

GPU Performance Comparison of Preconditioned Krylov Solvers



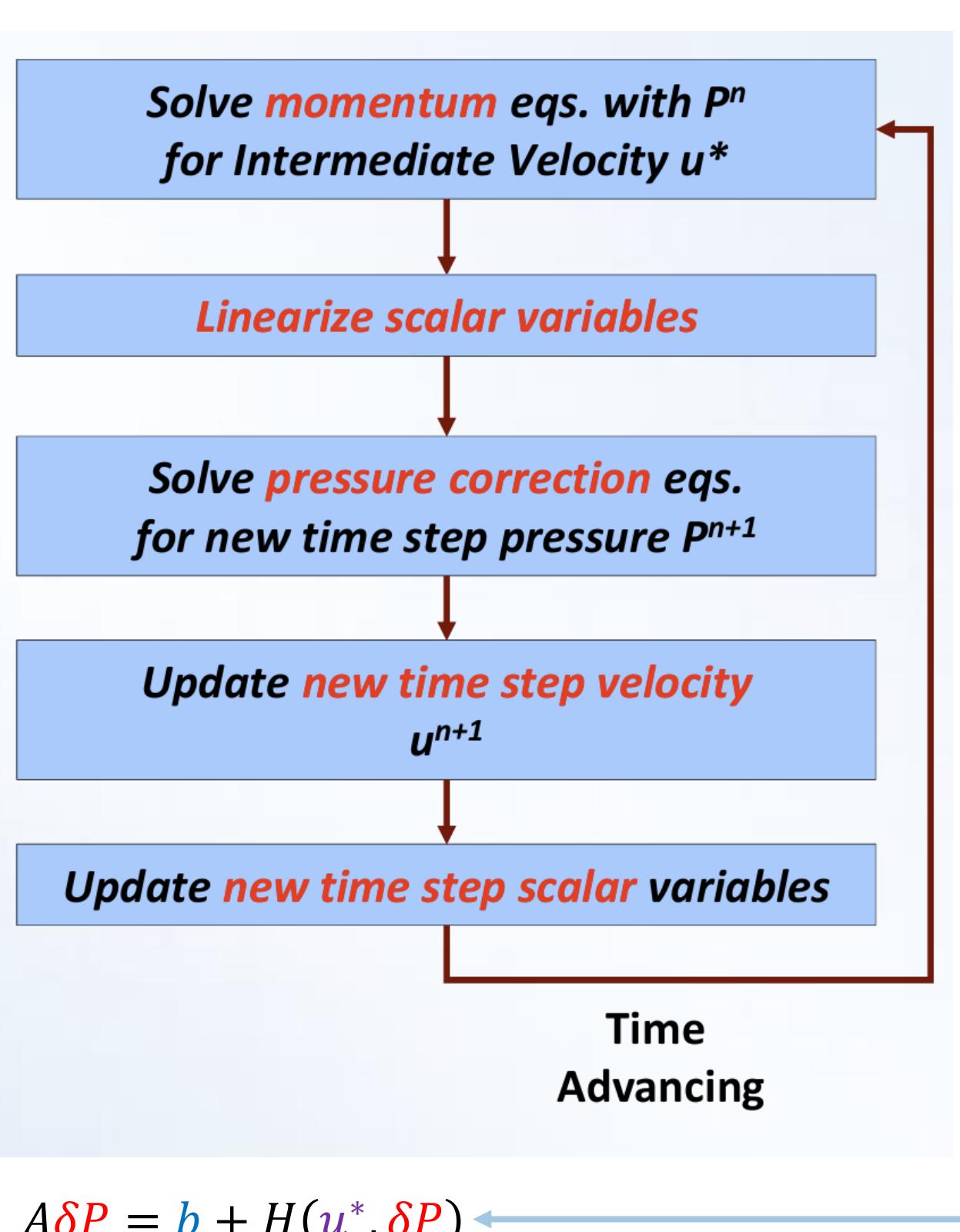
Ki-Ha Kim, Young Kwang Hwang, Jong-Chan Kim

Dept. of Supercomputing Acceleration Research, Korea Institute of Science and Technology Information

Abstract

A three-dimensional two-phase thermal-hydraulic simulation code is employed for nuclear reactor safety analysis. In this framework, the pressure field is obtained from a pressure-correction equation that forms a globally coupled sparse linear system. As the pressure solution accounts for over 80% of the total runtime, it represents the major computational bottleneck. To address this limitation, the present work accelerates the pressure solver on a single-GPU architecture by implementing and comparing three preconditioning strategies: diagonal, ILU(0), and AMG. The BiCGStab solver components were ported to CUDA using cuSPARSE and cuBLAS libraries. Performance was evaluated on KISTI's supercomputer comparing a 64-core Xeon Phi 7250 against an NVIDIA GH200 GPU. Results show that the diagonal preconditioner achieves the best balance between convergence and efficiency due to its fully parallel structure, while ILU(0) is limited by sparse triangular solves. The AMG solver converges rapidly but requires costly setup at each time step, reducing its overall benefit.

Introduction



Momentum (calc_momentum)

$$\begin{bmatrix} \vec{u}_g^* \\ \vec{u}_l^* \\ \vec{u}_d^* \end{bmatrix} = \begin{bmatrix} m_{gg} & m_{gl} & m_{gd} \\ m_{lg} & m_{ll} & m_{ld} \\ m_{dg} & m_{dl} & m_{dd} \end{bmatrix}^{-1} \begin{bmatrix} \vec{s}_g \\ \vec{s}_l \\ \vec{s}_d \end{bmatrix} - \begin{bmatrix} V\alpha_g \\ V\alpha_l \\ V\alpha_d \end{bmatrix} \nabla P^*$$

Scalar (calc_scalar): mass+energy

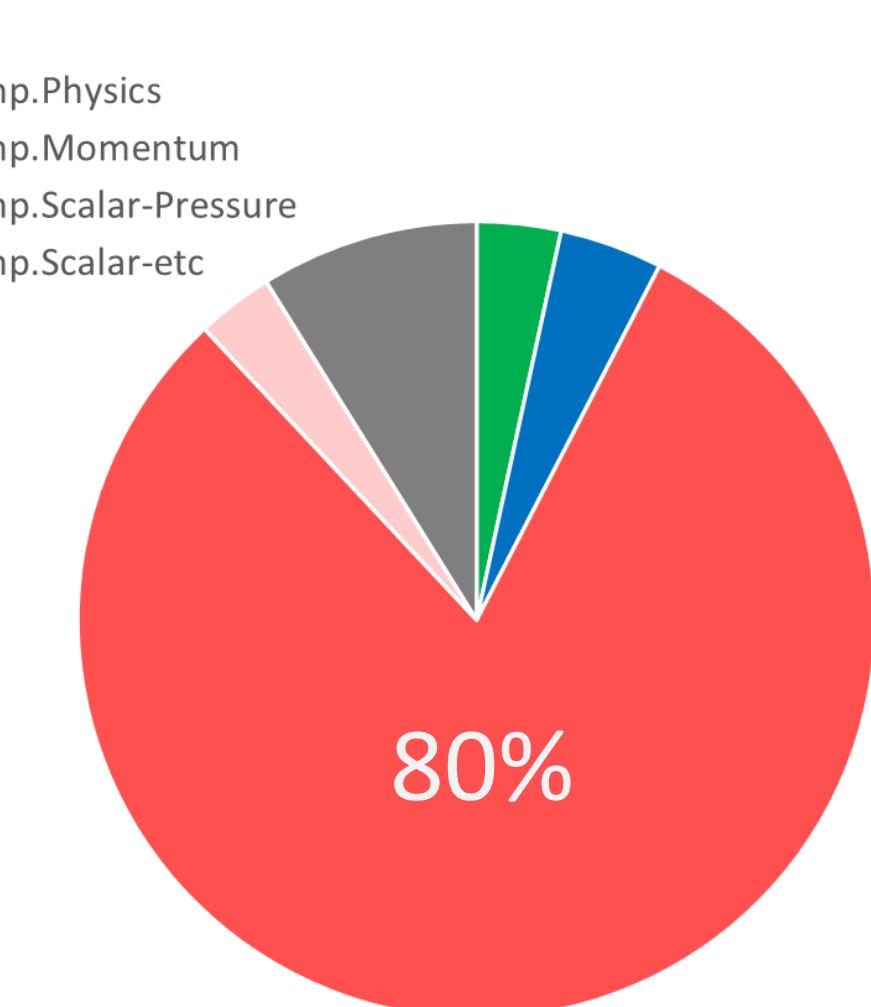
$$AX = B + \sum_f G \psi_g^{n+1} + \sum_f L \psi_l^{n+1} + \sum_f D \psi_d^{n+1}$$

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} & a_{15} & a_{16} \\ a_{21} & a_{22} & a_{23} & a_{24} & a_{25} & a_{26} \\ a_{31} & a_{32} & a_{33} & a_{34} & a_{35} & a_{36} \\ a_{41} & a_{42} & a_{43} & a_{44} & a_{45} & a_{46} \\ a_{51} & a_{52} & a_{53} & a_{54} & a_{55} & a_{56} \\ a_{61} & a_{62} & a_{63} & a_{64} & a_{65} & a_{66} \end{bmatrix} \begin{bmatrix} \delta e_g \\ \delta e_l \\ \delta e_d \\ \delta X_g \\ \delta X_l \\ \delta X_d \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \\ b_6 \end{bmatrix}, \quad \begin{bmatrix} g_{1,f} \\ g_{2,f} \\ g_{3,f} \\ g_{4,f} \\ g_{5,f} \\ g_{6,f} \end{bmatrix}, \quad \begin{bmatrix} l_{1,f} \\ l_{2,f} \\ l_{3,f} \\ l_{4,f} \\ l_{5,f} \\ l_{6,f} \end{bmatrix}, \quad \begin{bmatrix} d_{1,f} \\ d_{2,f} \\ d_{3,f} \\ d_{4,f} \\ d_{5,f} \\ d_{6,f} \end{bmatrix}$$

Pressure correction

$$\begin{bmatrix} \vec{u}_g^{n+1} \\ \vec{u}_l^{n+1} \\ \vec{u}_d^{n+1} \end{bmatrix} = \begin{bmatrix} \vec{u}_g^* \\ \vec{u}_l^* \\ \vec{u}_d^* \end{bmatrix} - \begin{bmatrix} m_{gg} & m_{gl} & m_{gd} \\ m_{lg} & m_{ll} & m_{ld} \\ m_{dg} & m_{dl} & m_{dd} \end{bmatrix}^{-1} \begin{bmatrix} V\alpha_g \\ V\alpha_l \\ V\alpha_d \end{bmatrix} \nabla \delta P$$

$$A\delta P = b + H(u^*, \delta P)$$



A three-dimensional two-phase thermal-hydraulic code is employed for nuclear reactor safety analysis. It employs unstructured meshes for complex geometries and MPI-based domain decomposition for parallel execution. The simulation proceeds through momentum calculation, scalar transport (mass and energy), and pressure correction stages. Among these, the pressure correction step dominates the computational cost, accounting for over 80% of the total runtime. The pressure field is obtained by solving a globally coupled sparse linear system derived from the pressure-correction equation. This system must be solved at every time step, making efficient solver design essential for large-scale transient simulations.

Result

(1) Convergence History Comparison

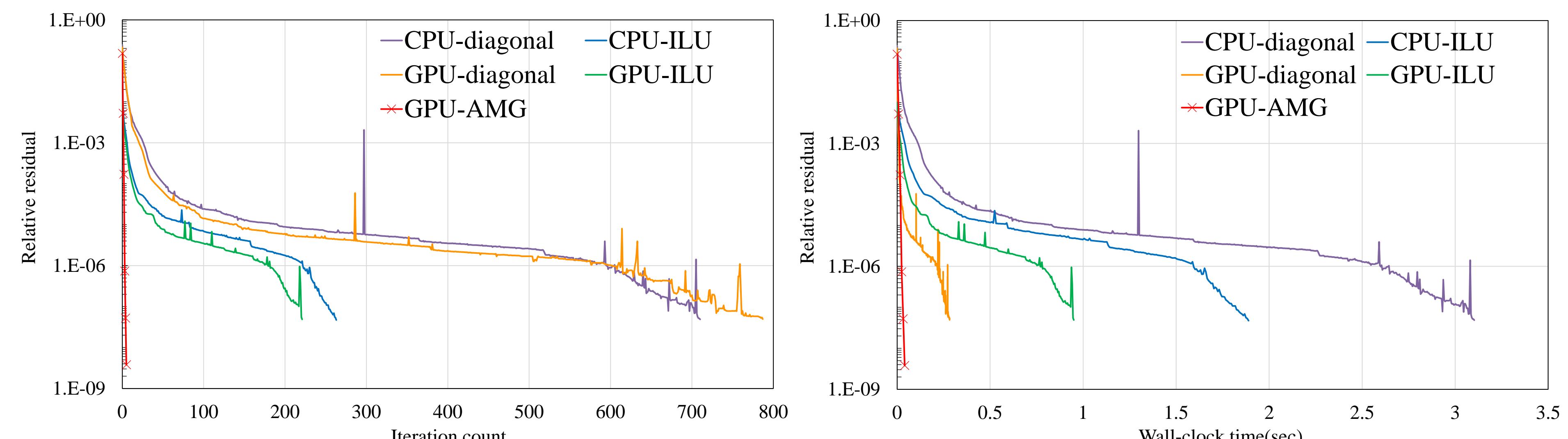


Fig1. Convergence history of preconditioned solvers: residual versus iteration count (left) and residual versus wall-clock time (right).

Figure 1 compares residual convergence histories by iteration count (left) and wall-clock time (right). In terms of iterations, ILU(0) converges in approximately 200 iterations while diagonal requires about 700, and AMG achieves convergence within only 5 iterations. However, when measured by wall time, the GPU diagonal solver reaches the convergence threshold faster than GPU ILU, despite requiring more iterations. This is because the diagonal preconditioner has much lower per-iteration cost due to its simple parallel structure. The AMG solver demonstrates the fastest convergence in both metrics, completing in under 0.1 seconds.

(3) Residual Distribution at Convergence

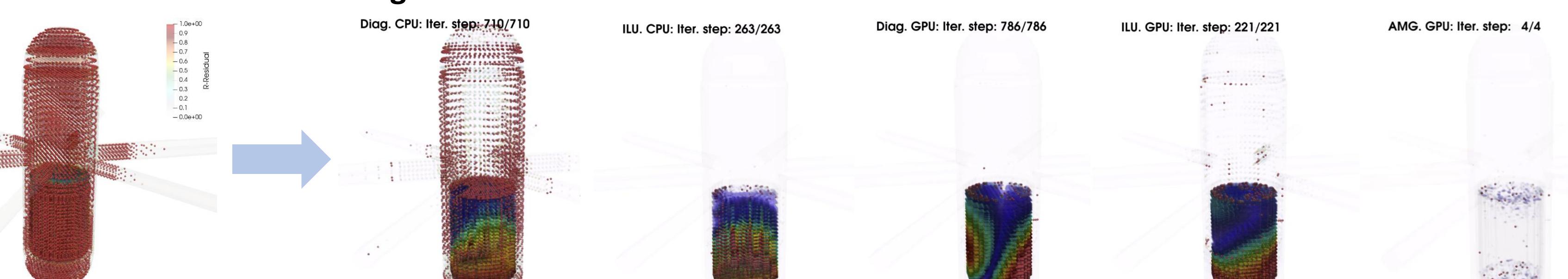


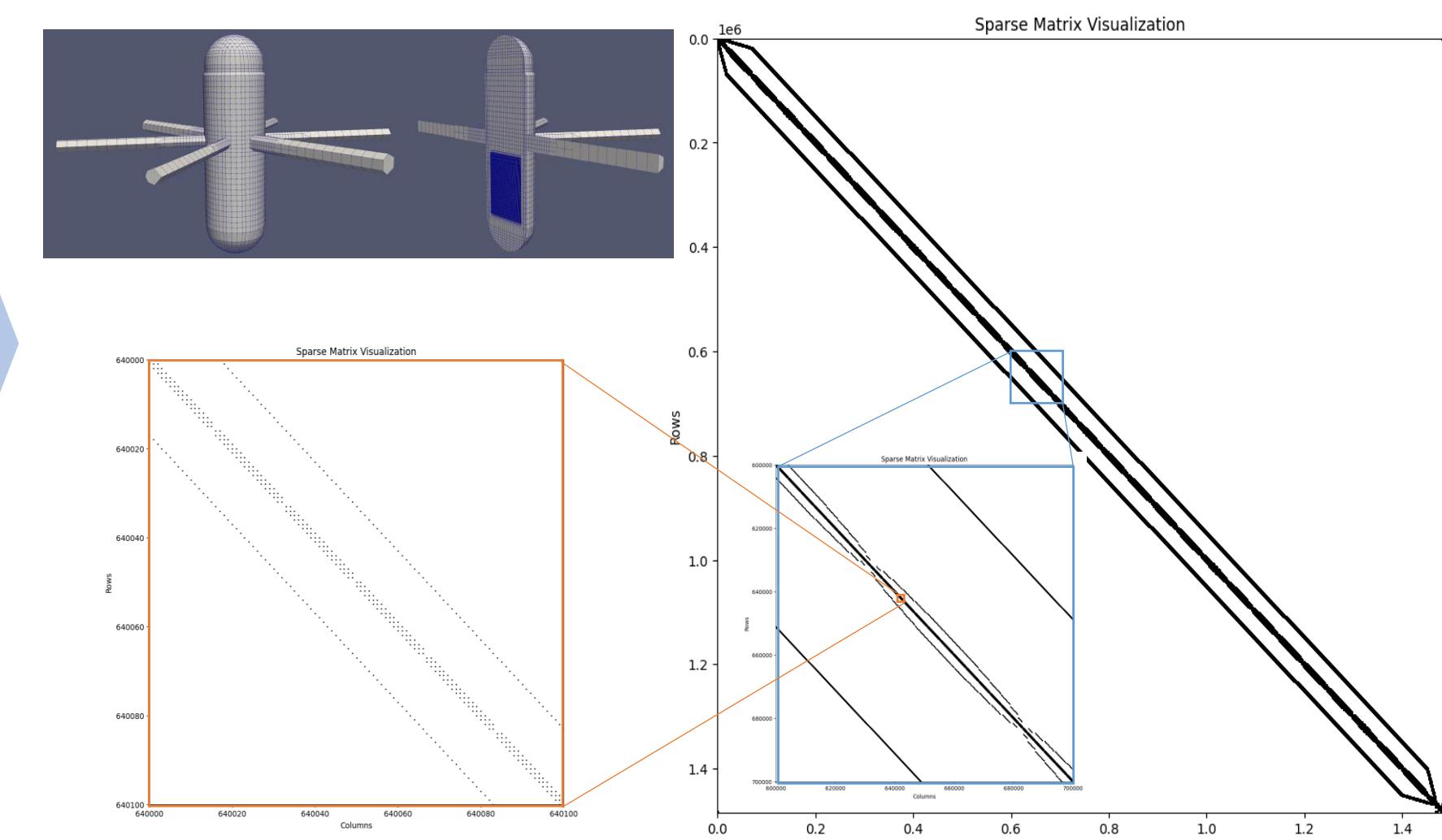
Fig3. Spatial distribution of relative residual at convergence for each solver (left to right: CPU diagonal, CPU ILU, GPU diagonal, GPU ILU, GPU AMG). Red/opaque indicates high residual; blue/transparent indicates low residual.

Method

Pressure Equation as Sparse Linear System

$$A\delta P = b + H(u^*, \delta P)$$

$$\begin{bmatrix} 1 + \sum_i CC_{g,i} \\ \vdots \\ (CC_{g,i})_{1,j} \\ \vdots \\ (CC_{g,i})_{N,j} \\ \vdots \\ (CC_{g,i})_{N,N} \end{bmatrix} \dots \begin{bmatrix} (CC_{g,i})_{1,j} \\ \vdots \\ (CC_{g,i})_{N,N} \end{bmatrix} \dots \begin{bmatrix} (CC_{g,i})_{1,N} \\ \vdots \\ (CC_{g,i})_{N,N} \end{bmatrix} \begin{bmatrix} \delta P_g \\ \vdots \\ \delta P_l \\ \vdots \\ \delta P_d \end{bmatrix} = \begin{bmatrix} (BB_g)_1 \\ \vdots \\ (BB_g)_i \\ \vdots \\ (BB_g)_N \end{bmatrix}$$



GPU Acceleration of Preconditioned Krylov Solver

comm(x₀)

- $r_0 = b - Ax_0$ Matmul
- Choose an arbitrary vector r_0 such that $(r_0, r_0) \neq 0$, e.g., $r_0 = r_0$
- $\rho_0 = (r_0, r_0)$ InnerProduct
- $p_0 = r_0$
- For $i = 1, 2, 3, \dots$

comm(y)

- $y = K_1^{-1} K_1^{-1} p_{i-1}$ Preconditioner
- $v = Ay$ Matmul
- $\alpha = p_{i-1} \cdot (r_0, v)$ InnerProduct
- $h = x_{i-1} + \alpha y$ Elementwise vector add
- $s = r_{i-1} - \alpha v$
- If s is accurate enough then $x_i = h$ and quit
- $z = K_2^{-1} K_1^{-1} s$ Preconditioner
- $t = Az$ Matmul
- $\omega = (K_1^{-1} K_1^{-1} s) \cdot (K_1^{-1} t, K_1^{-1} t)$ InnerProduct
- $x_i = h + \alpha z$
- $r_i = s - \alpha t$ Elementwise vector add
- If r_i is accurate enough then quit
- $\rho_i = (r_i, r_i)$ InnerProduct
- $\beta = (\rho_i / \rho_{i-1}) (\alpha / \omega)$ Elementwise vector add
- $p_i = r_i + \beta (p_{i-1} - \alpha v)$ Projection

Matmul

Sparse matrix vector mul.

`khh_cuspmv(Ax, matA, x, handle_spmv, dBuffer, n);
↳ cusparseSpMV(handle_spmv, ..., matA, ..., dBuffer);`

InnerProduct

Reduce sum

`cublasDdot(handle_dotp, n, r, 1, r, 1, &rho_a);`

Projection

Elementwise vector add

`khh_cuvecadd<<blocksPerGrid, threadsPerBlock>>();`

Preconditioner

(1) Diagonal, (2) ILU, (3) AMG

(1) Diagonal

```
global void khh_cuvecmul(double *c, double *a, double *b, int n)
{
    int idx = threadIdx.x + blockDim.x * blockIdx.x;
    if (idx < n) {
        c[idx] = a[idx] * b[idx];
    }
}
```

The diagonal (Jacobi) preconditioner scales each residual component by the inverse of the corresponding diagonal entry. This **element-wise operation** is embarrassingly parallel, ideal for GPU execution with a simple CUDA kernel.

(2) ILU

```
khh_cuspsv_init(handleL, matL, dBufferL, spsvDescrL,
    CUSPARSE_FILL_MODE_LOWER, CUSPARSE_DIAG_TYPE_UNIT,
    L_rpntr, L_cinx, L_val, Ux, x, nnzL, n);
khh_cuspsv_init(handleU, matU, dBufferU, spsvDescrU,
    CUSPARSE_FILL_MODE_UPPER, CUSPARSE_DIAG_TYPE_NON_UNIT,
    U_rpntr, U_cinx, U_val, Ux, x, nnzU, n);

khh_cuspsv(Ux, matL, b, handleL, dBufferL, spsvDescrL, n);
khh_cuspsv(X, matU, Ux, handleU, dBufferU, spsvDescrU, n);
↳ cusparseSpSV_solve(handle, ..., spsvDescr);
```

The ILU(0) preconditioner requires incomplete LU factorization followed by forward/backward substitutions. The factorization is performed on CPU due to sequential dependencies, while **triangular solves run on GPU via cusparseSpSV**.

(3) AMG

```
AMGX_solver_create(&solver, rsrc, AMGX_mode_dDDI, cfg);
AMGX_solver_setup(solver, Amatrix);
AMGX_solver_solve(solver, rhs, sol);
```

The algebraic multigrid (AMG) preconditioner constructs a hierarchy of coarser grids to accelerate convergence. **NVIDIA's AmGX library** provides GPU-optimized AMG coupled with GMRES solver.

(2) Execution Time Breakdown and Speedup

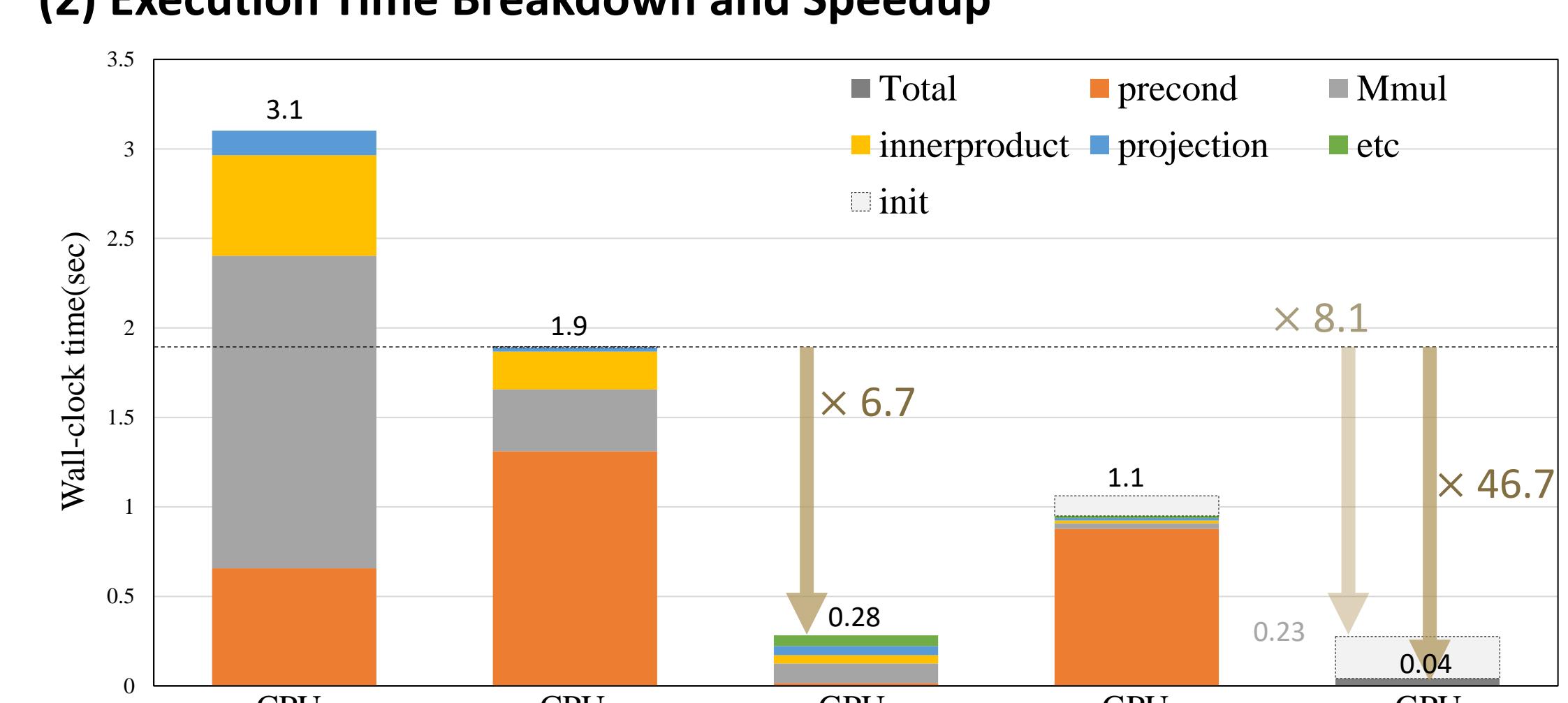


Fig2. Breakdown of execution time by solver component on CPU and GPU, with speedup ratios indicated.

Figure 2 shows the execution time breakdown. The **GPU diagonal solver achieves 6.7x speedup** with high parallel efficiency. The GPU ILU shows limited speedup due to sequential dependencies in sparse triangular solves (cusparseSpSV). The **AMG solver is extremely fast** (0.04 sec solve time), but since the pressure matrix changes every time step, the **hierarchy must be rebuilt each time step**, making total runtime comparable to diagonal.

Conclusion

- The pressure solver in the thermal-hydraulic code was successfully accelerated on a single GPU using cuSPARSE and cuBLAS libraries, achieving significant speedup over the 64-core CPU implementation.
- The **diagonal preconditioner** demonstrated the best balance between convergence rate and computational efficiency, providing the highest speedup due to its simple and fully parallel structure.
- The **ILU(0) preconditioner** reduces iteration count but is limited by sparse triangular solve efficiency. The **AMG solver** converges extremely fast but requires costly setup at each time step when matrix coefficients change.

This research is supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MIST) (RS-2023-00282764). This work is also supported by Global top research lab funded by National research council of science and technology (GTL24031-000).