

Objective: Investigate the performance and scalability of a multigrid pressure Poisson equation solver running on a GPU cluster.

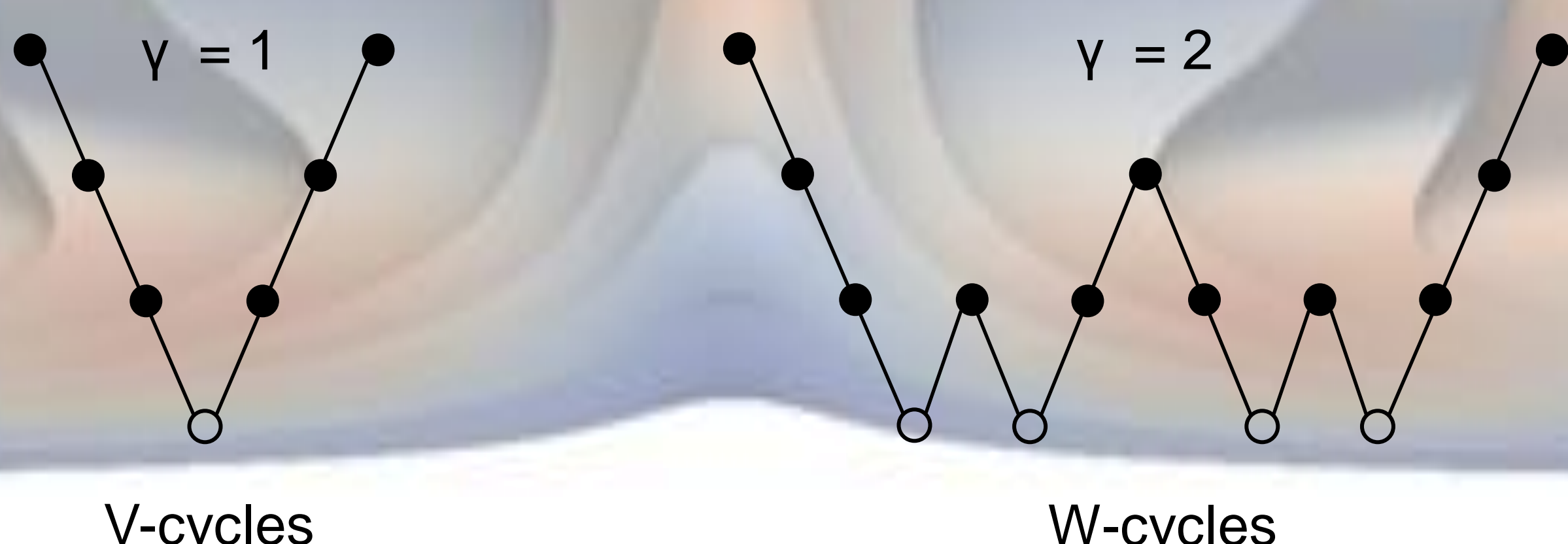
Method: A parallel 3D multigrid pressure solver was written for GIN3D, a 3D incompressible Navier-Stokes flow solver which runs on GPU clusters. Tests were performed using the well-known lid-driven cavity and natural convection in a cavity problems.

Solver:

- For pressure Poisson equation in incompressible Navier-Stokes flow solver.
- Typically consumes a majority of the time.
- Simple iterative solvers such as Jacobi and Gauss-Seidel can run very efficiently on the GPU, but solutions converge very slowly.
- Multigrid methods can converge rapidly and allow convergence to be relatively independent from the grid size, which is important for large computational meshes.

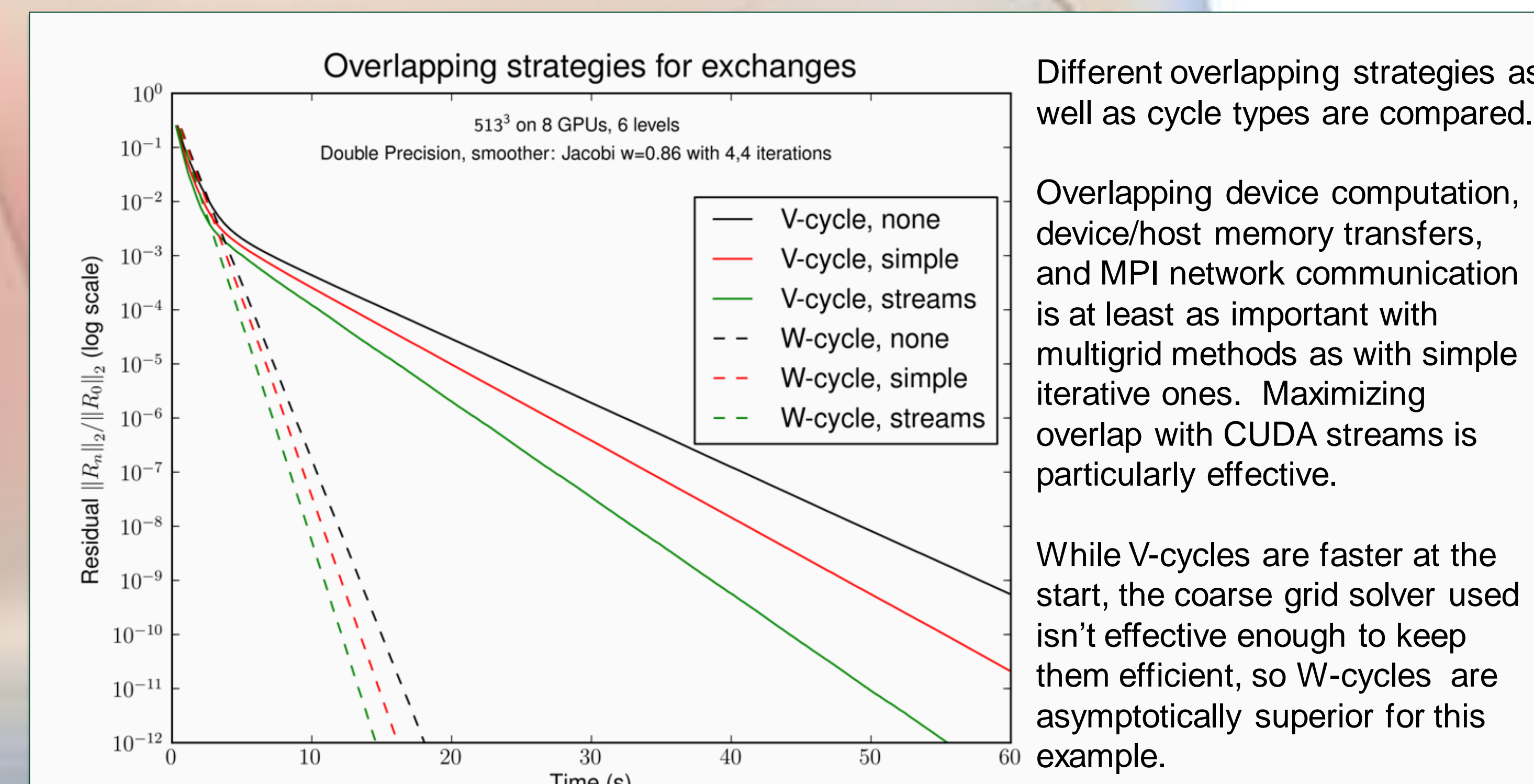
Multigrid Algorithm

- $\text{cycle}(\gamma, u_k, f, v_1, v_2)$
 - Smooth v_1 times
 - Compute residual $r_k = f - \mathcal{L} u_k$
 - Restrict residual $r_{k-1} = \mathcal{R} r_k$
 - Compute approximate solution v_{k-1} :
 - “Direct Solve” if lowest level, otherwise
 - Repeat γ times: $\text{cycle}(\gamma, 0, r_{k-1}, v_1, v_2)$
 - Prolongate correction: $u_k = u_k + \mathcal{P} v_{k-1}$
 - Smooth v_2 times



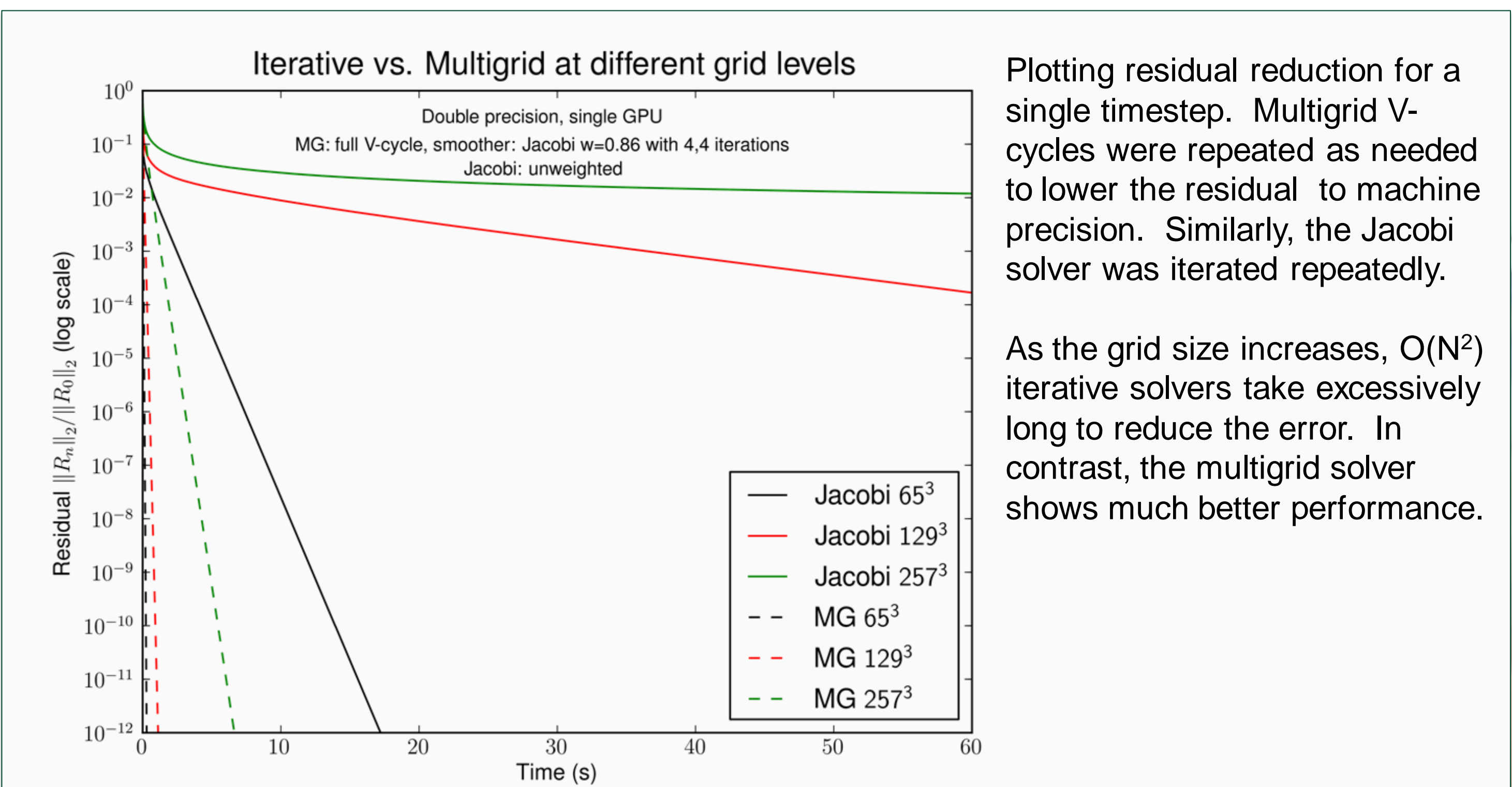
Question:

Which of these parameters are most important to tune? *Truncation level, cycle type, smoother type, method for coarsest grid solve.*



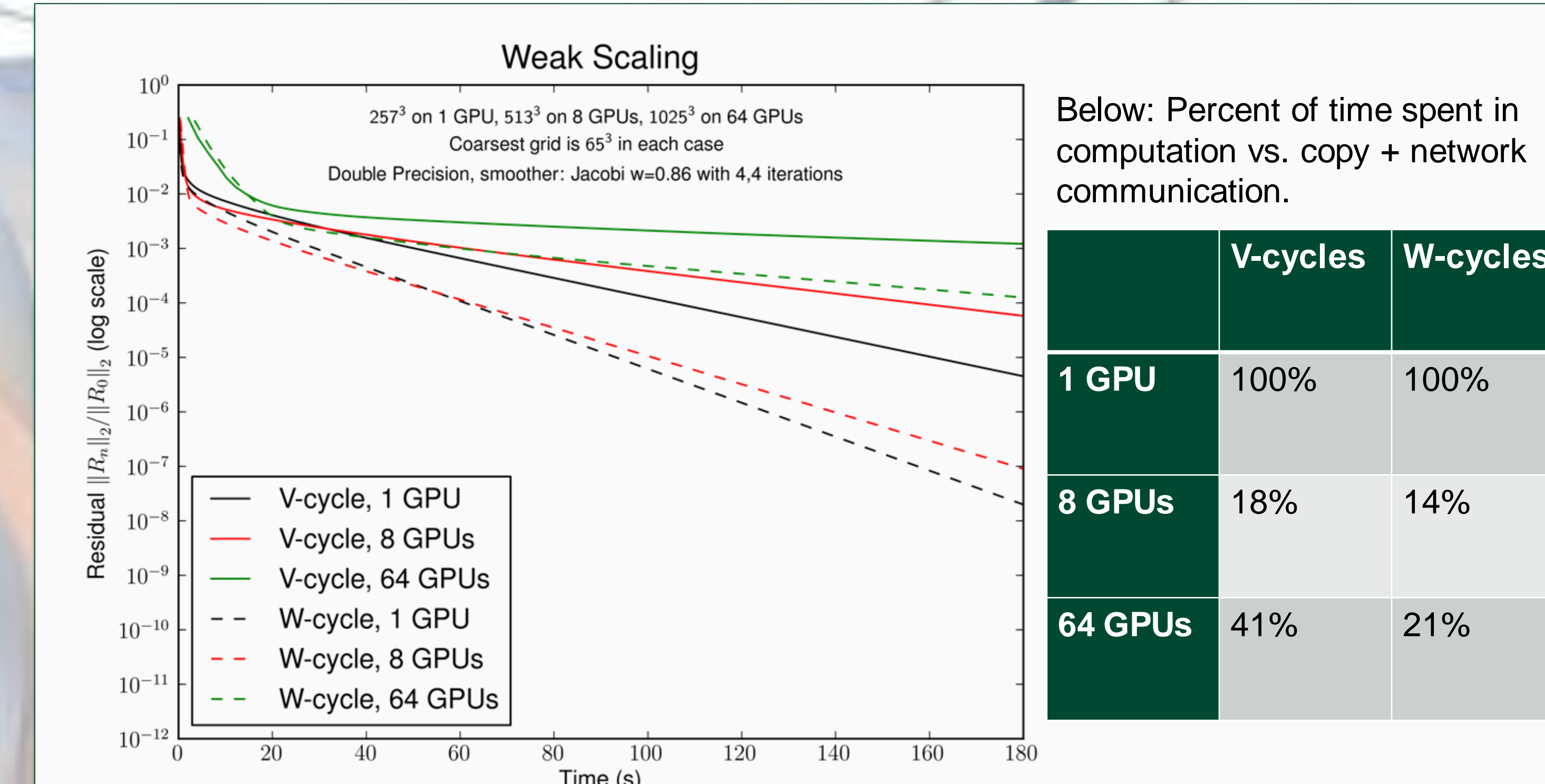
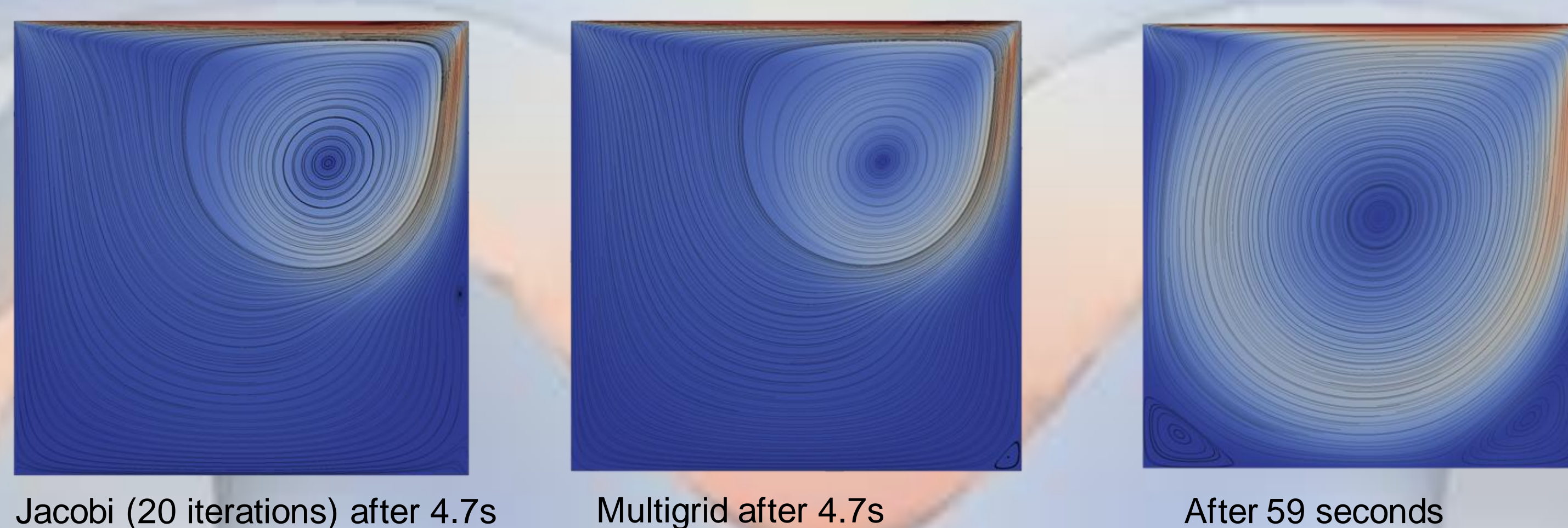
Details of the overlapping strategies are shown in the 2009 NVIDIA GTC poster and 2010 AIAA paper by the authors.

- No overlap: computation followed by synchronous exchange.
- Simple overlap: Compute edges, do async MPI exchange while computing middle.
- Streams overlap: Use CUDA Streams to asynchronously overlap computation, host / device memory transfers, and network communication.

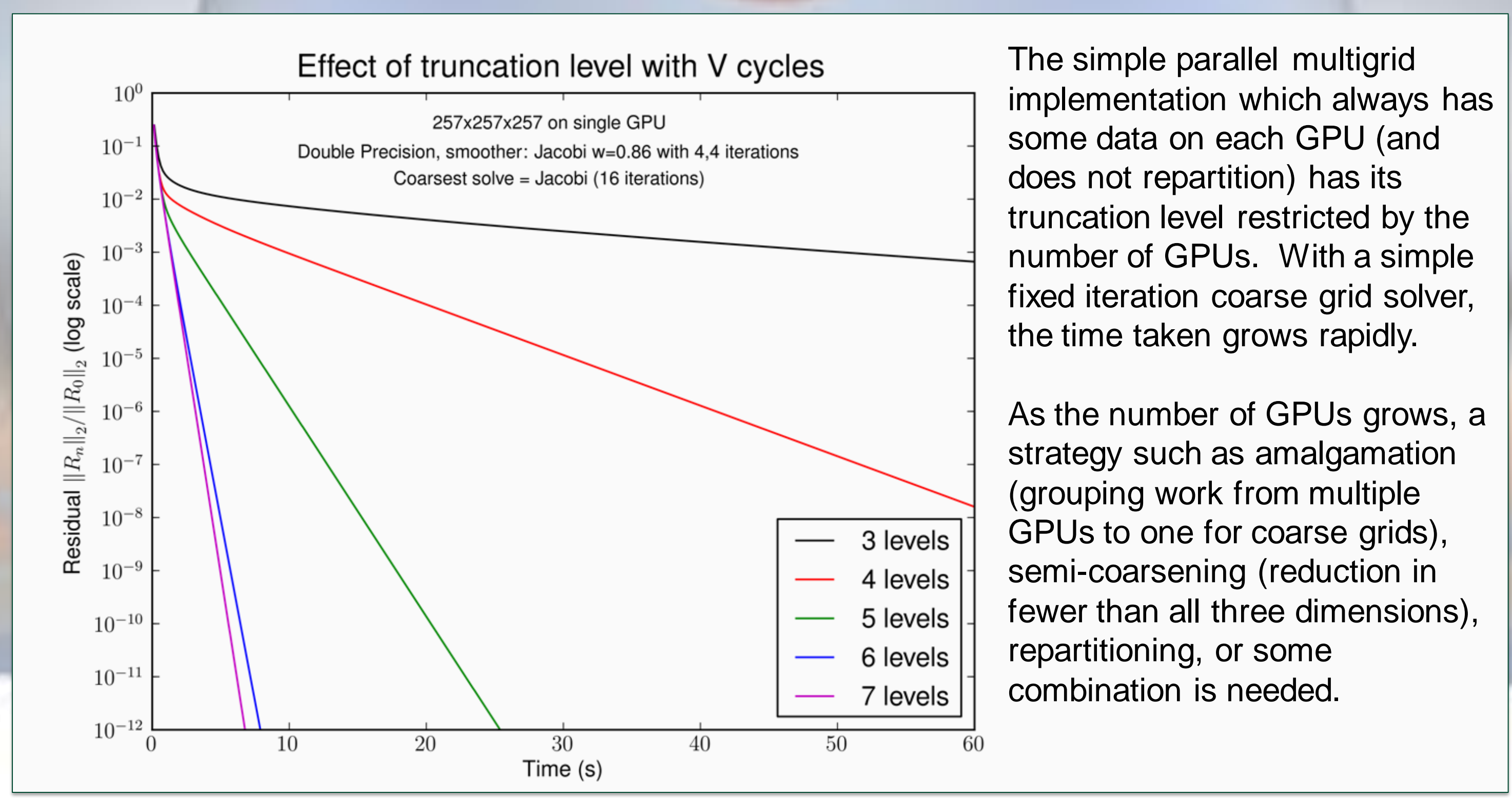
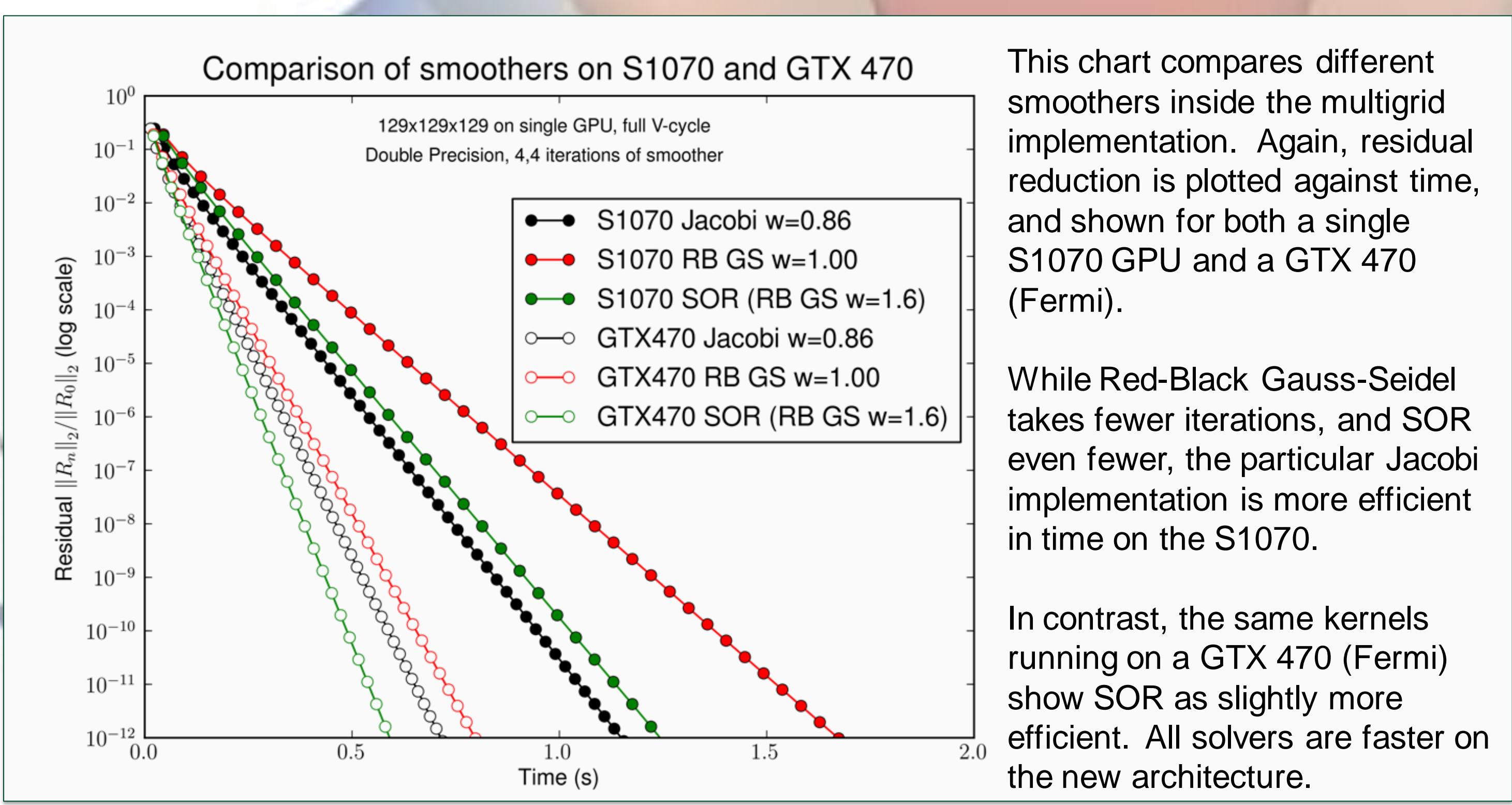


Multigrid Implementation

Four kernels: Laplacian, Restriction (full-weighting), Smooth (Jacobi or Red-Black Gauss-Seidel), and Prolongation.



Changing the problem size for weak scaling results in significantly different convergence behavior with truncated multigrid. The weak scaling analysis of multigrid is much more complex than for iterative solvers.



Conclusions:

- Multigrid on GPUs can be effective and efficient.
- Multigrid performance significantly benefits from deeper truncation levels.
- The Fermi architecture makes a difference in relative smoother performance.
- Coarsest level operations have a large impact on cluster scalability.