

A Scalable Barotropic Solver for the Parallel Ocean Program

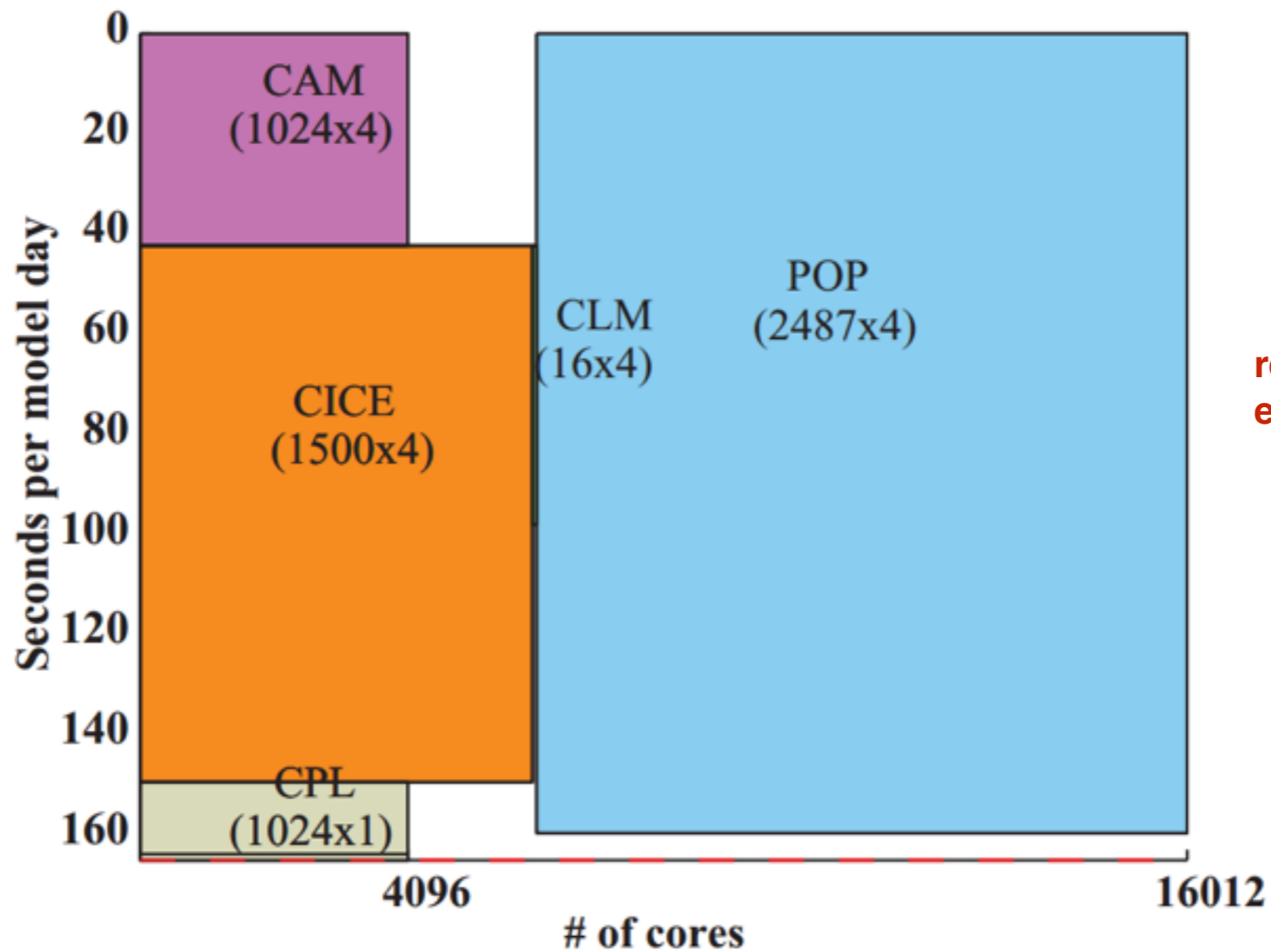
Yong Hu¹, Xiaomeng Huang¹, Yu-Heng Tseng², Allison Baker³,
Frank Bryan², and John Dennis³

1 Ministry of Education Key Laboratory for Earth System Modeling
Tsinghua University, China

2 Oceanography Section, Climate and Global Dynamics Division

3 Computational and Information Systems Laboratory
National Center for Atmospheric Research, USA

POP dominates in CESM

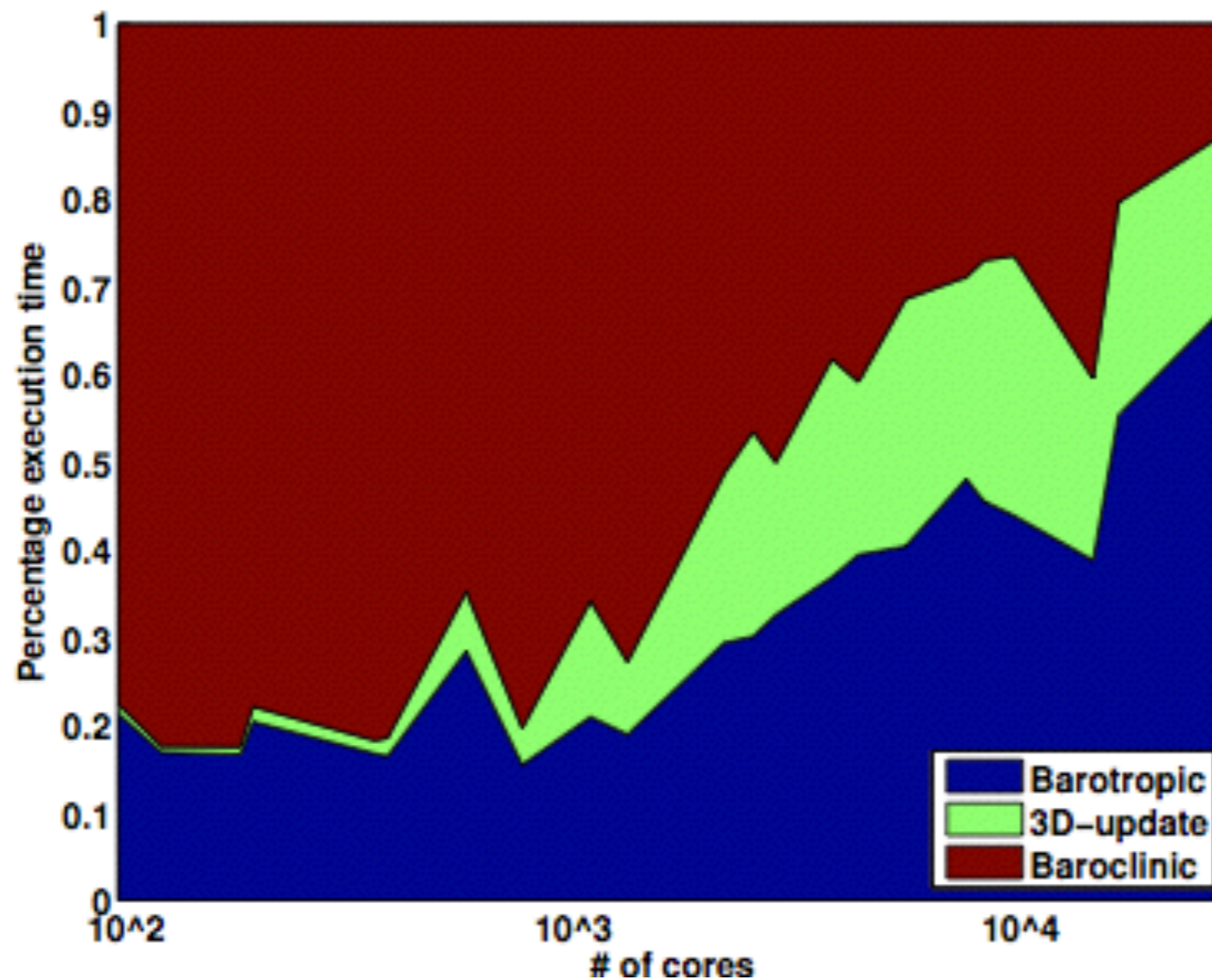


resources : 58% to 66%
efficiency : 42% to 58%

Computational performance of ultra-high-resolution capability in the Community Earth System Model. John Dennis, et al. HPCA'12

Scalability of 0.1° POP

Baroclinic mode time dominates on a small number of cores,
Barotropic mode time dominates in a large number of cores (>5000)



Kraken, 99,072 core Cray XT5 system at the National Institute for Computational Science (NICS). Andrew Stone. 2011

Barotropic mode in POP

- Require to solve an elliptic system for SSH in the barotropic mode.
- The elliptic problem is approximated as a linear system.

$$[\nabla \cdot H \nabla - \phi(\tau)]\eta^{n+1} = \psi(\eta^n, \eta^{n-1}, \tau) \longrightarrow Ax = b$$

- PCG (ChronGear) is used to solve this linear system.

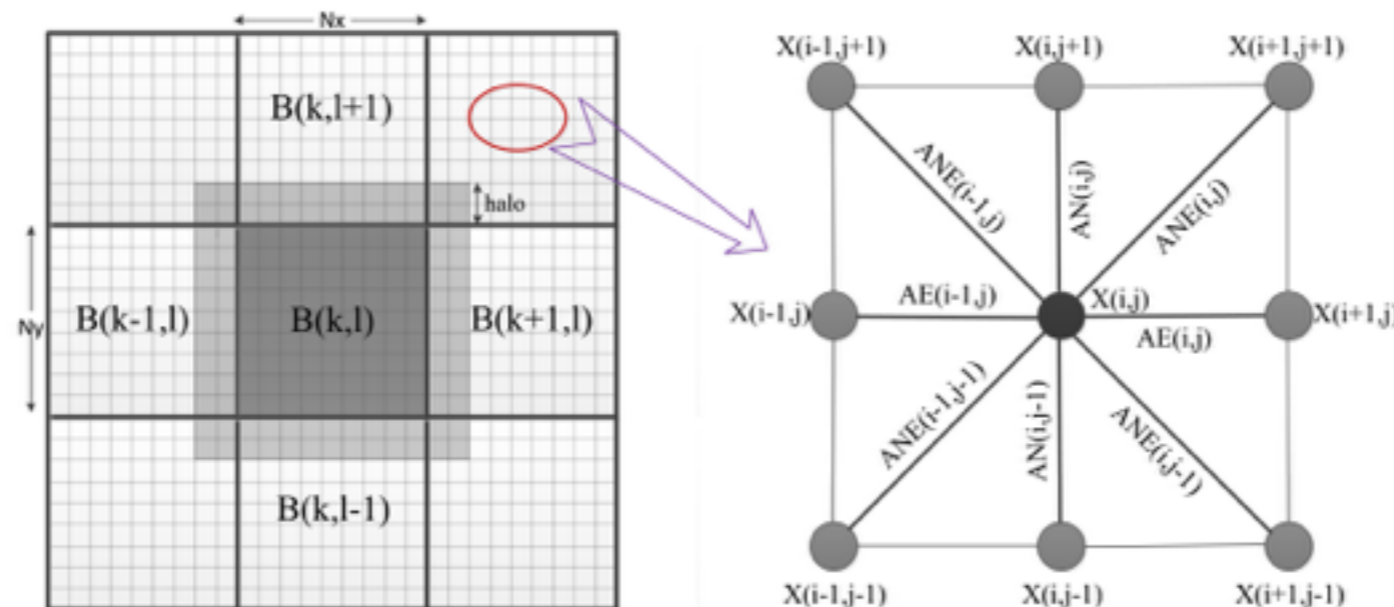
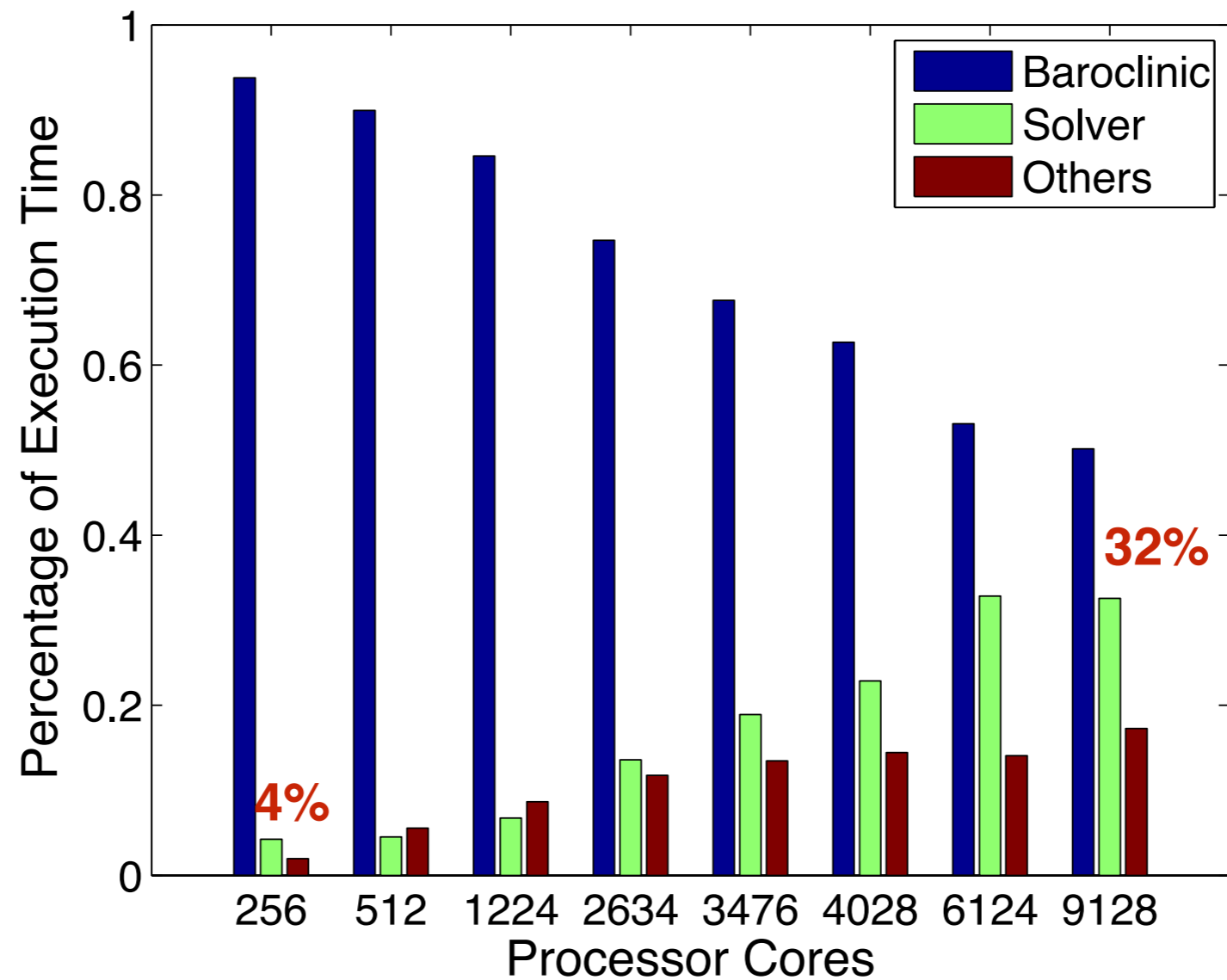


Fig. 1. Grid domain decomposition of POP

Scalability of 0.1° POP on Yellowstone

ChronGear Solver does not scale well on large number of processes.



ChronGear Solver

Algorithm 1 Chronopoulos-Gear Solver

Require: Coefficient matrix \mathbf{B} , preconditioner \mathbf{M} , initial guess \mathbf{x}_0 and \mathbf{b} associated with grid block $B_{i,j}$

// do in parallel with all processes

```

1:  $\mathbf{r}_0 = \mathbf{b} - \mathbf{B}\mathbf{x}_0$ ,  $\mathbf{s}_0 = \mathbf{0}$ ,  $\mathbf{p}_0 = \mathbf{0}$ ;  $\rho_0 = 1, \sigma_0 = 0, k = 0$ ;
2: while  $k \leq k_{max}$  do
3:    $k = k + 1$ ;
4:    $\mathbf{r}'_k = \mathbf{M}^{-1}\mathbf{r}_{k-1}$ ; /* preconditioning */
5:   update_halo( $\mathbf{r}'_{k-1}$ ); /* boundary communication */
6:    $\mathbf{z}_k = \mathbf{B}\mathbf{r}'_k$ ; /* matrix-vector multiplication */
7:   update_halo( $\mathbf{z}_k$ ); /* boundary communication */
8:    $\tilde{\rho}_k = \mathbf{r}'_{k-1}\mathbf{r}'_k$ ;
9:    $\tilde{\delta}_k = \mathbf{z}_k\mathbf{r}'_k$ ;
10:  ( $\rho_k, \delta_k$ ) = global_sum( $\tilde{\rho}_k, \tilde{\delta}_k$ ); /* global reduction */
11:   $\beta_k = \rho_k / \rho_{k-1}$ ;
12:   $\sigma_k = \delta_k - \beta_k^2 \sigma_{k-1}$ ;
13:   $\alpha_k = \rho_k / \sigma_k$ ;
14:   $\mathbf{s}_k = \mathbf{r}'_{k-1} + \beta_k \mathbf{s}_{k-1}$ ;
15:   $\mathbf{p}_k = \mathbf{z}_k + \beta_k \mathbf{p}_{k-1}$ ;
16:   $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{s}_k$ ;
17:   $\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_k \mathbf{p}_k$ ;
18:  if  $k \% n_c == 0$  then
19:    check convergence;
20:  end if
21: end while

```

Computation

$$\mathcal{T}_p + 15 \frac{\mathcal{N}^2}{p} \theta$$

θ : time unit per floating-point operation

\mathcal{N}^2 : number of grids

p : number of processors

ChronGear Solver

Algorithm 1 Chronopoulos-Gear Solver

Require: Coefficient matrix \mathbf{B} , preconditioner \mathbf{M} , initial guess \mathbf{x}_0 and \mathbf{b} associated with grid block $B_{i,j}$
// do in parallel with all processes

- 1: $\mathbf{r}_0 = \mathbf{b} - \mathbf{B}\mathbf{x}_0$, $\mathbf{s}_0 = 0$, $\mathbf{p}_0 = 0$; $\rho_0 = 1, \sigma_0 = 0$, $k = 0$;
- 2: **while** $k \leq k_{max}$ **do**
- 3: $k = k + 1$;
- 4: $\mathbf{r}'_k = \mathbf{M}^{-1}\mathbf{r}_{k-1}$; */* preconditioning */*
- 5: *update_halo*(\mathbf{r}'_{k-1}); */* boundary communication */*
- 6: $\mathbf{z}_k = \mathbf{B}\mathbf{r}'_k$; */* matrix-vector multiplication */*
- 7: *update_halo*(\mathbf{z}_k); */* boundary communication */*
- 8: $\tilde{\rho}_k = \mathbf{r}'_{k-1}{}^T \mathbf{r}'_k$;
- 9: $\tilde{\delta}_k = \mathbf{z}_k{}^T \mathbf{r}'_k$;
- 10: $(\rho_k, \delta_k) = \text{global_sum}(\tilde{\rho}_k, \tilde{\delta}_k)$; */* global reduction */*
- 11: $\beta_k = \rho_k / \rho_{k-1}$;
- 12: $\sigma_k = \delta_k - \beta_k^2 \sigma_{k-1}$;
- 13: $\alpha_k = \rho_k / \sigma_k$;
- 14: $\mathbf{s}_k = \mathbf{r}'_{k-1} + \beta_k \mathbf{s}_{k-1}$;
- 15: $\mathbf{p}_k = \mathbf{z}_k + \beta_k \mathbf{p}_{k-1}$;
- 16: $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{s}_k$;
- 17: $\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_k \mathbf{p}_k$;
- 18: **if** $k \% n_c == 0$ **then**
- 19: check convergence;
- 20: **end if**
- 21: **end while**

Boundary
Communication

$$4\alpha + \left(\frac{8\mathcal{N}}{\sqrt{p}}\right)\beta$$

α : communication latency

β : transfer time per byte

ChronGear Solver

Algorithm 1 Chronopoulos-Gear Solver

Require: Coefficient matrix \mathbf{B} , preconditioner \mathbf{M} , initial guess \mathbf{x}_0 and \mathbf{b} associated with grid block $B_{i,j}$
// do in parallel with all processes

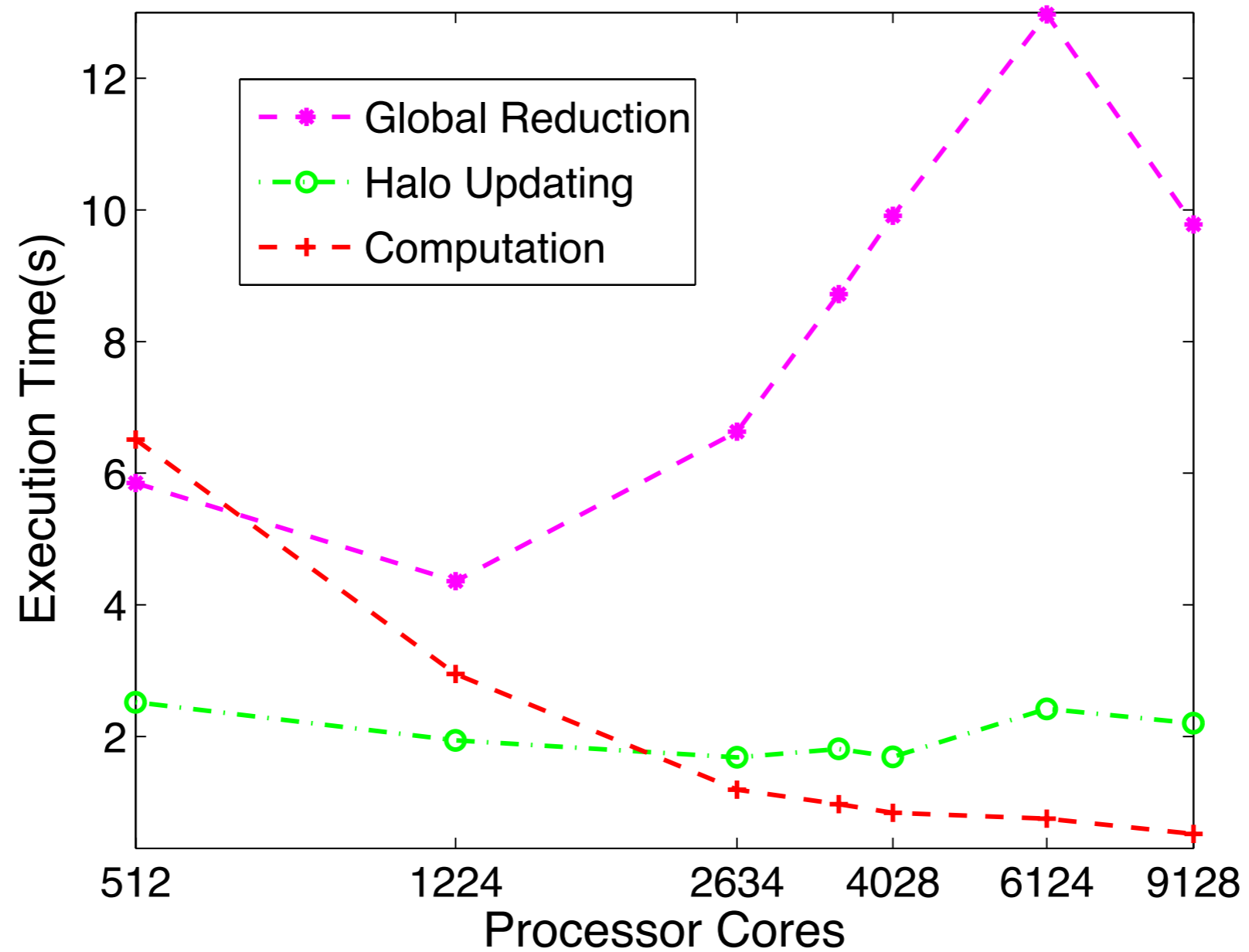
- 1: $\mathbf{r}_0 = \mathbf{b} - \mathbf{B}\mathbf{x}_0$, $\mathbf{s}_0 = 0$, $\mathbf{p}_0 = 0$; $\rho_0 = 1, \sigma_0 = 0, k = 0$;
- 2: **while** $k \leq k_{max}$ **do**
- 3: $k = k + 1$;
- 4: $\mathbf{r}'_k = \mathbf{M}^{-1}\mathbf{r}_{k-1}$; */* preconditioning */*
- 5: $update_halo(\mathbf{r}'_{k-1})$; */* boundary communication */*
- 6: $\mathbf{z}_k = \mathbf{B}\mathbf{r}'_k$; */* matrix-vector multiplication */*
- 7: $update_halo(\mathbf{z}_k)$; */* boundary communication */*
- 8: $\tilde{\rho}_k = \mathbf{r}'_{k-1T}\mathbf{r}'_k$;
- 9: $\tilde{\delta}_k = \mathbf{z}_k^T\mathbf{r}'_k$;
- 10: $(\rho_k, \delta_k) = global_sum(\tilde{\rho}_k, \tilde{\delta}_k)$; */* global reduction */*
- 11: $\beta_k = \rho_k / \rho_{k-1}$;
- 12: $\sigma_k = \delta_k - \beta_k^2\sigma_{k-1}$;
- 13: $\alpha_k = \rho_k / \sigma_k$;
- 14: $\mathbf{s}_k = \mathbf{r}'_{k-1} + \beta_k\mathbf{s}_{k-1}$;
- 15: $\mathbf{p}_k = \mathbf{z}_k + \beta_k\mathbf{p}_{k-1}$;
- 16: $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k\mathbf{s}_k$;
- 17: $\mathbf{r}_k = \mathbf{r}_{k-1} - \alpha_k\mathbf{p}_k$;
- 18: **if** $k \% n_c == 0$ **then**
- 19: check convergence;
- 20: **end if**
- 21: **end while**

Global
Reduction

$$\mathcal{T}_g = \log p\alpha$$

ChronGear Solver Components

Global reduction is the bottleneck of the ChronGear Solver in 0.1° POP.



P-CSI Solver

Algorithm 2 Preconditioned Stiefel Iteration solver

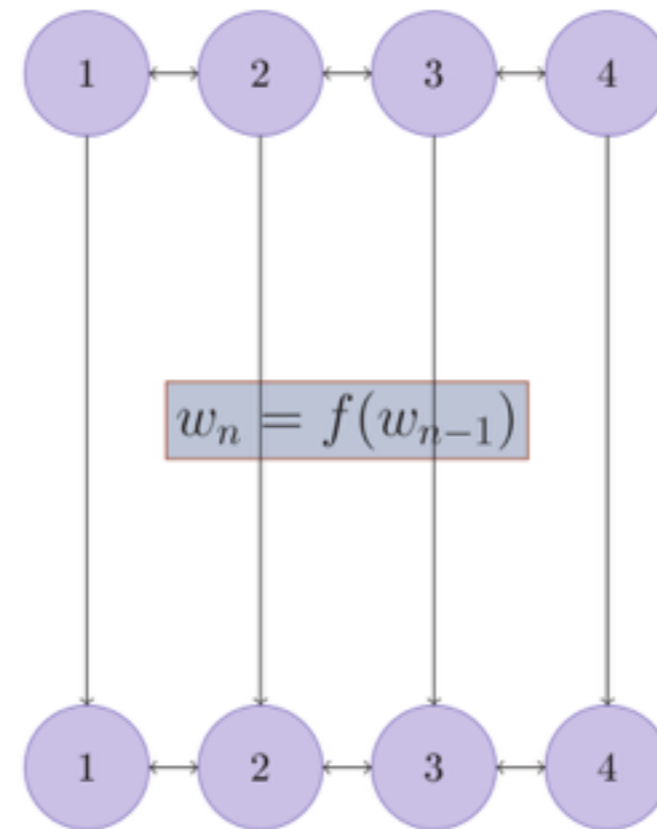
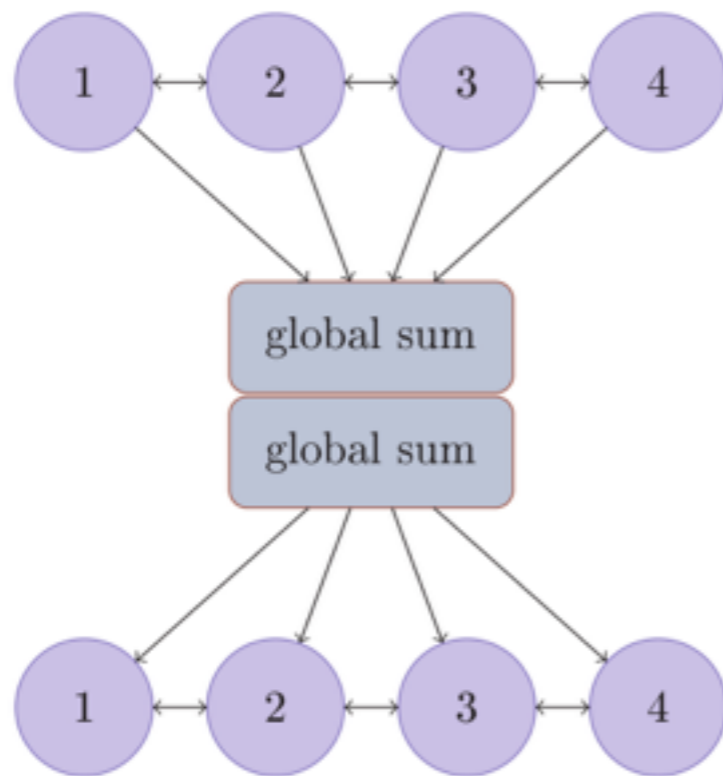
Require: Coefficient matrix \mathbf{B} , preconditioner \mathbf{M} , initial guess \mathbf{x}_0 and \mathbf{b} associated with grid block $B_{i,j}$; Estimated eigenvalue boundary $[\nu, \mu]$;

// do in parallel with all processes

- 1: $\alpha = \frac{2}{\mu - \nu}$, $\beta = \frac{\mu + \nu}{\mu - \nu}$, $\gamma = \frac{\beta}{\alpha}$, $\omega_0 = \frac{2}{\gamma}$; $k = 0$;
 - 2: $\mathbf{r}_0 = \mathbf{b} - \mathbf{B}\mathbf{x}_0$; $\mathbf{x}_1 = \mathbf{x}_0 - \gamma^{-1}\mathbf{M}^{-1}\mathbf{r}_0$; $\mathbf{r}_1 = \mathbf{b} - \mathbf{B}\mathbf{x}_1$;
 - 3: **while** $k \leq k_{max}$ **do**
 - 4: $k = k + 1$;
 - 5: $\omega_k = 1/(\gamma - \frac{1}{4\alpha^2}\omega_{k-1})$; */* the iterated function */*
 - 6: $\mathbf{r}'_{k-1} = \mathbf{M}^{-1}\mathbf{r}_{k-1}$; */* preconditioning */*
 - 7: $\Delta\mathbf{x}_k = \omega_k\mathbf{r}'_{k-1} + (\gamma\omega_k - 1)\Delta\mathbf{x}_{k-1}$;
 - 8: $\mathbf{x}_k = \mathbf{x}_{k-1} + \Delta\mathbf{x}_{k-1}$;
 - 9: $\mathbf{r}_k = \mathbf{b} - \mathbf{B}\mathbf{x}_k$; */* matrix-vector multiplication */*
 - 10: *update_halo*(\mathbf{r}_k); */* boundary communication */*
 - 11: **if** $k \% n_c == 0$ **then**
 - 12: check convergence;
 - 13: **end if**
 - 14: **end while**
-

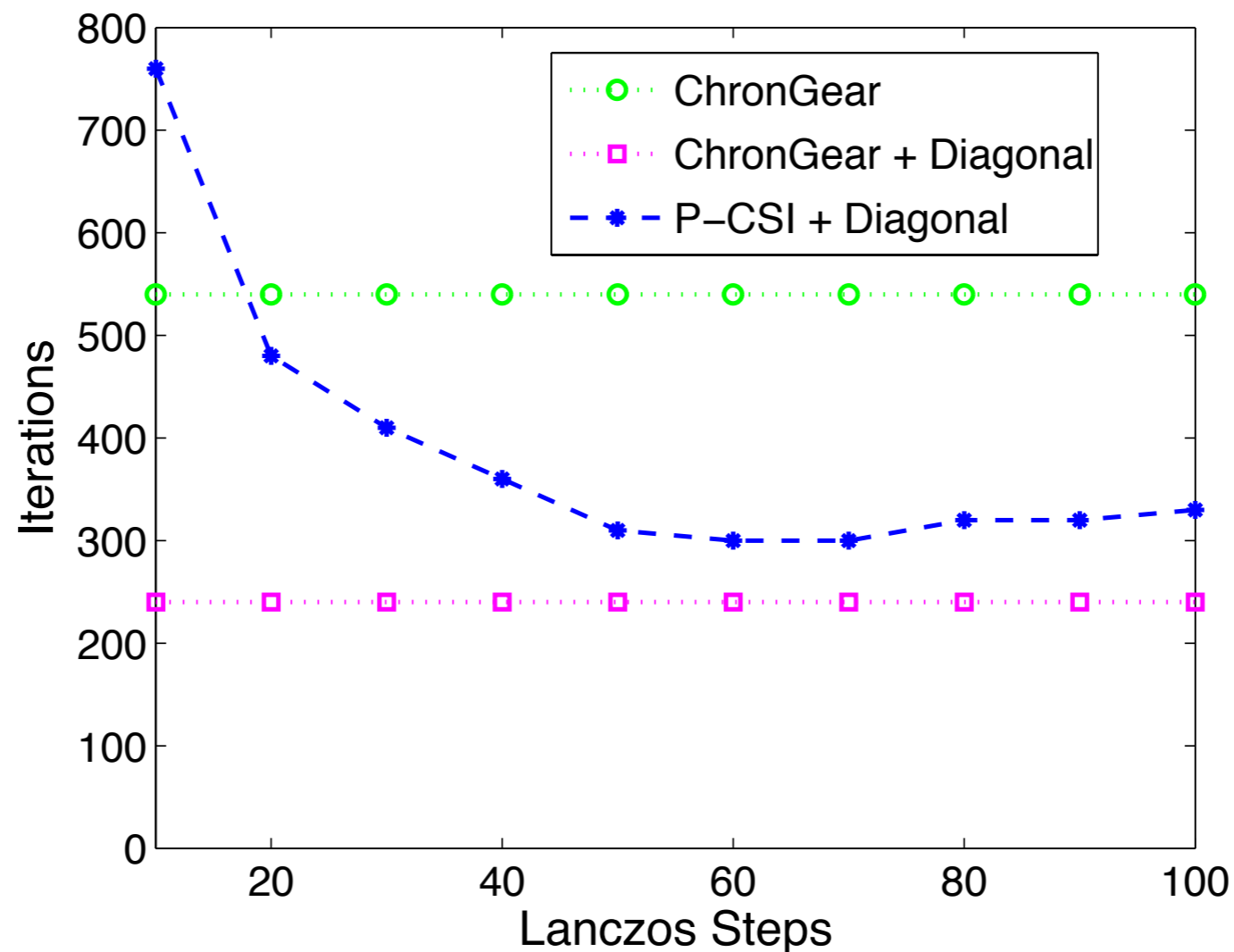
P-CSI .vs. ChronGear

- Unlike ChronGear, P-CSI has no global reduction
- P-CSI needs slightly more iterations, but less computation in each iteration
- P-CSI requires two extreme eigenvalues



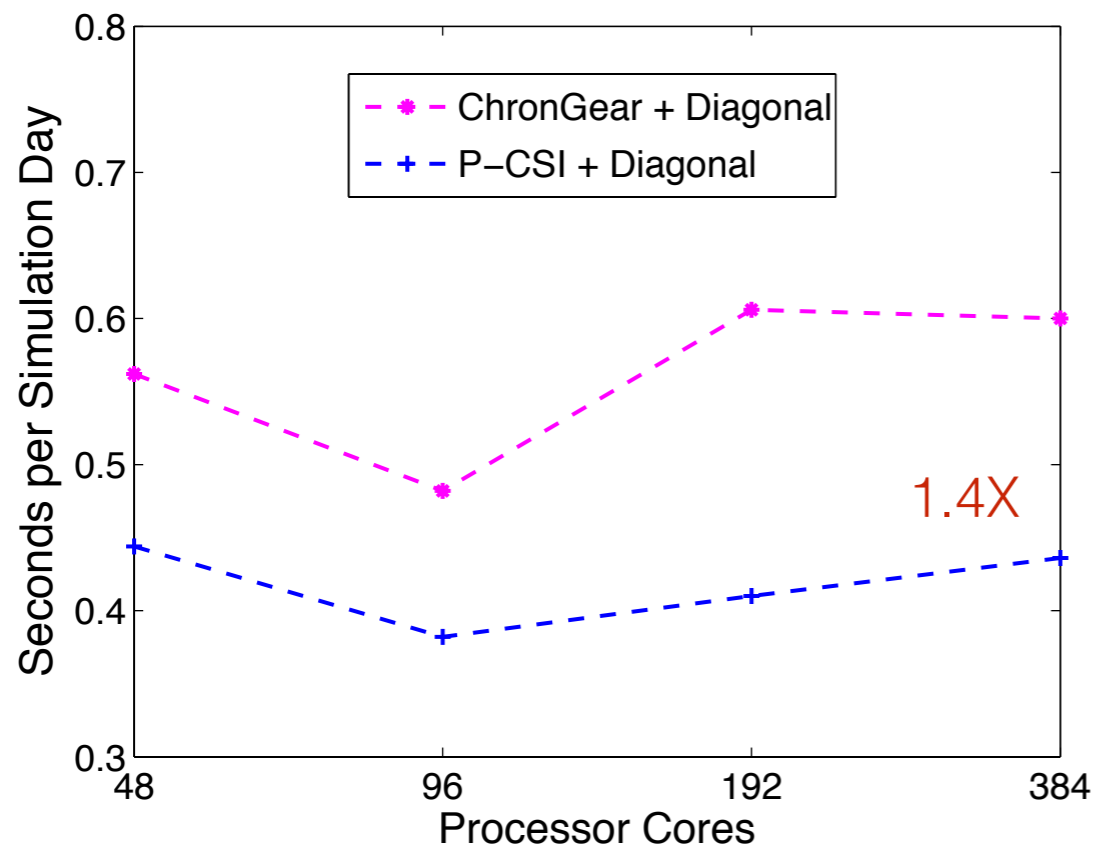
Eigenvalue Estimation

- Lanczos method to construct a tridiagonal matrix T
- T has eigenvalues close to $M^{-1}A$
- Estimation needed only once, extra overhead less than one barotropic step.

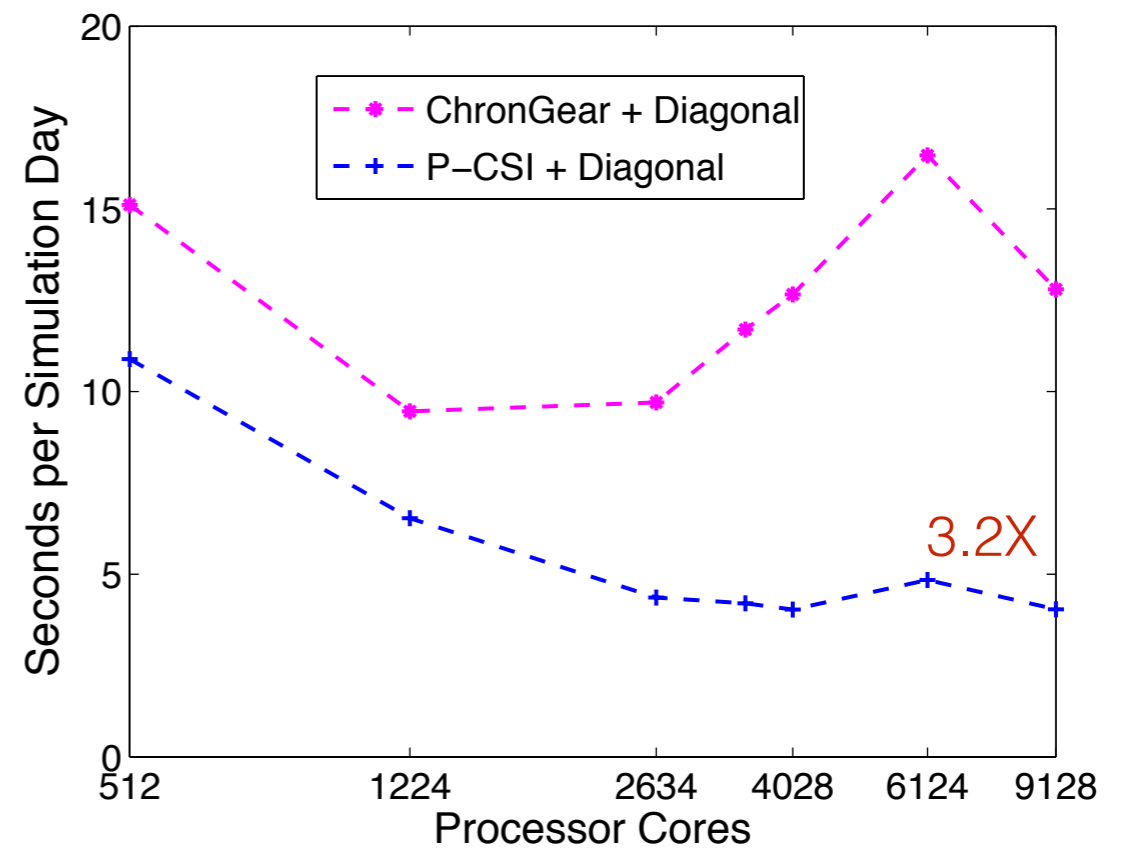


Scalability on Yellowstone

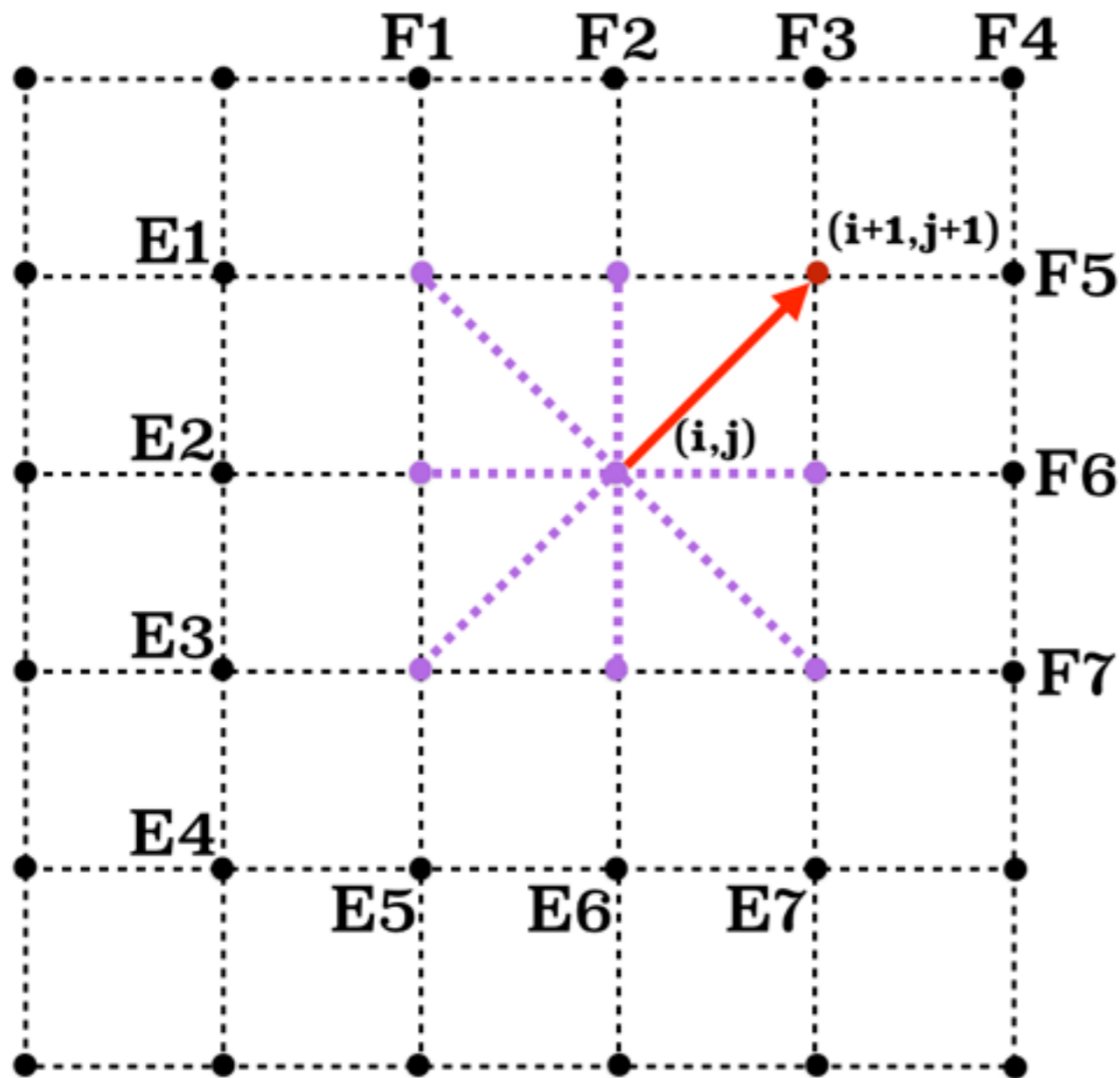
1 degree POP



0.1 degree POP



Error Vector Propagation (EVP) Preconditioning



$$\begin{aligned}
 &A_{i,j}^0 \eta_{i,j} + A_{i,j}^e \eta_{i+1,j} + A_{i,j}^n \eta_{i,j+1} + A_{i,j}^{ne} \eta_{i+1,j+1} \\
 &+ A_{i-1,j}^{ne} \eta_{i-1,j+1} + A_{i-1,j}^e \eta_{i-1,j} + A_{i-1,j-1}^{ne} \eta_{i-1,j-1} \\
 &+ A_{i,j-1}^n \eta_{i,j-1} + A_{i+1,j-1}^{ne} \eta_{i,j-1} = \psi_{i,j}
 \end{aligned}$$



$$\begin{aligned}
 \eta_{i+1,j+1} = &(1/A_{i,j}^{ne})(\psi_{i,j} - A_{i,j}^0 \eta_{i,j} - A_{i,j}^e \eta_{i+1,j} \\
 &- A_{i,j}^n \eta_{i,j+1} - A_{i-1,j}^{ne} \eta_{i-1,j+1} + A_{i-1,j}^e \eta_{i-1,j} \\
 &- A_{i-1,j-1}^{ne} \eta_{i-1,j-1} - A_{i,j-1}^n \eta_{i,j-1} - A_{i+1,j-1}^{ne} \eta_{i,j-1})
 \end{aligned}$$

Roache 1995, *Elliptic marching methods and domain decomposition*.

Preconditioning and Iterations

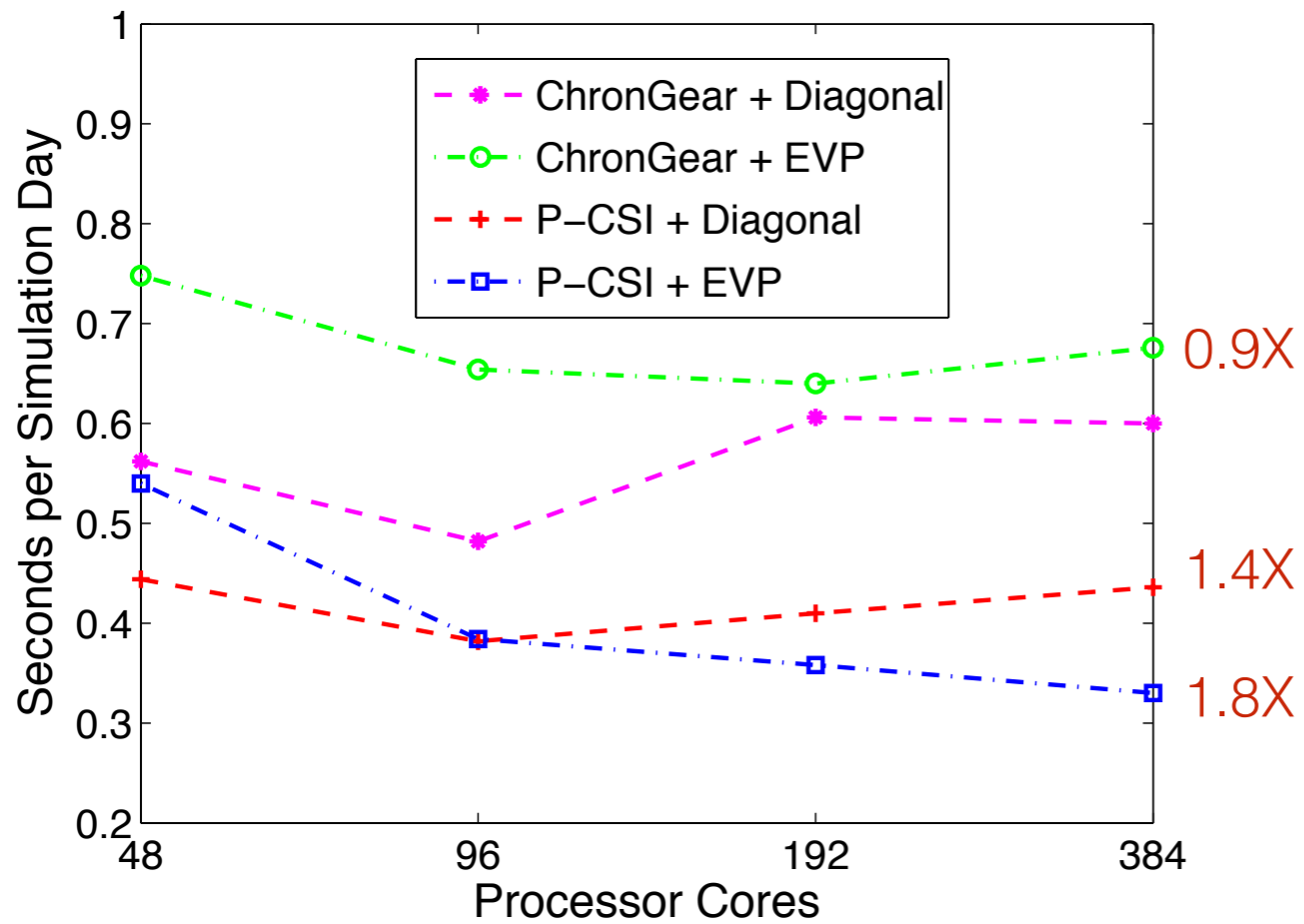
$$\begin{bmatrix} B_i^{nw} & B_i^n & B_i^{ne} \\ B_i^w & B_i & B_i^e \\ B_i^{sw} & B_i^s & B_i^{se} \end{bmatrix} M^{-1} = \begin{bmatrix} B_1^{-1} & & \\ & \ddots & \\ & & B_{m^2}^{-1} \end{bmatrix}$$

EVP preconditioning reduces the iteration number to about one-third

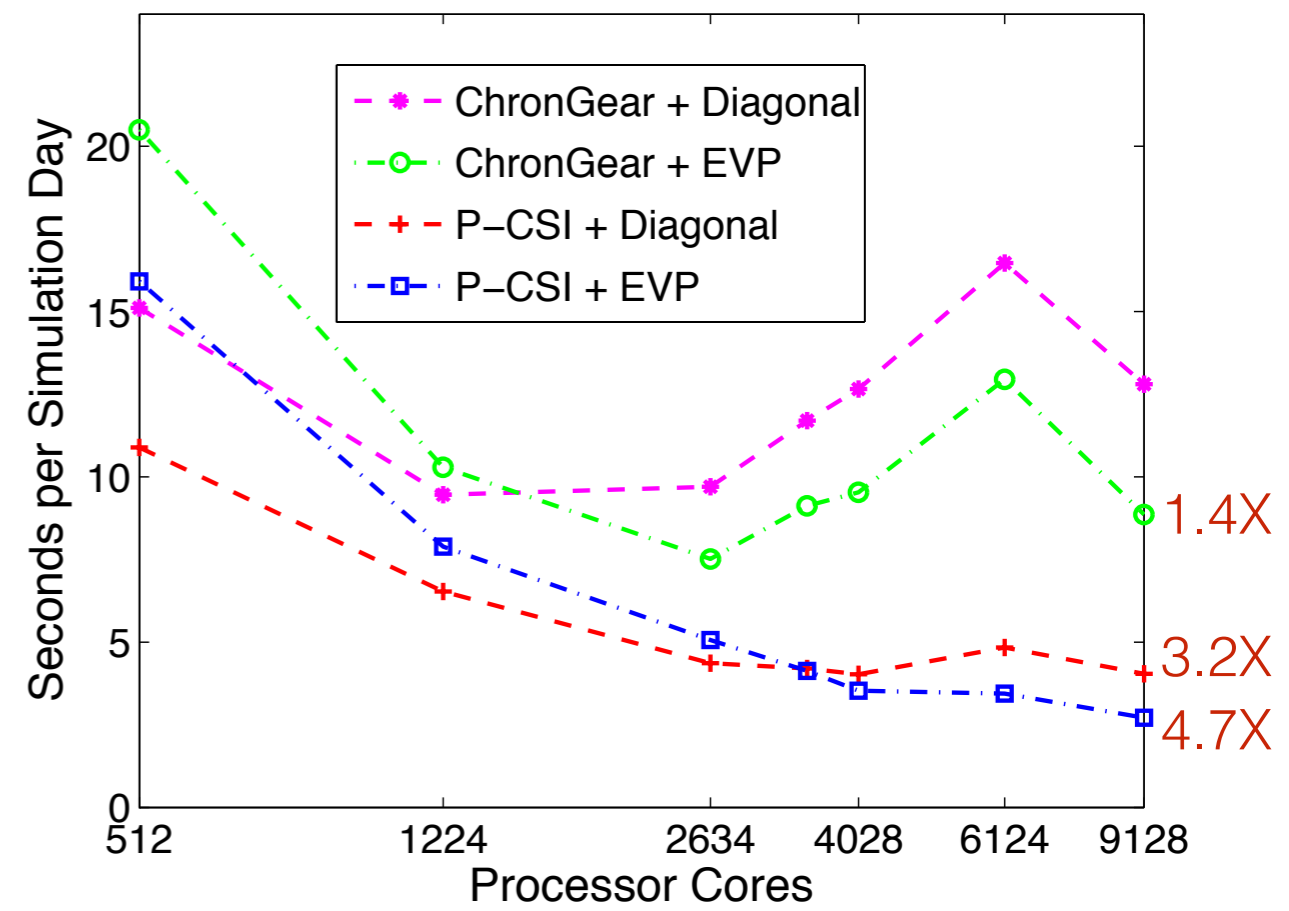
Resolution	1 degree		0.1 degree	
Solvers	ChronGear	P-CSI	ChronGear	P-CSI
NONE	540	570	150	160
DIAG	240	300	100	120
EVP	160	200	50	60

Scalability on Yellowstone

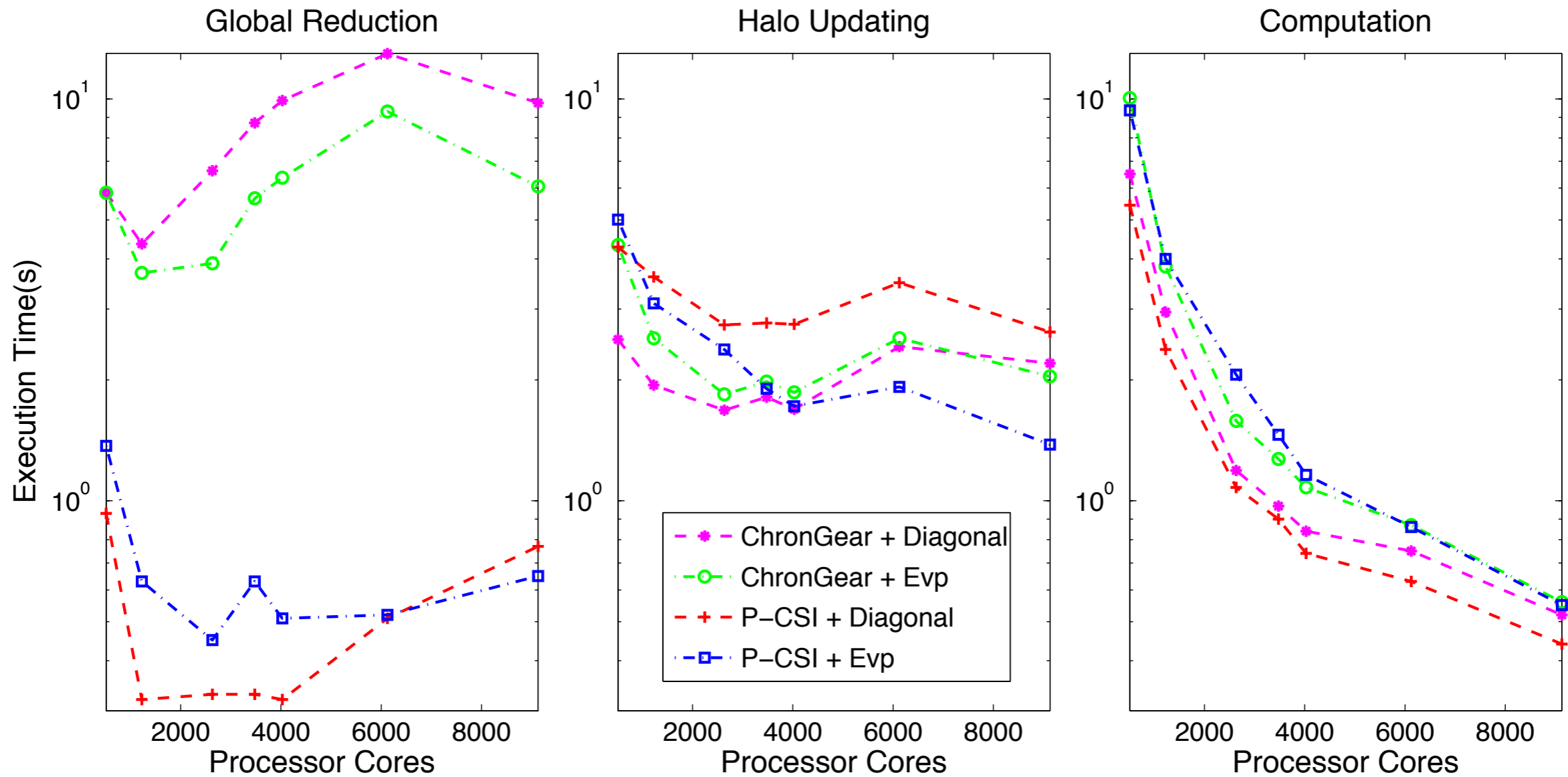
1 degree POP



0.1 degree POP



Solver Component



Simulated Years for Solver Per wall-clock Day

Simulation rate without considering the I/O and initialization in 0.1 degree POP.

Number of cores	512	1224	2634	3476	4028	6124	9128
ChronGear	0.71	1.69	3.31	3.83	4.28	4.72	6.02
ChronGear+Evp	0.70	1.68	3.42	3.99	4.53	5.08	6.69
P-CSI+Diagonal	0.72	1.72	3.58	4.36	5.07	6.15	7.75
P-CSI+Evp	0.71	1.70	3.54	4.36	5.12	6.38	8.10



35%

Q & A

Thanks!