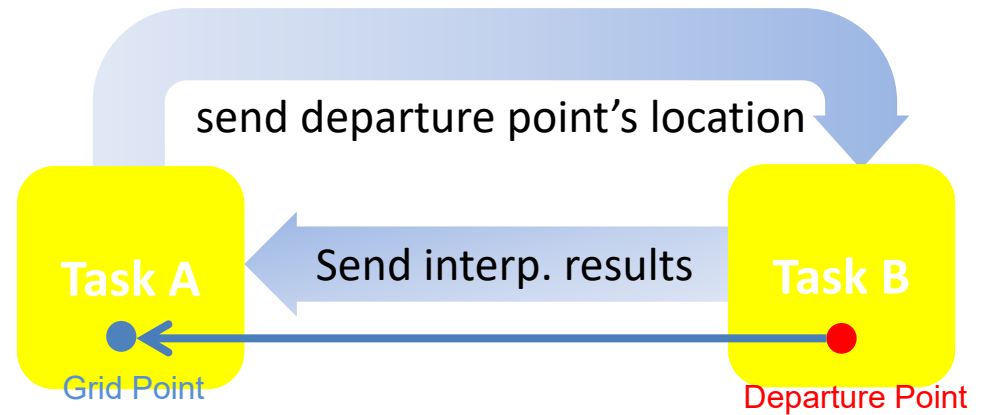
Performance at poles is very poor



# Semi Lagrange departure point interpolation



MPI task partition for grid and departure points



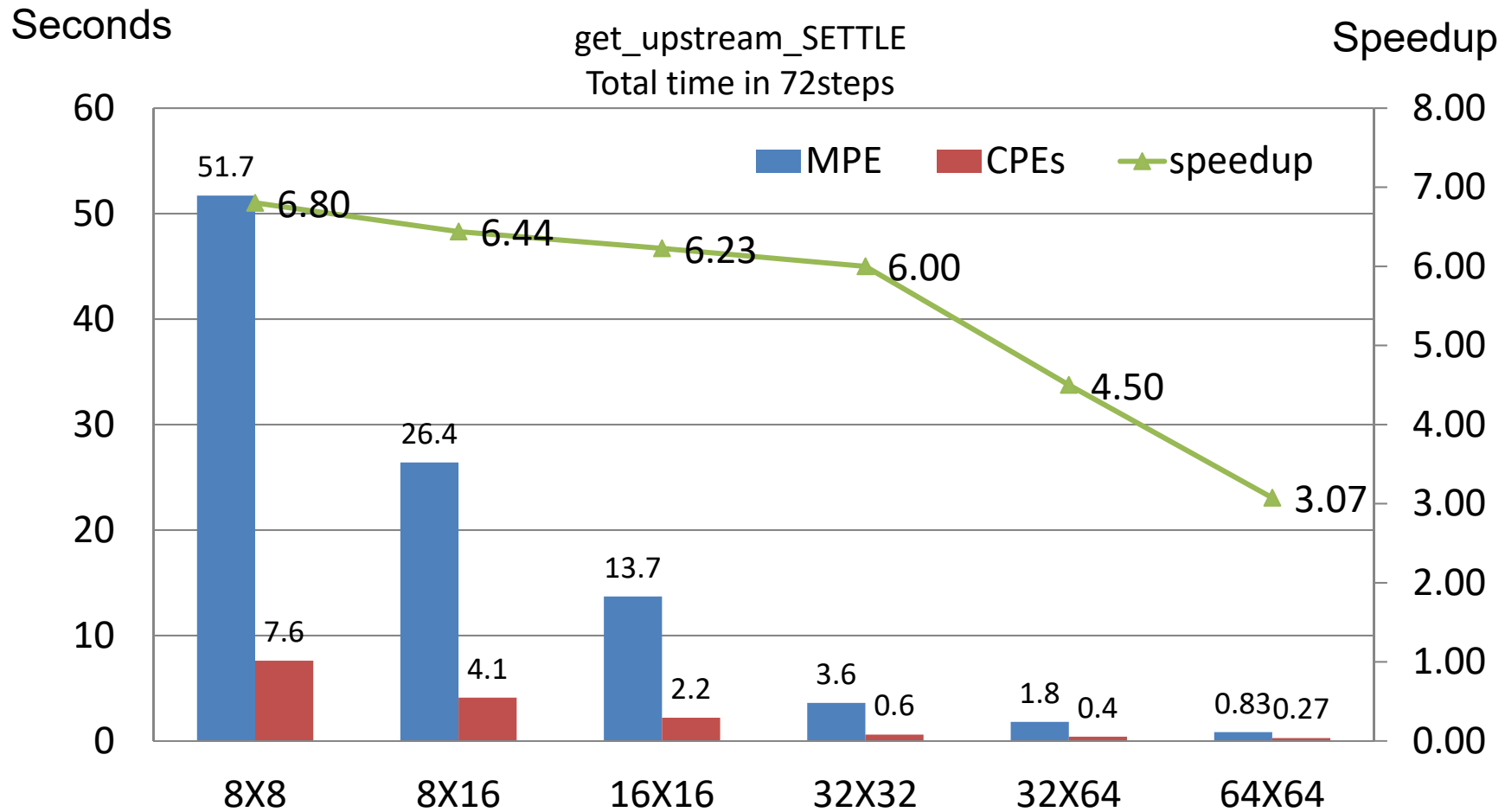
● Grid Point ● Departure Point

thread partition in a MPI task

The grid points and departure points in a MPI task is partitioned both in horizontal & vertical dimensions, small enough to fit into the 64k bytes SPM of a CPE. Subdomain have small halo needed by interpolation.



# Semi-Lagrange Interpolation



GRAPES global, 0.5°, DP, Taihu light, subroutine "get\_upstream\_SETTLE", total time in 72 steps



# Helmholtz Solver



# Characteristics of Helmholtz Eq.

- Math. Model & Matrix Characteristics

- Large Scale non-symmetrical Linear Equations for Globe

- 25km H-resolution,  $1440 \times 720 \times 36 = 37,324,800$

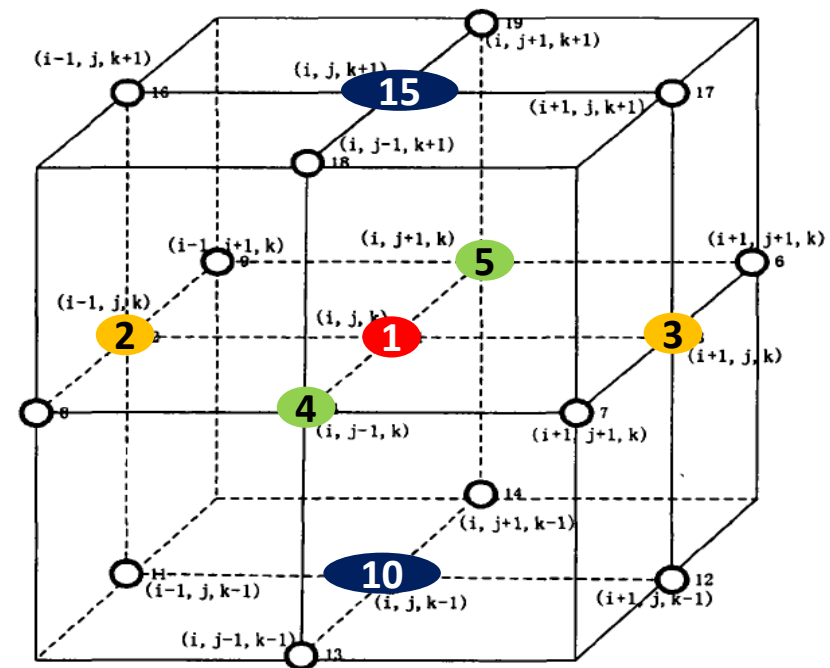
- 19-diagonal Coefficient Matrix

After scaling with diagonal elem.

- $C1 = 1.0$
- $C10/C15 \sim 10^{-1}$
- $C2/C3 \sim 10^{-2}$
- $C4/C5 \sim 10^{-3}$
- Others  $\leq 10^{-5}$

- Not good distribution of eigenvalues

- 100km H-resolution, max/min  $\sim 3 \times 10^4$



# Improved pre-GCR algorithm (IGCR)

- Preconditioning with Restricted Additive-Schwarz domain decomposition scheme
  - Only one overlapping layer is enough
  - Additional halo update is introduced
- Improved GCR algorithm for strength reduction and communication reduction

01- Compute  $R_0 = b - Ax_0$ ,  $\hat{R}_0 = M^{-1}R_0$ ,  $p_0 = \hat{R}_0$

02- Do  $i = 1, 2, \dots$ , until convergence

03-  $\alpha_{i-1} = \langle R_{i-1}, Ap_{i-1} \rangle / \langle Ap_{i-1}, Ap_{i-1} \rangle$

04-  $x_i = x_{i-1} + \alpha_{i-1}p_{i-1}$

05-  $R_i = R_{i-1} - \alpha_{i-1}Ap_{i-1}$

06-  $\hat{R}_i = M^{-1}R_i$

07- Do  $j = \text{int}[(i-1)/k]k, \dots, i-1$

08-  $\beta_{ij} = -\langle A\hat{R}_i, Ap_j \rangle / \langle Ap_j, Ap_j \rangle$

09- EndDo

10-  $p_i = \hat{R}_i + \sum_{j=\text{int}[(i-1)/k]k}^{i-1} \beta_{ij}p_j$

11-  $Ap_i = A\hat{R}_i + \sum_{j=\text{int}[(i-1)/k]k}^{i-1} \beta_{ij}Ap_j$

12- EndDo

$$(r_{j+1}, Ap_i) = (r_j, Ap_i) - \alpha_j(Ap_j, Ap_i) =$$

$$(r_{j-1}, Ap_i) - \alpha_{j-1}(Ap_{j-1}, Ap_i) =$$

$$(r_i, Ap_i) - \alpha_i(Ap_i, Ap_i) = 0 \quad i \leq j$$

Compare to Baseline Algorithm  
(k=10)

✓ (-) 1 Allreduce

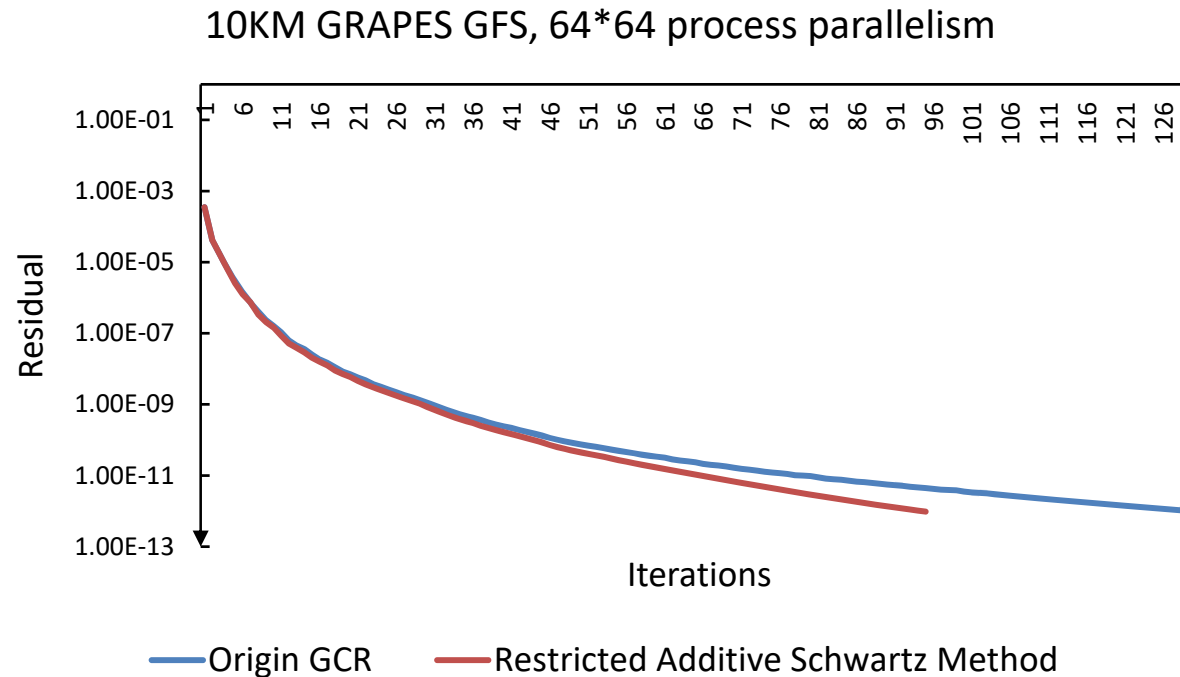
✓ (-) 1 SpMV

✓ (+) k BLAS1



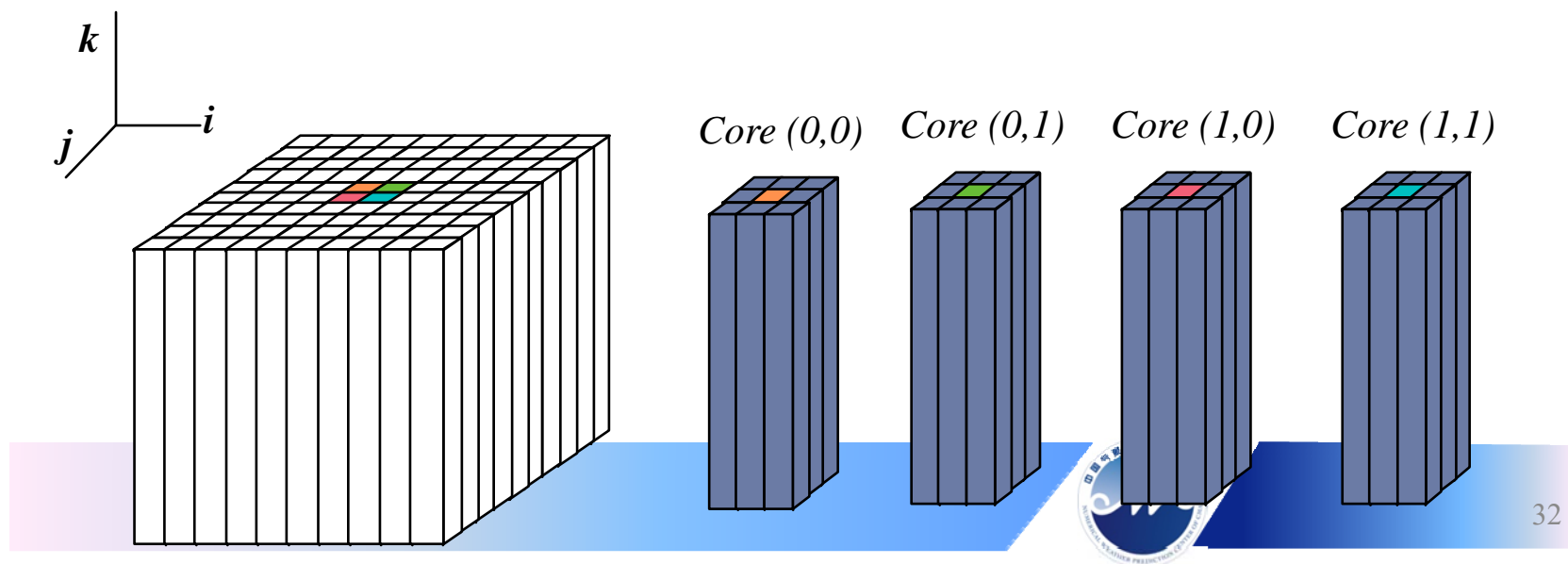
# Convergence improvement

- Convergence of RAS method and origin method



# Sparse Matrix-vector Multiplication for SW26010 processor

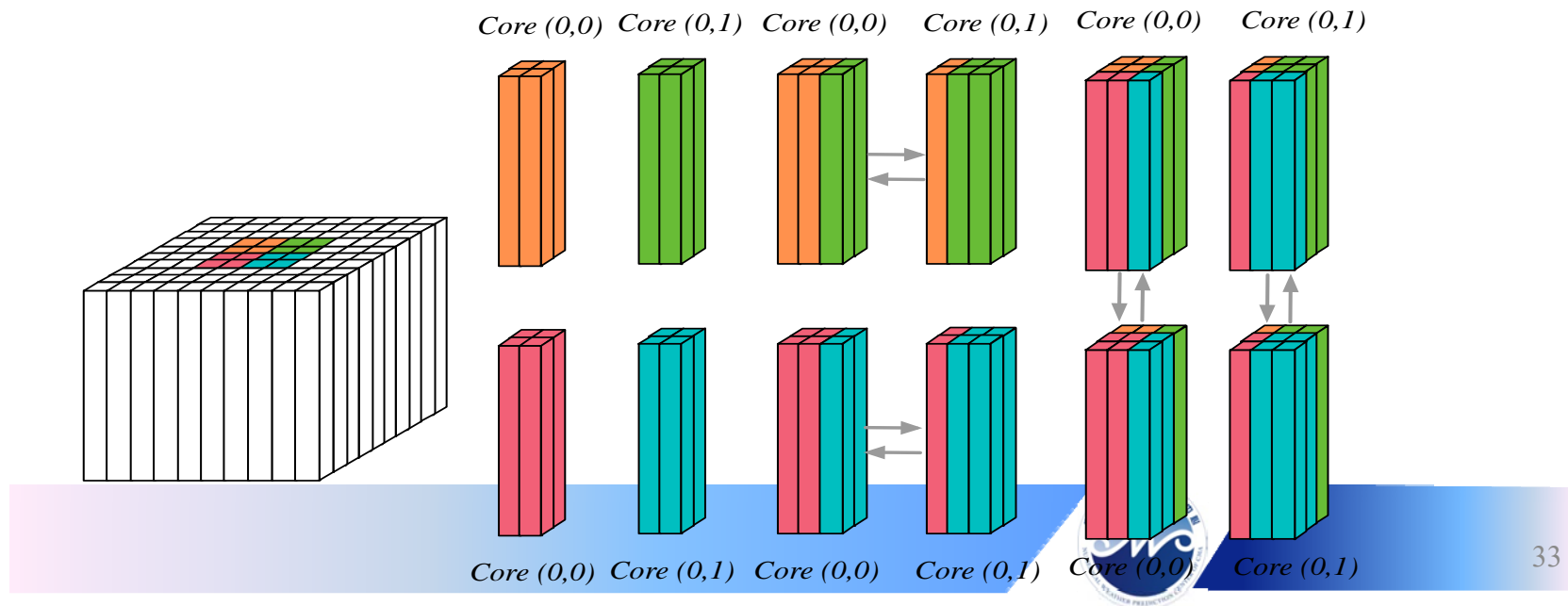
- Chunk access of a  $k$  column
  - Matrix: continuous storage with  $(19, k, i, j)$  order and bulk access of 19 nonzeros of each row of matrix
  - Vector: stride access following the geometry
- For vector side, each core should read 9 column data to compute one column





# Sparse Matrix-vector Multiplication for SW26010 processor

- For computation of each column, 9 column data has to be access for each core
- Collaborated data access by multiple cores reduces the data access volume in total but introduces more synchronizations
- Good trade-off: 2\*2 Cores share data with each other by register communication, each core only need to read 4 column data



# Sparse Matrix-vector Multiplication for SW26010 processor

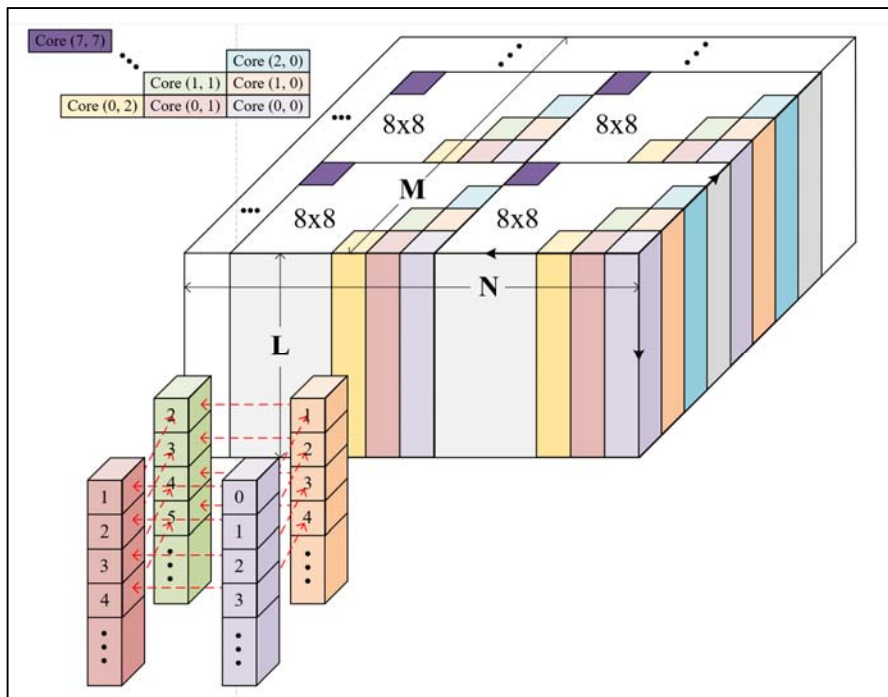
- More Cores share data with each other, less x need to read.
- Cores which share data with each other need synchronization, more Cores, more expensive for synchronization.

Groups	Average Columns of data to read per Core	Cores to synchronization
1*1	9	0
2*2	4	4
4*4	9/4	16
8*8	25/16	64



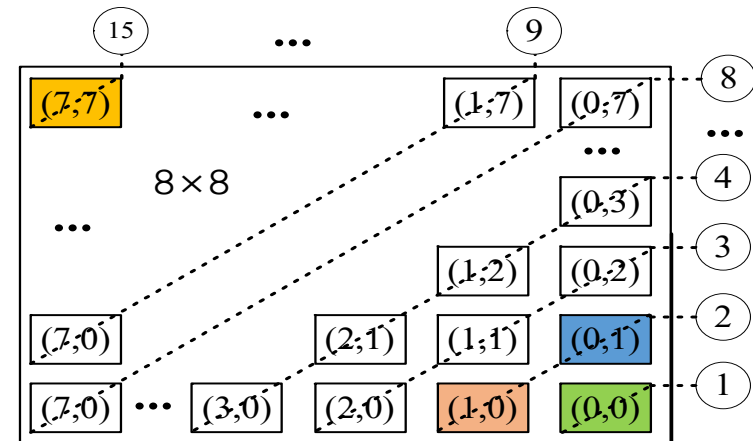
# Fine-grained parallel Incomplete LU factorization for SW26010 processor

- Example of 7 points ILU



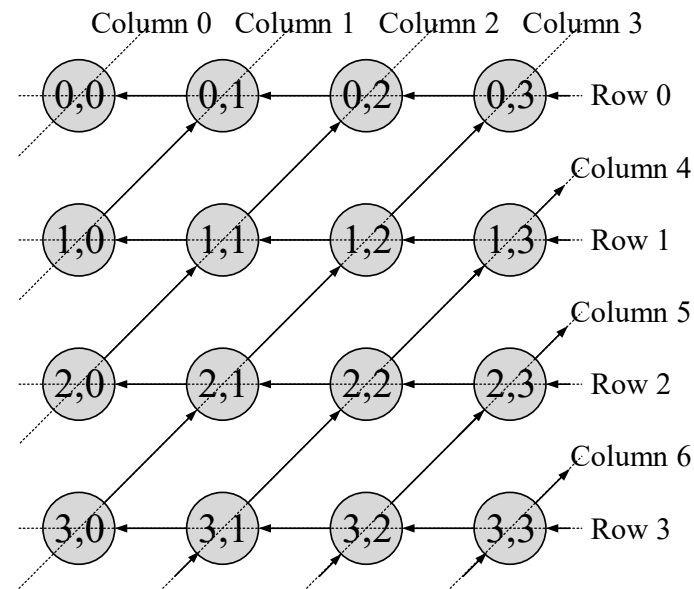
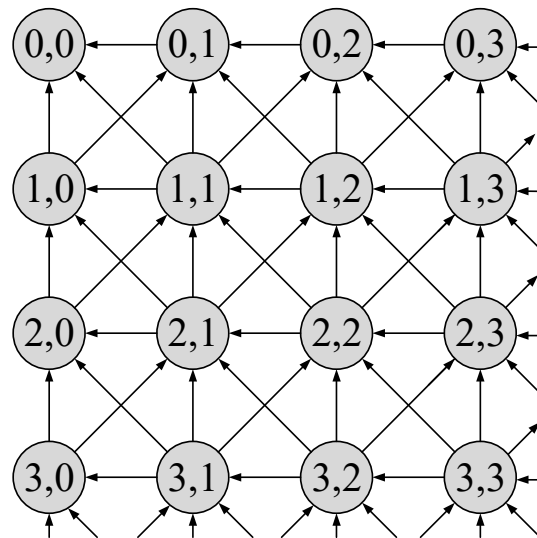
① ~ ⑮ From step 1 to step 15

- Computed by Core (0,0)
- Computed by Core (0,1)
- Computed by Core (1,0)
- Computed by Core (7,7)



# Fine-grained parallel Incomplete LU factorization for SW26010 processor

- How to handle diagonal communication on SW26010?
  - Reduction in communication path graph
  - Diagonal 2D partition



# Communication-Avoiding GCR for Sunway TaihuLight

Algorithm : communication avoiding preconditioned generalized conjugate residual method

$$\mathbf{x}_0, \mathbf{r}_0 = \mathbf{b} - \mathbf{A}\mathbf{x}_0, \mathbf{z}_0 = \mathbf{M}^{-1}\mathbf{r}_0, \mathbf{p}_0 = \mathbf{z}_0$$

while not converged do

- 1, Calculate  $\mathbf{V} = [\mathbf{p}_0, \mathbf{M}^{-1}\mathbf{A}\mathbf{p}_0, (\mathbf{M}^{-1}\mathbf{A})^2\mathbf{p}_0, \dots, (\mathbf{M}^{-1}\mathbf{A})^s\mathbf{p}_0, \mathbf{z}_0, \mathbf{M}^{-1}\mathbf{A}\mathbf{z}_0, \dots, (\mathbf{M}^{-1}\mathbf{A})^{s-1}\mathbf{z}_0]$

- 2, Let  $\mathbf{G} = \mathbf{V}^T\mathbf{A}^T\mathbf{A}\mathbf{V}, \mathbf{G}_m = \mathbf{V}^T\mathbf{M}^T\mathbf{A}\mathbf{V}$

- 3,  $\mathbf{m}_0 = [0_{s+1}, 1, 0_{s-1}]^T, \mathbf{n}_0 = [1, 0_{2s}]^T, \mathbf{l}_0 = [0_{2s+1}]^T$

- 4, for  $k = 0 : s-1$ , do

- 5,  $\alpha_k = \mathbf{m}_k^T \mathbf{G}_m \mathbf{n}_k / \mathbf{n}_k^T \mathbf{G} \mathbf{n}_k$

- 6,  $\mathbf{l}_{k+1} = \mathbf{l}_k + \alpha_k \mathbf{n}_k$

- 7,  $\mathbf{m}_{k+1} = \mathbf{m}_k - \alpha_k \mathbf{T} \mathbf{n}_k$

- 8,  $\beta_j^k = -\mathbf{m}_{k+1}^T \mathbf{G} \mathbf{n}_j / \mathbf{n}_j^T \mathbf{G} \mathbf{n}_j$

- 9,  $\mathbf{n}_{k+1} = \mathbf{m}_{k+1} + \sum_{j=0}^k \beta_j^k \mathbf{n}_j$

- 10, end for

- 11,  $\mathbf{z}_s = \mathbf{V} \mathbf{m}_s, \mathbf{p}_s = \mathbf{V} \mathbf{n}_s, \mathbf{x}_s = \mathbf{x}_0 + \mathbf{V} \mathbf{l}_s$

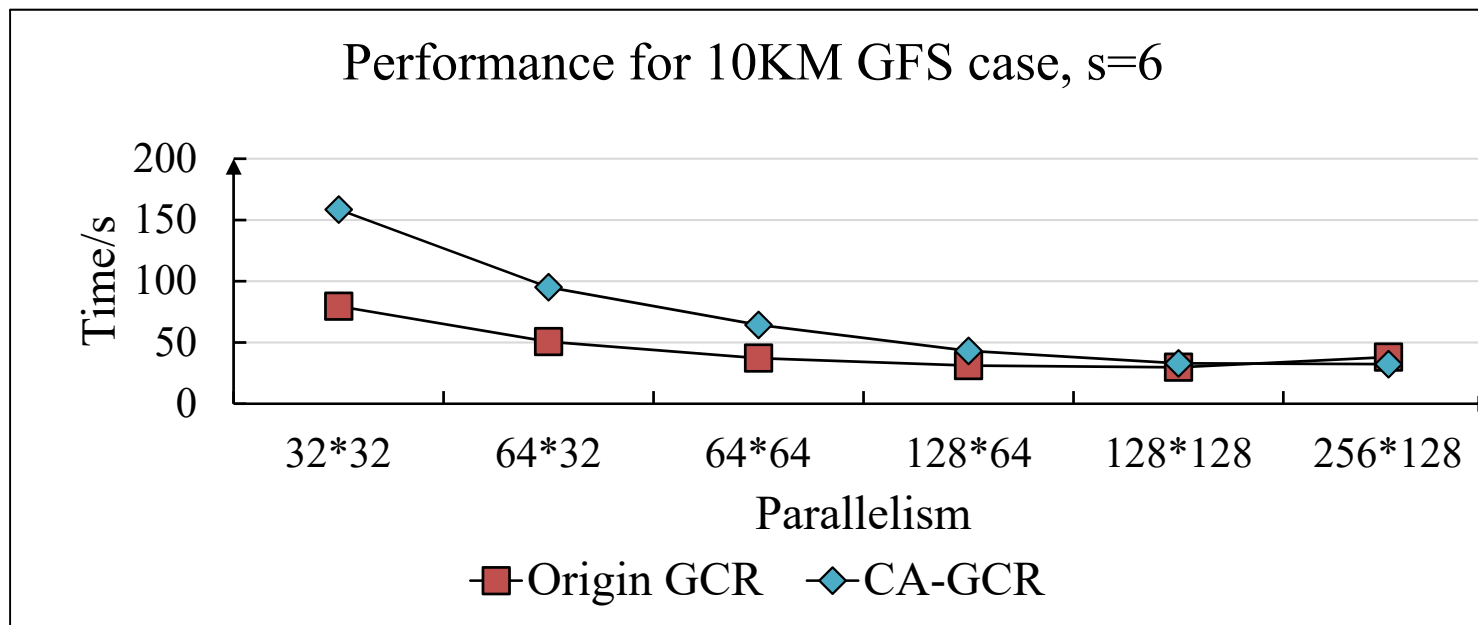
- 12,  $\mathbf{z}_0 = \mathbf{z}_s, \mathbf{p}_0 = \mathbf{p}_s$

end while



# Communication-Avoiding GCR for Sunway TaihuLight

- Some computing kernels still have some room for optimization
- If the scale continues to expand, CA-GCR may beat origin GCR.



# Halo Communication Optimization



# Halo communication optimization

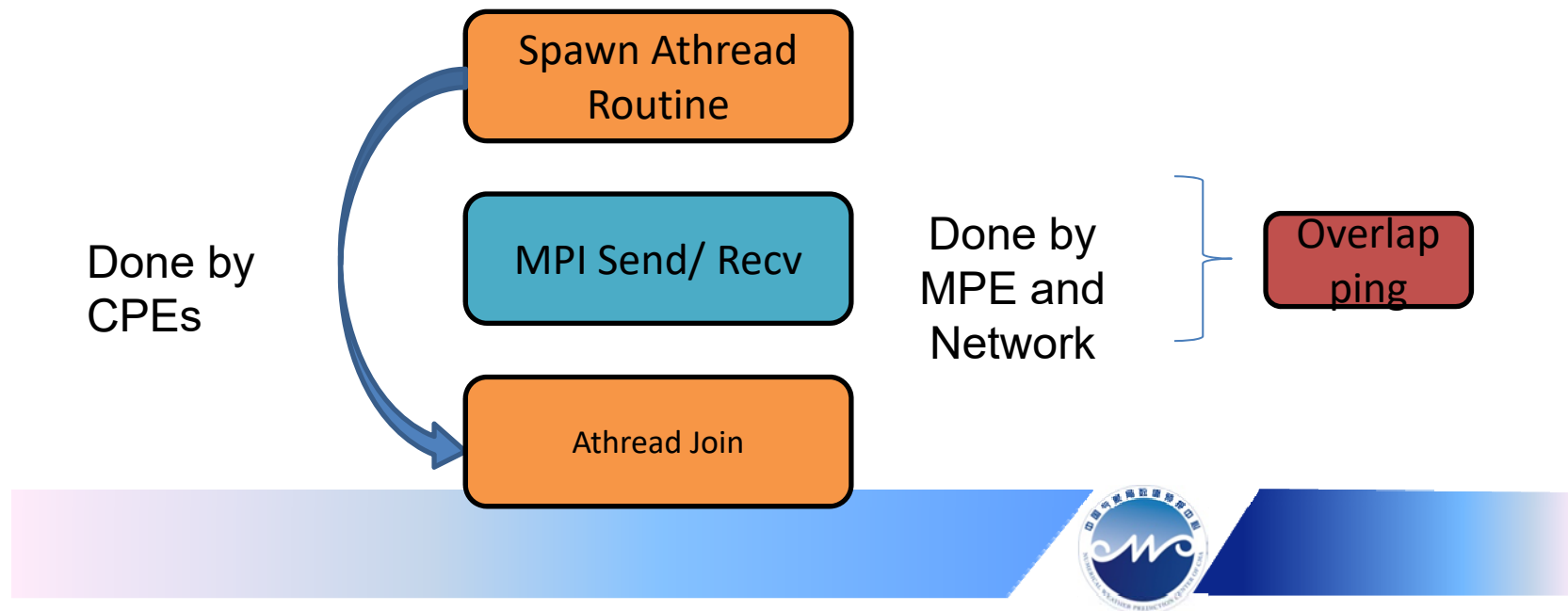
- Array transpose
  - Domain partition on i and j dimensions
  - (i , k, j) order -> (k, i, j) order for DMA friendly data access
  - The conversion can be conducted at the beginning and at the end of loop code to minimize memory access overhead.
- Assign task to partial CPEs by column
  - 64-core simultaneous access is over-provisioning for memory subsystem of SW26010 if good access pattern
  - Fewer core access makes larger chunk data access, thus better bandwidth utilization

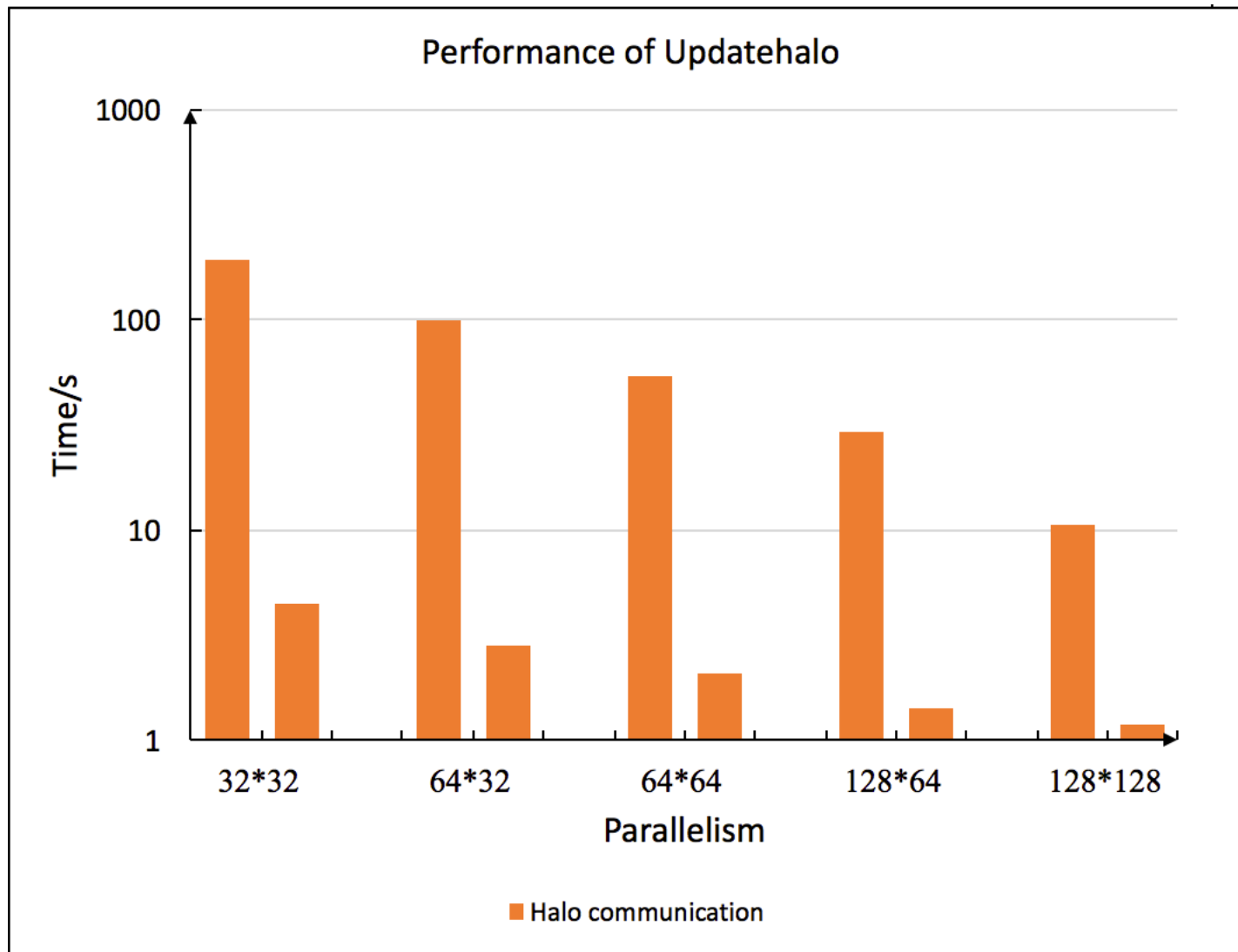




# Halo communication optimization

- Quite a few halo communications may hurt performance
- Solution:
  - Offloading the data package to CPE cluster
  - Performing data package and communication overlapping



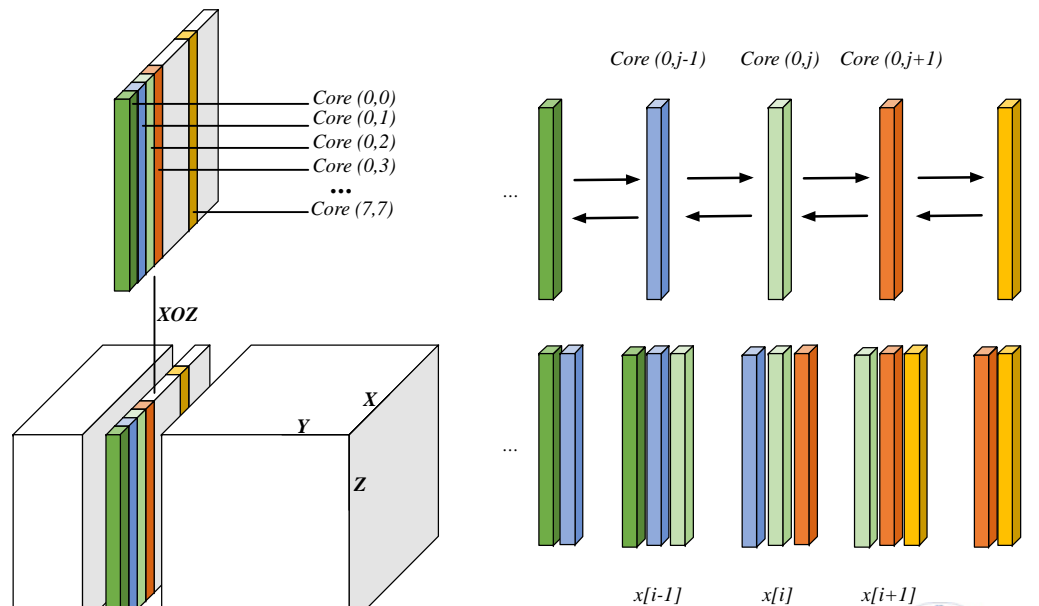


# Stencil-like kernel Optimization

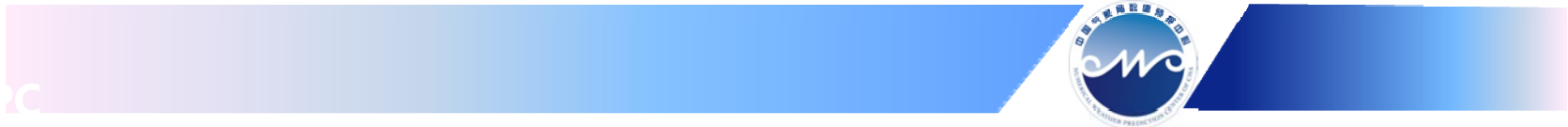


# Stencil kernel optimization (x-axis stencil case)

- Neighbour CPEs can share the dependent data by register communication rather than reading from memory



# Exploring thread-level parallelism of GRAPES with OpenACC\*



# OpenACC\* optimization on TaihuLight

```

!$acc parallel loop &
!$acc copyout(c)
!$acc collapse(2) tile(k:16)
DO j = 1, 45
  DO k= 1,60
    DO i =1, 90
      c(i,k,j) = 0.
    END DO
  END DO
END DO
!$acc end parallel loop

```

X7

```

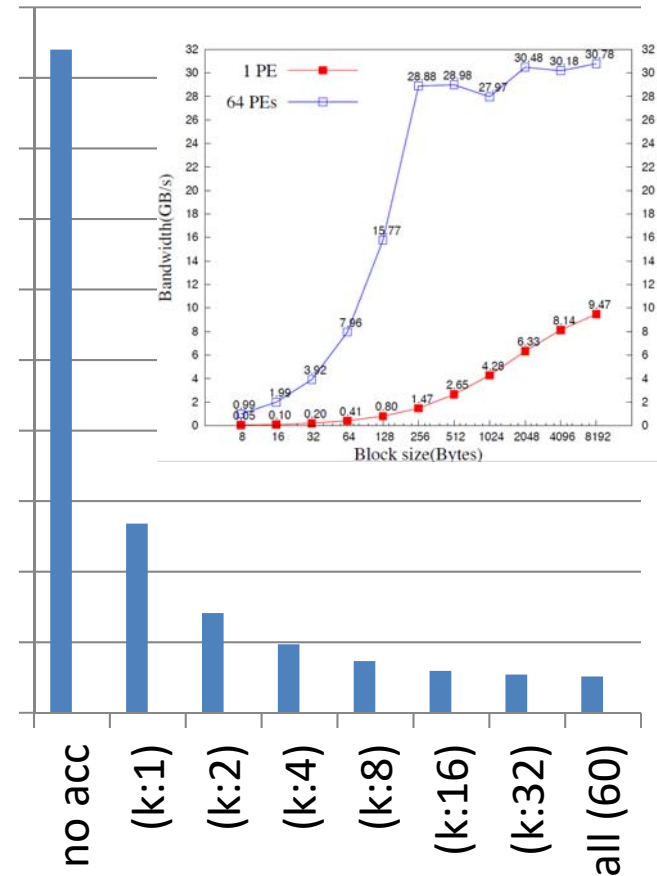
....
!$acc parallel loop copyin(A) &
!$acc copyout(b)
!$acc collapse(2) tile(k:16)
DO j = 1, 45
  DO k= 1,60
    DO i =1, 90
      A(i,k,j) = B(i,k,j)
    END DO
  END DO
END DO
!$acc end parallel loop

```

X13

Tuning tiling configuration

second  
5.00E-02  
4.50E-02  
4.00E-02  
3.50E-02  
3.00E-02  
2.50E-02  
2.00E-02  
1.50E-02  
1.00E-02  
5.00E-03  
0.00E+00

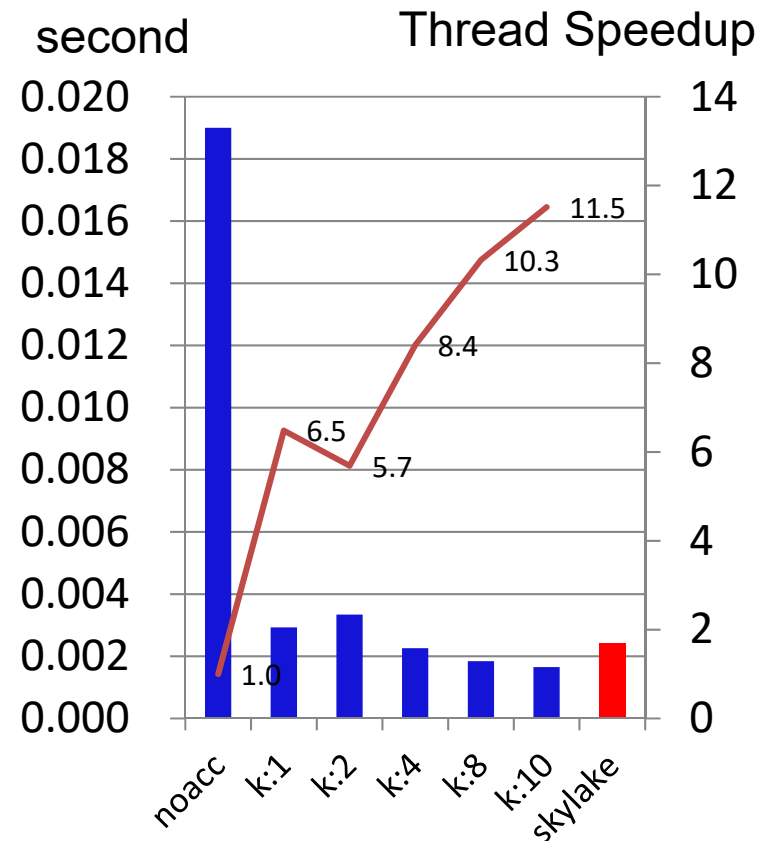


# OpenACC\* optimization on Taihu Light

```

Start_time = mpi_wtime()
!$acc parallel loop packin(its,ite,jts,jte,kts,kte,dc025,cp25)&
!$acc copyin(u,v,pi,d2k,th,thref,thv,fv,zsy,rkrf,dy)&
!$acc copyout(vl,vn) local(i,k,j,t,zdpdz,uv) &
!$acc collapse(2) tile(k:10) annotate(entire(dy,d2k,rkrf,fv))
DO j=jts,jte
  DO k=kts,kte
    DO i=its,ite
      uv=dc025*(u(i,k,j-1)+u(i-1,k,j-1)+u(i,k,j)+u(i-1,k,j))
      t=(pi(i,k,j)-pi(i,k,j-1))/dy(j)
      zdpdz=t+zsy(i,k,j)*(((pi(i,k+1,j)-pi(i,k,j))/d2k(k)+(pi(i,k,j)-pi(i,k-1,j))/d2k(k-1))+ &
        ((pi(i,k+1,j-1)-pi(i,k,j-1))/d2k(k)+(pi(i,k,j-1)-pi(i,k-1,j-1))/d2k(k-1)))*dc025)
      vl(i,k,j)=-fv(j)*uv-cp25*(thref(i,k,j-1)+thref(i,k+1,j-1)+thref(i,k,j)+thref(i,k+1,j))*zdpdz
      vn(i,k,j)=-((th(i,k,j-1)+th(i,k+1,j-1)+th(i,k,j)+th(i,k+1,j))*cp25*zdpdz &
        -cp25*(thv(i,k,j-1)+thv(i,k+1,j-1)+thv(i,k,j)+thv(i,k+1,j))*zdpdz-krf(k)*v(i,k,j)*3.
    ENDDO
  ENDDO
ENDDO
!$acc end parallel loop
end_time = mpi_wtime()
if ( myprcid == 0 ) write(*,*) 'vl,vn use ', end_time - start_time

```



# OpenACC\* optimization on TaihuLight

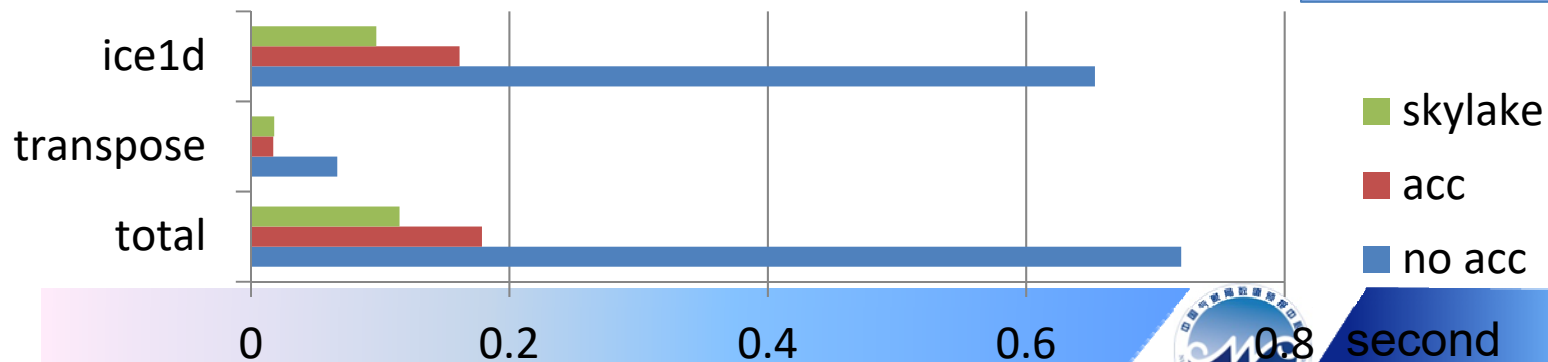
$A(i,k,j)$  **Transpose**  $\rightarrow$   $A\_c(k,i,j)$

```
!$acc parallel loop &
!$acc copy(A_c,...14 3D arrays) &
!$acc copyin(20 3D arrays) &
!$acc collapse(2) tile(i:1)
  DO j=1, 45
    DO i=1, 90
      CALL ice1D(A_c(:,i,j)... 34 1D arrays and some others)
    ENDDO
  ENDDO
!$acc end parallel loop
```

$A\_c(k,i,j)$  **Transpose**  $\rightarrow$   $A(i,k,j)$

```
subroutine ice1d(A...)
real :: A(1:60),...
...
do k=kts, kte
...
end do
do k=kts, kte
...
end do
...
end
```

1500 lines





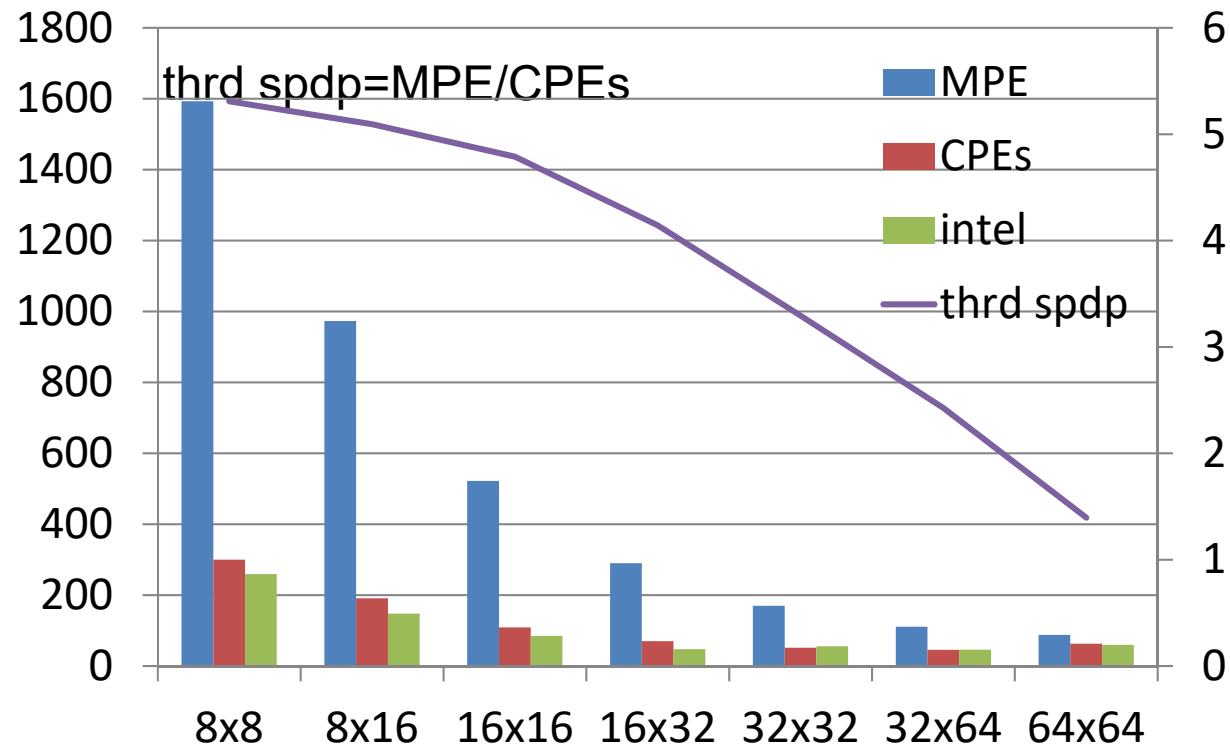
# Test Results

- Dynamic core of GRAPES global model in Double precision
  - 0.5° resolution
  - 0.25° resolution
- Physics of of GRAPES in single precision
  - 0.5° resolution



# GRAPES-GLB Dynamic core

0.5, Double Precision



GRAPES global model, 0.5° , no physics, 72 steps

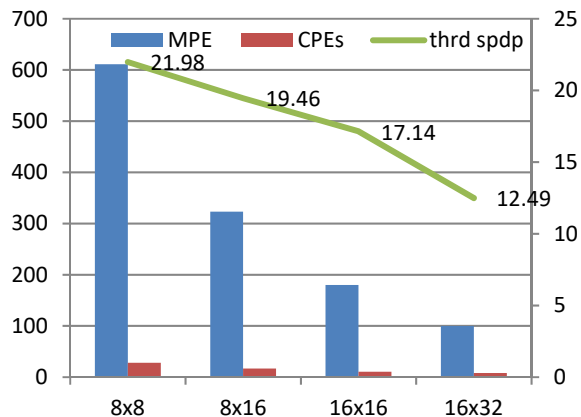
MPE: “Taihu Light”, no thread parallel processing

CPEs: “Taihu Light” with CPEs parallel processing

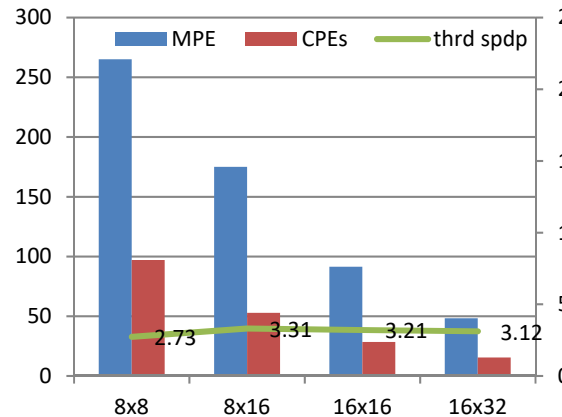
Intel: Node/CPU: Intel Xeon Gold 6124 (2.66GHz 16 cores) processor, 2 CPUs/node (16 cores/CPU) 100Gb/s EDR InfiniBand inter connection



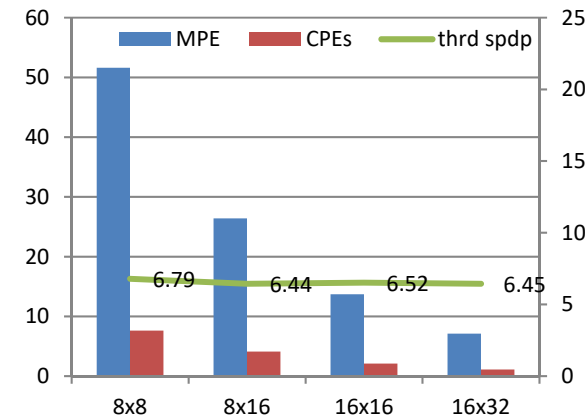
# 6 routines of dynamic core



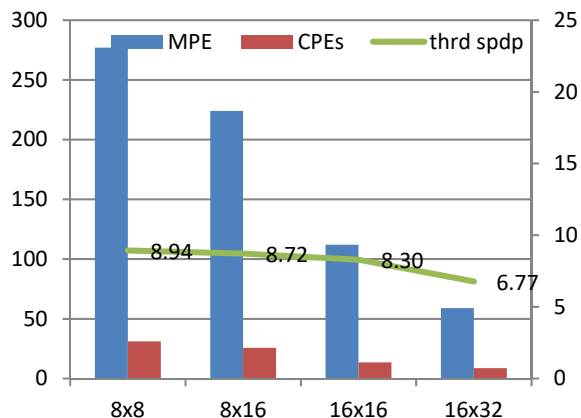
Helmholtz solver



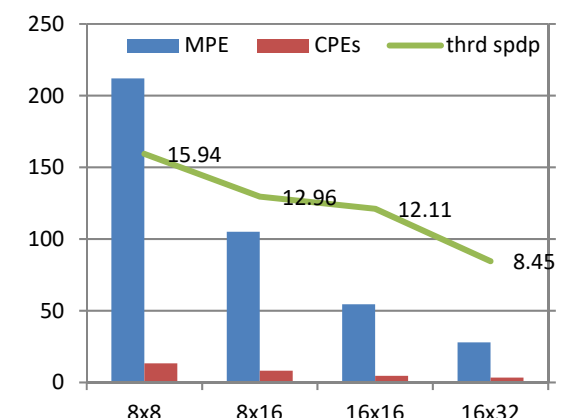
Departure points



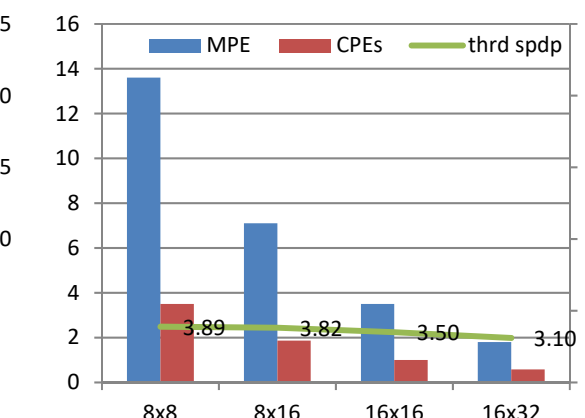
SL interp uvwthpi



SL interp trace substance



PRM



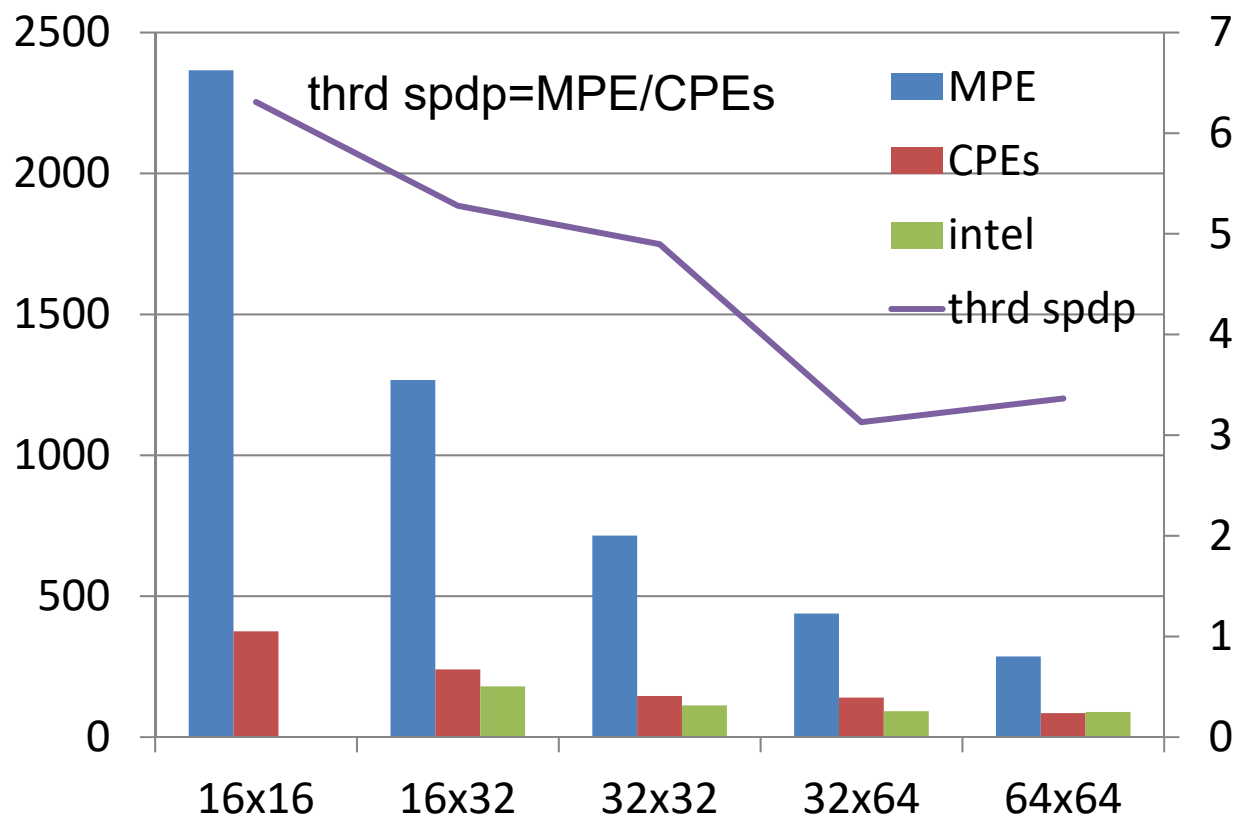
Linear & none linear

GRAPES global model, 0.5°, no physics, 72 steps  
 MPE: TaihuLight, only on MPE  
 CPEs: TaihuLight, on MPE& CPEs



# GRAPES-GLB Dynamic core

0.25, Double Precision



GRAPES global model, 0.25° , no physics, 72 steps

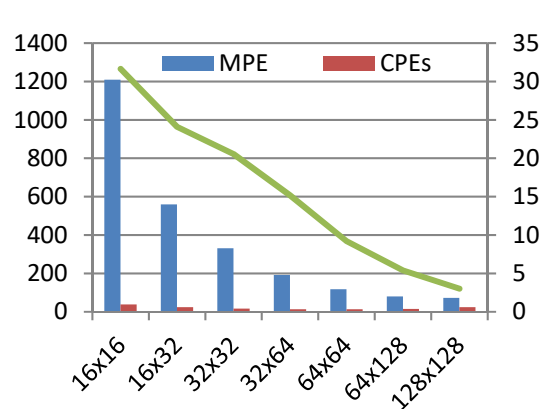
MPE: “Taihu Light”, no thread parallel processing

CPEs: “Taihu Light” with CPEs parallel processing

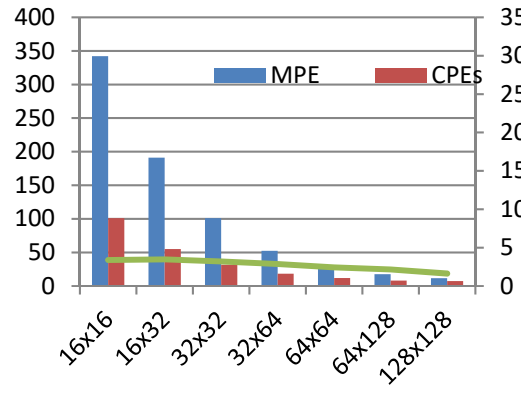
Intel: Node/CPU: Intel Xeon Gold 6124 (2.66GHz 16 cores) processor, 2 CPUs/node\* (16 cores/CPU) 100Gb/s EDR InfiniBand inter connection



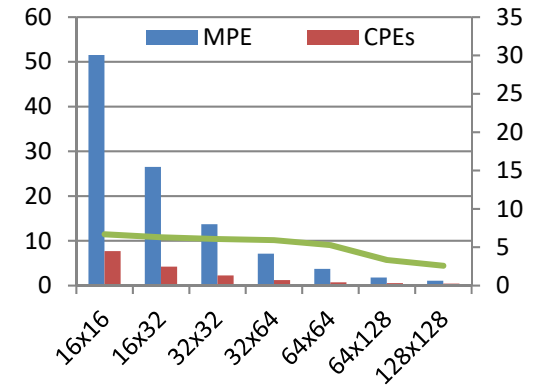
# 6 routines of dynamic core



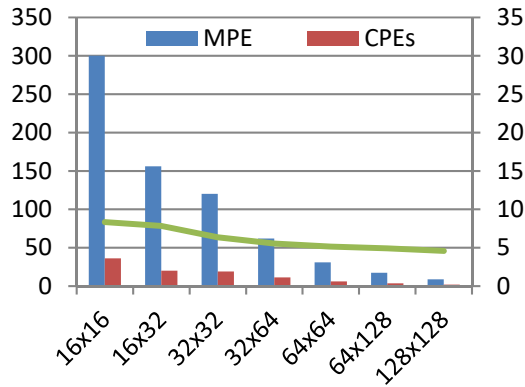
Helmholtz solver



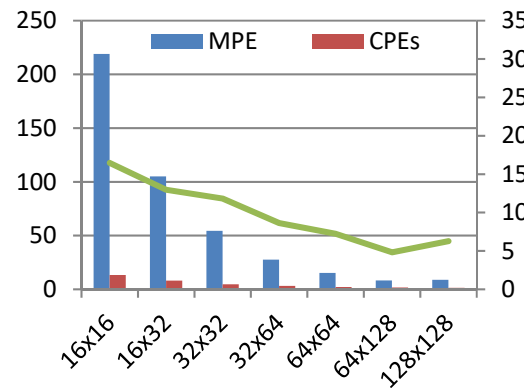
Departure points



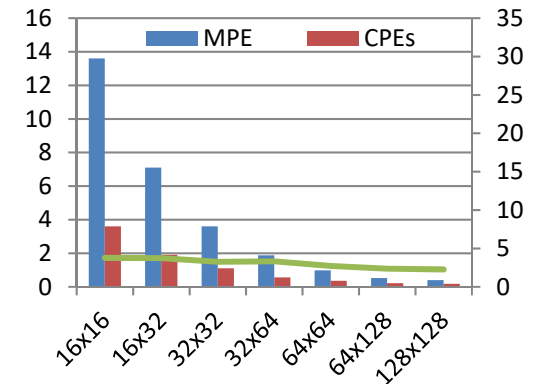
SL interp uvwthpi



SL interp trace substances



PRM

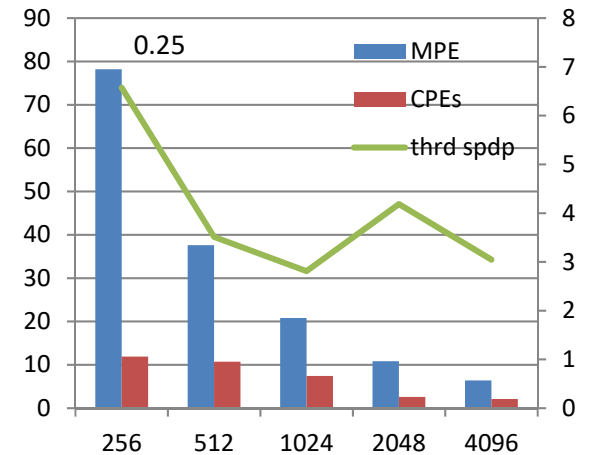
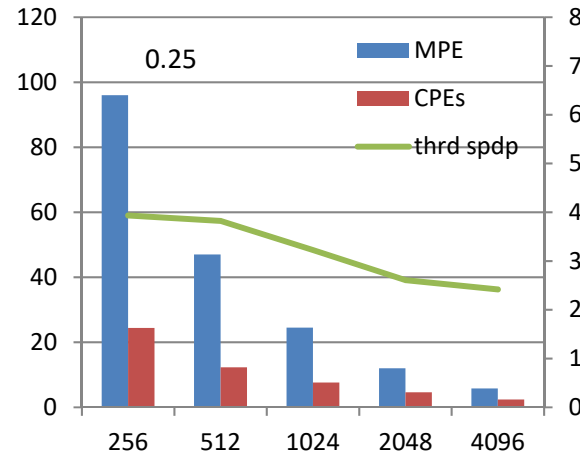
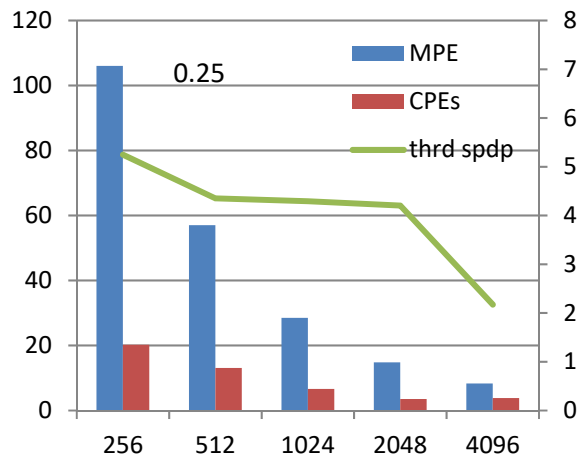
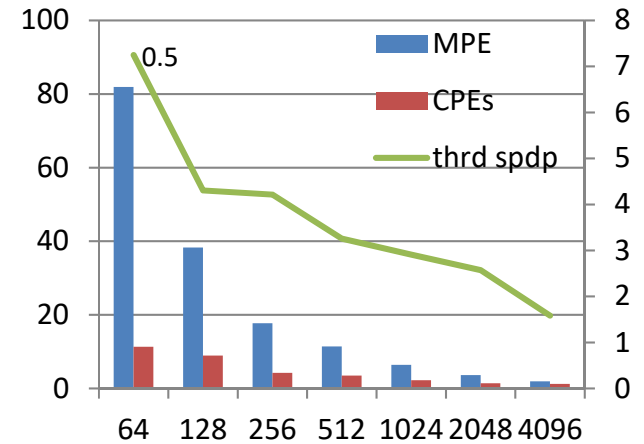
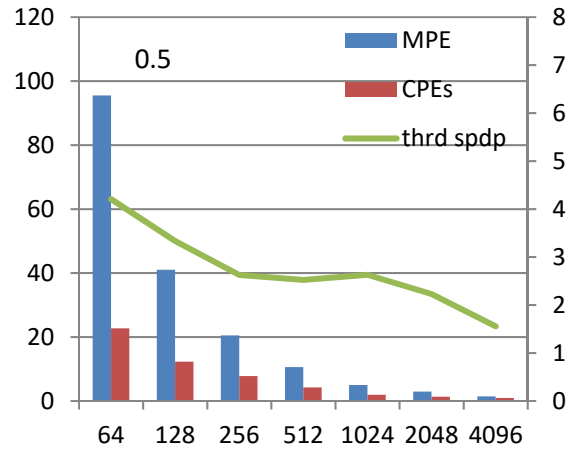
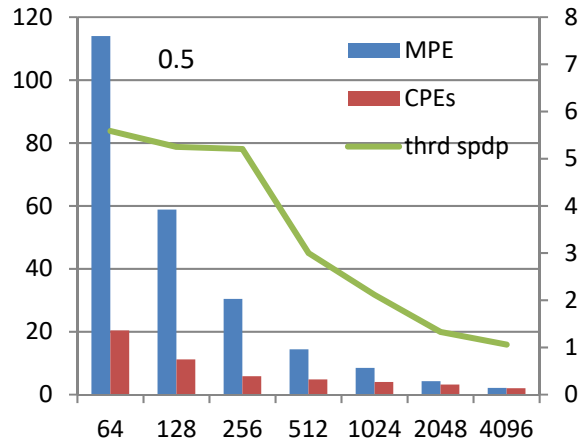


Linear & none linear

GRAPES global model, 0.25°, no physics, 72 steps  
 MPE: "Taihu Light", no thread parallel processing  
 CPEs: "Taihu Light" with CPEs parallel processing



# 3 routines in physics package



Phy\_prep

Microphysics

Phy\_post\_back

GRAPES global model, no physics, 72 steps  
MPE: "Taihu Light", only on MPE  
CPEs: "Taihu Light" on MPE & CPEs



# Summary

- Current swGRAPES dynamic core has comparable performance with Pi system (intel skylake processors system), suggesting that the refactorization of GRAPES having potential to use heterogeneous many-core platforms
- Fine-grained parallel algorithm for SpMV and ILU can better utilization of SW26010 processor, encouraging future many-core oriented algorithm design.
- Computation and communication overlapping can well fit to the heterogeneous supercomputing system
- OpenACC\* has moderate performance improvement and better portability for regular stencil-like computation loops, The key point of the performance is the best utilization of memory bandwidth
- Programming with Athread can get better performance due to full control of cache, accessing data, communication while it is lack of portability
- The work is far from end and further improvements is required



# Future Work

- GCR algorithm
  - Geometric Multigrid Preconditioning
  - Two communication optimization algorithms: Pipelined GCR and Chebyshev Methods
- Further Optimization on Sunway TaihuLight
- Improving portability by using OpenACC\* enhancement, high performance libraries, and code generation frameworks





**THANK YOU  
and questions?**

