

Addressing the Communication Bottleneck: **Towards a Modular Precision Ecosystem for High Performance Computing**

A Structured Approach for Programming Extreme-Scale and Heterogeneous Systems (EPIGRAM-HS project)
Zurich, September 9th, 2019

Hartwig Anzt, Terry Cojean, Goran Flegar, Thomas Grützmacher, Pratik Nayak, Tobias Ribizel
Steinbuch Centre for Computing (SCC)



Terry Cojean



Goran Flegar



Thomas
Grützmacher



Pratik Nayak



Tobias Ribizel

Where do we stand?



- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s** (**10^{18}**) for 16 bit floating point used in data analytics and artificial intelligence



Copyright@ORNL

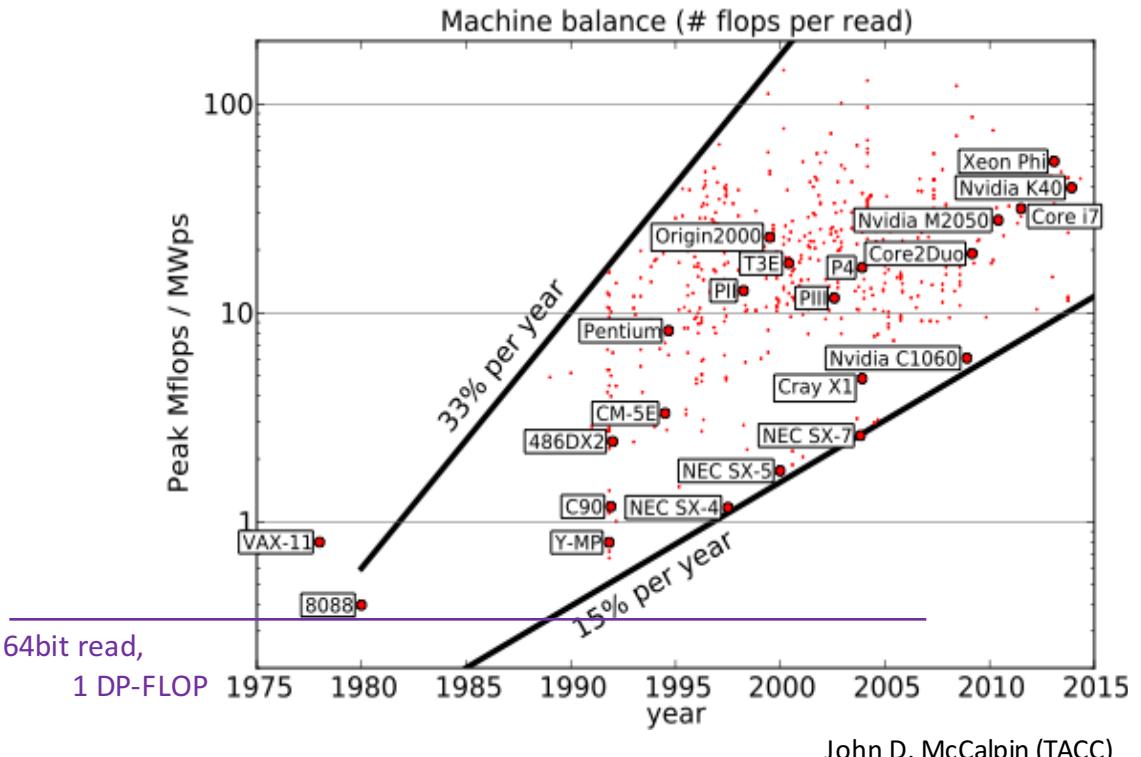
Where do we stand?



- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s** (10^{18}) for 16 bit floating point used in data analytics and artificial intelligence



Copyright@ORNL



- Compute power (#FLOPs) grows much faster than bandwidth.

"Operations are free, mem access and comm is what counts."

Where do we stand?



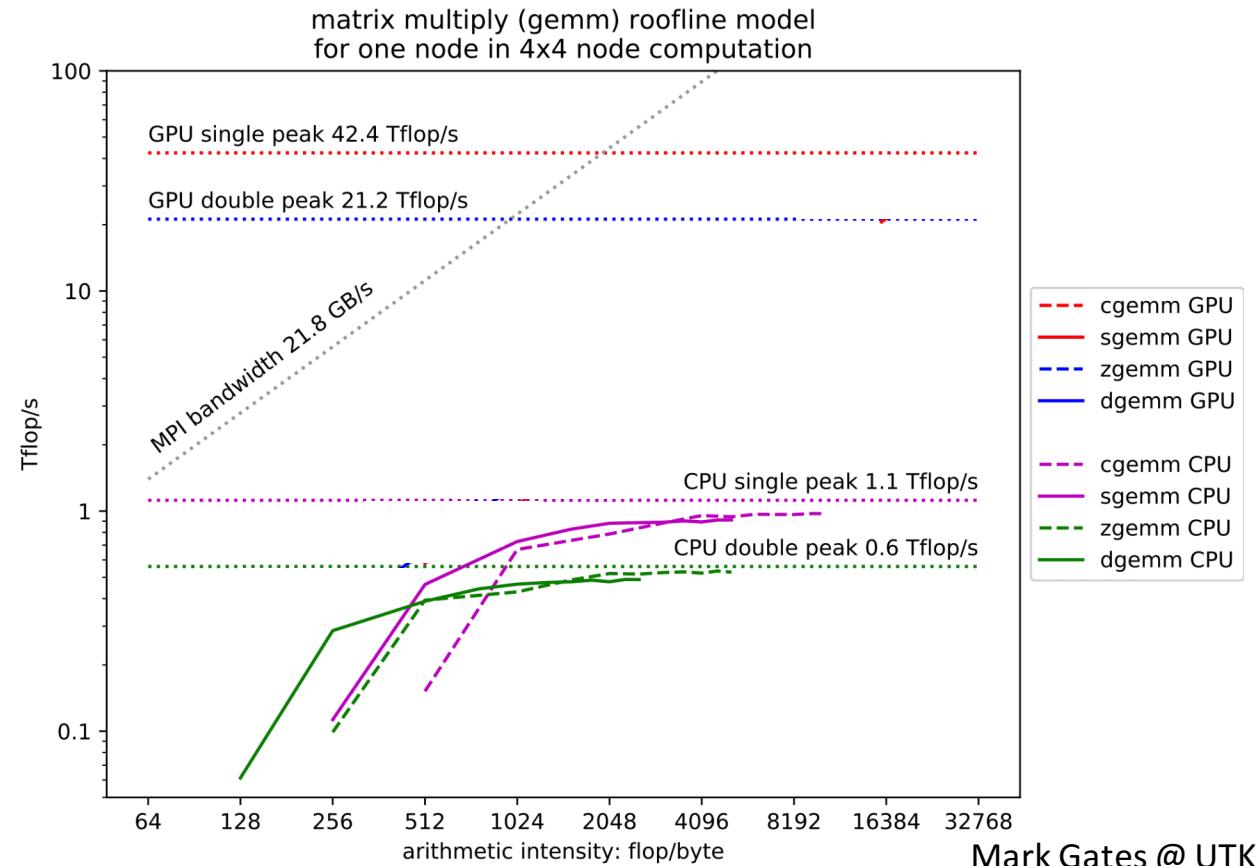
- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s (10^18)** for 16 bit floating point used in data analytics and artificial intelligence



Copyright@ORNL

Dense Matrix Operations?

- The inter-node communication is the limiting resource;
- Each node has more computational power than what we can leverage;



Where do we stand?



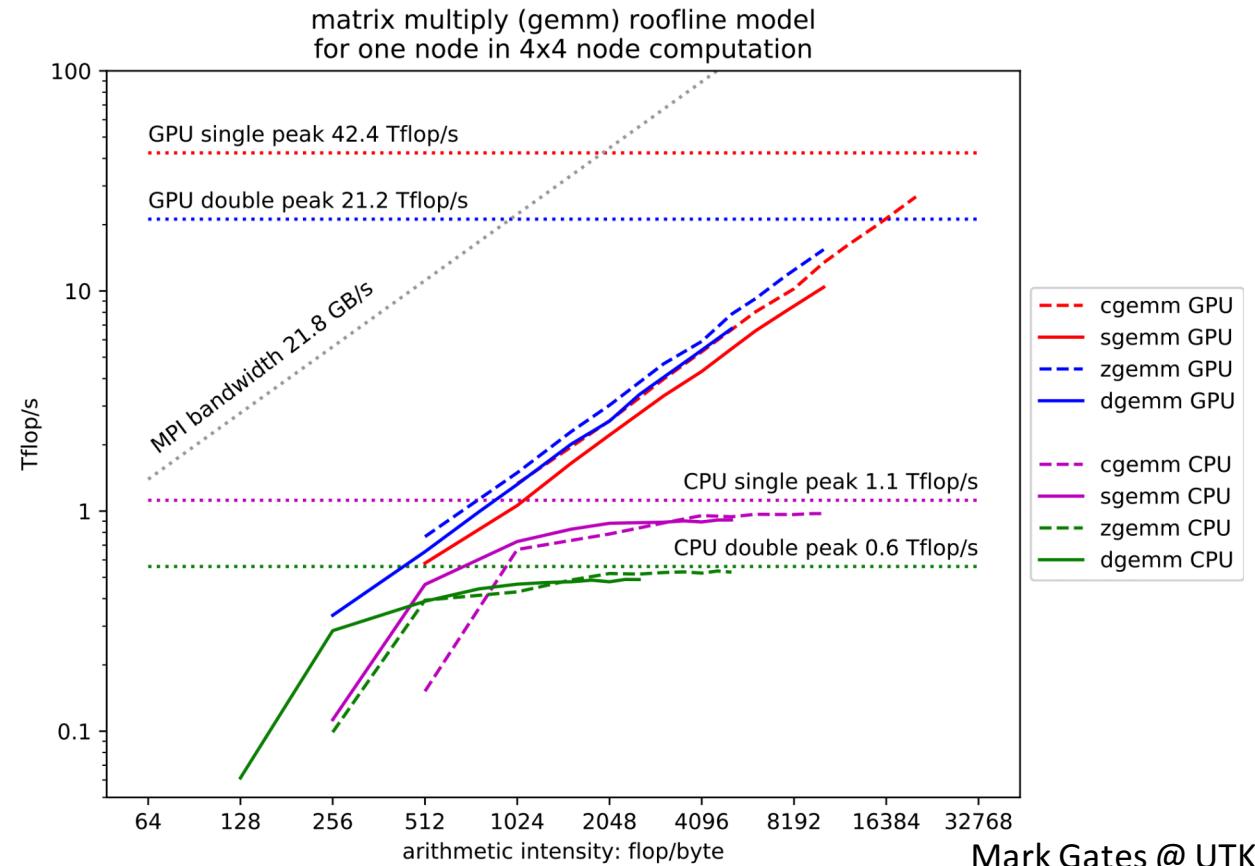
- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s (10^18)** for 16 bit floating point used in data analytics and artificial intelligence



Copyright@ORNL

Dense Matrix Operations?

- The inter-node communication is the limiting resource;
- Each node has more computational power than what we can leverage;



Where do we stand?



- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s** (**10^{18}**) for 16 bit floating point used in data analytics and artificial intelligence



Copyright@ORNL

Dense Matrix Operations?

- The inter-node communication is the limiting resource;
- Each node has more computational power than what we can leverage;

Sparse / Graph Problems?

- *Sparse Matrix Vector Product* (SpMV) is a central building block;

$$\begin{pmatrix} \times & \times & \times & \times & & \times & \times \\ \times & \times & \times & & & & \\ \times & \times & \times & \times & & & \\ \times & & \times & \times & \times & & \\ \times & & \times & \times & \times & & \\ \times & & \times & \times & \times & \times & \\ \times & & & \times & \times & & \\ \times & & & \times & \times & & \end{pmatrix} \cdot \begin{pmatrix} \times \\ \times \end{pmatrix}$$

Where do we stand?



- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s** (**10^{18}**) for 16 bit floating point used in data analytics and artificial intelligence



Copyright@ORNL

Dense Matrix Operations?

- The inter-node communication is the limiting resource;
- Each node has more computational power than what we can leverage;

Sparse / Graph Problems?

- ***Sparse Matrix Vector Product*** (SpMV) is a central building block;
- It looks like for the test problems at the SuiteSparse Matrix Collection¹, a Multi-node SpMV is slower than a Single-node SpMV;
- The inter-node communication is an order of magnitude slower than the local computations.

$$\begin{pmatrix} \times & \times & \times & \times & & \times & \times \\ \times & \times & \times & & & & \\ \times & \times & \times & \times & & & \\ \times & & \times & \times & \times & & \times \\ & & & \times & \times & \times & \\ \times & & \times & \times & \times & \times & \times \\ \times & & & \times & \times & & \times \\ & & & & \times & \times & \\ \times & & & & & \times & \times \end{pmatrix} \cdot \begin{pmatrix} \times \\ \times \end{pmatrix}$$

¹SuiteSparse Matrix Collection: <https://sparse.tamu.edu/>

Where do we stand?



- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s** (**10^{18}**) for 16 bit floating point used in data analytics and artificial intelligence



Radically decouple storage format from arithmetic format.

Dense Matrix Operations?

- The inter-node communication is the limiting resource;
- Each node has more computational power than what we can leverage;

Sparse / Graph Problems?

- *Sparse Matrix Vector Product* (SpMV) is a central building block;
- Like for the test problems at the SuiteSparse Matrix Collection¹, a Multi-node SpMV is slower than a Single-node SpMV;

The inter-node communication is an order of magnitude slower than the local computations.

$$\begin{pmatrix} \times & \times & \times & \times & & \times & \times \\ \times & \times & \times & & & & \\ \times & \times & \times & \times & & & \\ \times & & \times & \times & \times & & \\ \times & & \times & \times & \times & & \\ \times & & \times & \times & \times & \times & \\ \times & & & \times & \times & & \\ \times & & & \times & \times & & \\ \end{pmatrix} \cdot \begin{pmatrix} \times \\ \times \end{pmatrix}$$

¹SuiteSparse Matrix Collection: <https://sparse.tamu.edu/>

Where do we stand?



- Node: 2 IBM POWER9 nodes
- 4,608 nodes, 9,216 IBM POWER9 cores
- 27,648 V100 GPUs (NVIDIA)
- Peak performance of 1 exaflop for simulation & simulation
- Peak performance of 1 exaflop for 16 bit floating point for analytics and artificial intelligence



Radically decouple storage format from arithmetic format.

- The arithmetic operations should use high precision formats natively supported by hardware.
- Data access should be as cheap as possible, reduced precision.

Challenges when using double precision
+ IEEE single/half precision:

- Need explicit conversion.
- Data range reduction: protect against under- / overflow.
- Need to duplicate data in memory (half/single/double).
- Better: Customized Precision Format

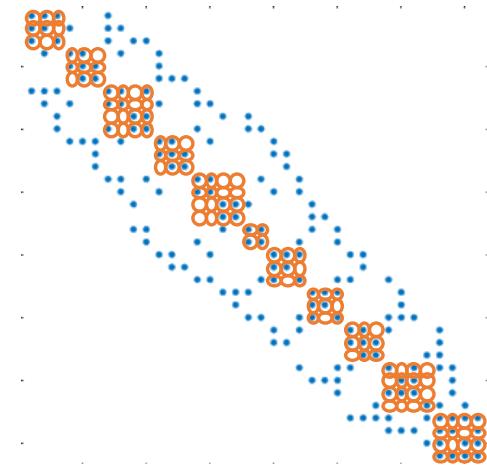
$$\begin{pmatrix} \times \\ \times \\ \times \\ \times \\ \times \\ \times \end{pmatrix} \cdot \begin{pmatrix} \times \\ \times \\ \times \\ \times \\ \times \\ \times \end{pmatrix}$$

¹SuiteSparse Matrix Collection: <https://sparse.tamu.edu/>

Use reduced precision for “approximate Operators”

- **Preconditioners for iterative solvers.**
 - Idea: Approximate inverse of system matrix to make the system “easier to solve”: $P^{-1} \approx A^{-1}$
 $\tilde{A} = P^{-1}A$, $\tilde{b} = P^{-1}b$, and we solve $Ax = b \Leftrightarrow \tilde{A}x = \tilde{b}$.
 - **Why should we store the preconditioner matrix P^{-1} in full (high) precision?**
 - **We have to ensure regularity!** (Reducing precision can turn matrix singular)
- **Jacobi method** based on **diagonal scaling**: $P = \text{diag}(A)$
- **Block-Jacobi** is based on **block-diagonal scaling**: $P = \text{diag}_B(A)$
 - Large set of small diagonal blocks.
 - Each block corresponds to one (small) linear system.
 - *Larger* blocks typically **improve convergence**.
 - *Larger* blocks make block-Jacobi **more expensive**.

Extreme case: one block of matrix size.



Block-Jacobi Preconditioning

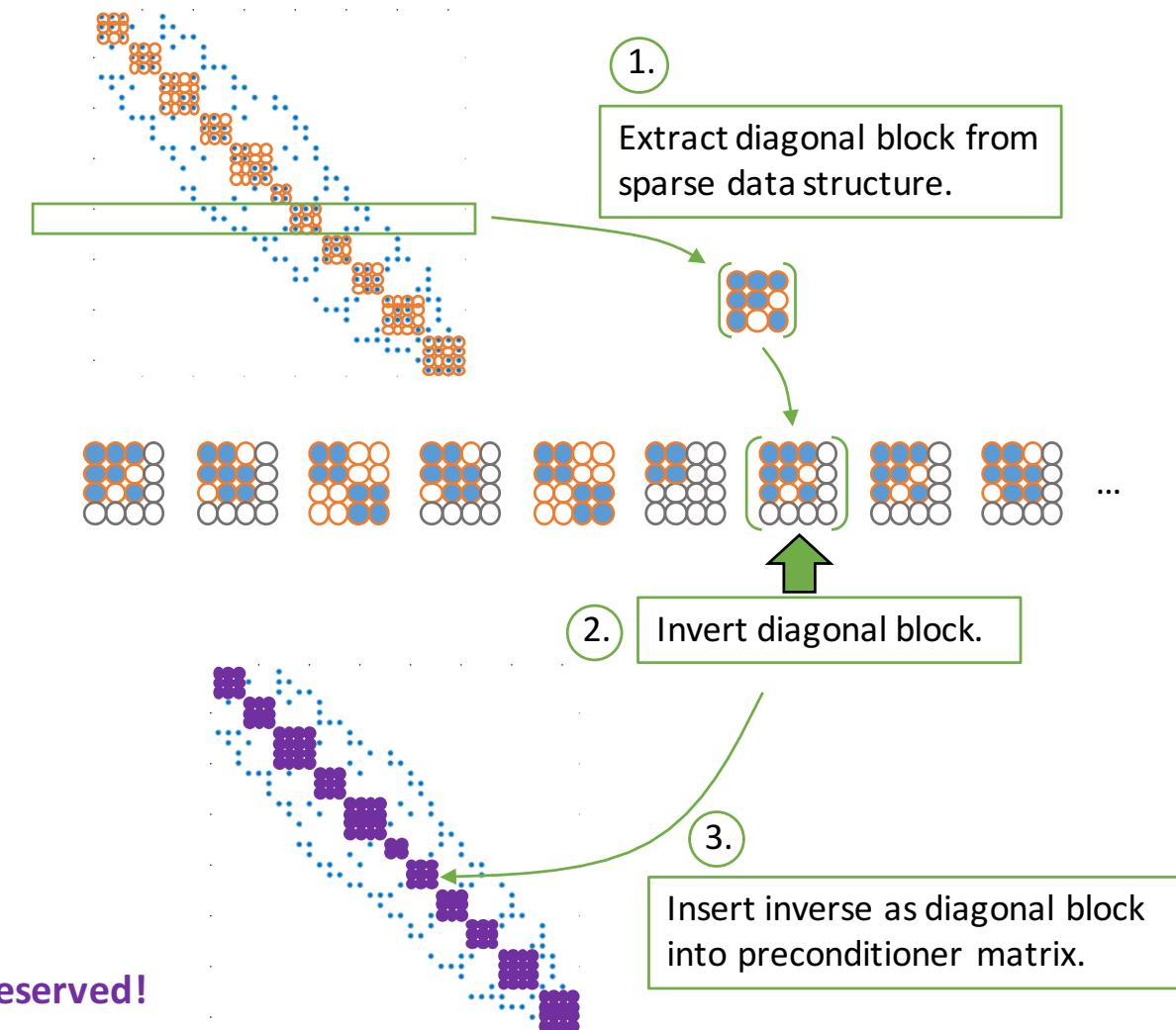
Preconditioner Setup:

- Identify the diagonal blocks $P = \text{diag}_B(A)$
- Form the block-Inverse $P^{-1} \approx A^{-1}$

Preconditioner Application:

- Apply the preconditioner in every solver iteration via:

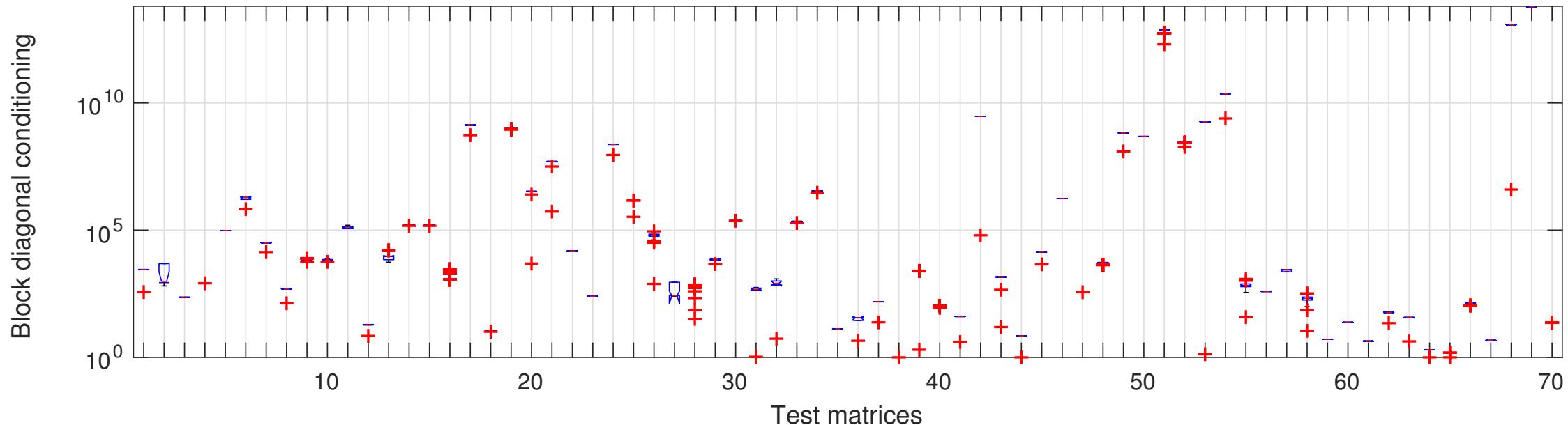
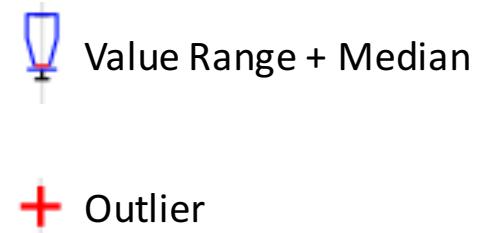
$$y := P^{-1}x$$



We can store diagonal blocks in lower precision, if regularity is preserved!

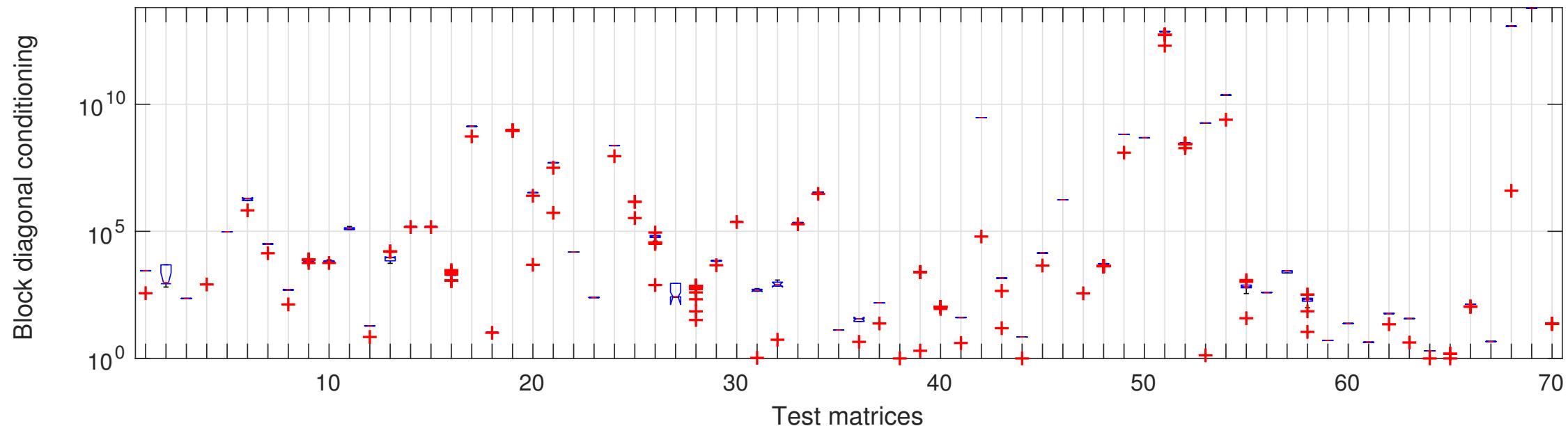
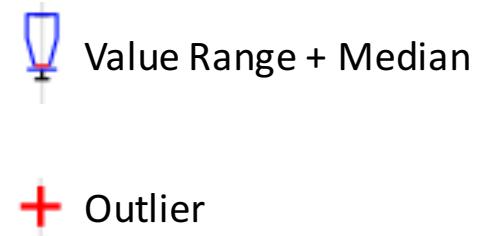
Adaptive Precision Block-Jacobi Preconditioning

- 70 matrices from the SuiteSparse Matrix Collection
- Use block-size 24 with Super-Variable agglomeration (24 is upper bound for size of blocks)
- Report conditioning of all arising diagonal blocks

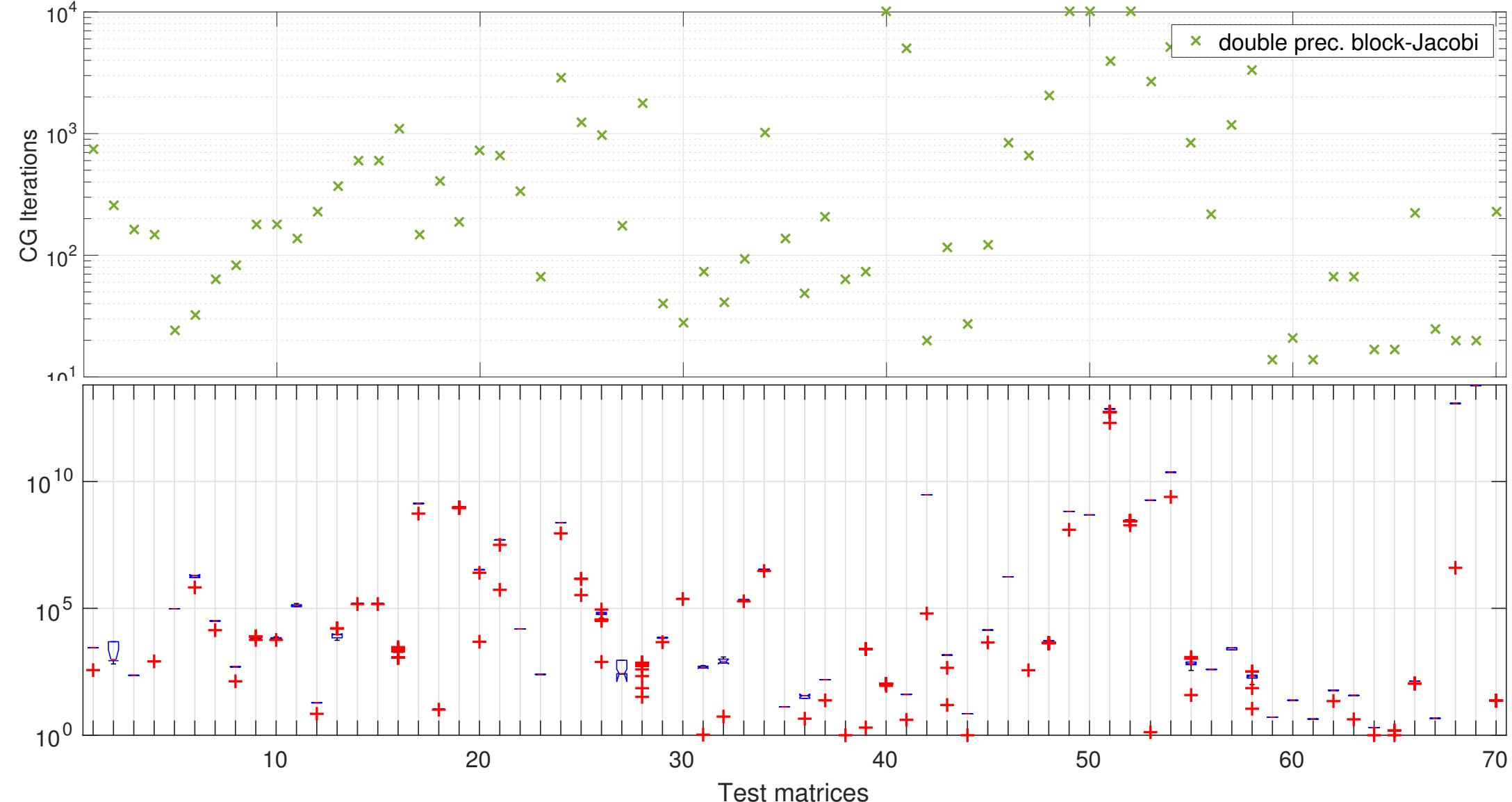


Adaptive Precision Block-Jacobi Preconditioning

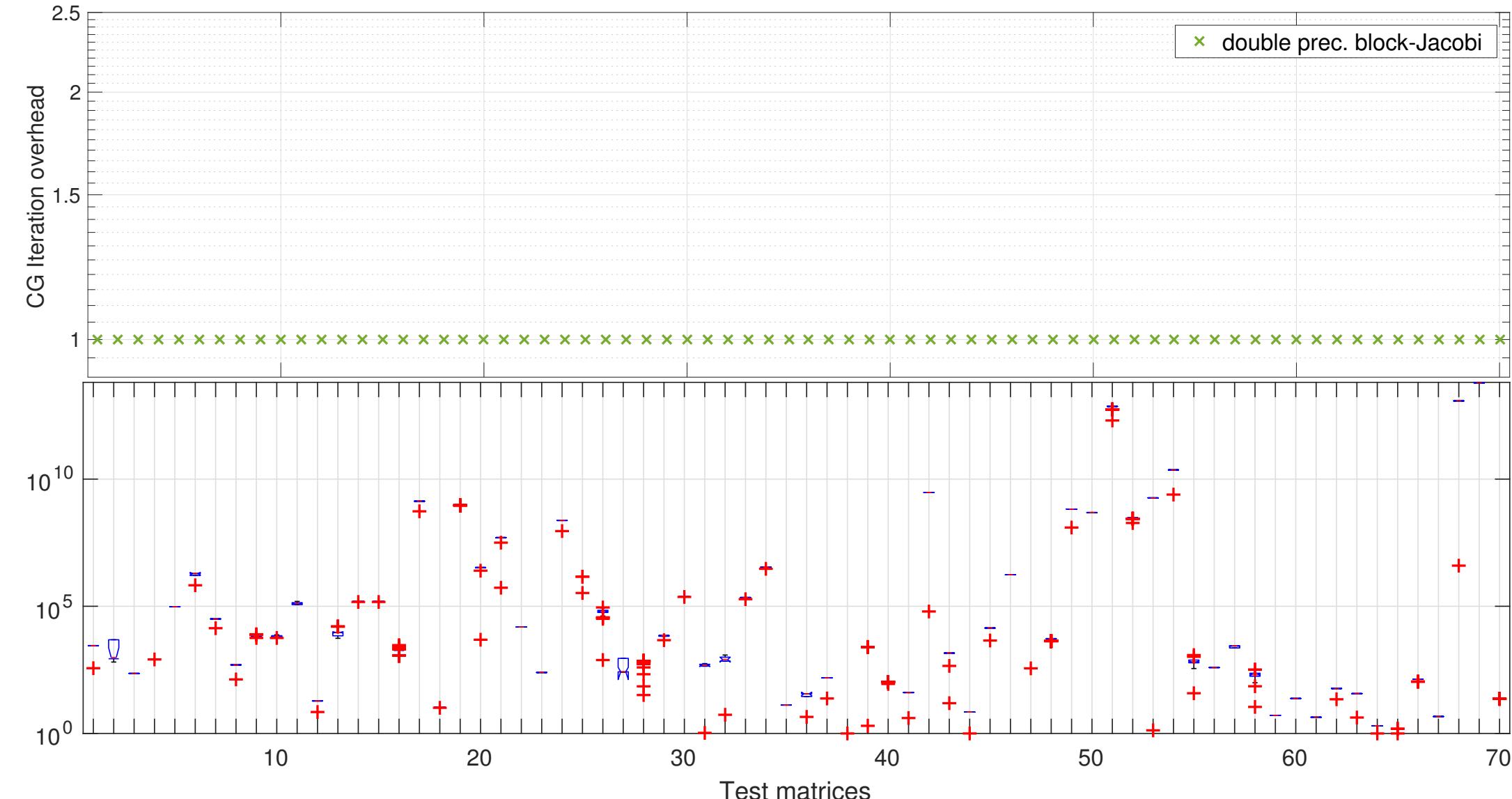
- 70 matrices from the SuiteSparse Matrix Collection
- Use block-size 24 with Super-Variable agglomeration (24 is upper bound for size of blocks)
- Report conditioning of all arising diagonal blocks
- Analyze the impact of storing block-Jacobi in lower precision a top-level Conjugate Gradient solver (CG)



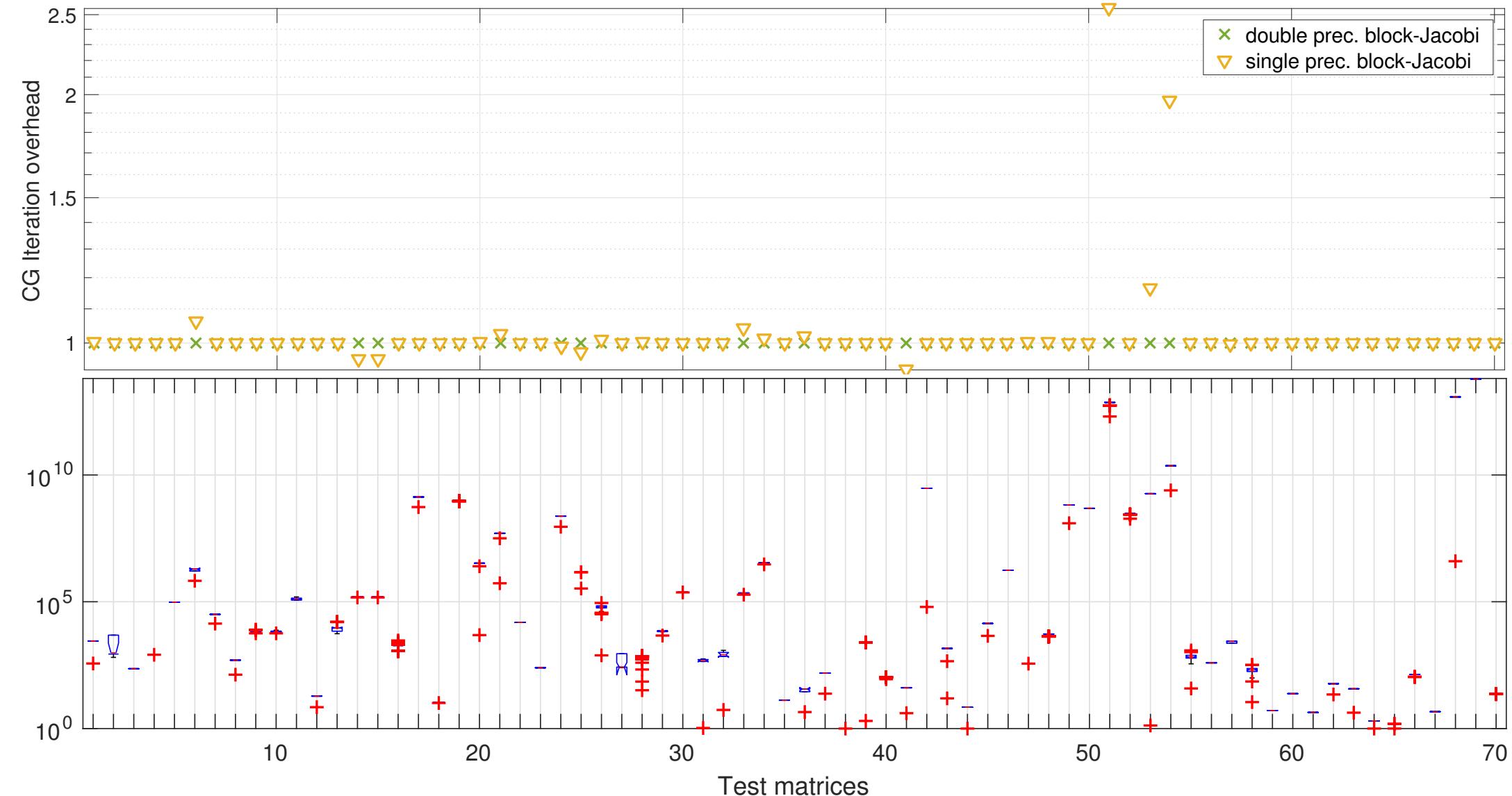
Adaptive Precision Block-Jacobi Preconditioning



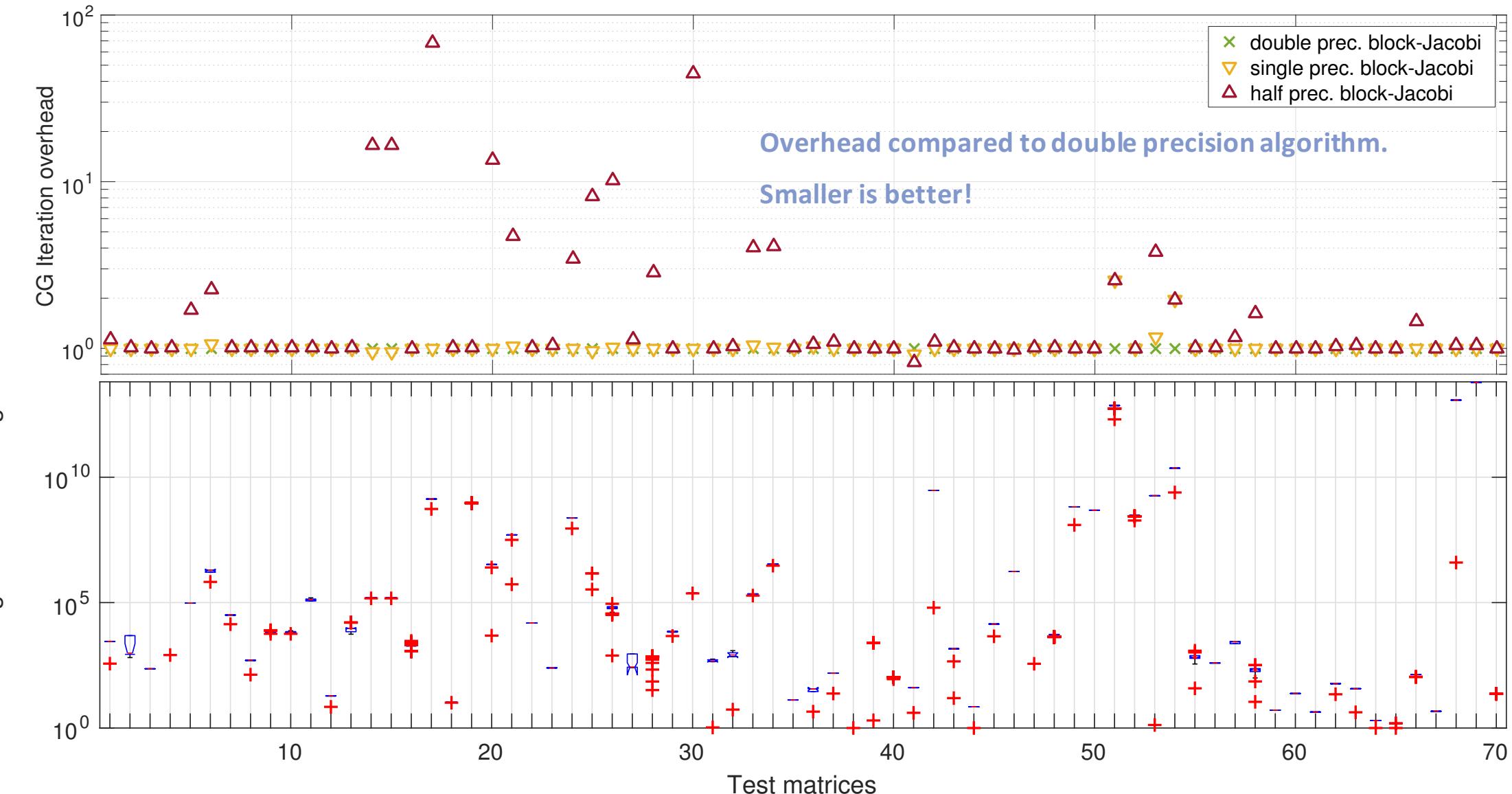
Adaptive Precision Block-Jacobi Preconditioning



Adaptive Precision Block-Jacobi Preconditioning



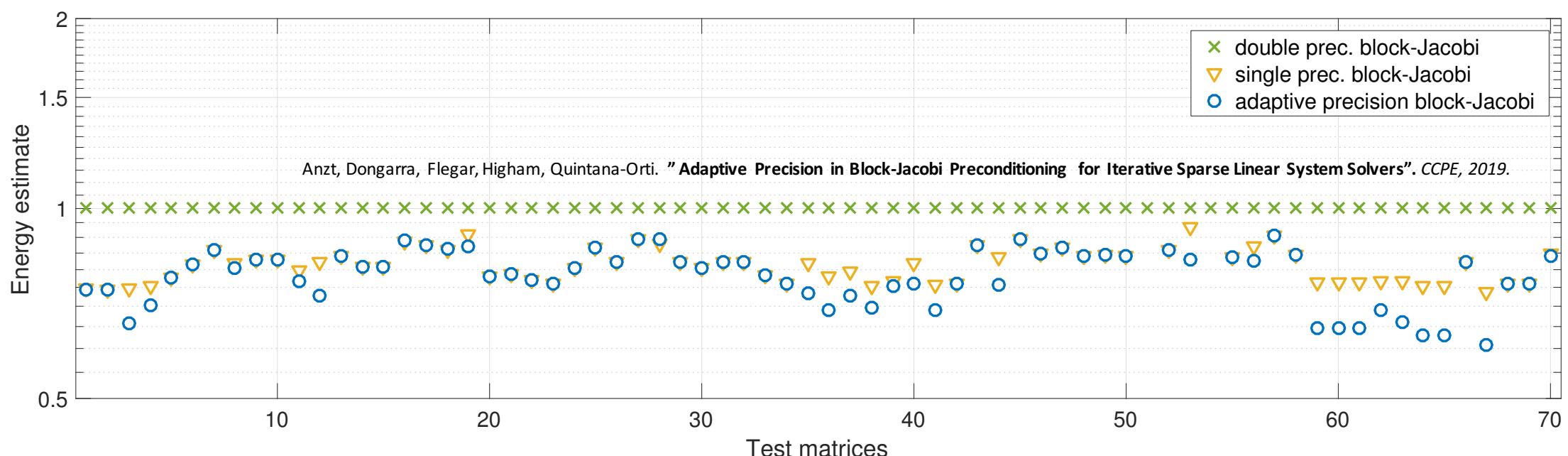
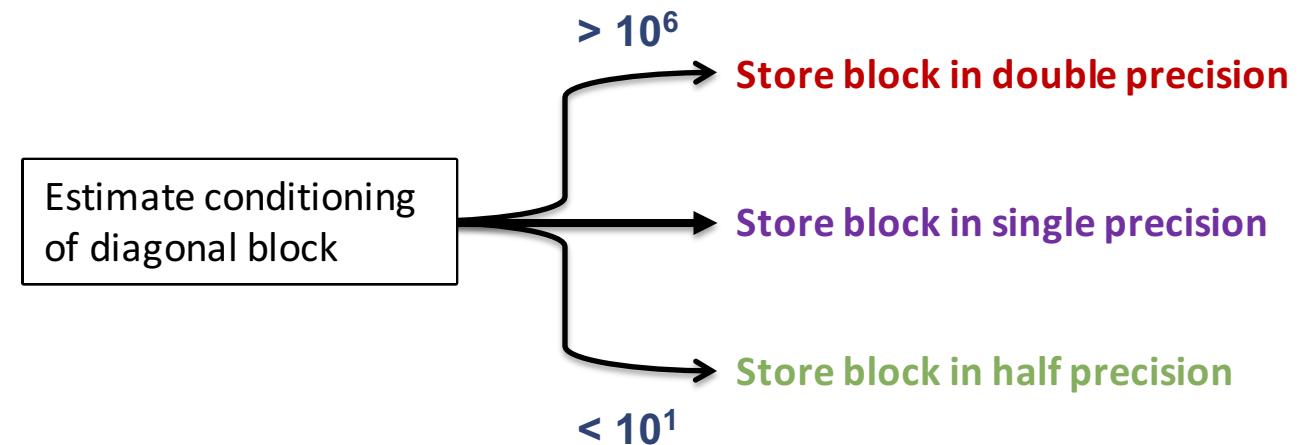
Adaptive Precision Block-Jacobi Preconditioning



Adaptive Precision Block-Jacobi Preconditioning

Multi-Precision Idea:

- All computations use double precision!
- Store distinct blocks in different formats
- Use single precision as standard storage format
- Where necessary: switch to double
- For well-conditioned blocks use half precision



Adaptive Precision Block-Jacobi Preconditioning

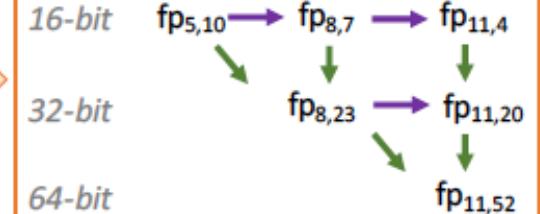
Multi-Precision Idea:

- All computations use double precision!
- Depart from the rigid IEEE precision formats!
- Preserve either 1 or 2 digits accuracy of the inverted diagonal blocks.

Invert the diagonal block using Gauss-Jordan elimination.

Compute condition number and exponent range.

Select storage format:

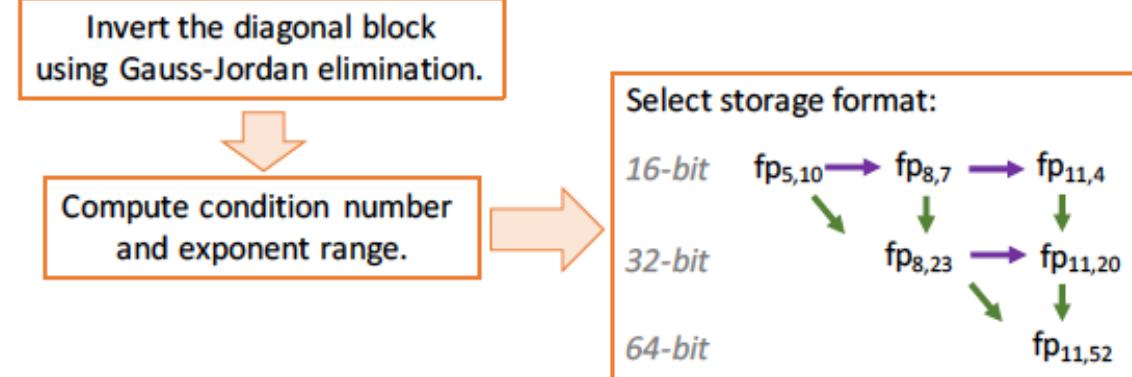


Flegar, Anzt, Quintana-Orti. "Customized-Precision Block-Jacobi Preconditioning for Krylov Iterative Solvers on Data-Parallel Manycore Processors". TOMS, submitted.

Adaptive Precision Block-Jacobi Preconditioning

Multi-Precision Idea:

- All computations use double precision!
- Depart from the rigid IEEE precision formats!
- Preserve either 1 or 2 digits accuracy of the inverted diagonal blocks.



- ✓ Regularity preserved;
- ✓ No flexible Krylov solver needed
 - (Preconditioner constant operator);
- ✓ Can handle non-spd problems
 - (inversion features pivoting);
- ✓ Preconditioner for any iterative preconditionable solver;
- Speedups / preconditioner quality **problem-dependent**;
- **Overhead** of the **precision detection**
 - (condition number calculation);
- **Overhead** from storing **precision information**
 - (need to additionally store/retrieve flag);

Flegar, Anzt, Quintana-Orti. "Customized-Precision Block-Jacobi Preconditioning for Krylov Iterative Solvers on Data-Parallel Manycore Processors". *TOMS, submitted.*

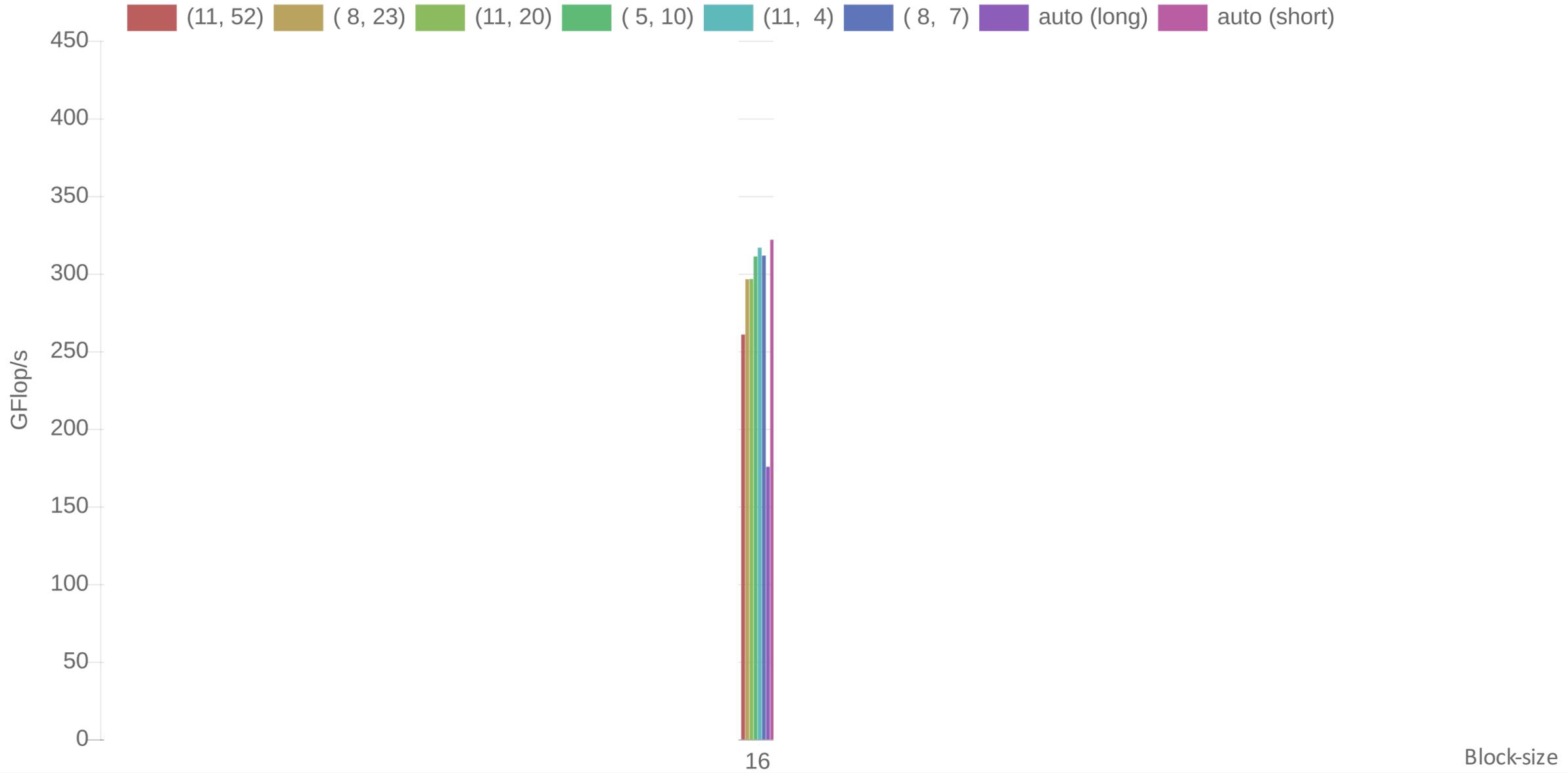
Precision distribution in Adaptive Block-Jacobi



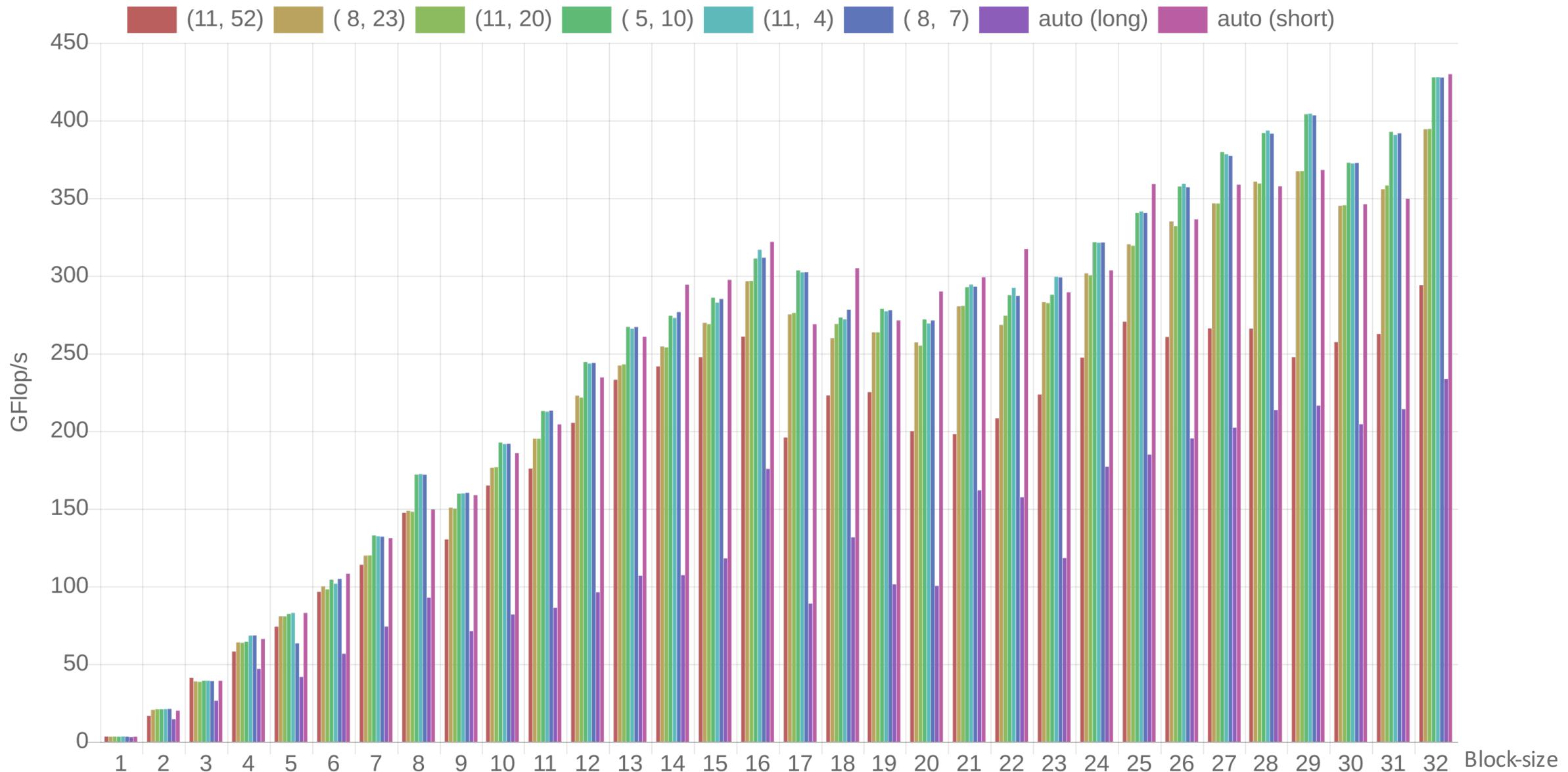
Precision distribution in Adaptive Block-Jacobi



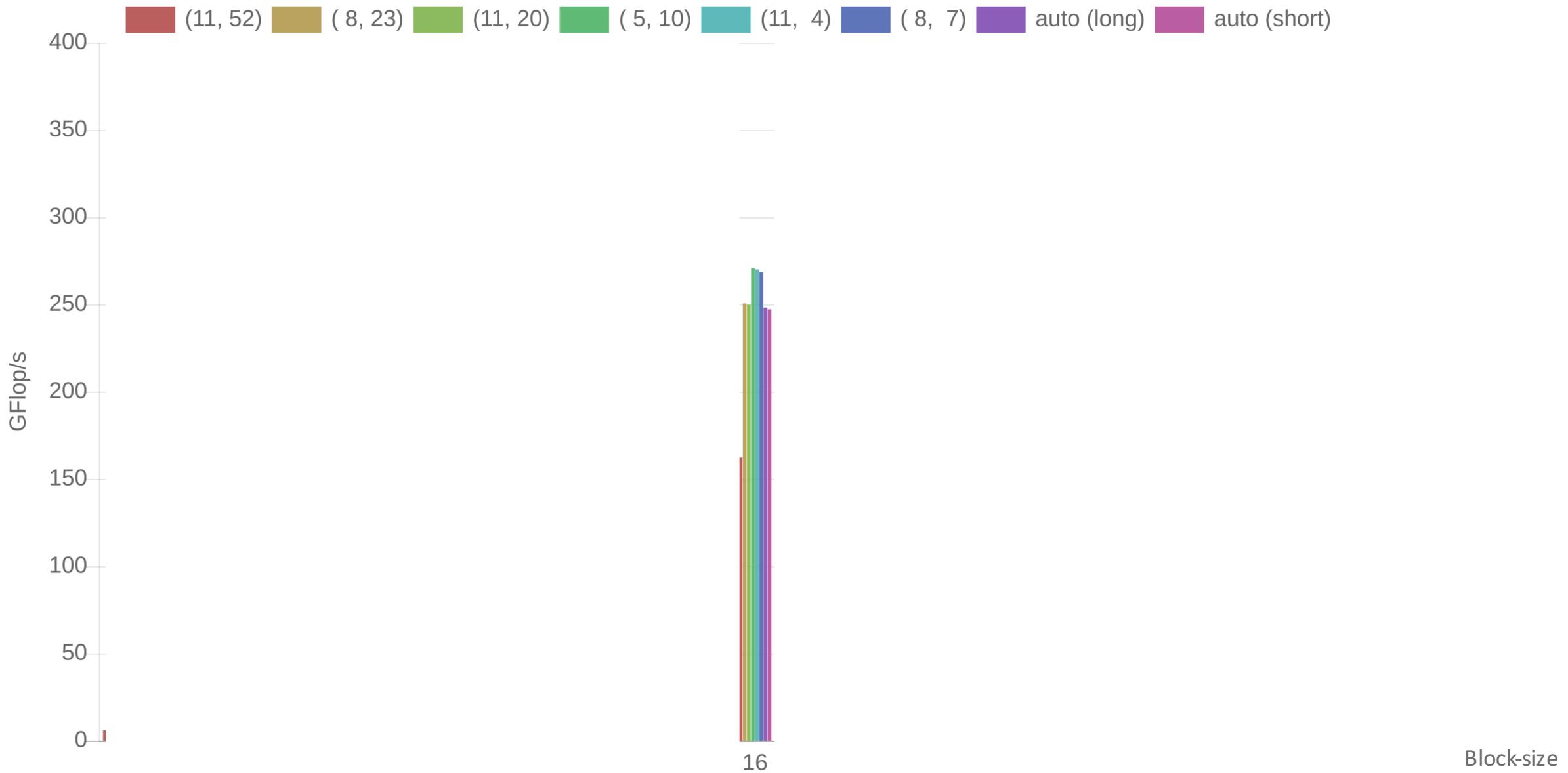
Adaptive Block-Jacobi Generation



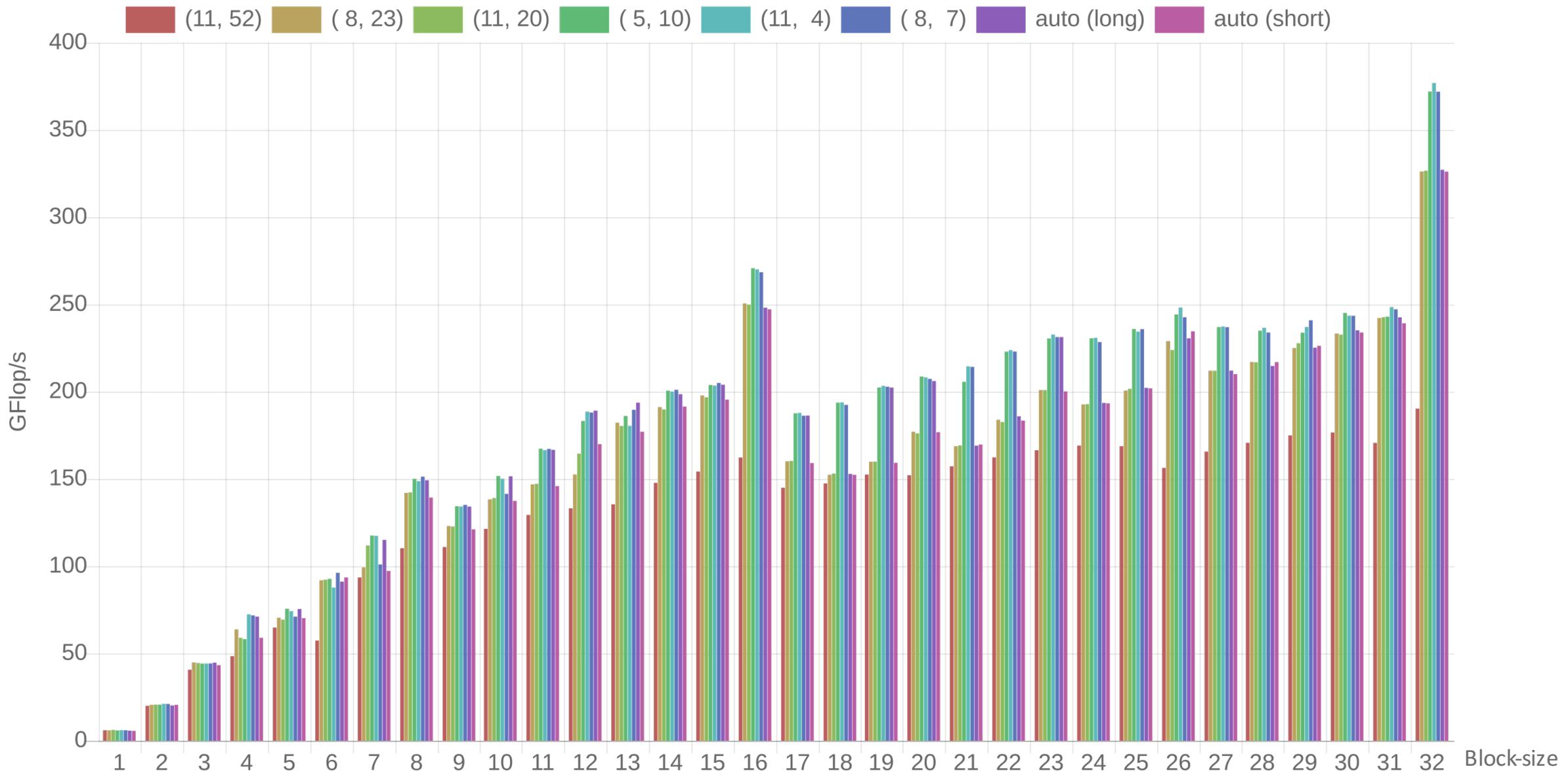
Adaptive Block-Jacobi Generation



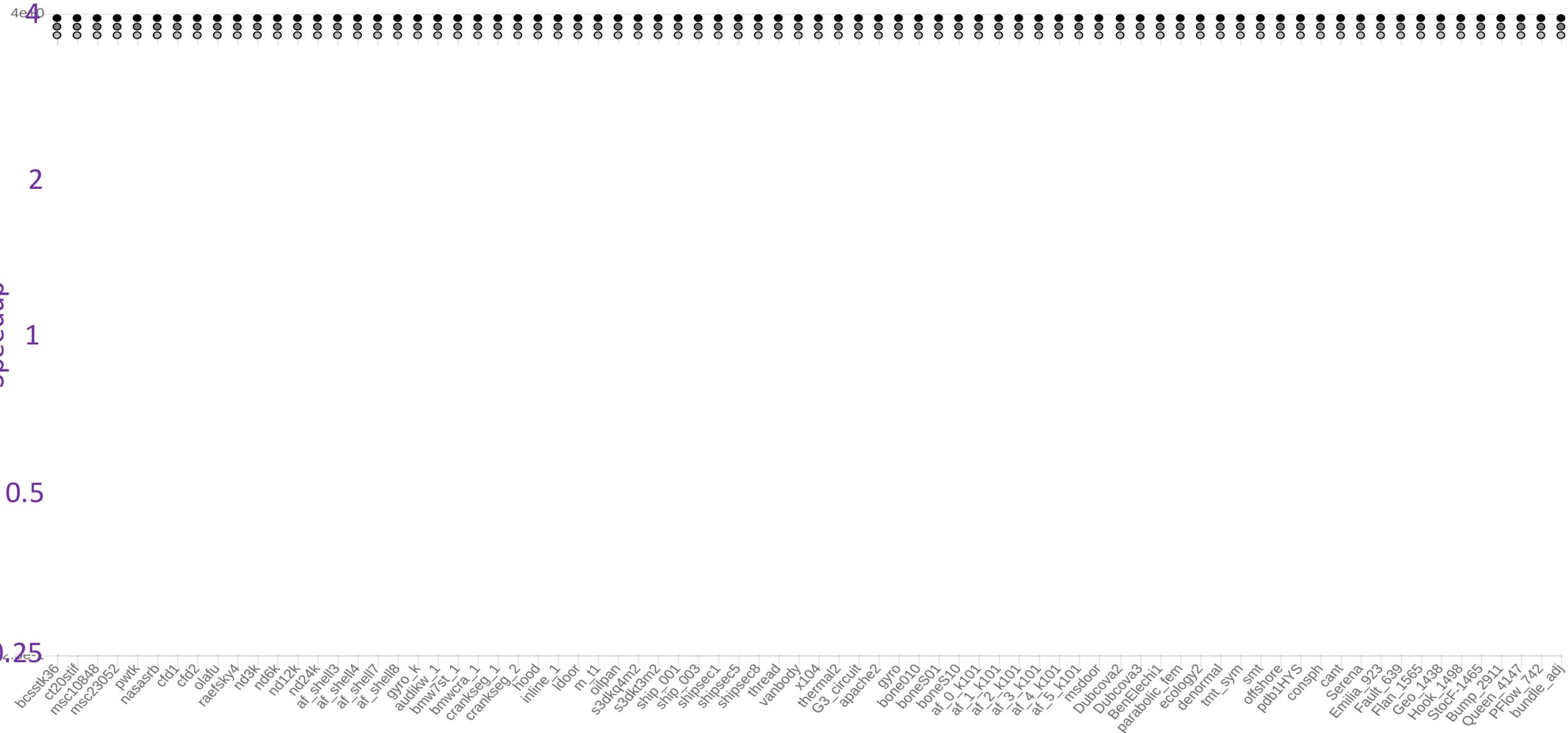
Adaptive Block-Jacobi Application



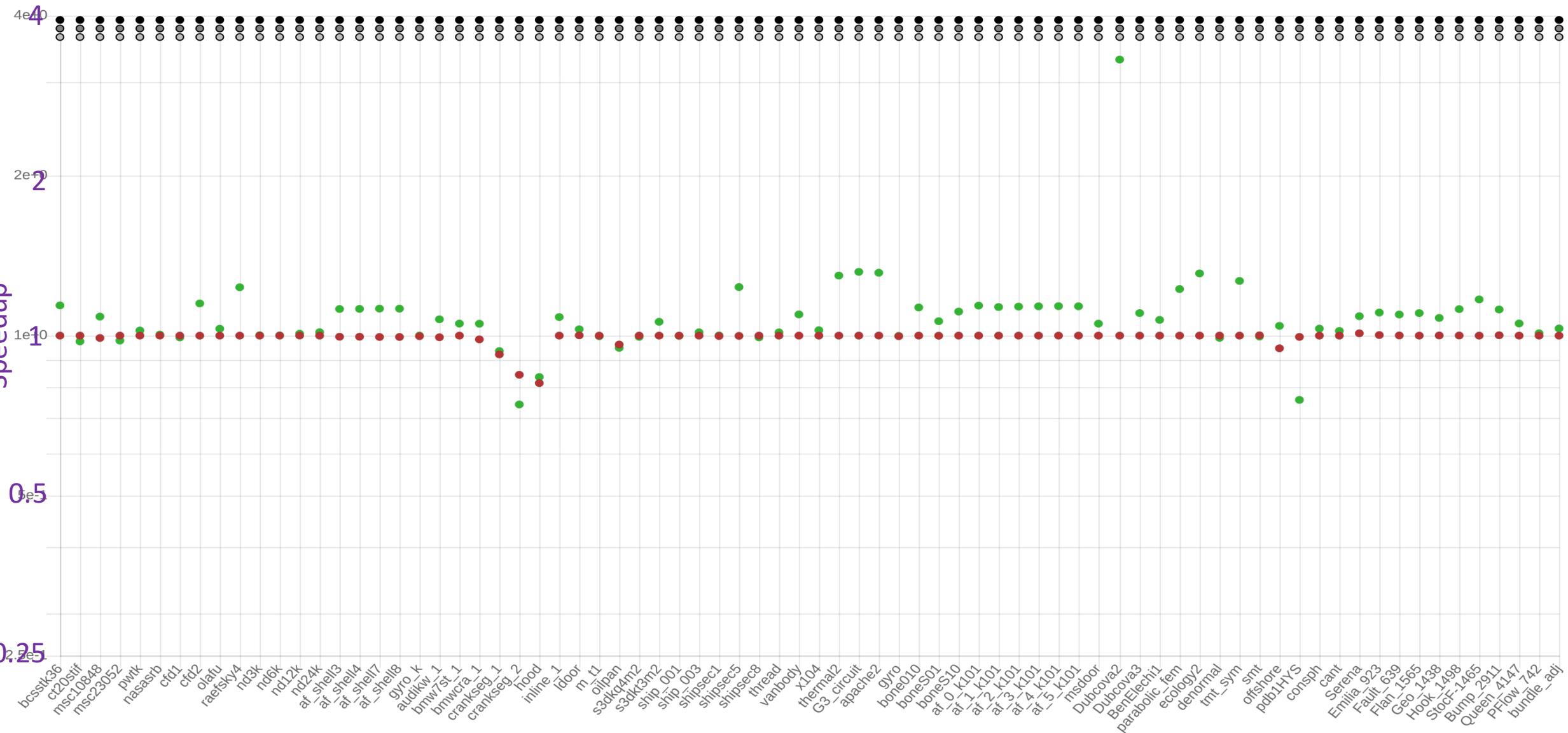
Adaptive Block-Jacobi Application

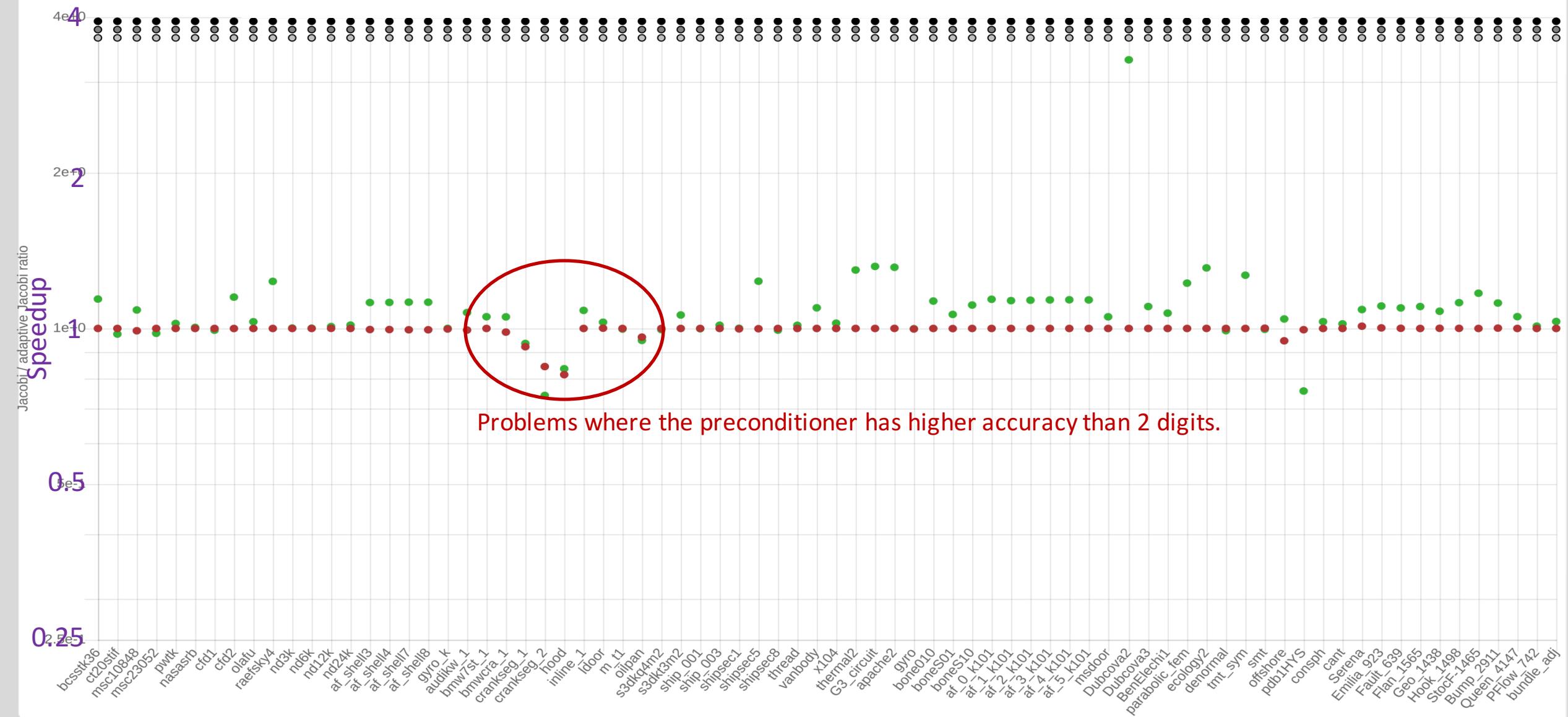


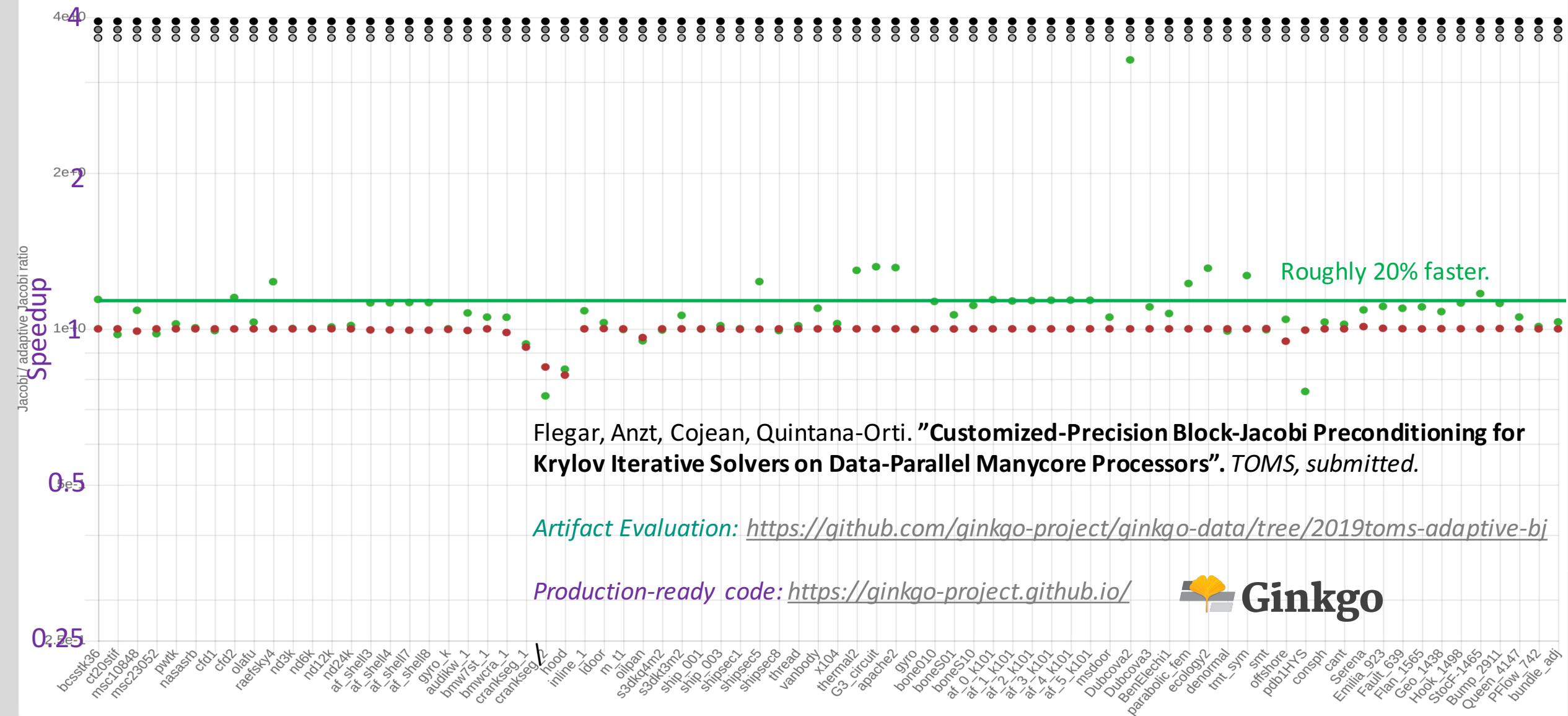
Iterations (adaptive) Time (adaptive) CG converged? CG + Jacobi converged? CG + adaptive Jacobi converged?



Iterations (adaptive) Time (adaptive) CG converged? CG + Jacobi converged? CG + adaptive Jacobi converged?

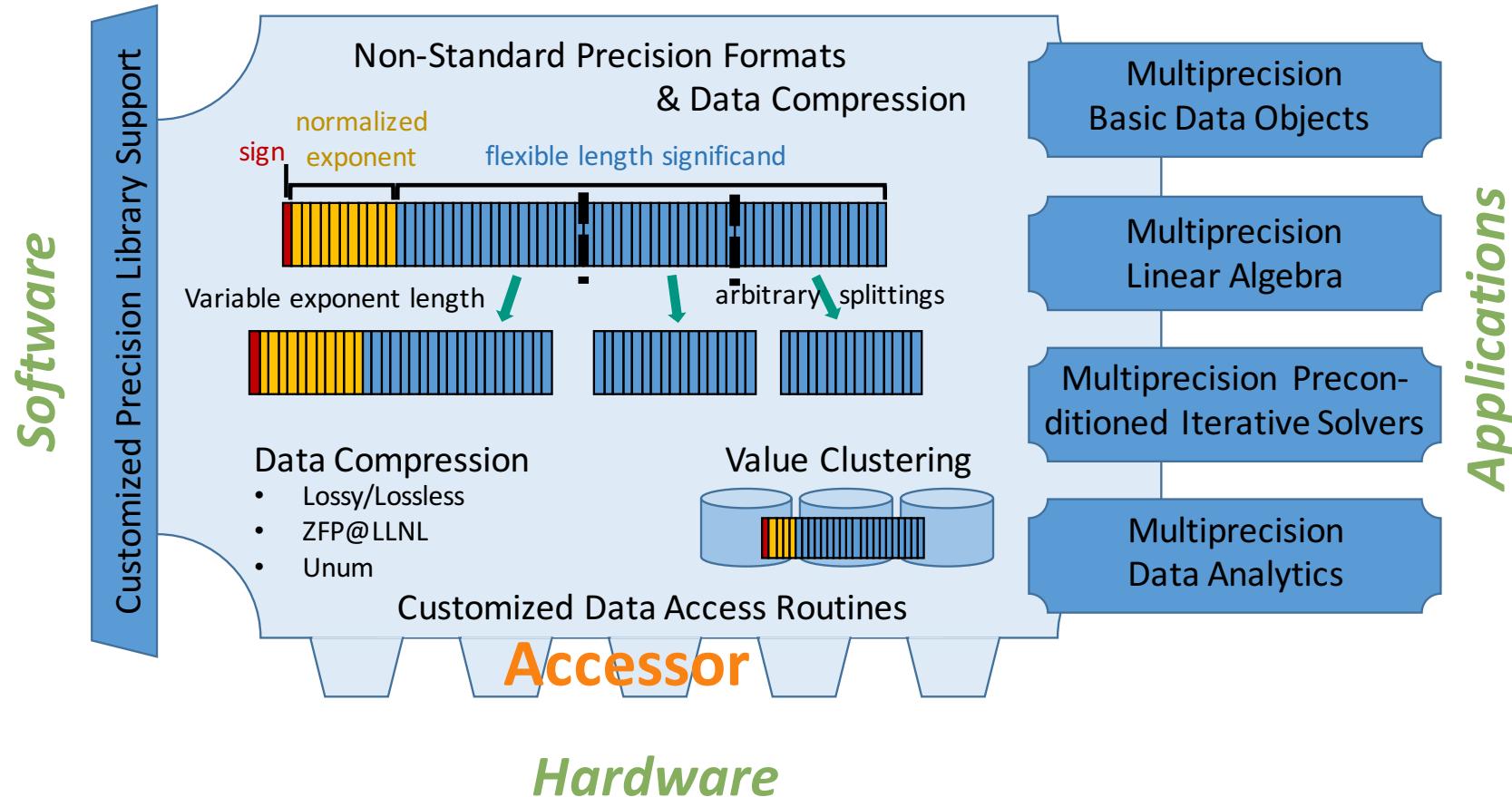






Towards a Modular Precision Ecosystem

- Creating a **Modular Precision Ecosystem** inside  **Ginkgo**.



Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

- Accessor
- SpMV
- Solver kernels
- Precond kernels
- ...

Core

Library Infrastructure
Algorithm Implementations

- Iterative Solvers
- Preconditioners
- ...

Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
- Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;

Core

Library Infrastructure
Algorithm Implementations

- Iterative Solvers
- Preconditioners
- ...

Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Runtime polymorphism selects the right kernel depending on the target architecture;

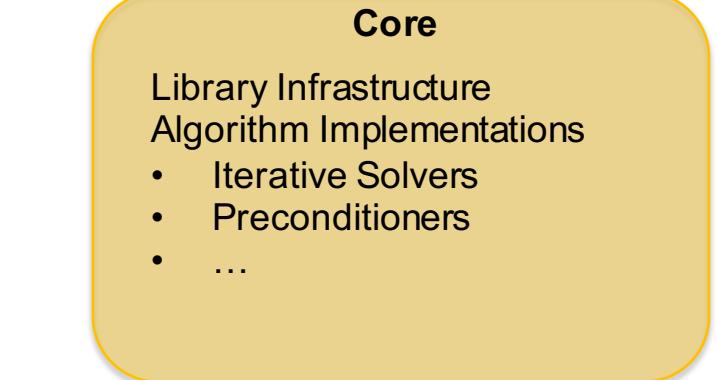
Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;



CUDA

- NVIDIA-GPU kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Optimized architecture-specific kernels;

OpenMP

- OpenMP-kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Runtime polymorphism selects the right kernel depending on the target architecture;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

CUDA

- NVIDIA-GPU kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

OpenMP

- OpenMP-kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;

Optimized architecture-specific kernels;



Core Concept: Separate Algorithm from Kernels



Library core contains architecture-agnostic algorithm implementation;

Runtime polymorphism selects the right kernel depending on the target architecture;

Architecture-specific kernels execute the algorithm on target architecture;

Kernels

Reference

- Reference kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

Reference are sequential kernels to check correctness of algorithm design and optimized kernels;

CUDA

- NVIDIA-GPU kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

OpenMP

- OpenMP-kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

HIP

- AMD-GPU kernels
 - Accessor
 - SpMV
 - ...
- Multi-GPU
 - NVIDIA-GPU kernels
 - Accessor
 - SpMV
 - Solver kernels
 - Precond kernels
 - ...

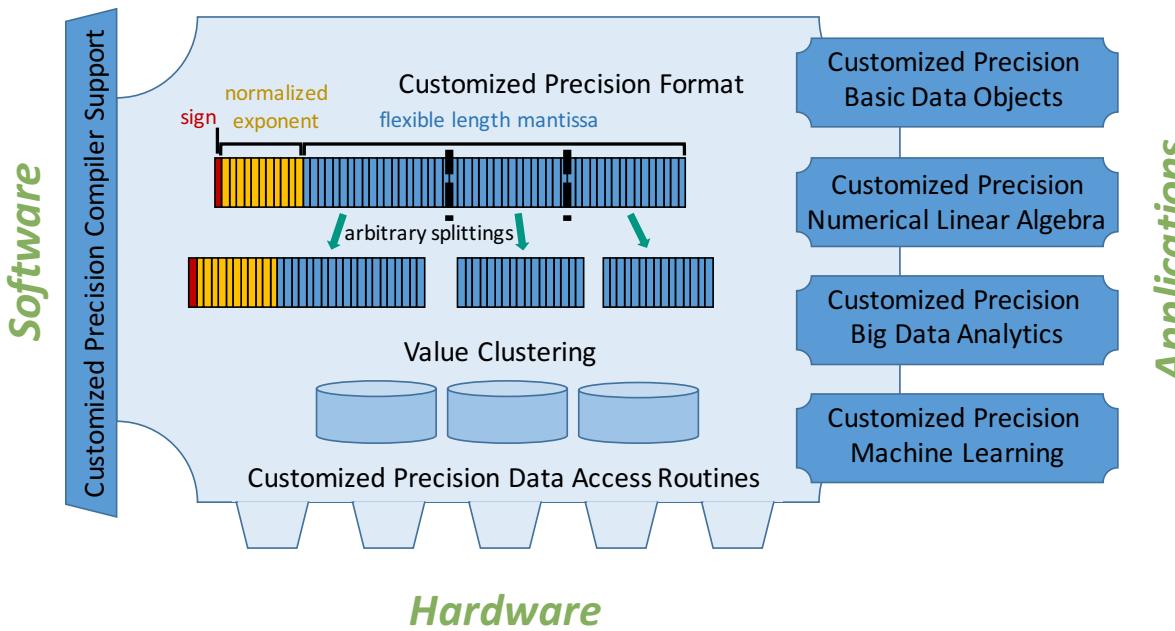
Optimized architecture-specific kernels;



Summary and next steps

- Decouple arithmetic precision from memory precision.
- Using **customized precisions** for memory operations.
- Speedup of up to 1.3x for adaptive precision block-Jacobi preconditioning.
- Creating a **Modular Precision Ecosystem** inside  **Ginkgo**.

<https://github.com/ginkgo-project/ginkgo>



HELMHOLTZ
RESEARCH FOR GRAND CHALLENGES

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration and the Helmholtz Impuls und Vernetzungsfond VH-NG-1241.