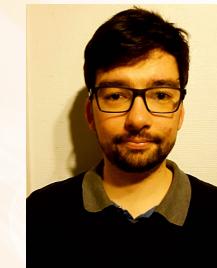


Addressing the Communication Bottleneck: **Towards a Modular Precision Ecosystem for High Performance Computing**

Focus Session: New Approaches, Algorithms Towards Exascale Computing
ISC High Performance 2019, June 19th, Frankfurt

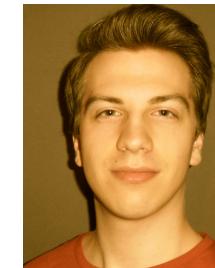
Hartwig Anzt, Terry Cojean, Goran Flegar, Thomas Grützmacher, Pratik Nayak, Tobias Ribizel
Steinbuch Centre for Computing (SCC)



Terry Cojean



Goran Flegar



Thomas
Grützmacher



Pratik Nayak



Tobias Ribizel

Where do we stand?

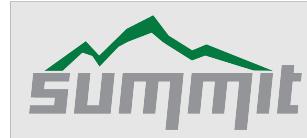


- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s** (10^{18}) for 16 bit floating point used in data analytics and artificial intelligence



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

Roofline Model

Given certain **hardware characteristics**:

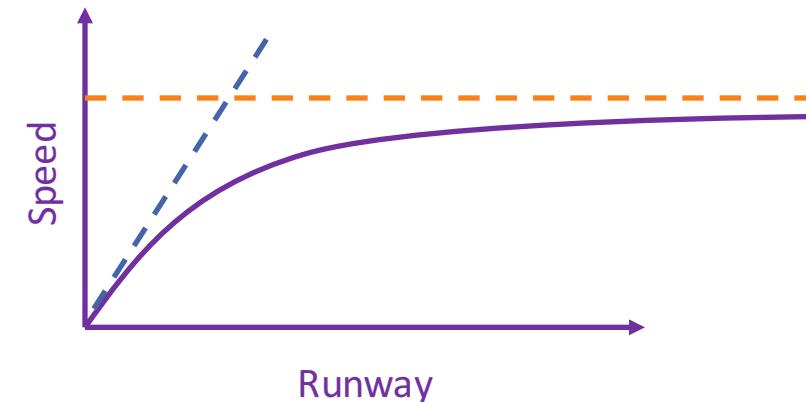
memory bandwidth,
arithmetic power,

Acceleration
Top Speed



the performance of any operation is

- either bound by the data access/communication (*memory bound*),
- or by the arithmetic operations (*compute bound*).



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

Roofline Model

Given certain **hardware characteristics**:

memory bandwidth,
arithmetic power,

Acceleration
Top Speed



the performance of any operation is

- either bound by the data access/communication (*memory bound*),
- or by the arithmetic operations (*compute bound*).

Matrix-Matrix Product (GEMM): $C = A \times B$ $A, B, C \in \mathbb{R}^{n \times n}$

$3n^2$ Memory operations

$2n^3$ Arithmetic operations

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

Roofline Model

Given certain **hardware characteristics**:

memory bandwidth,
arithmetic power,

**Acceleration
Top Speed**



the performance of any operation is

- either bound by the data access/communication (*memory bound*),
- or by the arithmetic operations (*compute bound*).

Matrix-Matrix Product (GEMM): $C = A \times B$ $A, B, C \in \mathbb{R}^{n \times n}$

$3n^2$ Memory operations
 $2n^3$ Arithmetic operations

We just need to increase the size, and at some point the operation becomes compute bound.

"we infinitely extend the acceleration runway"

Where do we stand?



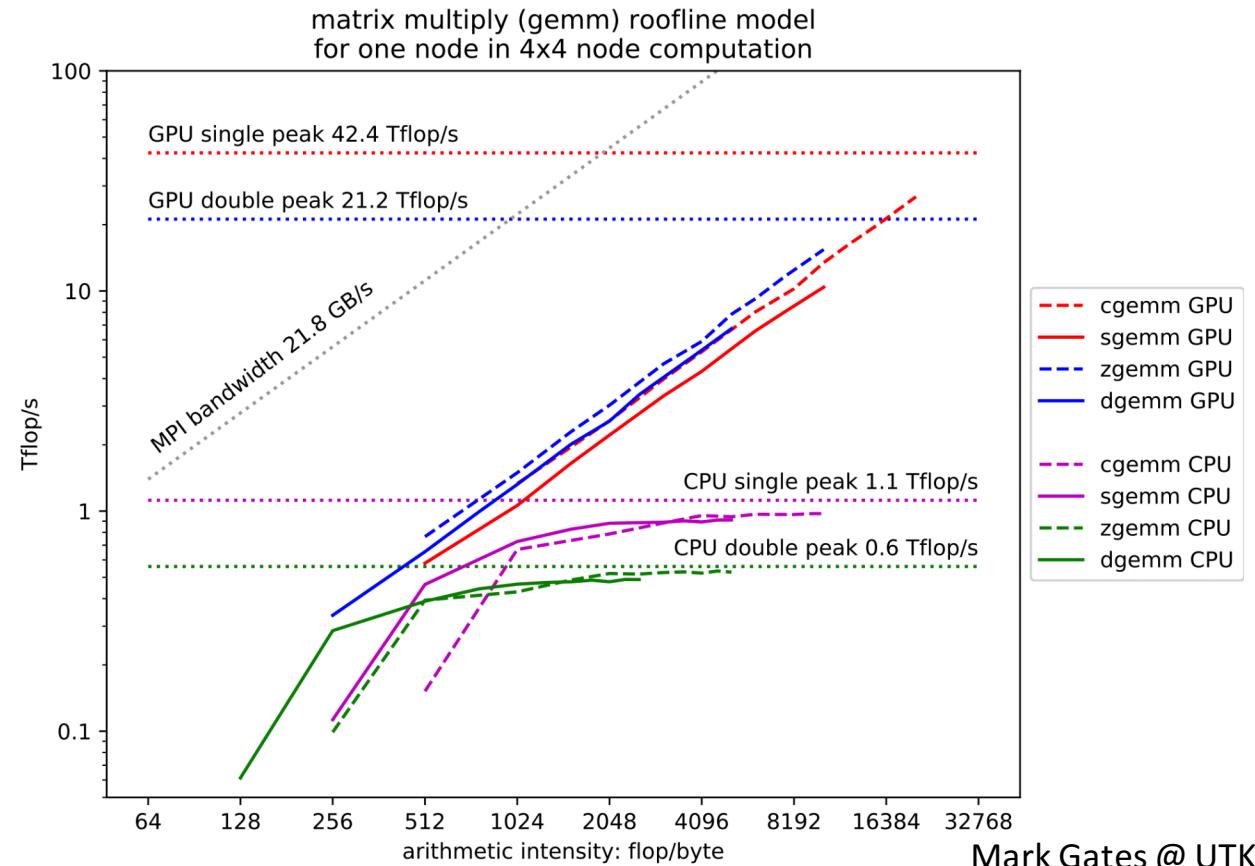
- Node: 2 IBM POWER9 + 6 NVIDIA V100 GPUs
- 4,608 nodes, 9,216 IBM Power9 CPUs
- 27,648 V100 GPUs (**8 TFLOPs / GPU**)
- Peak performance of **200 Pflop/s** for modeling & simulation
- Peak performance of **3.3 Eflop/s (10^18)** for 16 bit floating point used in data analytics and artificial intelligence



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;

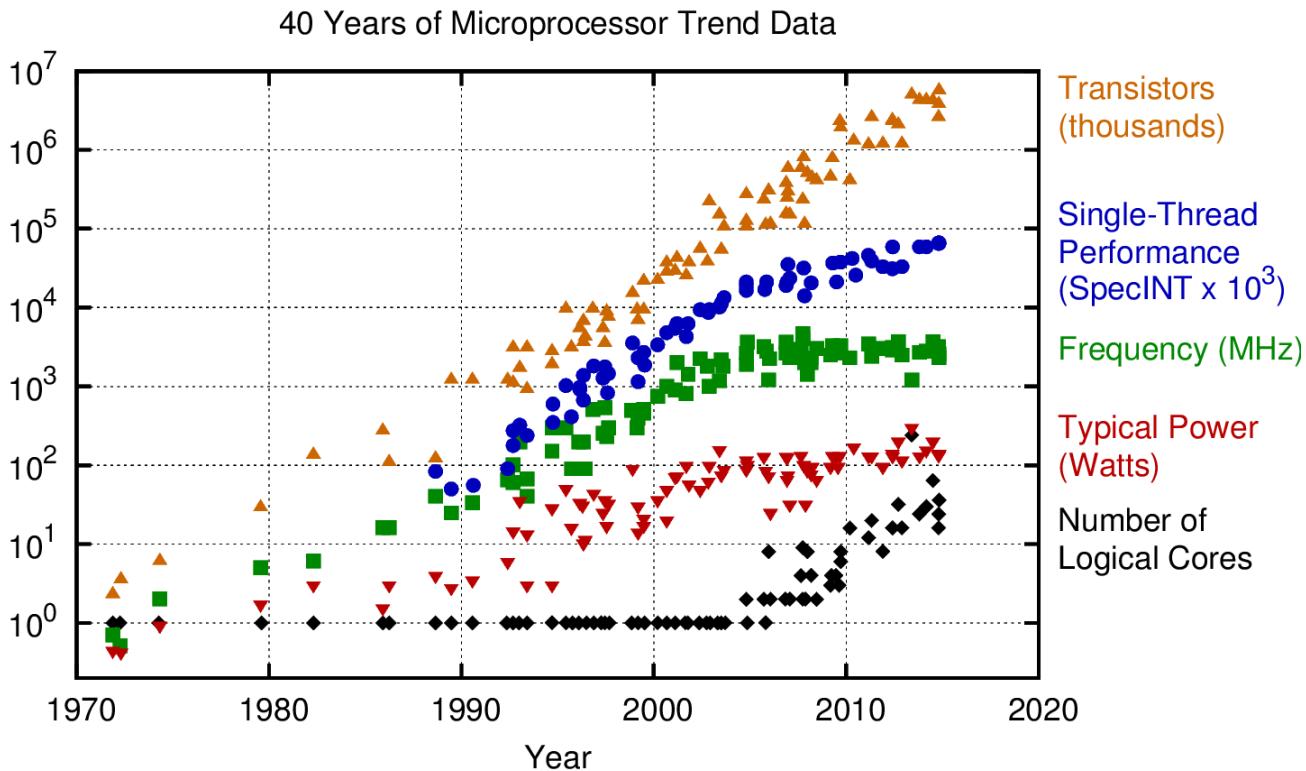
Sparse / Graph Problems?

- *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.

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

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

How did we get there?

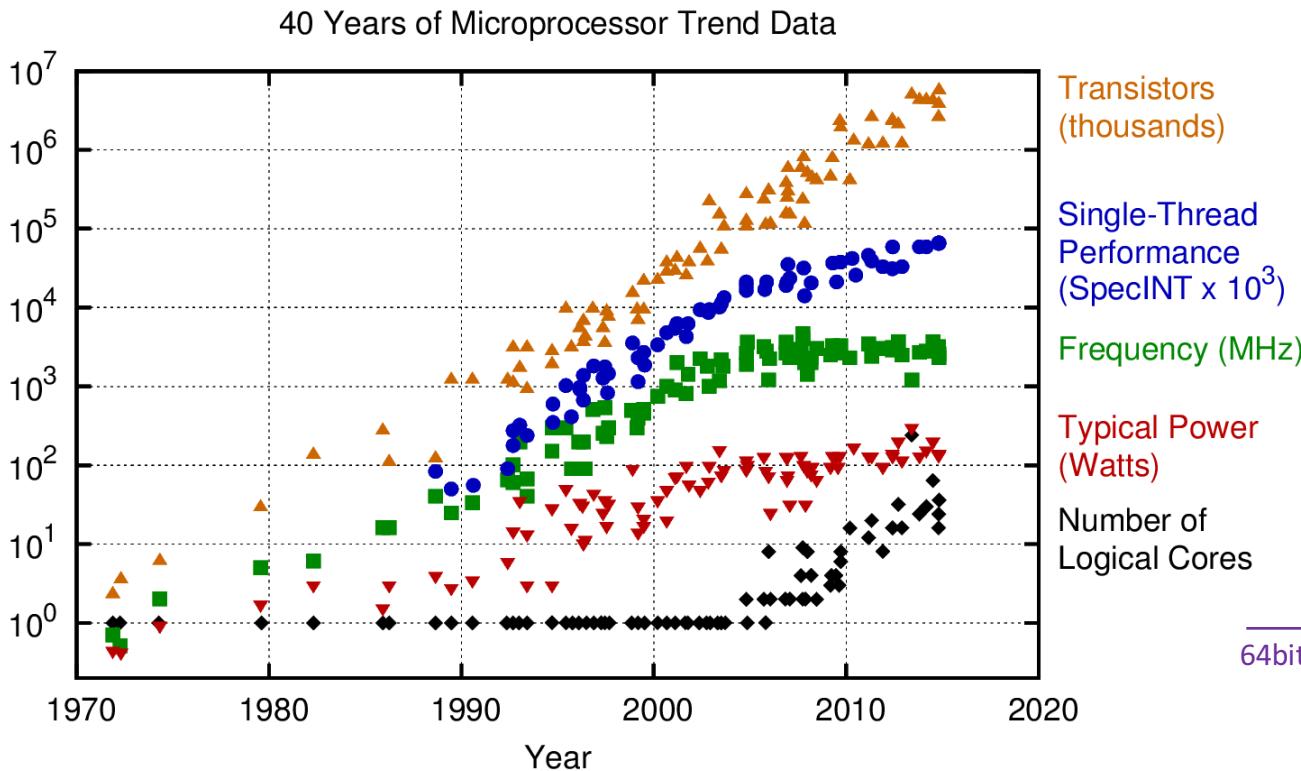


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

- Explosion in core count.

"Parallelism needed – synchronization kills performance."

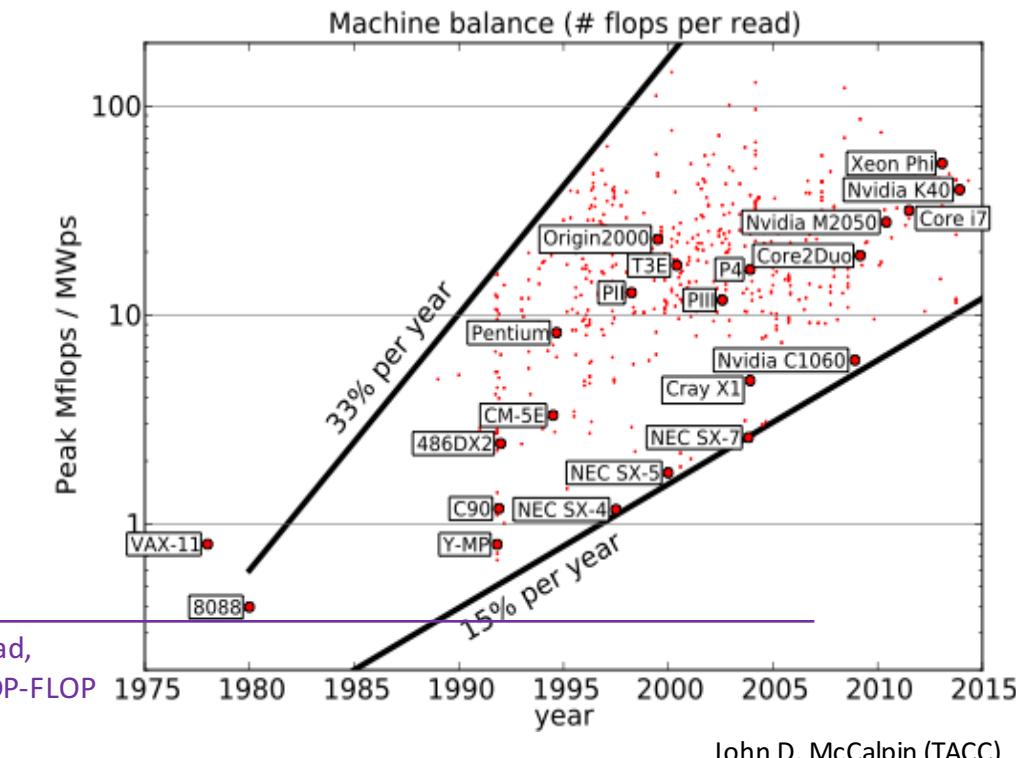
How did we get there?



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

- Explosion in core count.
- Compute power (#FLOPs) grows much faster than bandwidth.

"Parallelism needed – synchronization kills performance."



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

How did we get there?

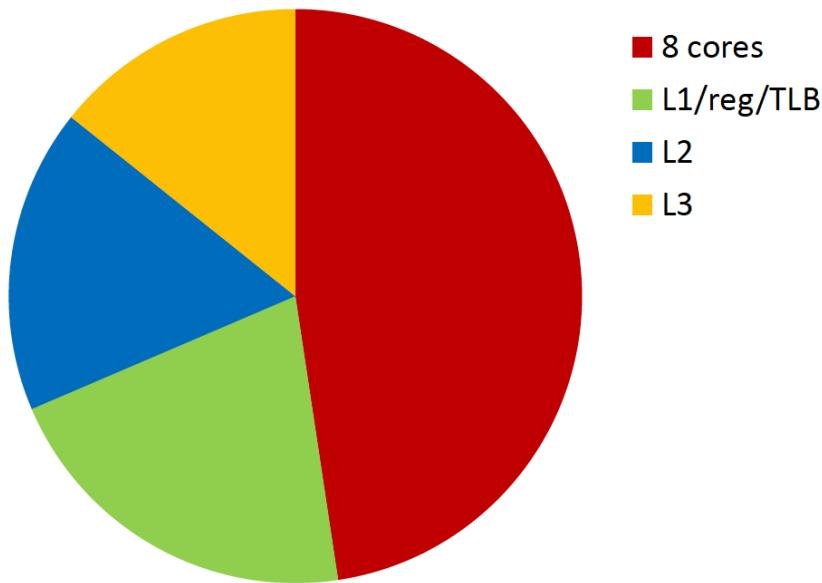
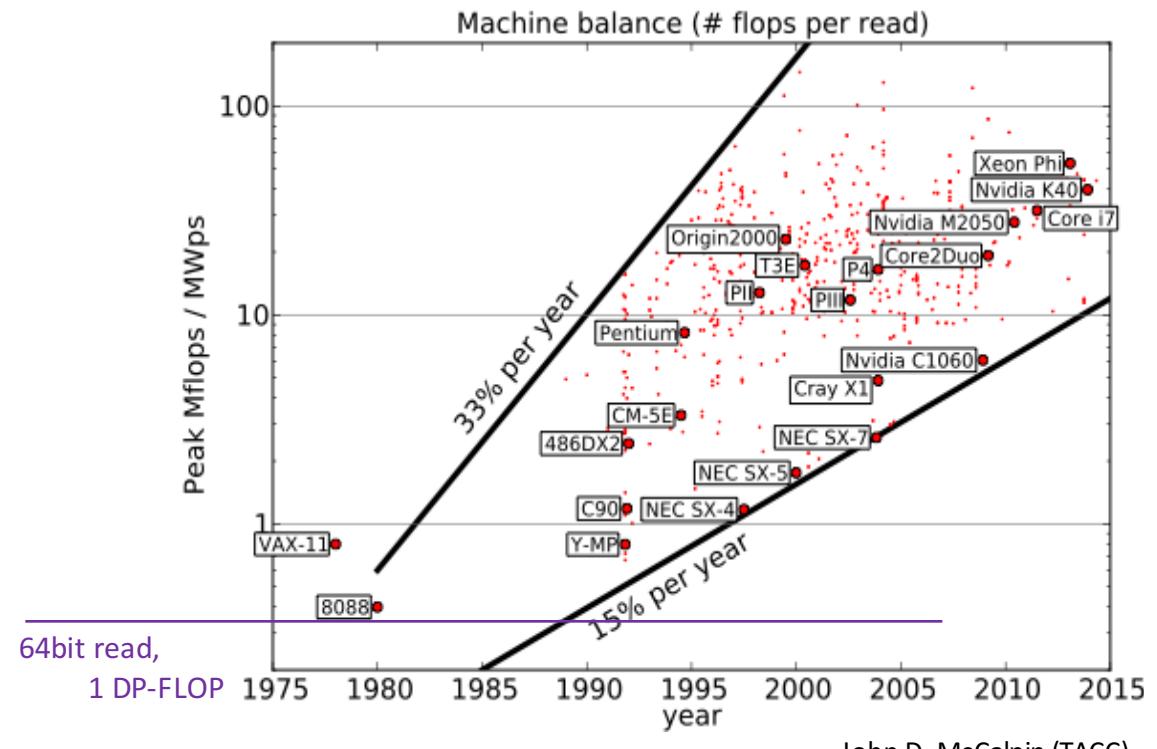


Figure 1.1.7: Power breakdown of an 8 core server chip.

Mark Horowitz (2014): Computing's energy problem (and what we can do about it)

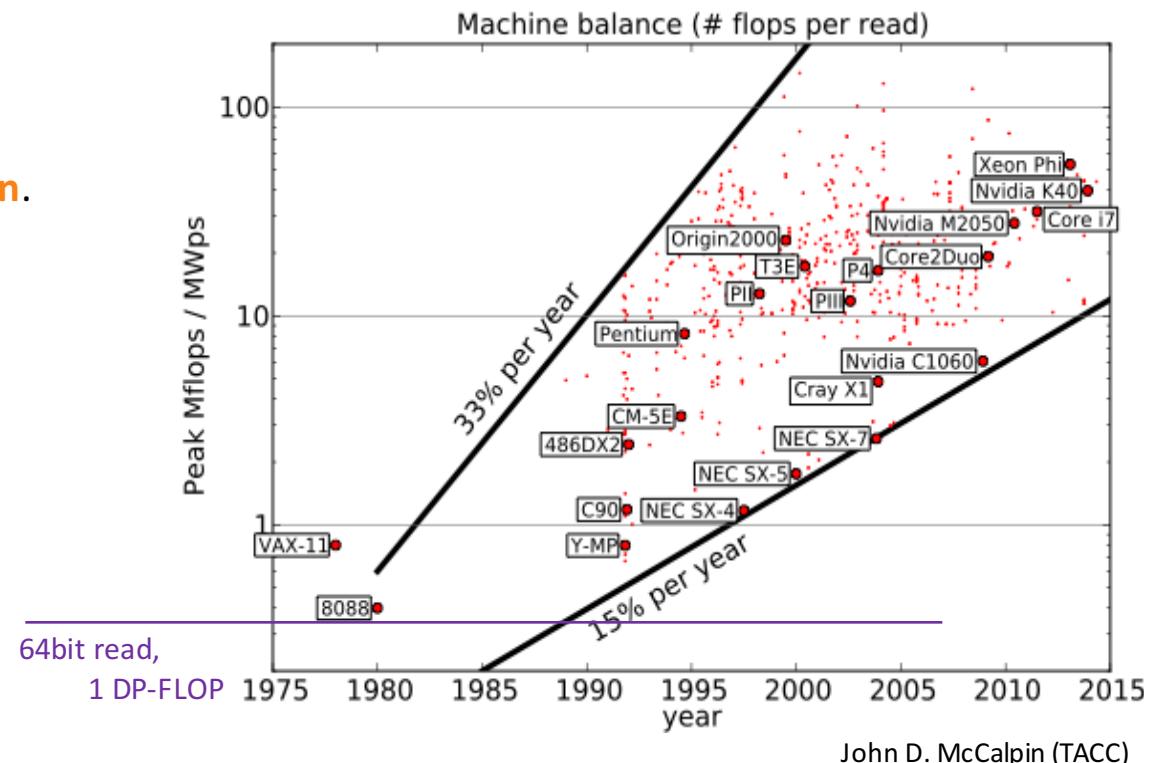


- Compute power (#FLOPs) grows much faster than bandwidth.

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

Algorithms reflecting Hardware Evolution

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



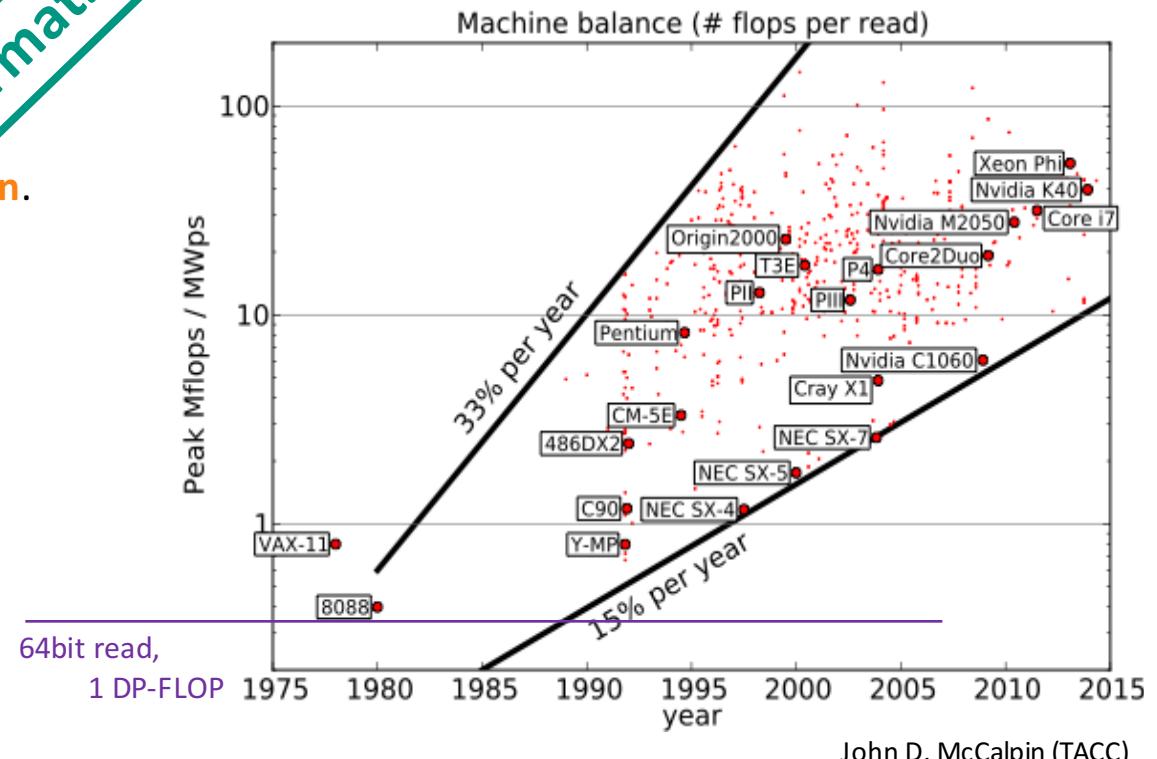
- Compute power (#FLOPs) grows much faster than bandwidth.

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

Algorithms reflecting Hardware Evolution

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

Radically decouple storage format from arithmetic format.

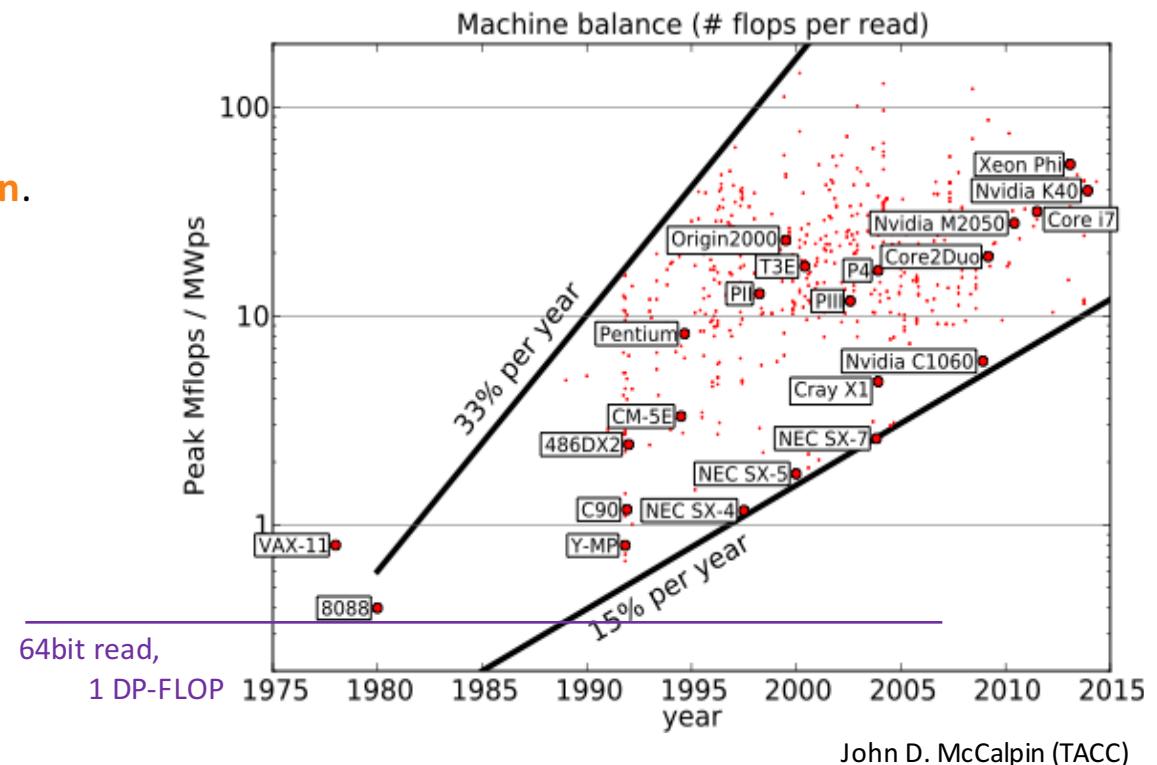


- Compute power (#FLOPs) grows much faster than bandwidth.

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

Algorithms reflecting Hardware Evolution

- The **arithmetic operations** should use **high precision formats** natively supported by hardware.
- **Data access** should be as cheap as possible, **reduced precision**.
- **Radically decouple storage format from arithmetic format.**
 - *Double precision in all arithmetic operations.*
 - *Algorithm-aware precision when accessing data.*



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

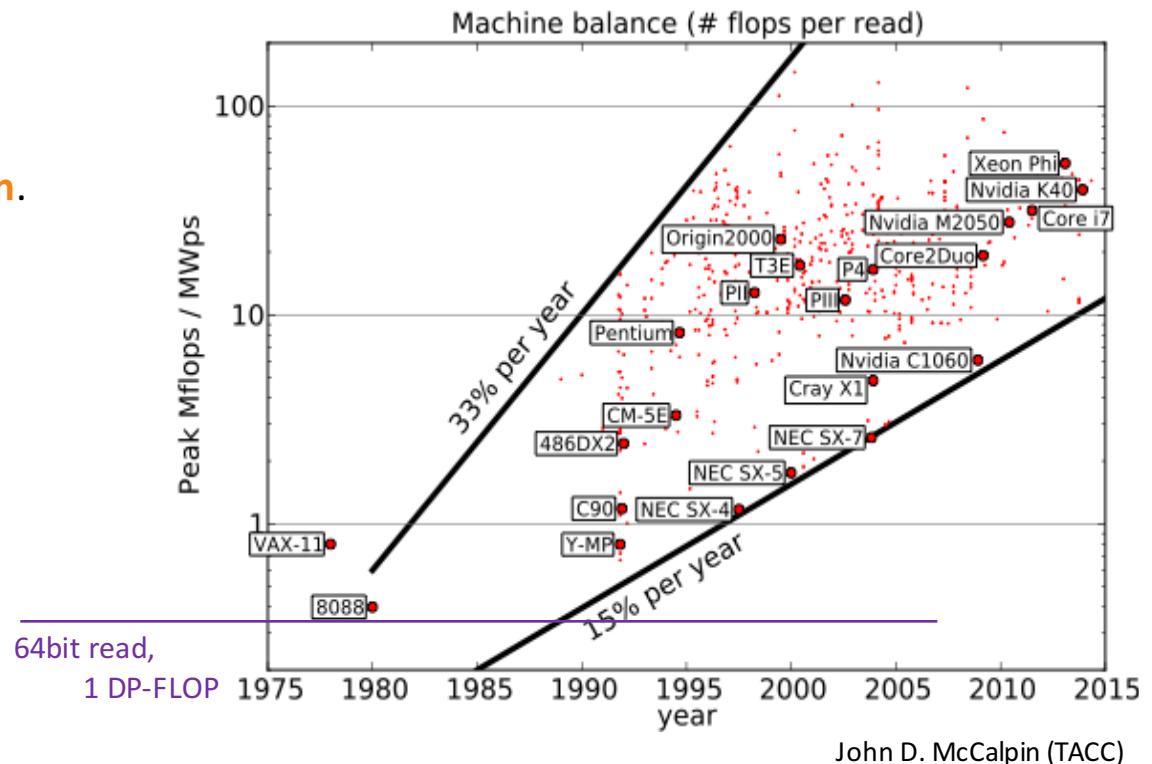
"Operations are free, mem access and comm is what counts."

Algorithms reflecting Hardware Evolution

- The **arithmetic operations** should use **high precision formats** natively supported by hardware.
- **Data access** should be as cheap as possible, **reduced precision**.
- **Radically decouple storage format from arithmetic format.**
 - *Double precision in all arithmetic operations.*
 - *Algorithm-aware precision when accessing data.*

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).



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

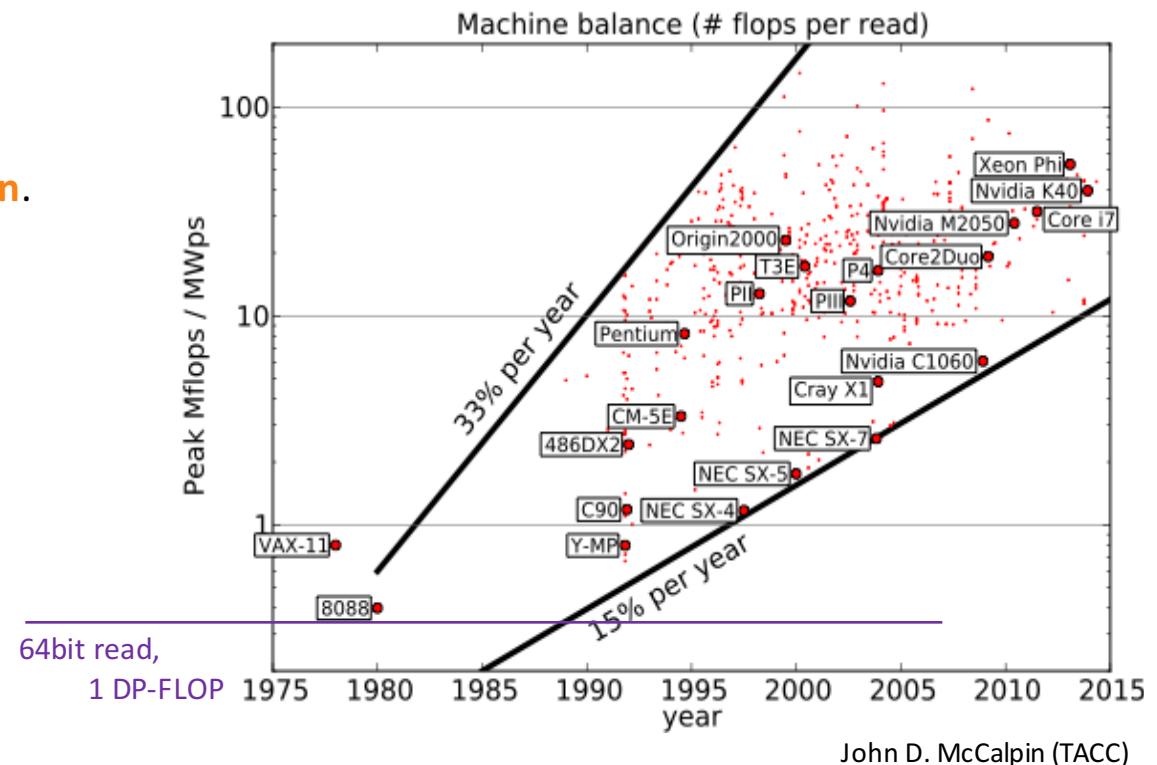
"Operations are free, mem access and comm is what counts."

Algorithms reflecting Hardware Evolution

- The **arithmetic operations** should use **high precision formats** natively supported by hardware.
- **Data access** should be as cheap as possible, **reduced precision**.
- **Radically decouple storage format from arithmetic format.**
 - *Double precision in all arithmetic operations.*
 - *Algorithm-aware precision when accessing data.*

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**



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

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

Use reduced precision for “approximate Operators”

- **Preconditioners for iterative solvers.**

Use reduced precision for “approximate Operators”

- **Preconditioners for iterative solvers.**

- Idea: Approximate inverse of system matrix to make the system “easier to solve”: $P^{-1} \approx A^{-1}$

$$Ax = b \Leftrightarrow \tilde{A}x = \tilde{b}$$

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 $Ax = b \Leftrightarrow \tilde{A}x = \tilde{b}$.

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?**

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!**

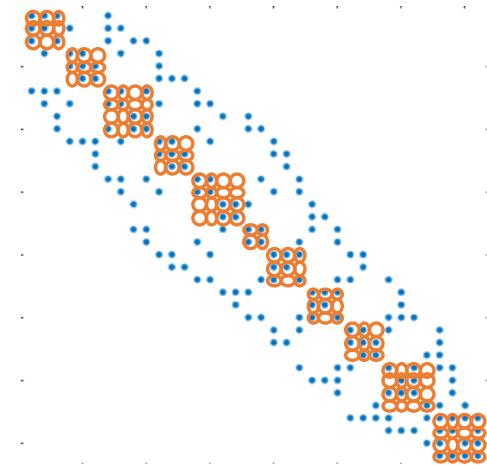
Use reduced precision for “approximate Operators”

- **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**: $P = \text{diag}(A)$

Use reduced precision for “approximate Operators”

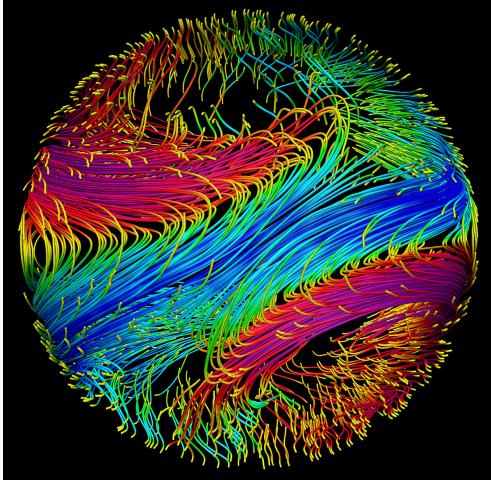
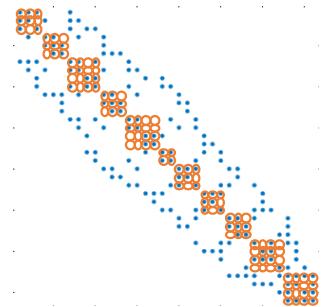
- **Preconditioners for iterative solvers.**
 - Idea: Approximate inverse of system matrix to make the system “easier to solve”: $P^{-1} \approx A^{-1}$
 $\tilde{A} = P^{-1}A$, $\tilde{b} = P^{-1}b$, and we solve $Ax = b \Leftrightarrow \tilde{A}x = \tilde{b}$.
 - **Why should we store the preconditioner 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

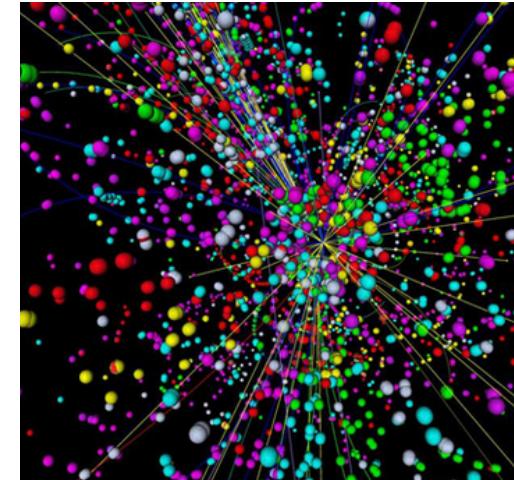
- **Block-Jacobi** method typically used as **preconditioner** inside **linear- / Eigenvalue solvers**.
- Target: large, sparse linear systems.
discretizations often carry a **block-structure** (multiple variables per node).
- “Natural blocks” of small size (8, 12..).
- System matrix often stored in **sparse data structure** (CSR).



<http://www.nas.nasa.gov/SC14/>



<https://science.nasa.gov/earth-science/focus-areas/earth-weather>



<http://i.livescience.com>

Block-Jacobi Preconditioning

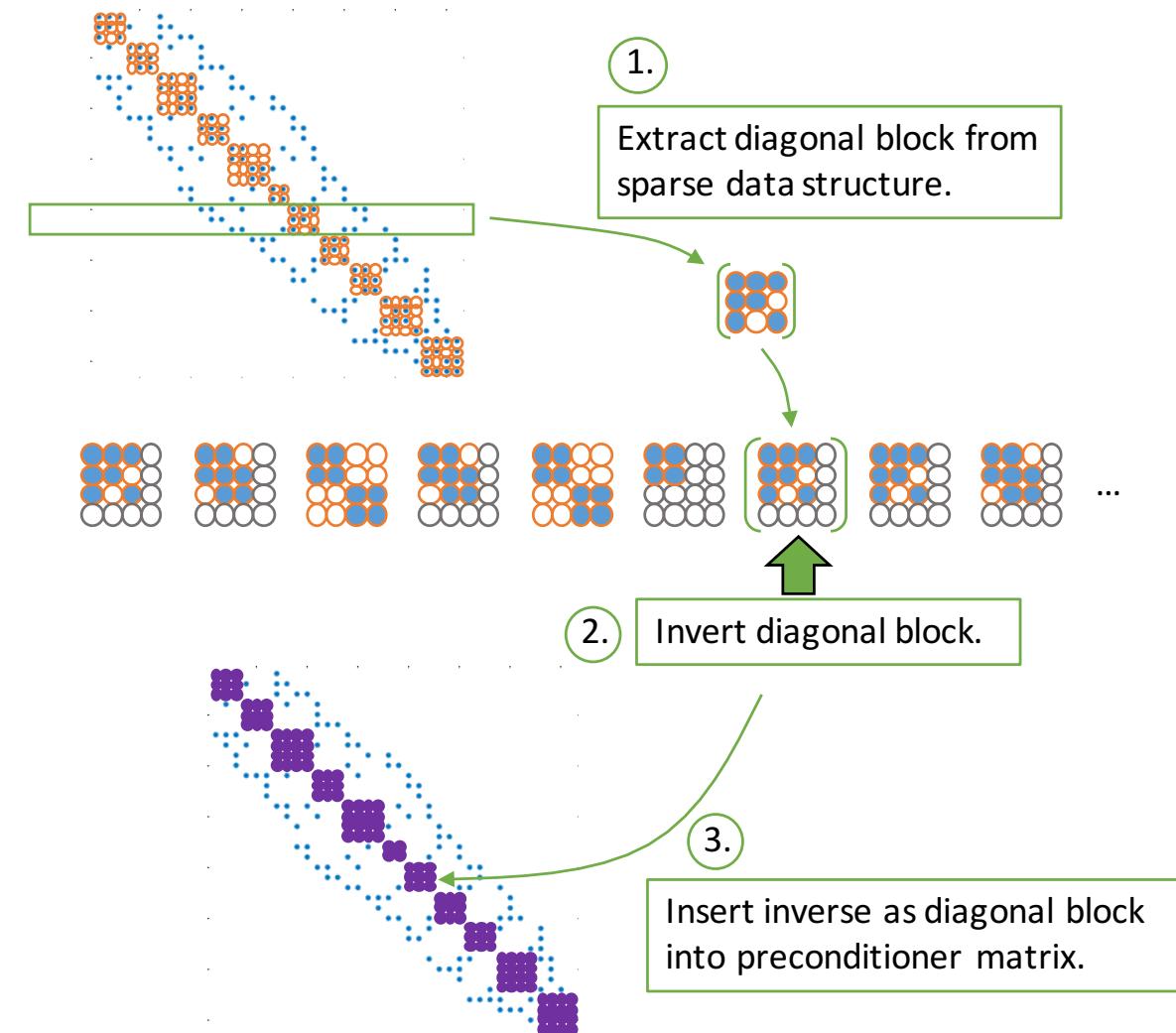
Preconditioner Setup:

- 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:

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



Block-Jacobi Preconditioning

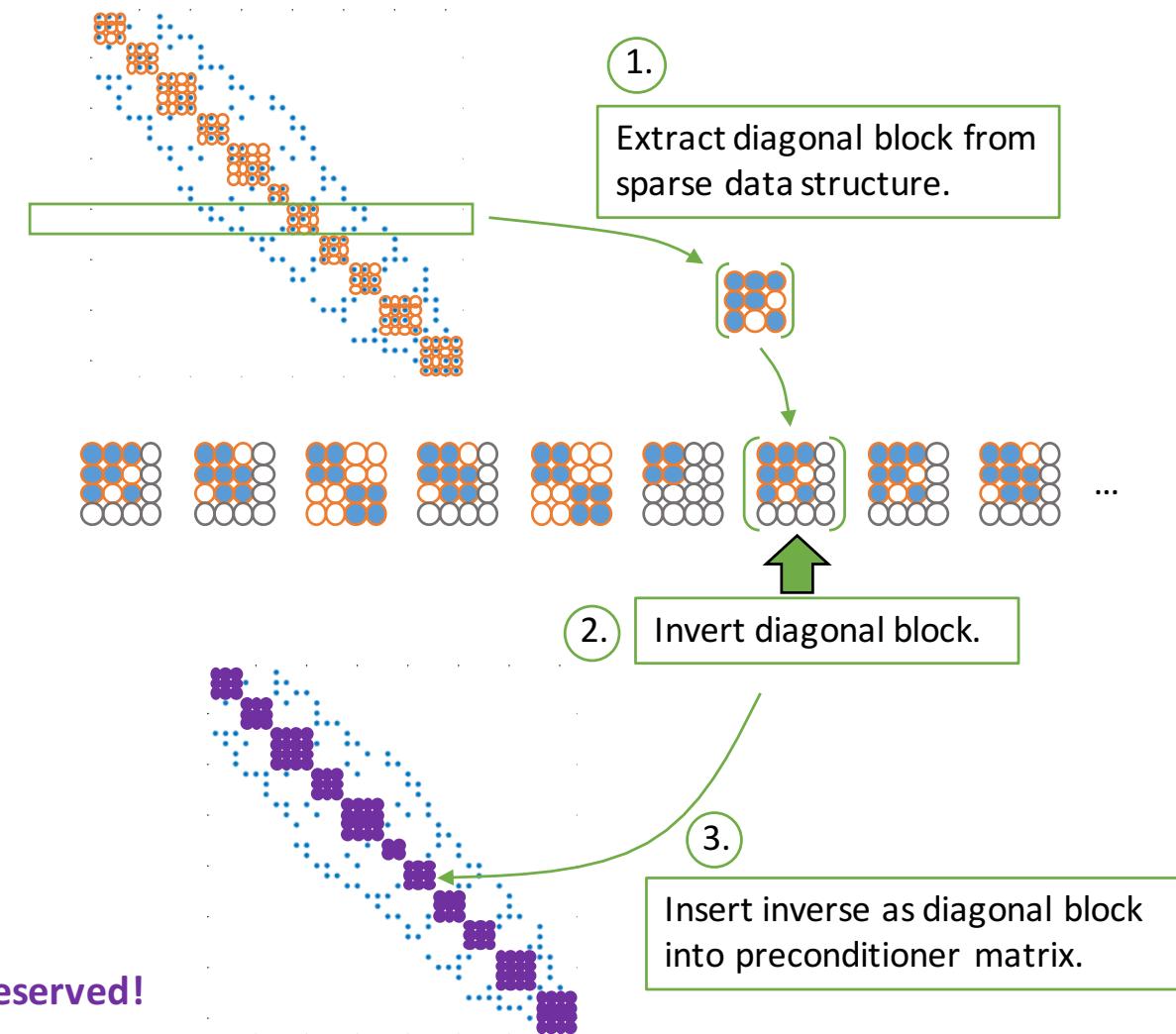
Preconditioner Setup:

- 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:

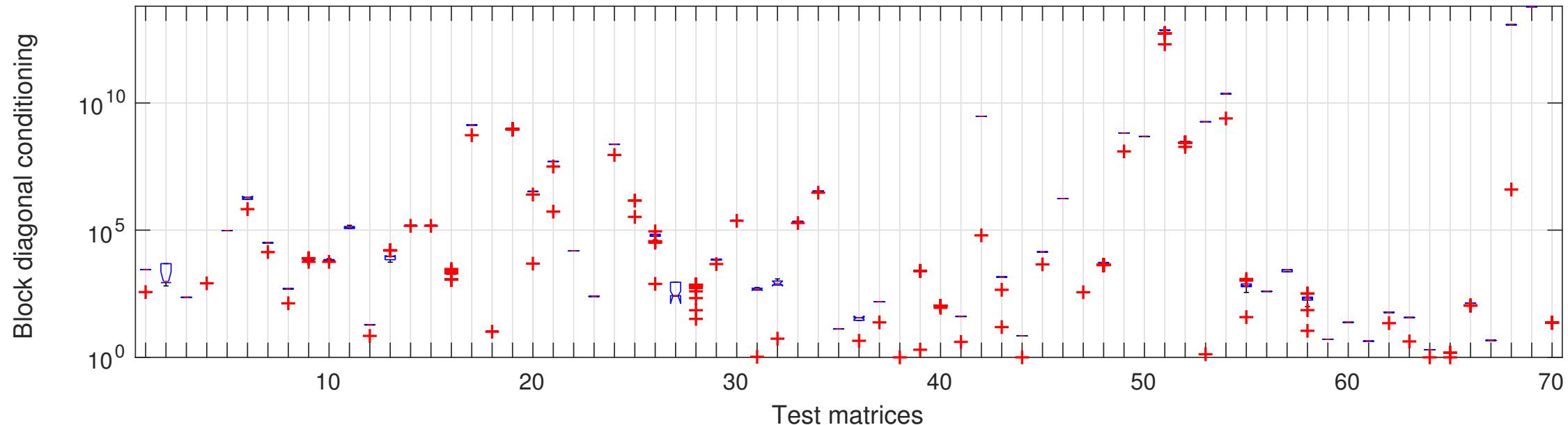
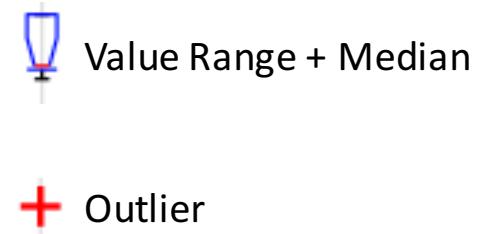
$$y := P^{-1}x$$



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

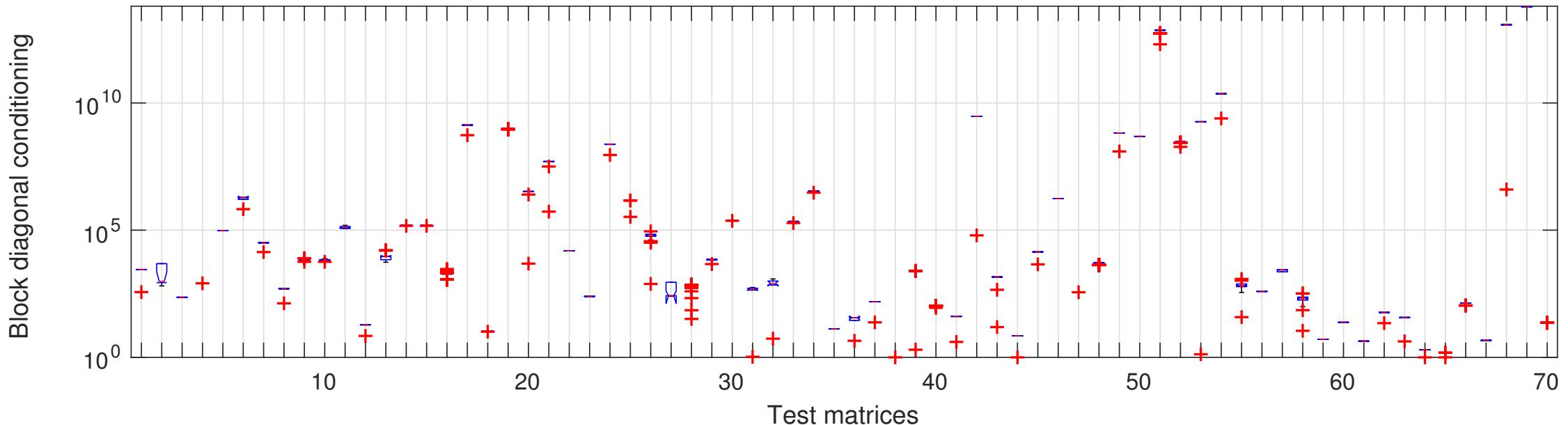
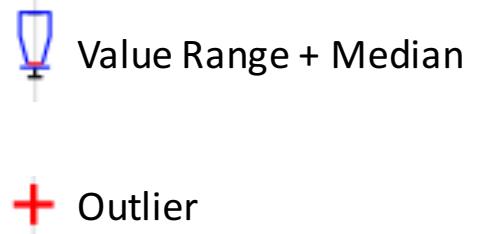
Adaptive Precision Block-Jacobi Preconditioning

- 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

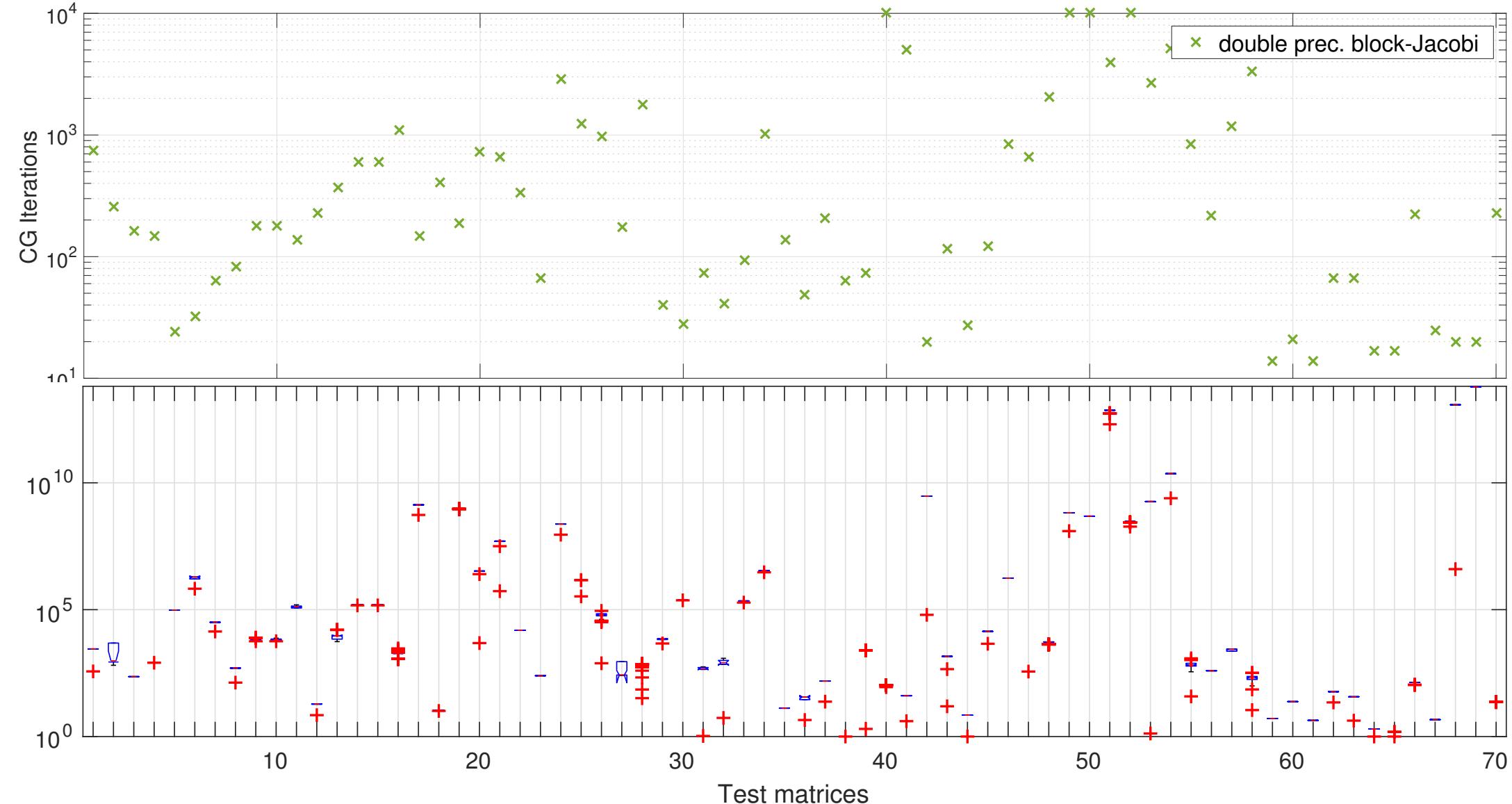


Adaptive Precision Block-Jacobi Preconditioning

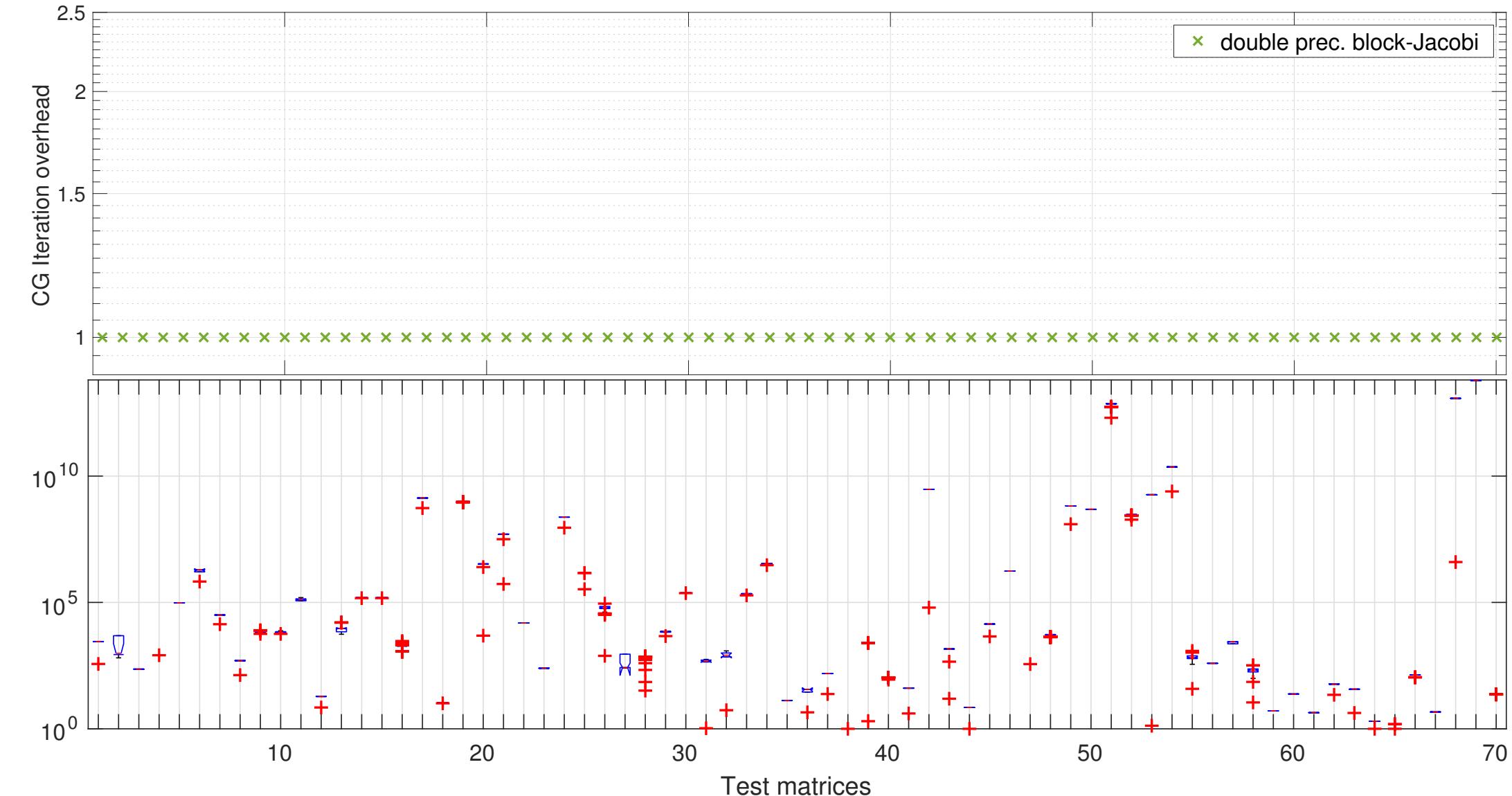
- 70 matrices from the SuiteSparse Matrix Collection
- Use block-size 24 with Super-Variable agglomeration (24 is upper bound for size of blocks)
- Report conditioning of all arising diagonal blocks
- Analyze the impact of storing block-Jacobi in lower precision a top-level Conjugate Gradient solver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



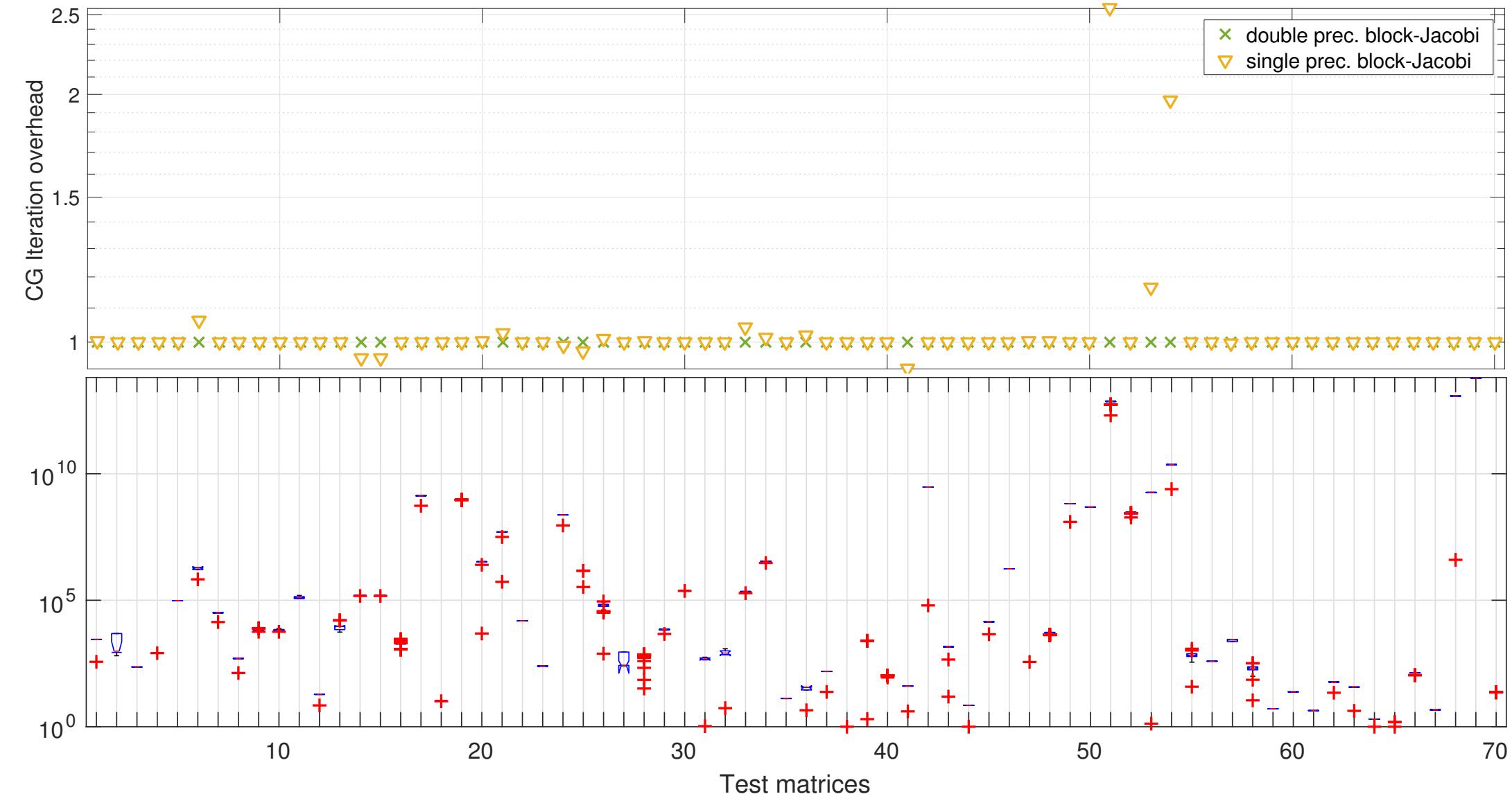
Adaptive Precision Block-Jacobi Preconditioning



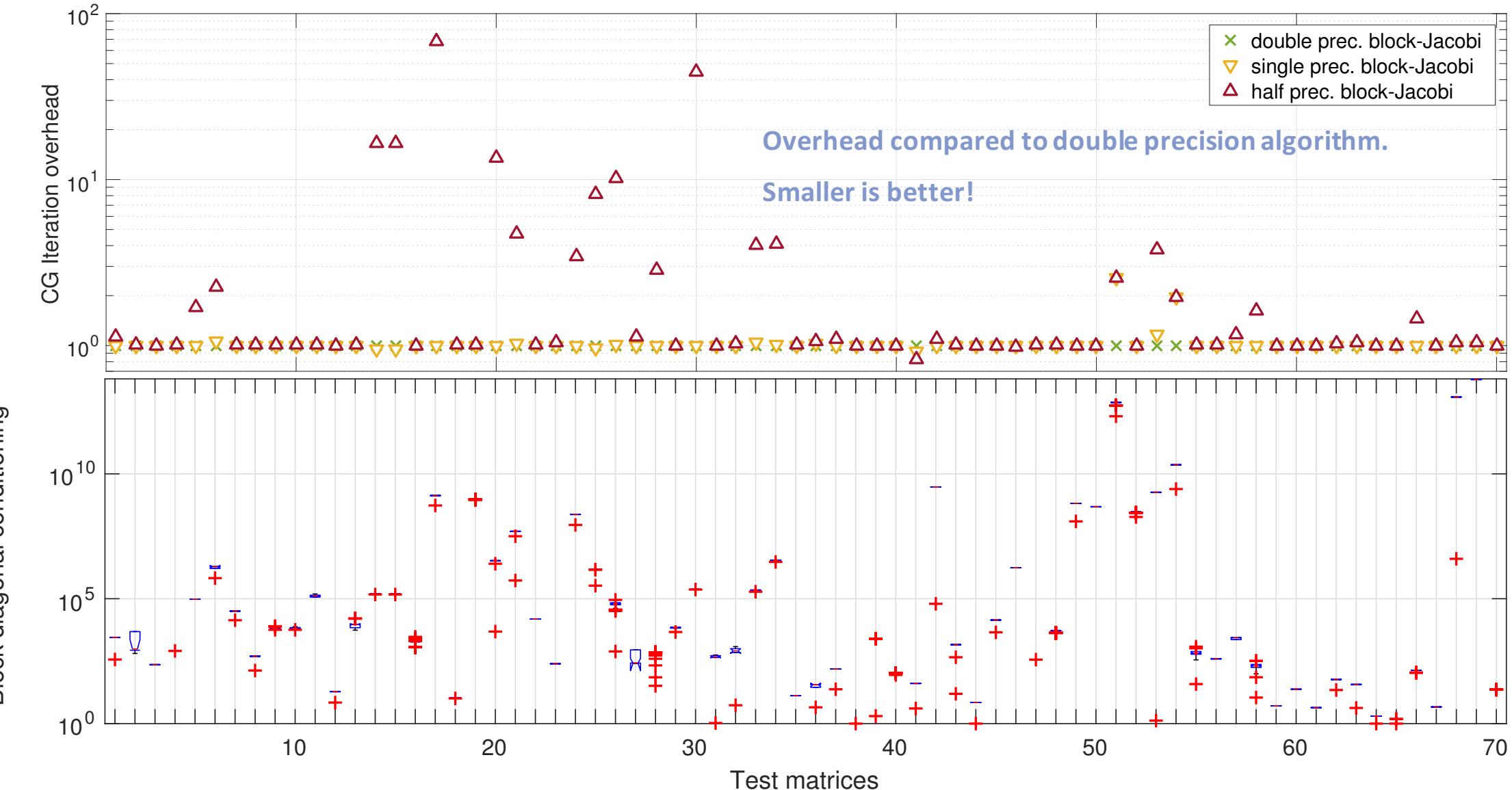
Adaptive Precision Block-Jacobi Preconditioning



Adaptive Precision Block-Jacobi Preconditioning



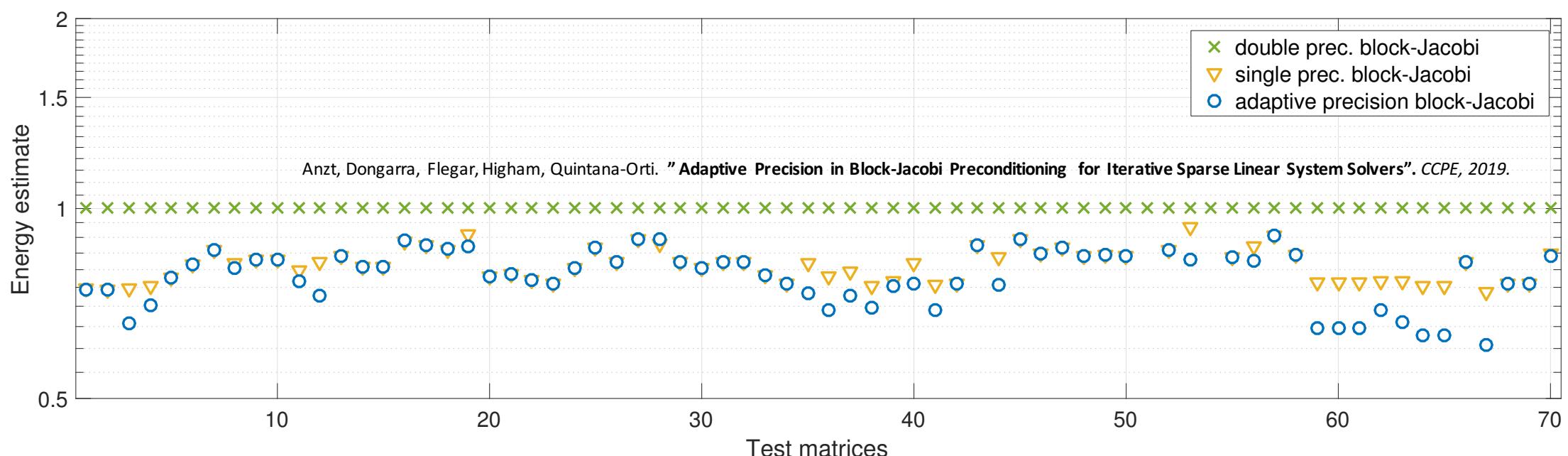
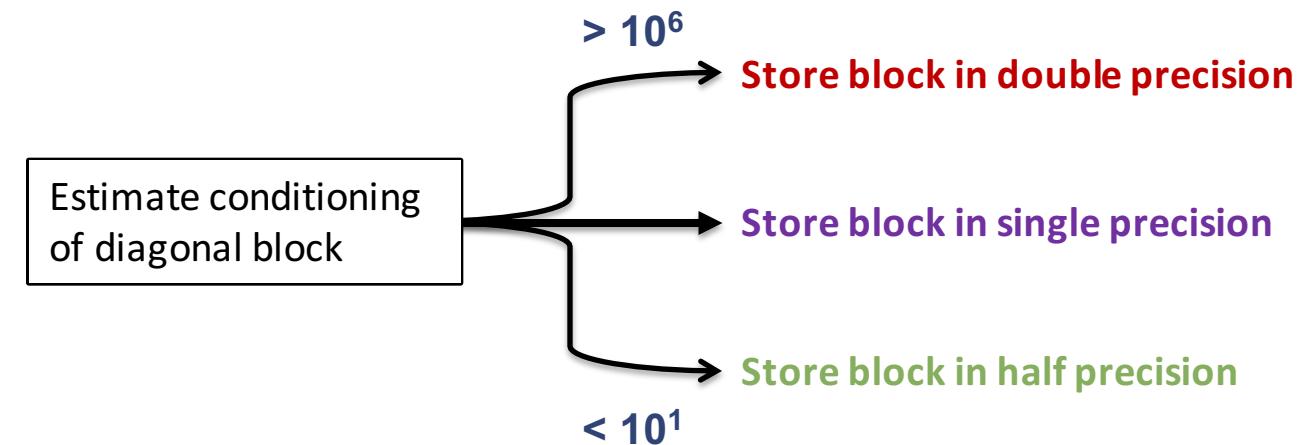
Adaptive Precision Block-Jacobi Preconditioning



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

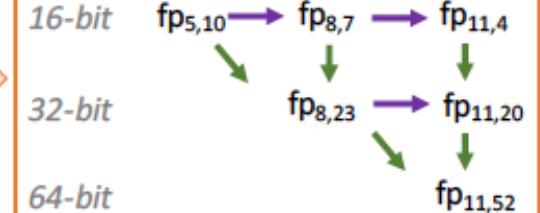
Multi-Precision Idea:

- 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:

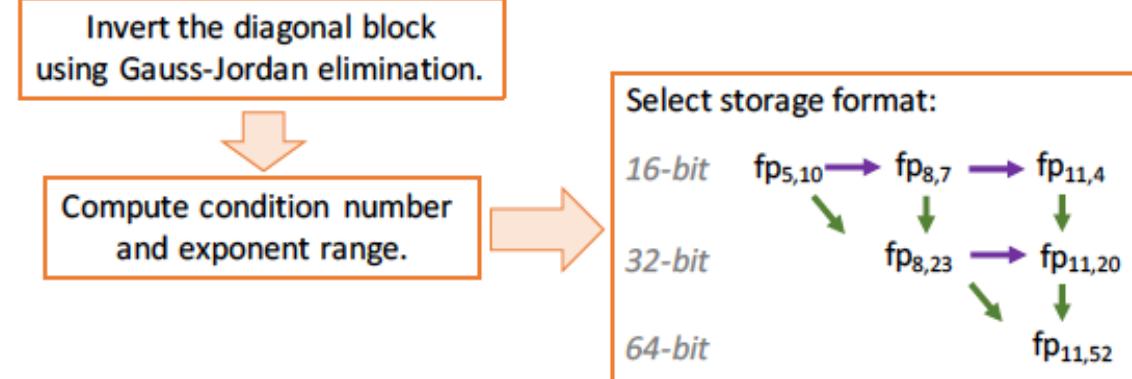


Flegar, Anzt, Quintana-Orti. "Customized-Precision Block-Jacobi Preconditioning for Krylov Iterative Solvers on Data-Parallel Manycore Processors". TOMS, submitted.

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.

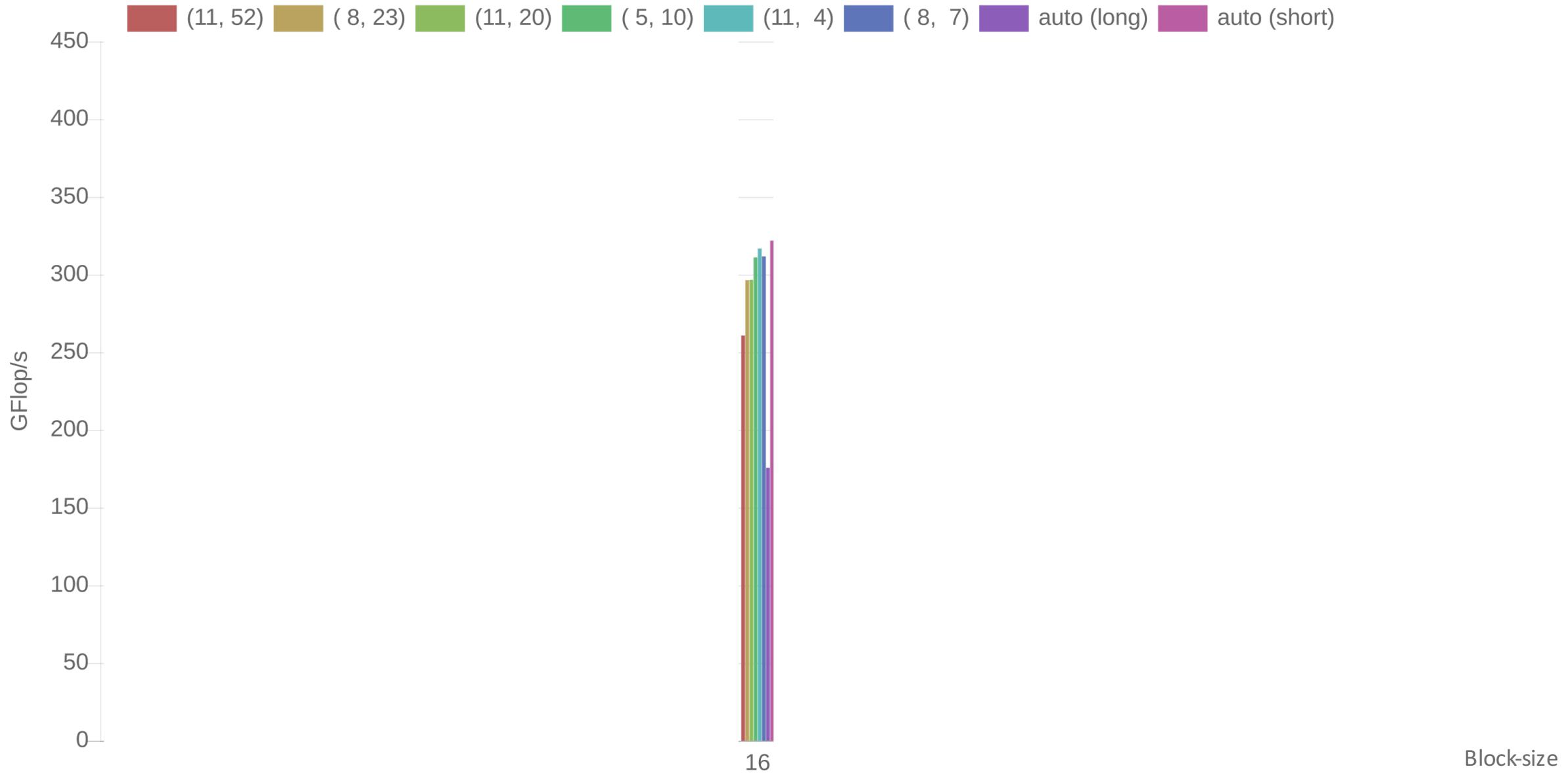
Precision distribution in Adaptive Block-Jacobi



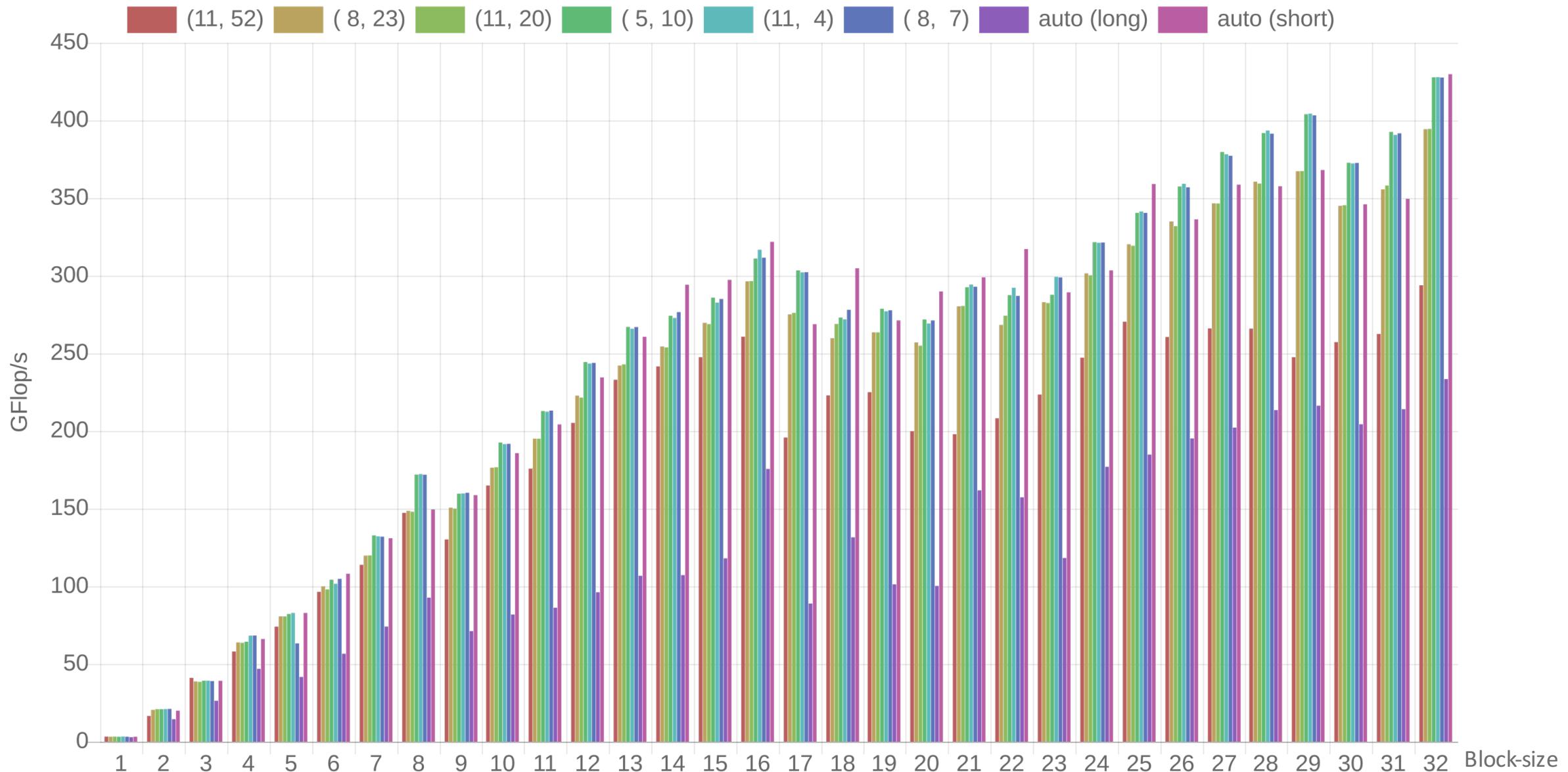
Precision distribution in Adaptive Block-Jacobi



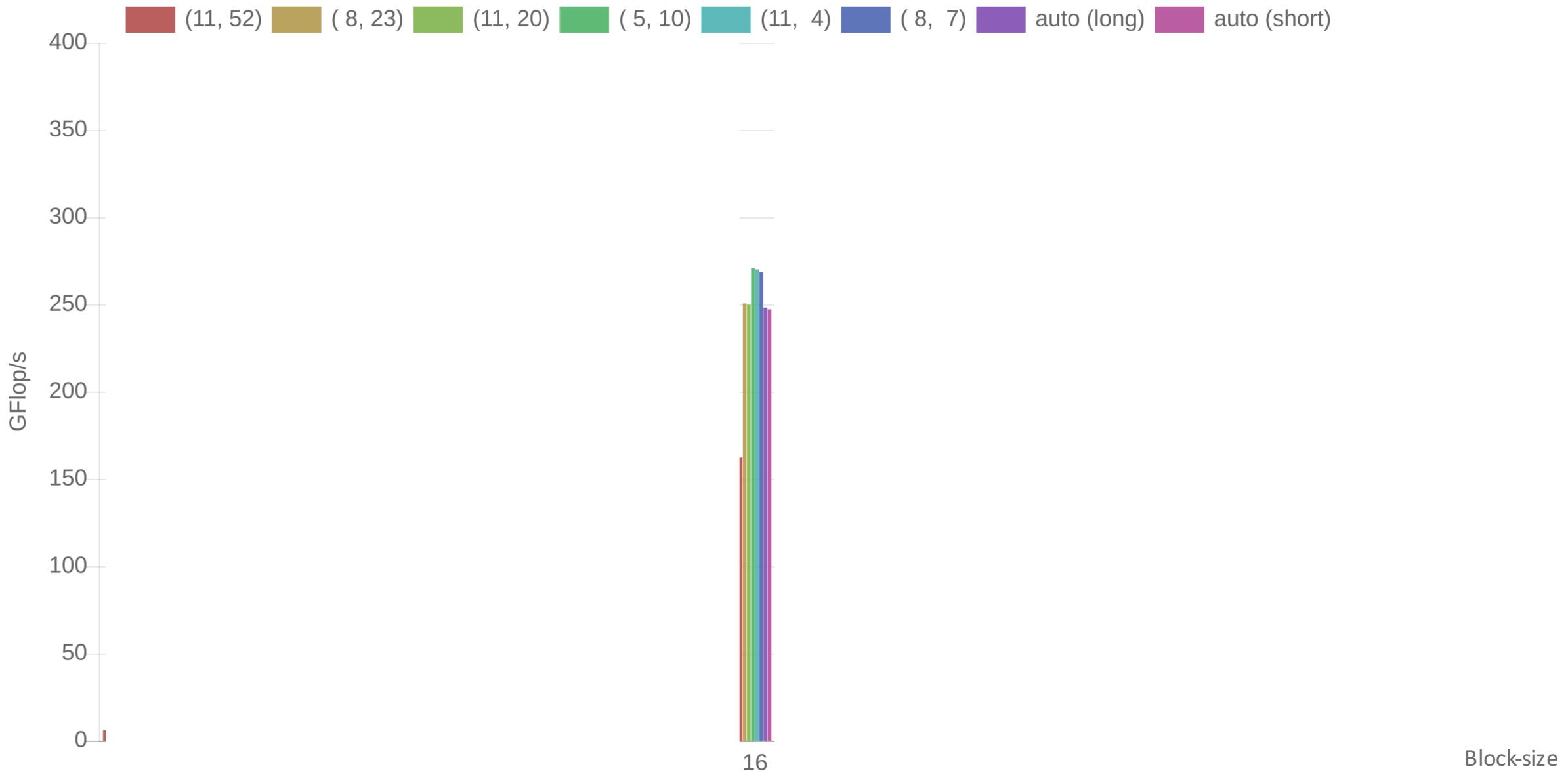
Adaptive Block-Jacobi Generation



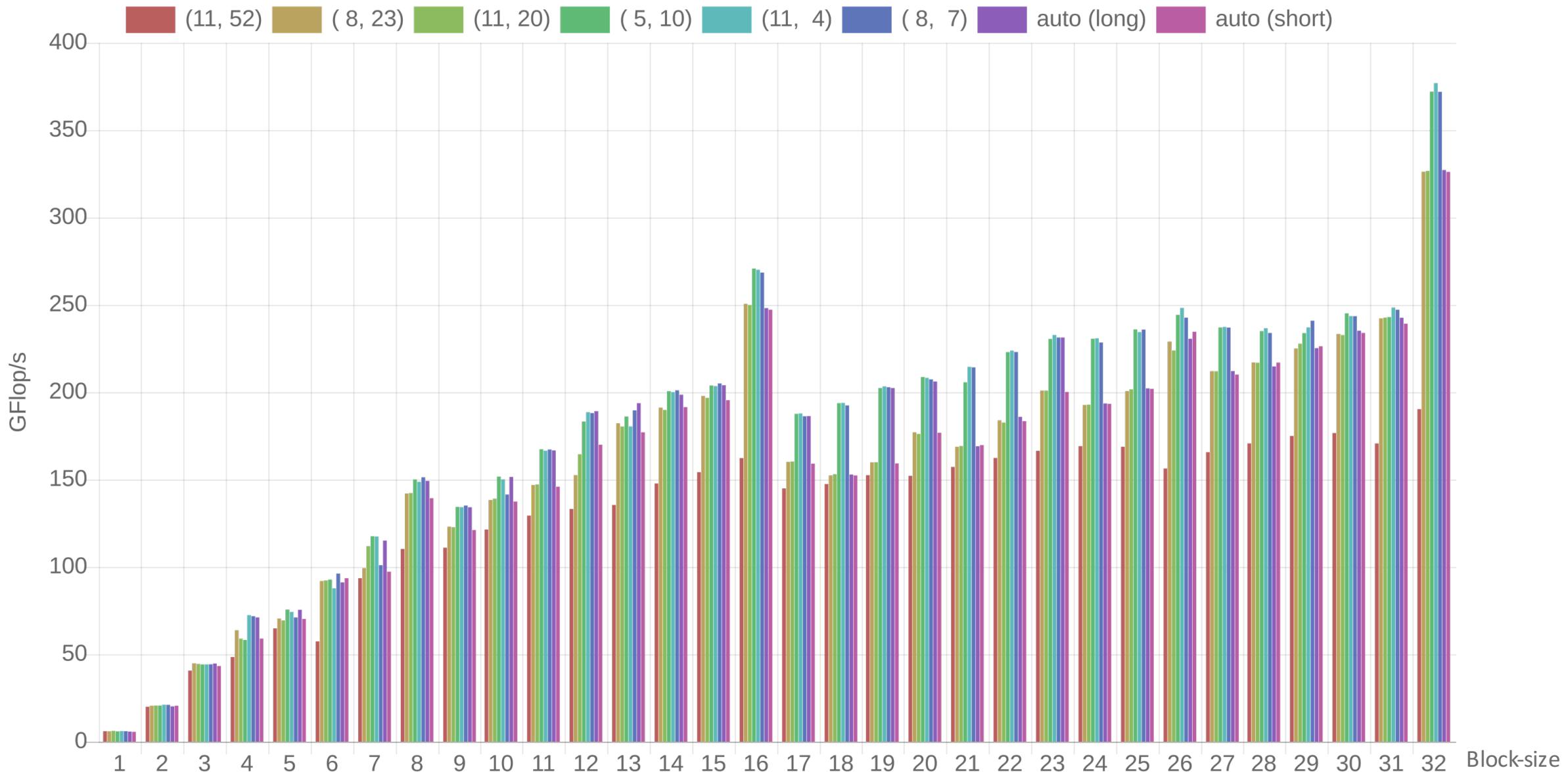
Adaptive Block-Jacobi Generation



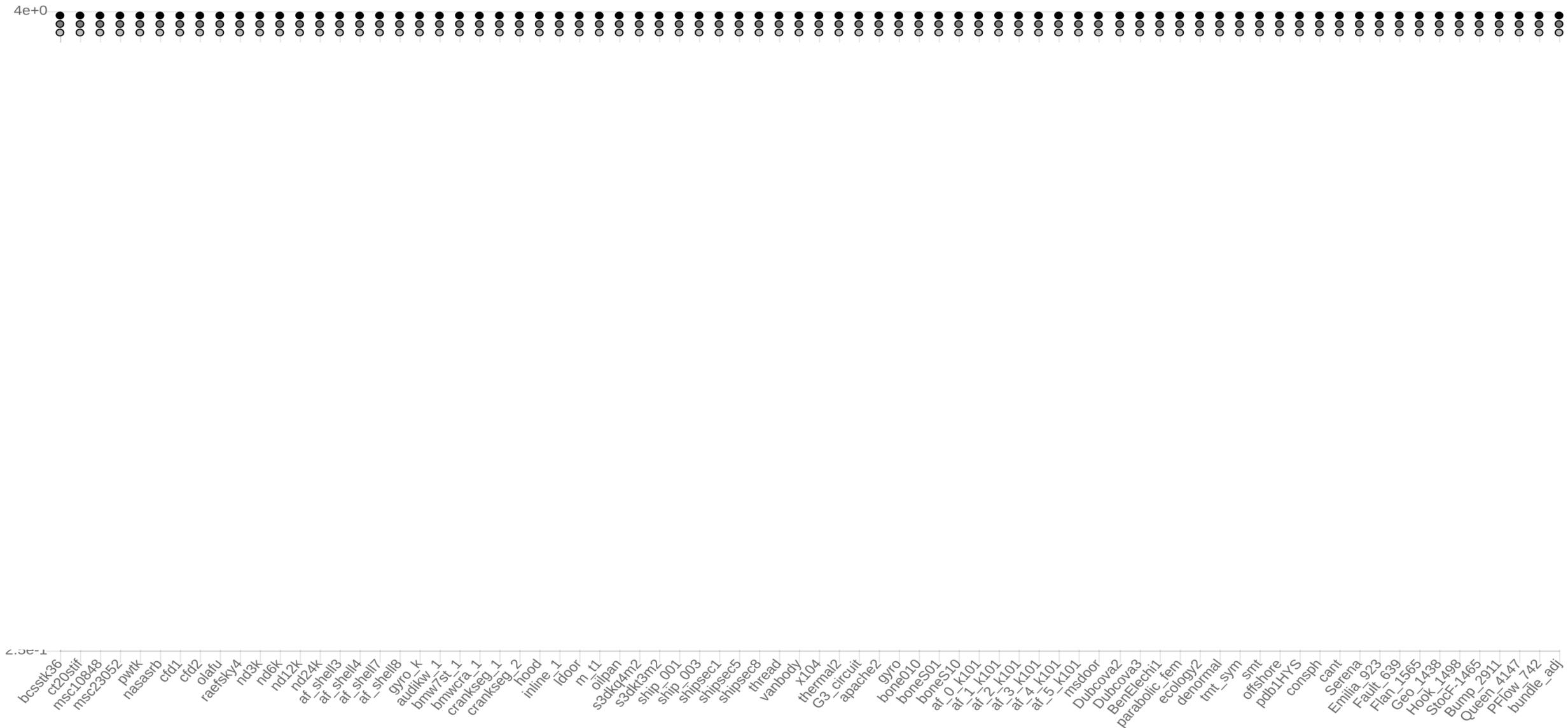
Adaptive Block-Jacobi Application



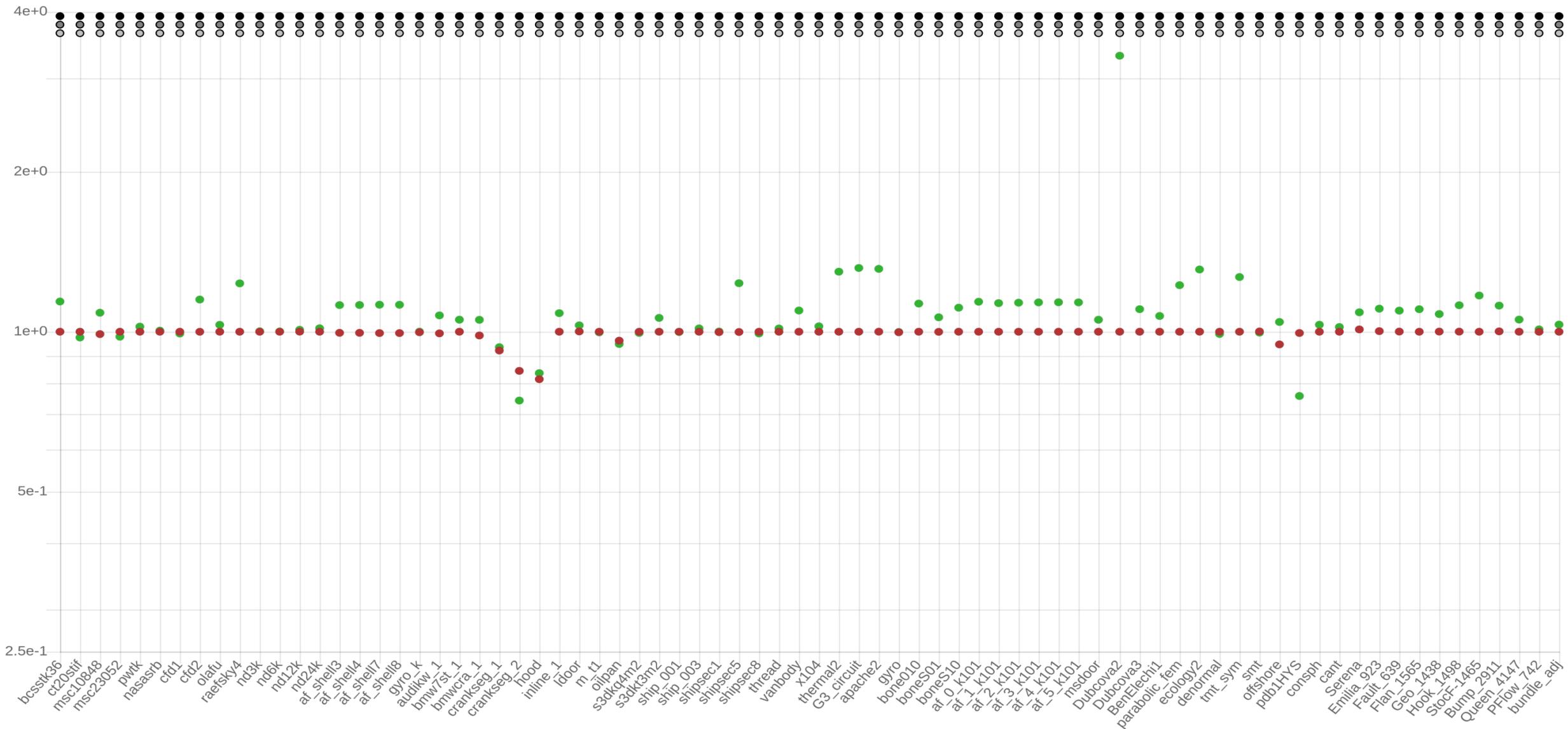
Adaptive Block-Jacobi Application

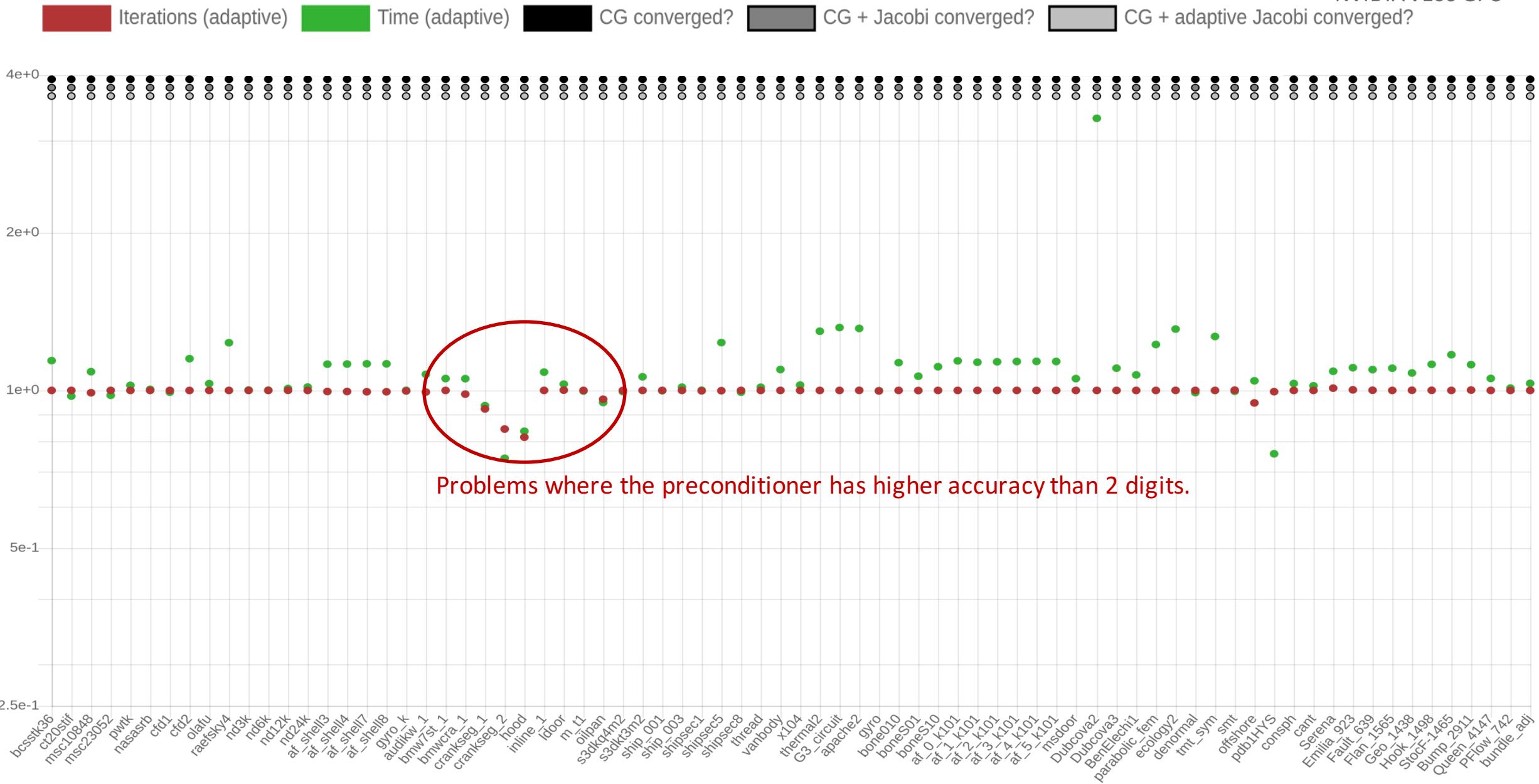


Iterations (adaptive) Time (adaptive) CG converged? CG + Jacobi converged? CG + adaptive Jacobi converged?

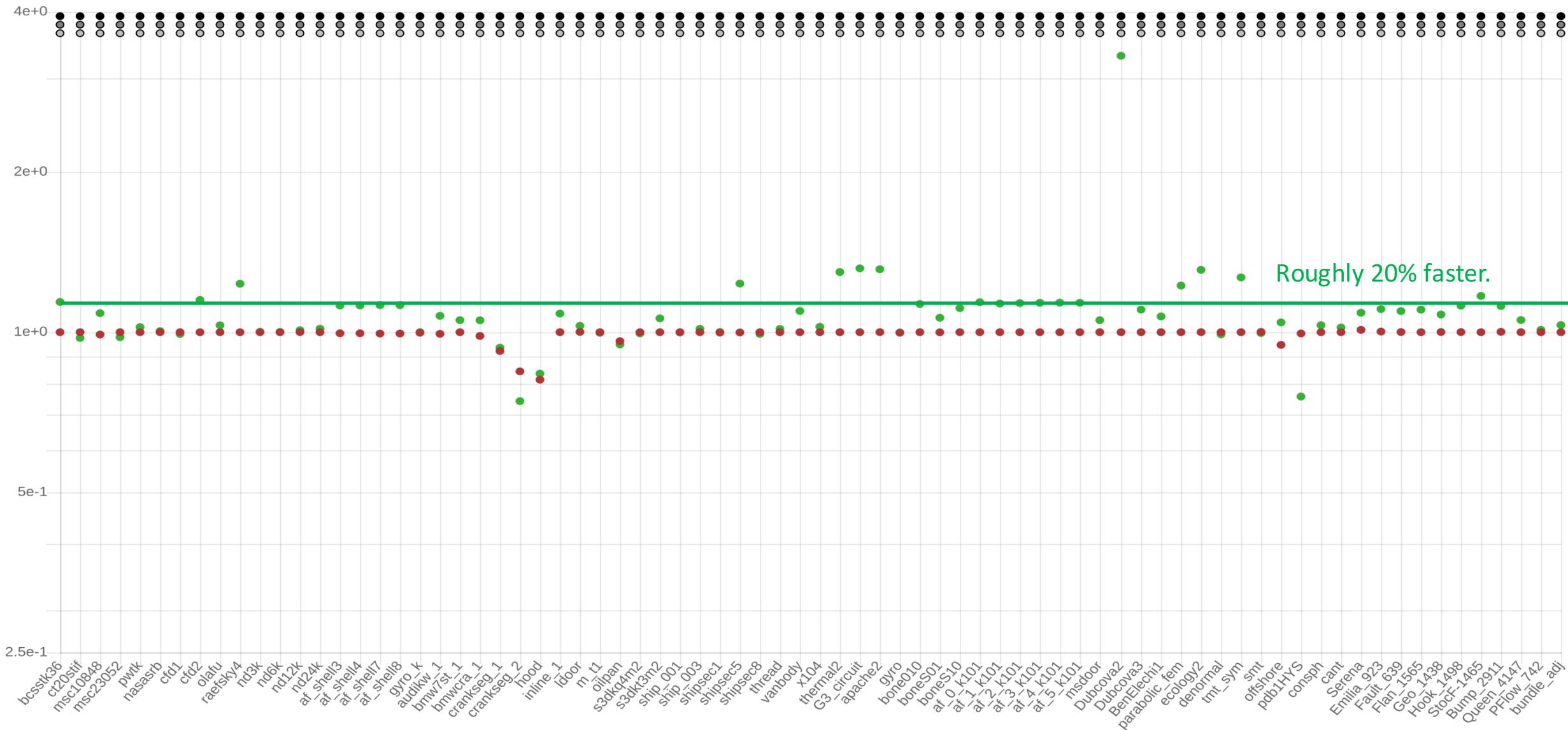


Iterations (adaptive) Time (adaptive) CG converged? CG + Jacobi converged? CG + adaptive Jacobi converged?

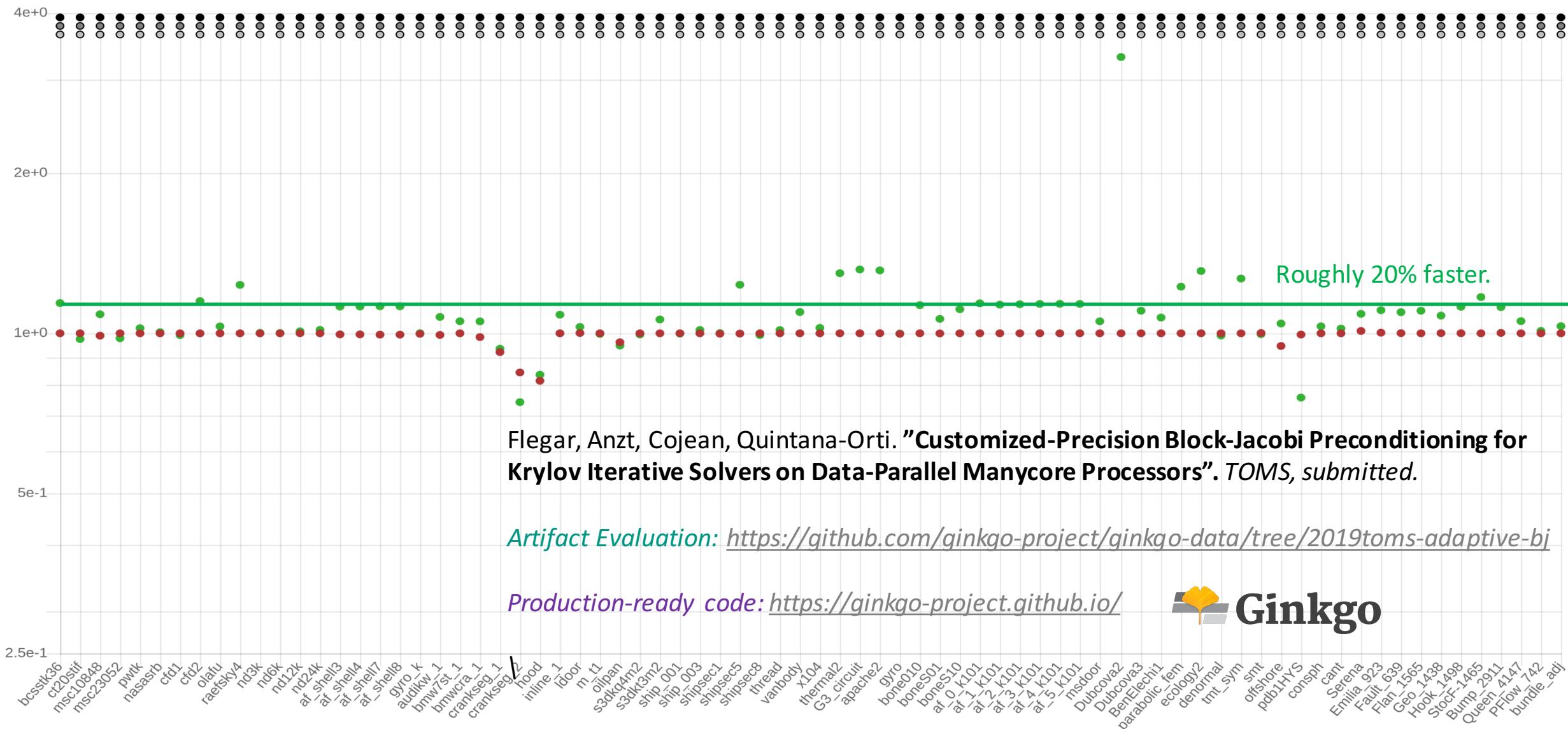




Iterations (adaptive) Time (adaptive) CG converged? CG + Jacobi converged? CG + adaptive Jacobi converged?



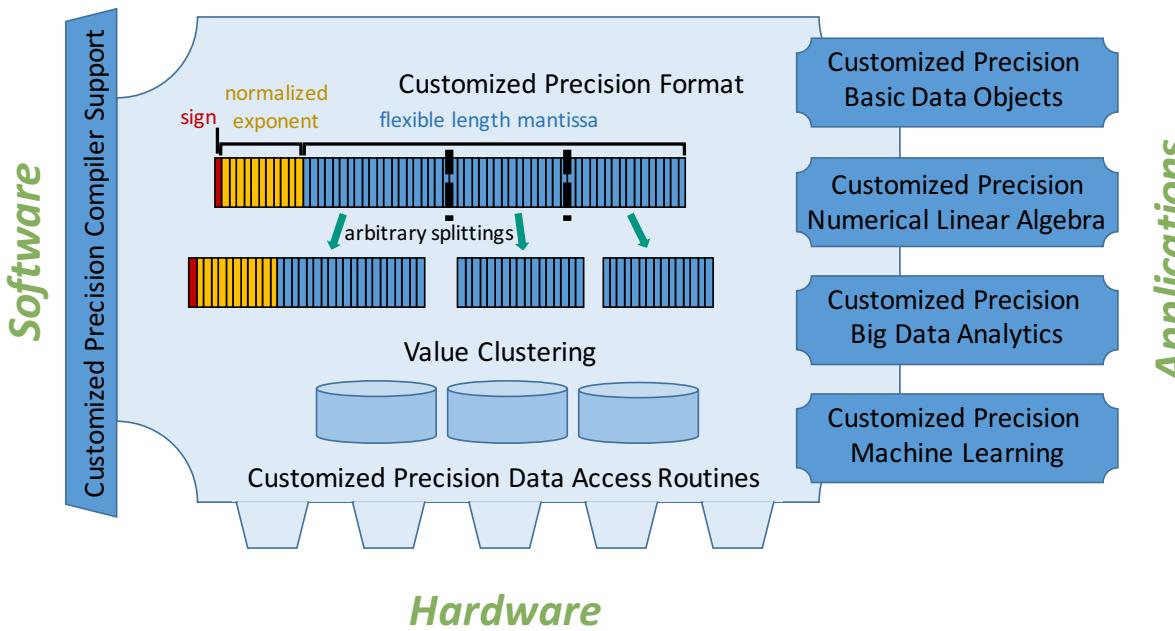
Iterations (adaptive) Time (adaptive) CG converged? CG + Jacobi converged? CG + adaptive Jacobi converged?



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>

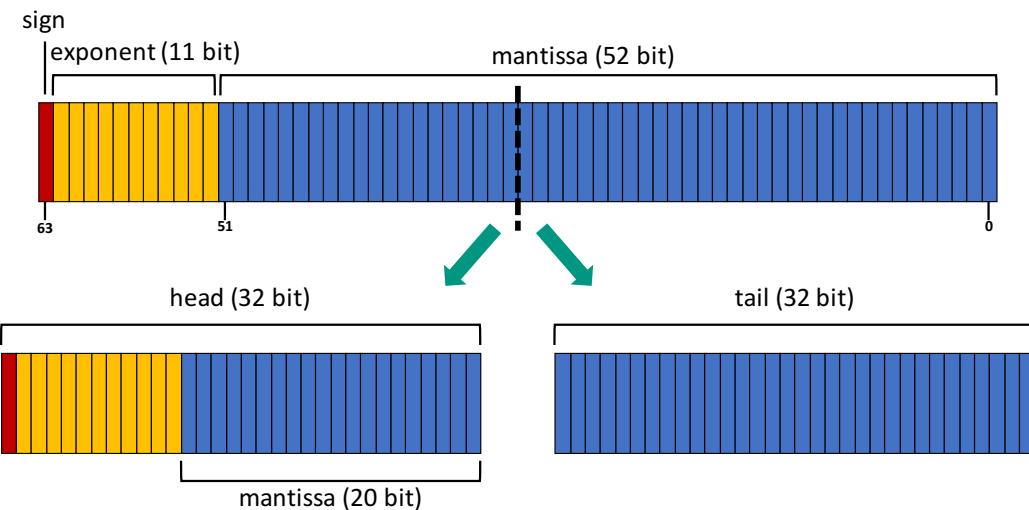


HELMHOLTZ
RESEARCH FOR GRAND CHALLENGES

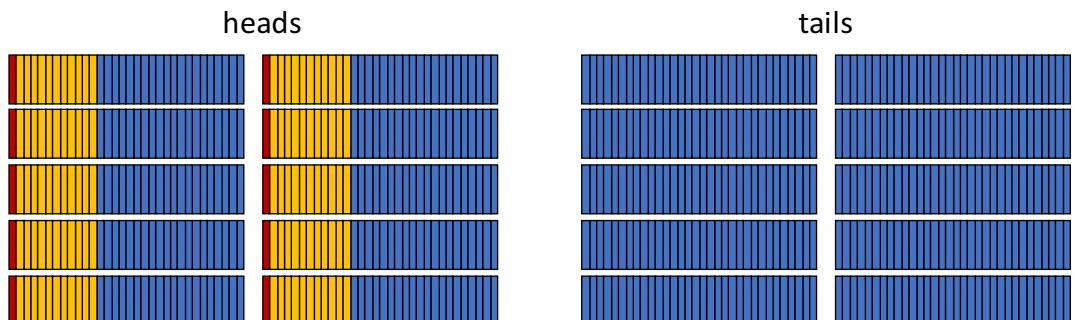
This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration and the Helmholtz Impuls und Vernetzungsfond VH-NG-1241.

Significand Segmentation

- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)

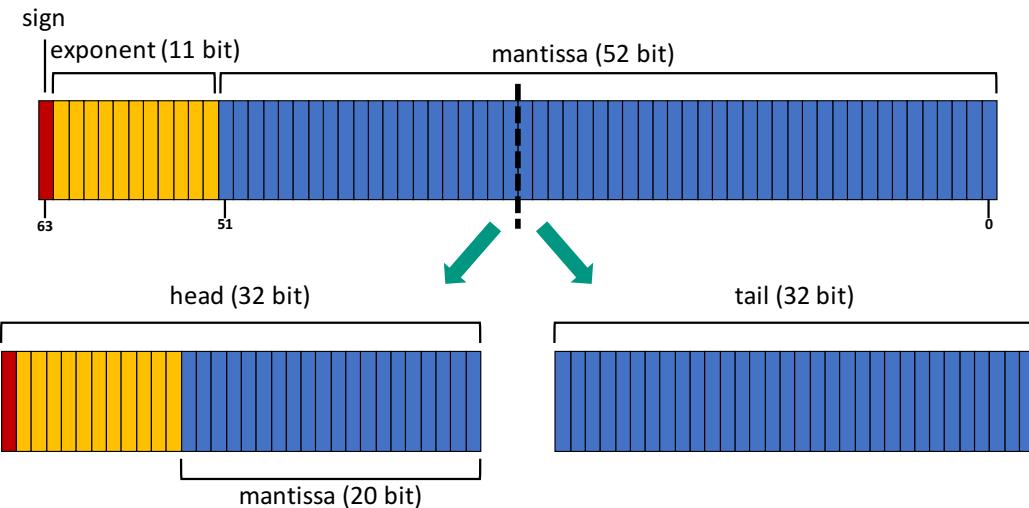


- Special “conversion” routines to double precision.
- Significand much shorter than IEEE single/half precision.
- No under- / overflow.

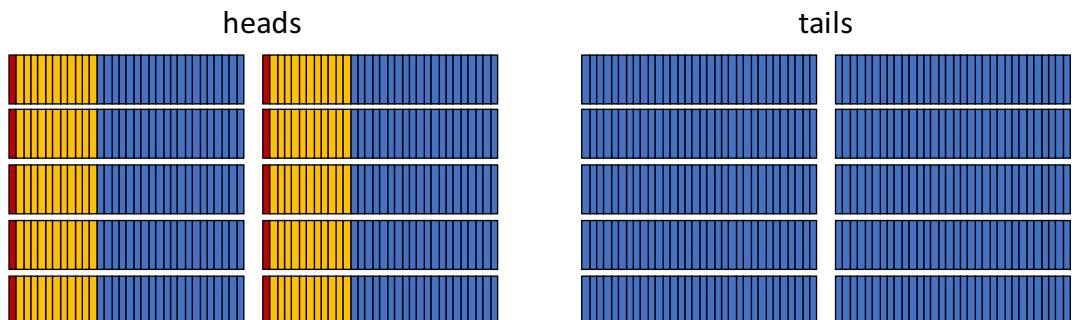


Significand Segmentation

- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.

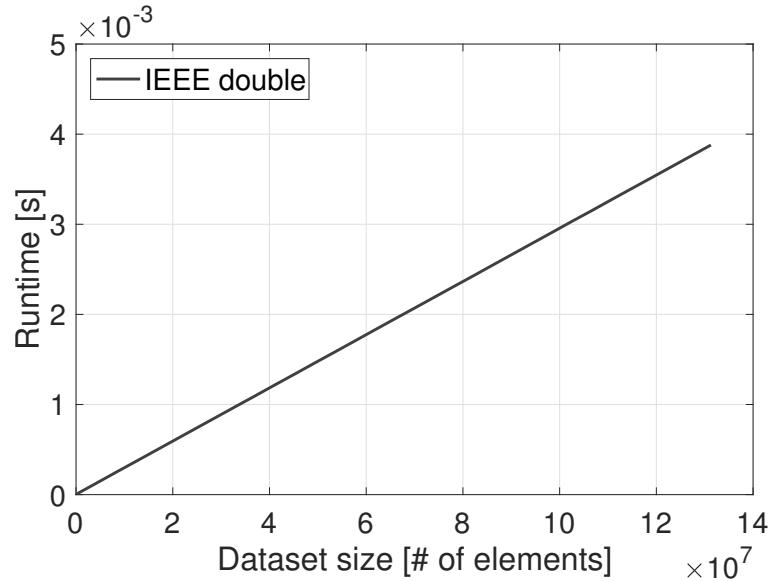


- Special “conversion” routines to double precision.
- Significand much shorter than IEEE single/half precision.
- No under- / overflow.

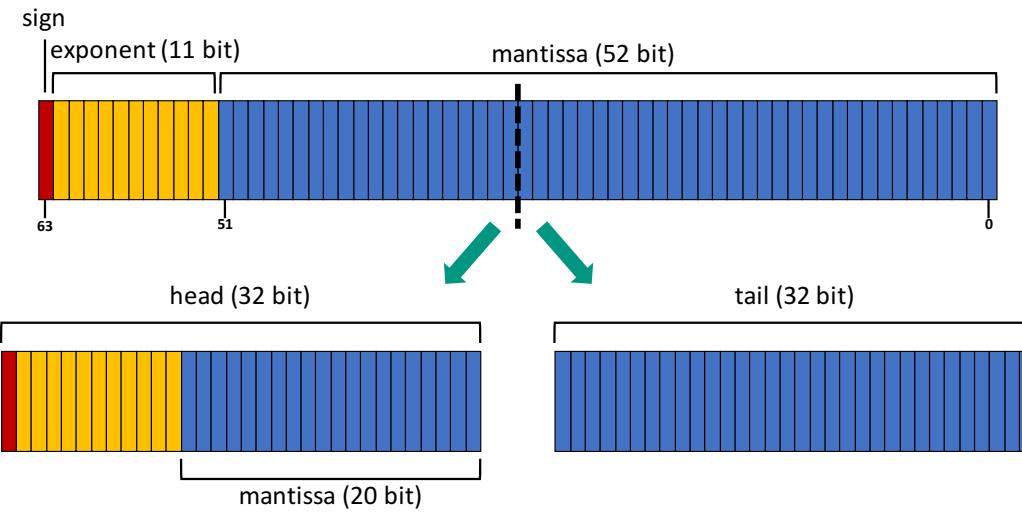


Significand Segmentation

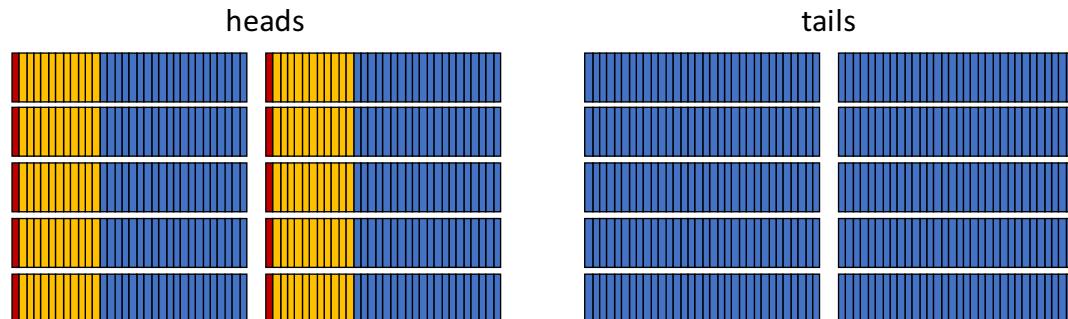
- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*



NVIDIA P100 "Pascal"
5.3 TFLOP/s DP
16GB RAM @720 GB/s

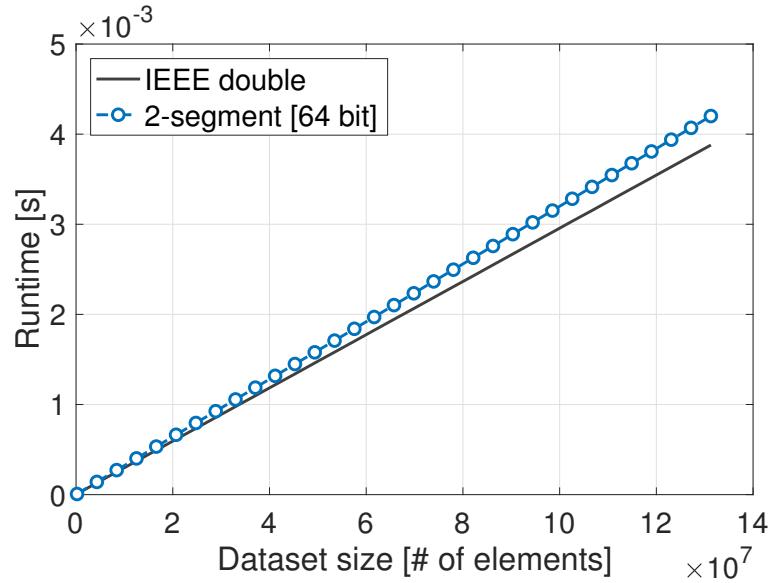


- *Special "conversion" routines to double precision.*
- *Significand much shorter than IEEE single/half precision.*
- *No under- / overflow.*

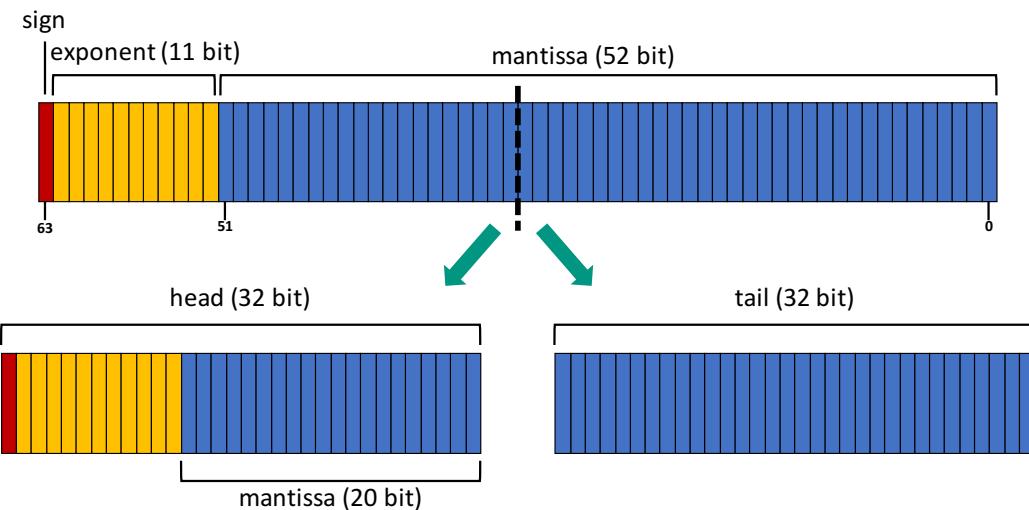


Significand Segmentation

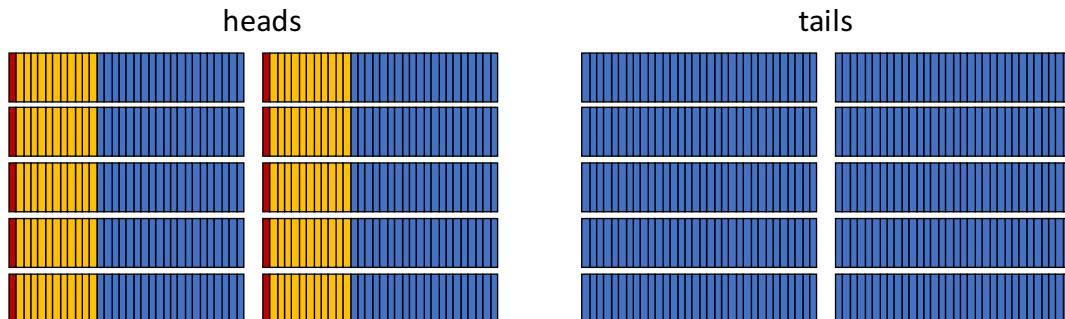
- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*



NVIDIA P100 "Pascal"
5.3 TFLOP/s DP
16GB RAM @720 GB/s

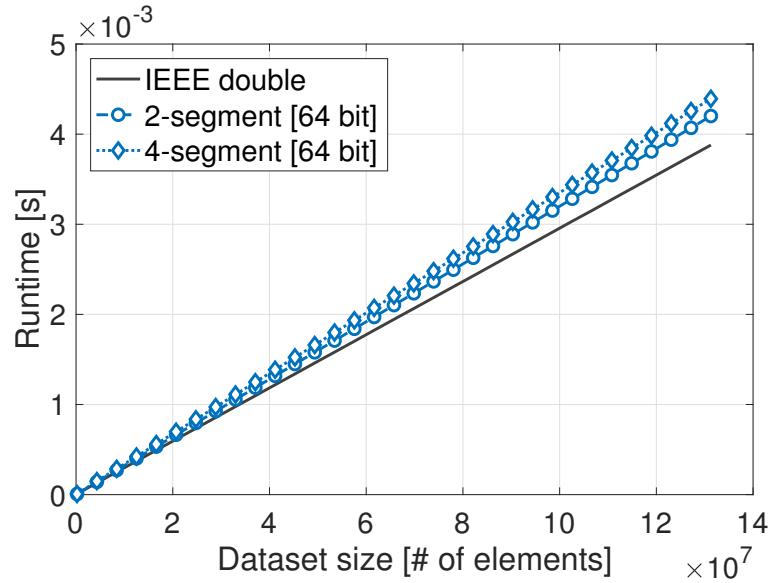


- *Special "conversion" routines to double precision.*
- *Significand much shorter than IEEE single/half precision.*
- *No under- / overflow.*

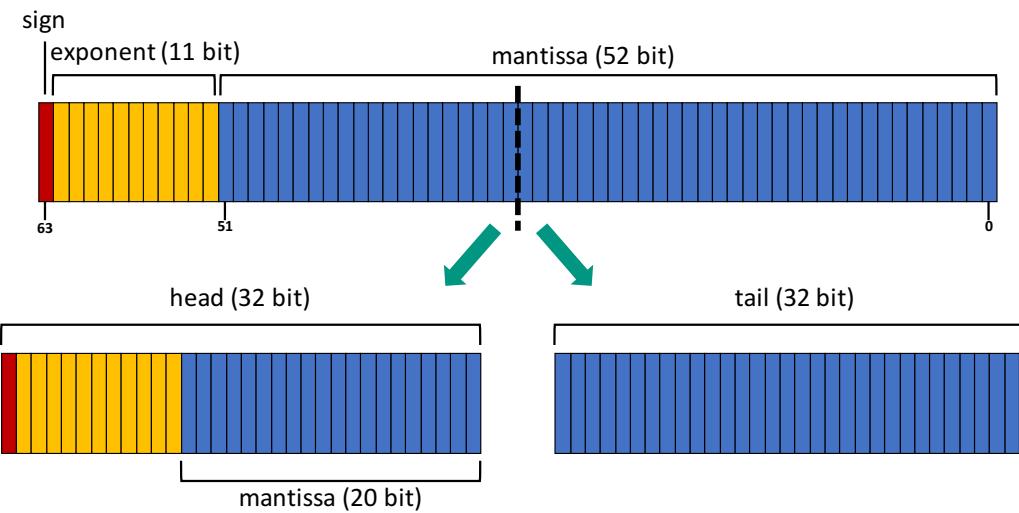


Significand Segmentation

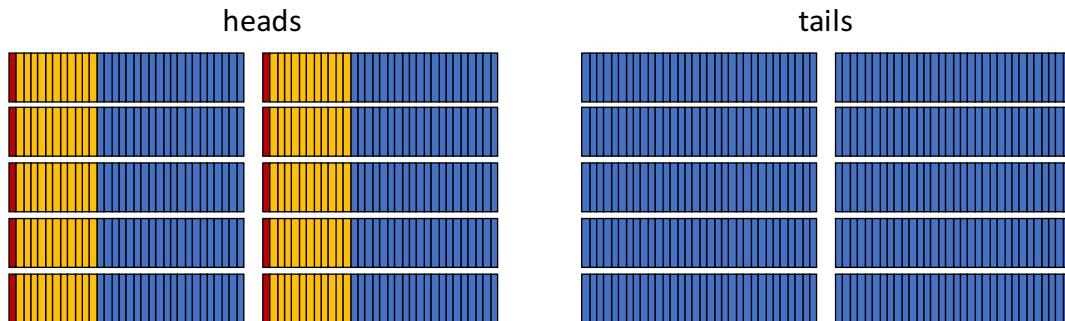
- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*



NVIDIA P100 "Pascal"
5.3 TFLOP/s DP
16GB RAM @720 GB/s

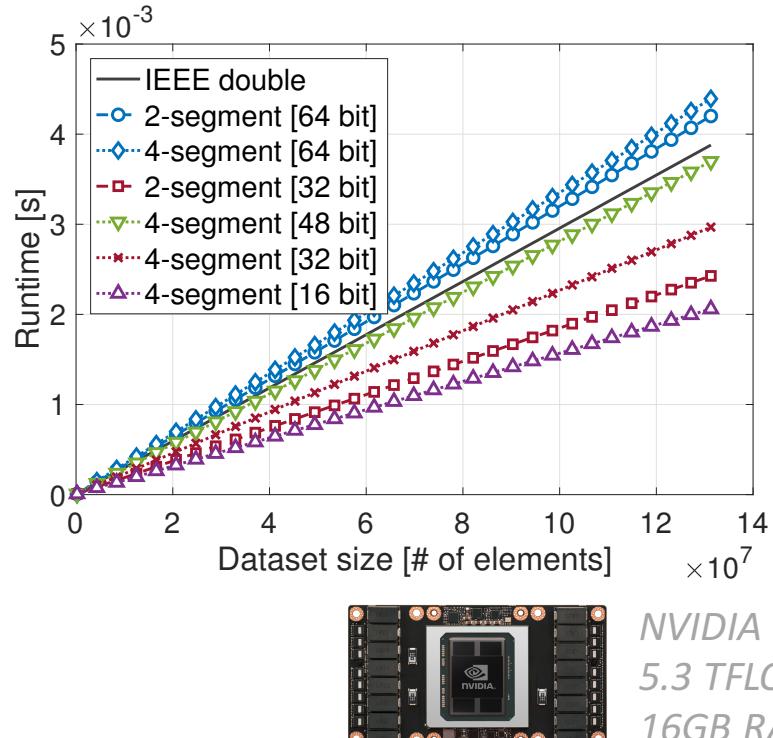


- *Special "conversion" routines to double precision.*
- *Significand much shorter than IEEE single/half precision.*
- *No under- / overflow.*

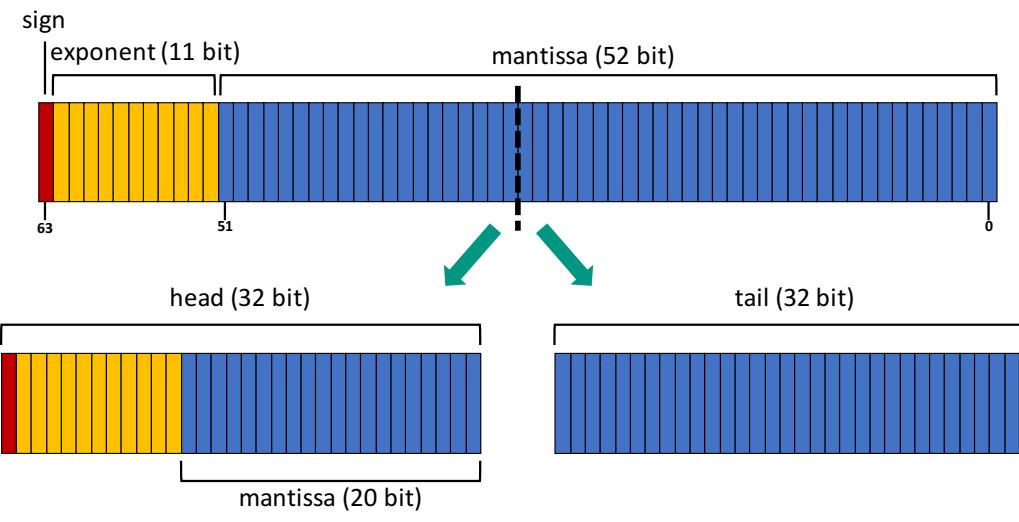


Significand Segmentation

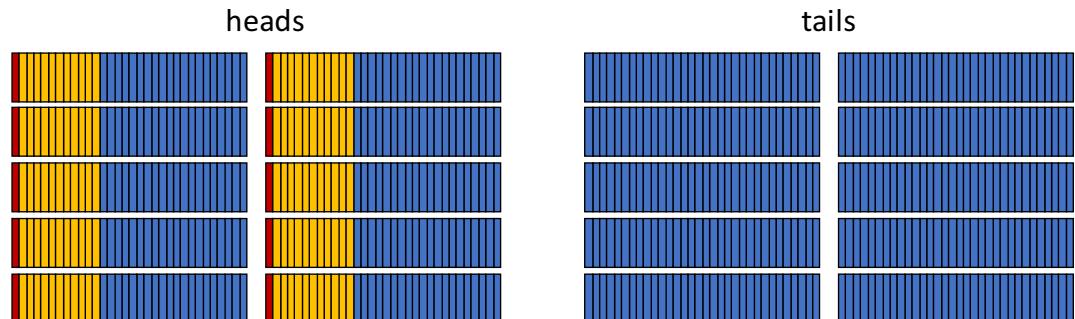
- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*



NVIDIA P100 "Pascal"
5.3 TFLOP/s DP
16GB RAM @720 GB/s

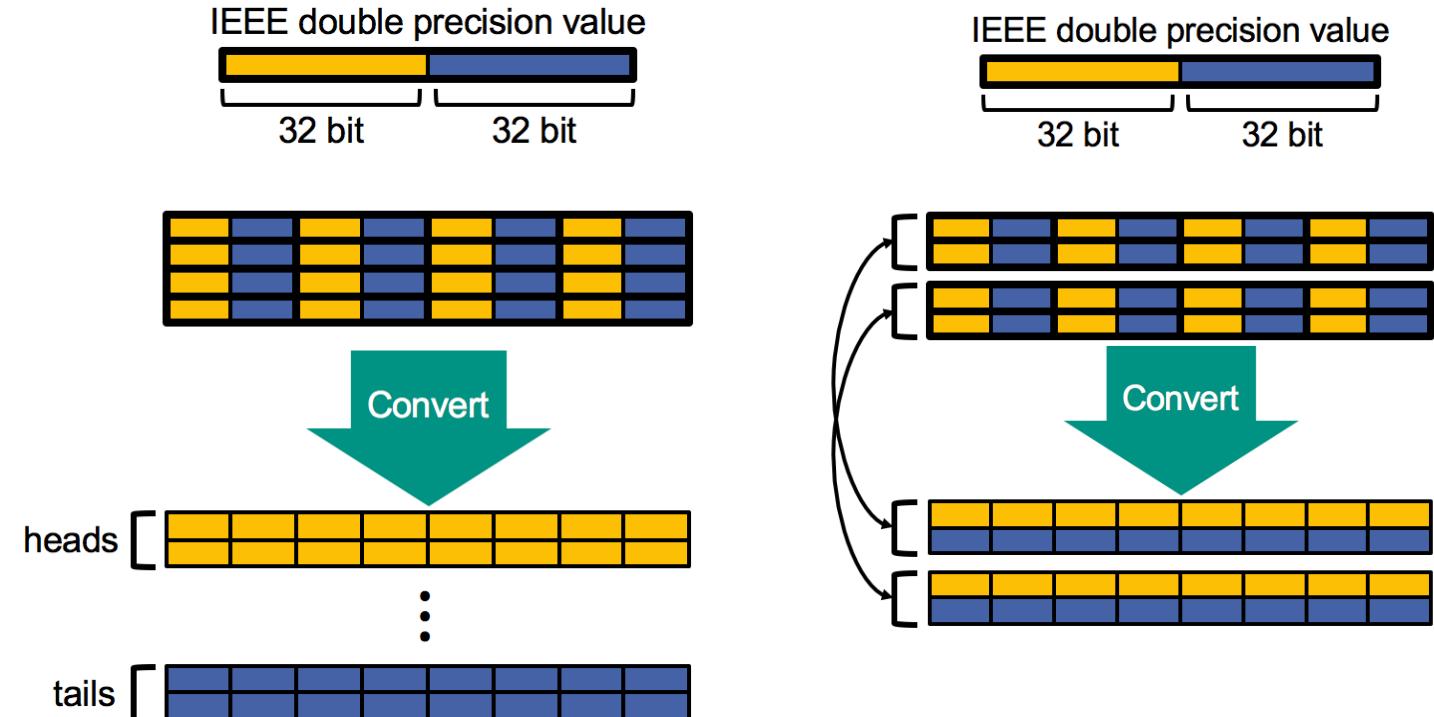


- *Special "conversion" routines to double precision.*
- *Significand much shorter than IEEE single/half precision.*
- *No under- / overflow.*



Significant Segmentation

- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*
- *Special “conversion” routines to double precision.*
- *Mantissa much shorter than IEEE single/half precision.*
- *No under- / overflow.*



Customized Precision Adaptive Jacobi

- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*
- Need precision-aware algorithms, e.g. Adaptive Precision Jacobi¹.
 - *Use low precision data reads whenever accuracy allows for it.*

Adaptive precision solvers for sparse linear systems

Anzt, Dongarra, Quintana-Orti, Published in SC-Workshops in 2015.

- *Start iterations with reading low precision.*
- *Increase accuracy over iteration process adapted to requirements.*
- *Adapt the accuracy on a component level.*

¹Anzt, Dongarra, Quintana-Orti. "Adaptive precision solvers for sparse linear systems". In: Proceedings of the 3rd International Workshop on Energy Efficient Supercomputing. (2015).

Customized Precision Adaptive Jacobi

- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*
- Need precision-aware algorithms, e.g. Adaptive Precision Jacobi¹.
 - *Use low precision data reads whenever accuracy allows for it.*

Adaptive precision solvers for sparse linear systems

Anzt, Dongarra, Quintana-Orti, Published in SC-Workshops in 2015.

- Start iterations with reading low precision.
- Increase accuracy over iteration process adapted to requirements.
- Adapt the accuracy on a component level.

Jacobi relaxation

Iterative method for solving $Ax = b$:

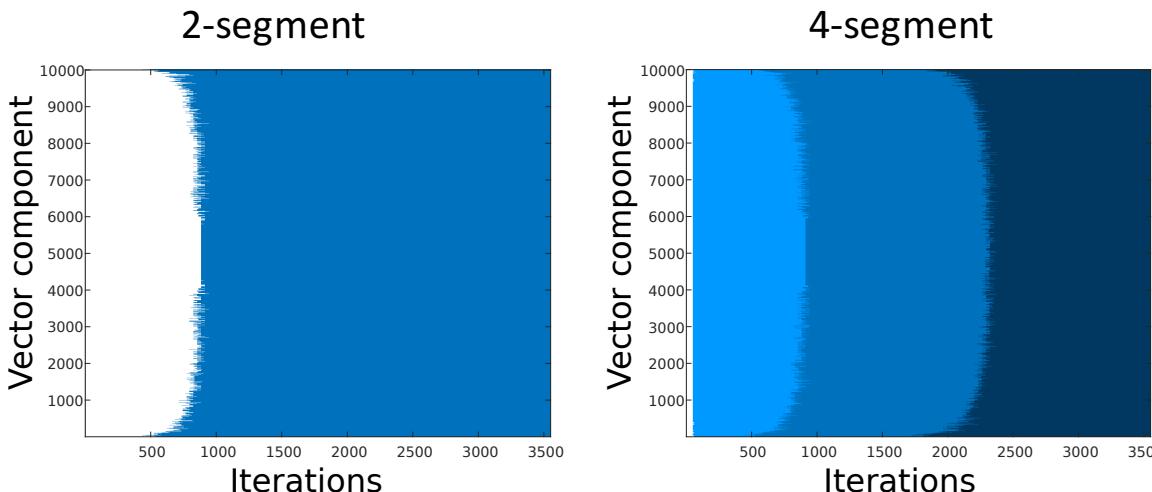
$$x^{k+1} = D^{-1} (b - (A - D)x^k)$$

$$x+ = -D^{-1} Ax + c$$

¹Anzt, Dongarra, Quintana-Orti. "Adaptive precision solvers for sparse linear systems". In: Proceedings of the 3rd International Workshop on Energy Efficient Supercomputing. (2015).

Customized Precision Adaptive Jacobi

- Split the IEEE double precision format into segments.
(2-segment customized precision, 4-segment customized precision...)
- For coalesced data access, interleave data in memory.
- *Data can be accessed much faster if low precision is acceptable.*
- Need precision-aware algorithms, e.g. Adaptive Precision Jacobi¹.
 - *Use low precision data reads whenever accuracy allows for it.*

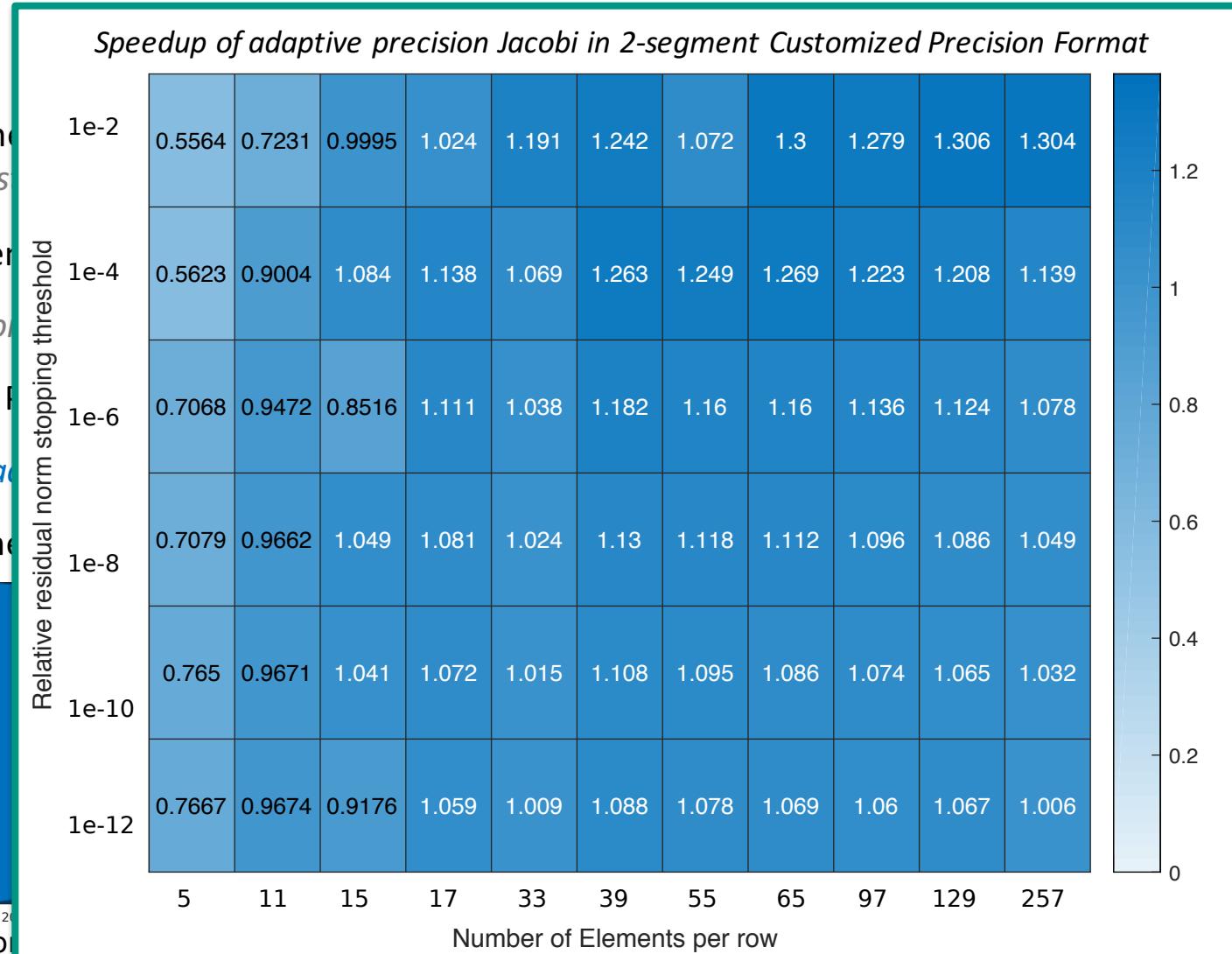
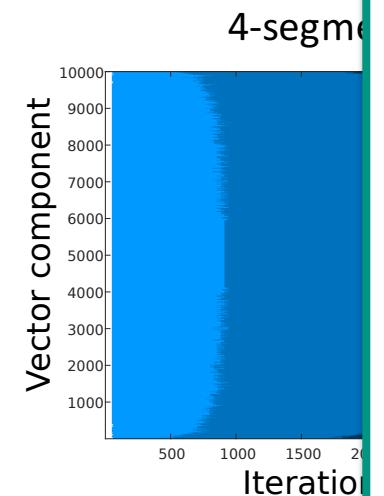
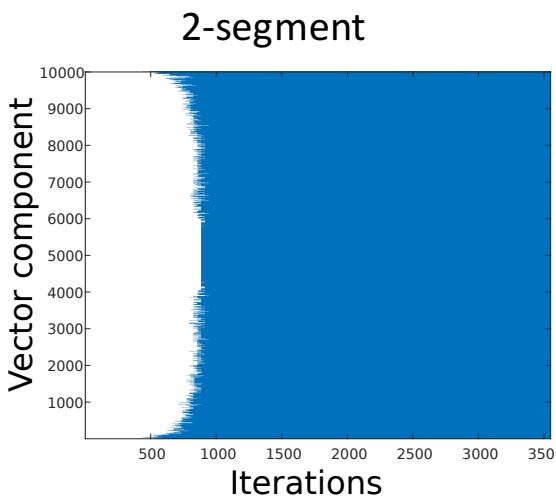


¹Anzt, Dongarra, Quintana-Orti. "Adaptive precision solvers for sparse linear systems". In: Proceedings of the 3rd International Workshop on Energy Efficient Supercomputing. (2015).

²Grützmacher, Anzt. "A Customized Precision Format for decoupling Arithmetic Format and Storage Format". Submitted to HeteroPar workshop 2018.

Customized Precision Adaptive Jacobi

- Split the IEEE double precision format into segments (2-segment customized precision, 4-segment customized precision)
- For coalesced data access, interleave data in memory
- *Data can be accessed much faster if low precision is used*
- Need precision-aware algorithms, e.g. Adaptive Jacobi
- *Use low precision data reads whenever appropriate*

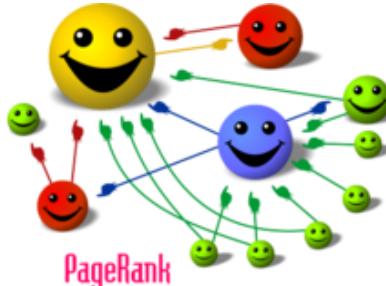


¹Anzt, Dongarra, Quintana-Orti. "Adaptive precision solvers for sparse linear systems". In: Proceedings of the 3rd International Workshop on Energy Efficient Supercomputing. (2015).

²Grützmacher, Anzt. "A Customized Precision Format for decoupling Arithmetic Format and Storage Format". Submitted to HeteroPar workshop 2018.

Customized Precision PageRank

- Algorithm ranking web-pages based on “importance.”
- Underlying algorithm is a eigenvalue solver based on the **Power Iteration**.
- **Fixed-point iteration type.**
- **Sparse matrix vector product** is the central and most expensive building block.



$\text{PageRank}(A, \varepsilon, \delta) :$

$\mathbb{S} = \text{All indices where } A \text{ has a empty row}$

$$p^{\{0\}} = \frac{1}{n} \cdot e$$

$k = 0$

repeat

$$k = k + 1$$

$$s = \sum_{i \in \mathbb{S}} p_i^{\{k-1\}}$$

$$p^{\{k\}} = \delta \cdot A^T \cdot p^{\{k-1\}} + \frac{(1 - \delta)}{n} \cdot e + \frac{s}{n} \cdot e$$

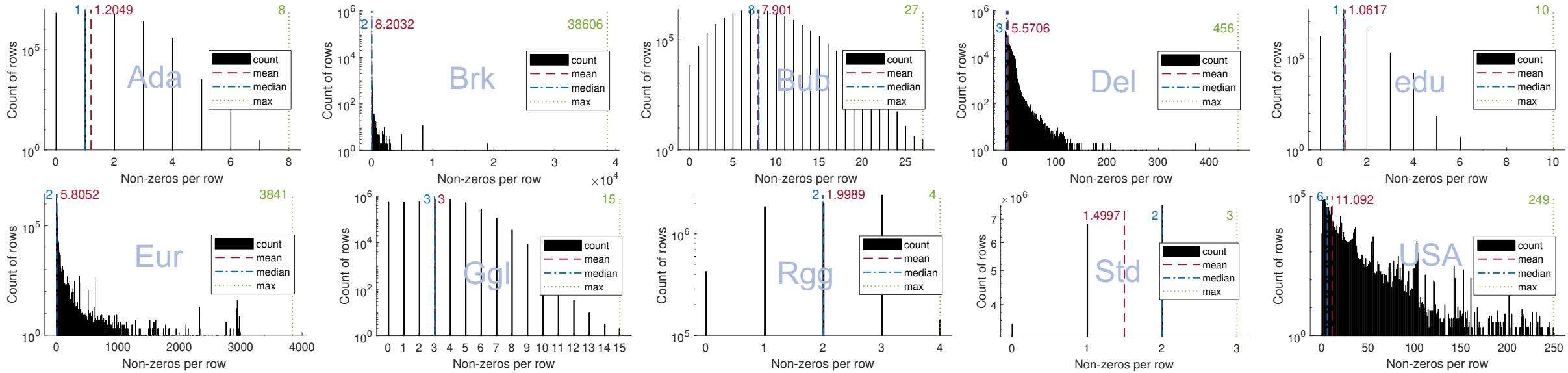
$$\gamma = \|p^{\{k\}} - p^{\{k-1\}}\|_1$$

until ($\gamma < \varepsilon$)

return $p^{\{k\}}$

The PageRank citation ranking: Bringing order to the Web L. Page, S. Brin, R. Motwani, and T. Winograd.
Proceedings of the 7th International World Wide Web Conference, page 161--172. Brisbane, Australia, (1998)

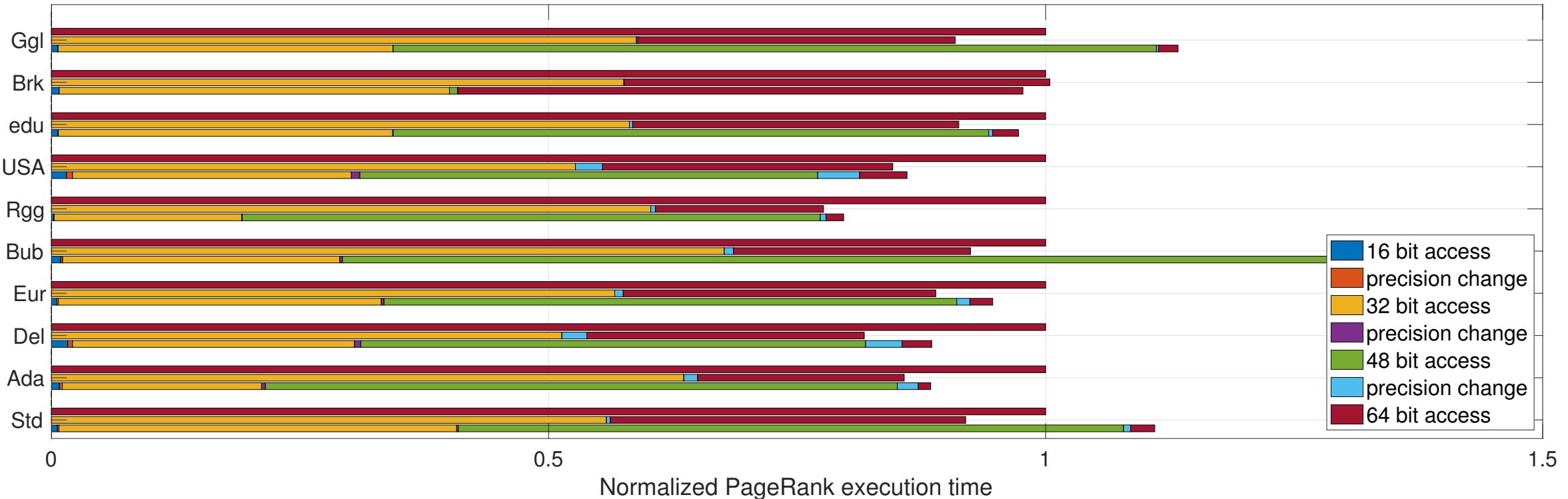
Customized Precision PageRank



10 test problems from Suite Sparse (community networks, etc.)

Customized Precision PageRank

Trace of 2-segment and 4-segment customized precision PageRank (NVIDIA V100 GPU, Relative residual stopping criterion 1e-12).



Benchmark results on NVIDIA V100 GPU

Submitting to SC Workshop on Irregular Applications: "Customized Precision PageRank on GPUs".

Customized Precision PageRank

PageRank speedup over reference implementation (NVIDIA V100 GPU). CSR format.

Stopping threshold	Std	Ada	Del	Eur	Bub	Rgg	USA	edu	Brk	Ggl
1e-2	1.164	1.469	1.373	1.27	1.176	2.042	1.311	1.218	0.9887	1.326
1e-4	1.209	1.551	1.469	1.318	1.24	2.194	1.387	1.23	0.977	1.266
1e-6	1.161	1.408	1.38	1.214	1.173	1.744	1.31	1.179	0.9798	1.197
1e-8	1.103	1.236	1.264	1.156	1.099	1.415	1.203	1.117	0.9857	1.133
1e-10	1.101	1.157	1.194	1.116	1.068	1.281	1.151	1.087	0.9902	1.1
1e-12	1.061	1.112	1.159	1.091	1.054	1.212	1.129	1.069	0.9916	1.073

Matrices
2-segment

Stopping threshold	Std	Ada	Del	Eur	Bub	Rgg	USA	edu	Brk	Ggl
1e-2	1.08	1.482	1.416	1.231	1.127	2.099	1.347	1.189	1.065	1.27
1e-4	1.076	1.597	1.402	1.261	1.068	1.997	1.387	1.182	1.041	1.186
1e-6	1.004	1.417	1.283	1.162	0.8915	1.682	1.324	1.127	1.03	1.05
1e-8	0.9283	1.199	1.184	1.104	0.7465	1.366	1.244	1.062	1.022	0.9356
1e-10	0.9118	1.129	1.105	1.05	0.6929	1.251	1.143	1.021	1.017	0.8828
1e-12	0.906	1.11	1.087	1.048	0.7395	1.191	1.121	1.025	1.013	0.8972

Matrices
4-segment

Customized Precision PageRank

PageRank speedup over reference implementation (NVIDIA V100 GPU). ELL format.

Stopping threshold	Ada	Del	Eur	Bub	Rgg	USA	Std
1e-2	1.442	1.486	1.465	1.293	1.477	1.505	1.514
1e-4	1.446	1.486	1.466	1.298	1.479	1.508	1.517
1e-6	1.266	1.267	1.287	1.177	1.295	1.3	1.35
1e-8	1.158	1.191	1.202	1.101	1.185	1.198	1.224
1e-10	1.105	1.142	1.148	1.07	1.132	1.149	1.165
1e-12	1.078	1.118	1.116	1.055	1.103	1.126	1.13

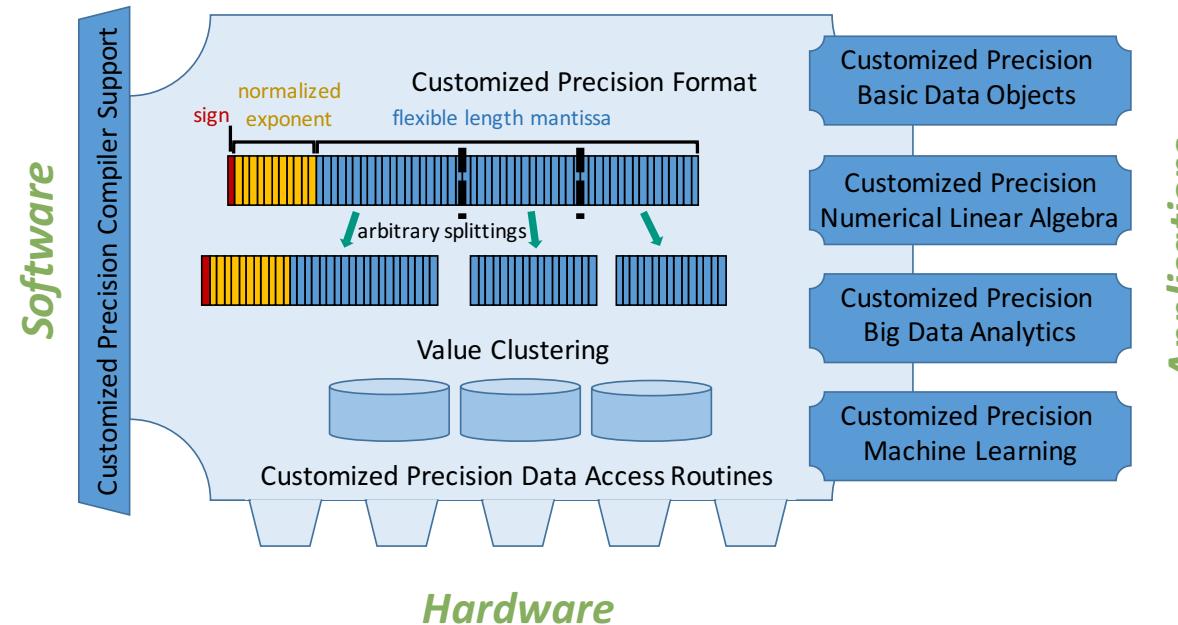
Matrices
2-segment

Stopping threshold	Ada	Del	Eur	Bub	Rgg	USA	Std
1e-2	1.515	1.693	1.496	1.308	1.583	1.7	1.284
1e-4	1.55	1.568	1.448	1.188	1.496	1.613	1.223
1e-6	1.369	1.451	1.329	0.9664	1.421	1.499	1.146
1e-8	1.175	1.351	1.262	0.7969	1.299	1.384	1.067
1e-10	1.091	1.219	1.186	0.7305	1.233	1.215	1.021
1e-12	1.084	1.179	1.154	0.7735	1.178	1.181	1.016

Matrices
4-segment

Summary and next steps

- Decouple arithmetic precision from memory precision.
- Using customized precisions for memory operations.
- Dynamically adapt the accuracy to the algorithm properties.
- Speedup of up to 1.3x for Jacobi iterations and PageRank algorithm.
- Speedup of up to 1.3x for adaptive precision block-Jacobi preconditioning.
- Creating a Modular Precision Ecosystem inside  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.

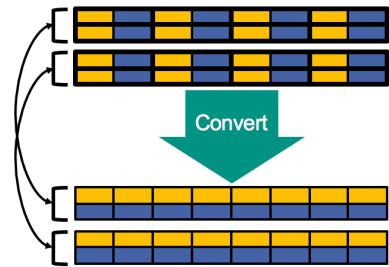
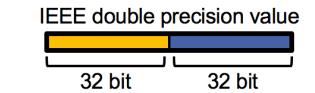
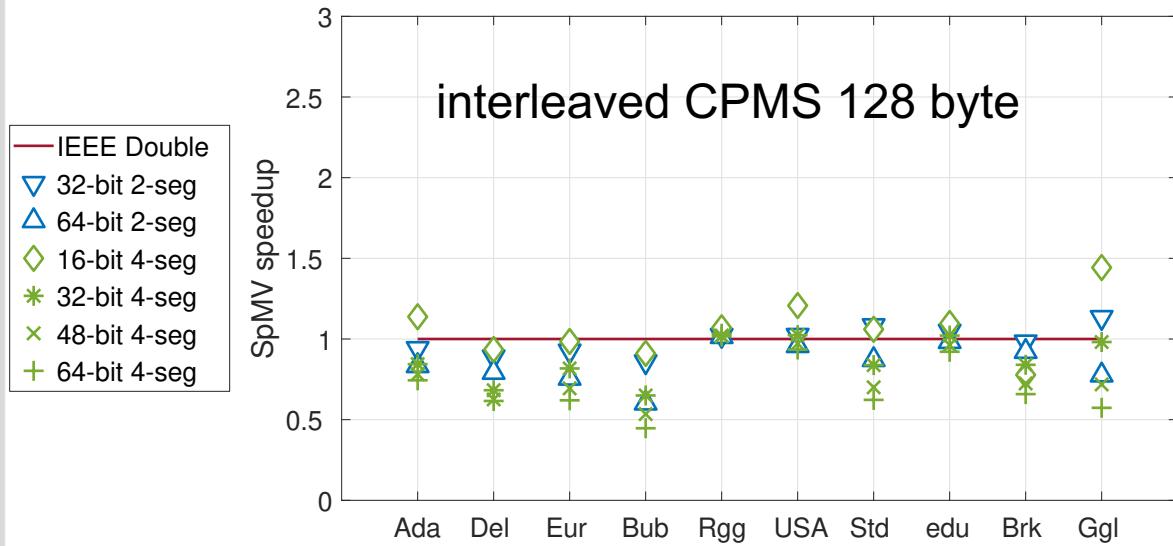
Customized Precision PageRank on CPUs?

Much harder!

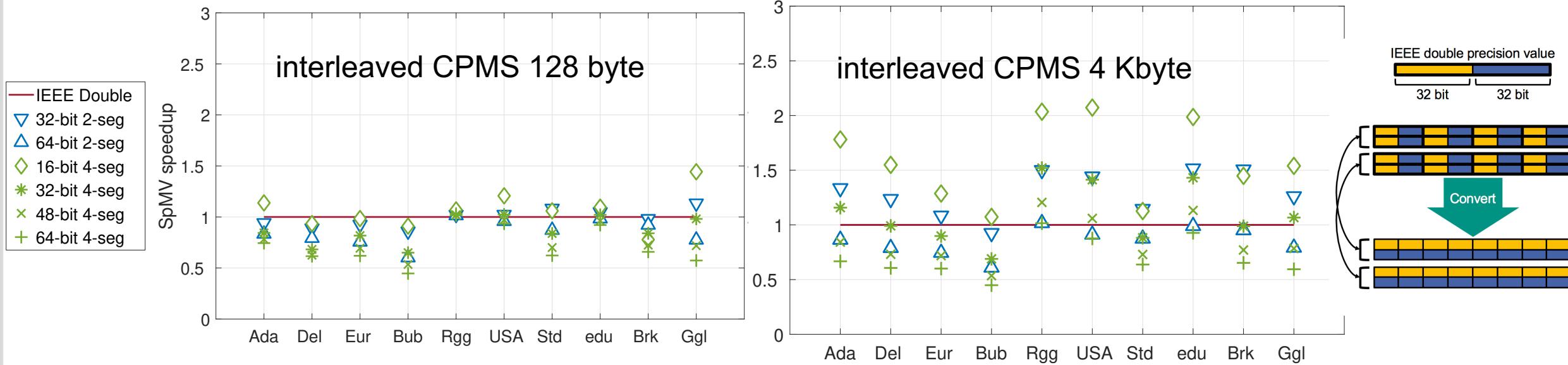
- Caches, prefetching, AVX...
- No easy-to-understand PTX
- Highly optimized compilers

```
#pragma omp parallel for schedule(guided, ompHelper::chunkSize)
for (std::size_t row = 0; row < m; ++row) {
    IndexType i = rowStartIndices[row];
    __m256d currentResPack = _mm256_setr_pd(0.0, 0.0, 0.0, 0.0);
    if (rowStartIndices[row + 1] > (4-1)) {
        for (i = rowStartIndices[row]; i < rowStartIndices[row + 1] - (4-1); i += 4) {
            const __m128i rankIndices32 = _mm_loadu_si128((__m128i *) (colIndices + i));
            __m256d matValPack =
                avx2Helper::readPack<partPointersToRead, NumberSegments, double>(valsSplit, i);
            __m256d rankPack =
                avx2Helper::readPack<partPointersToRead, NumberSegments, double>(oldRankSplit, rankIndices32);
            __m256d multResPack = _mm256_mul_pd(matValPack, rankPack);
            currentResPack = _mm256_add_pd(currentResPack, multResPack);
        }
    }
}
```

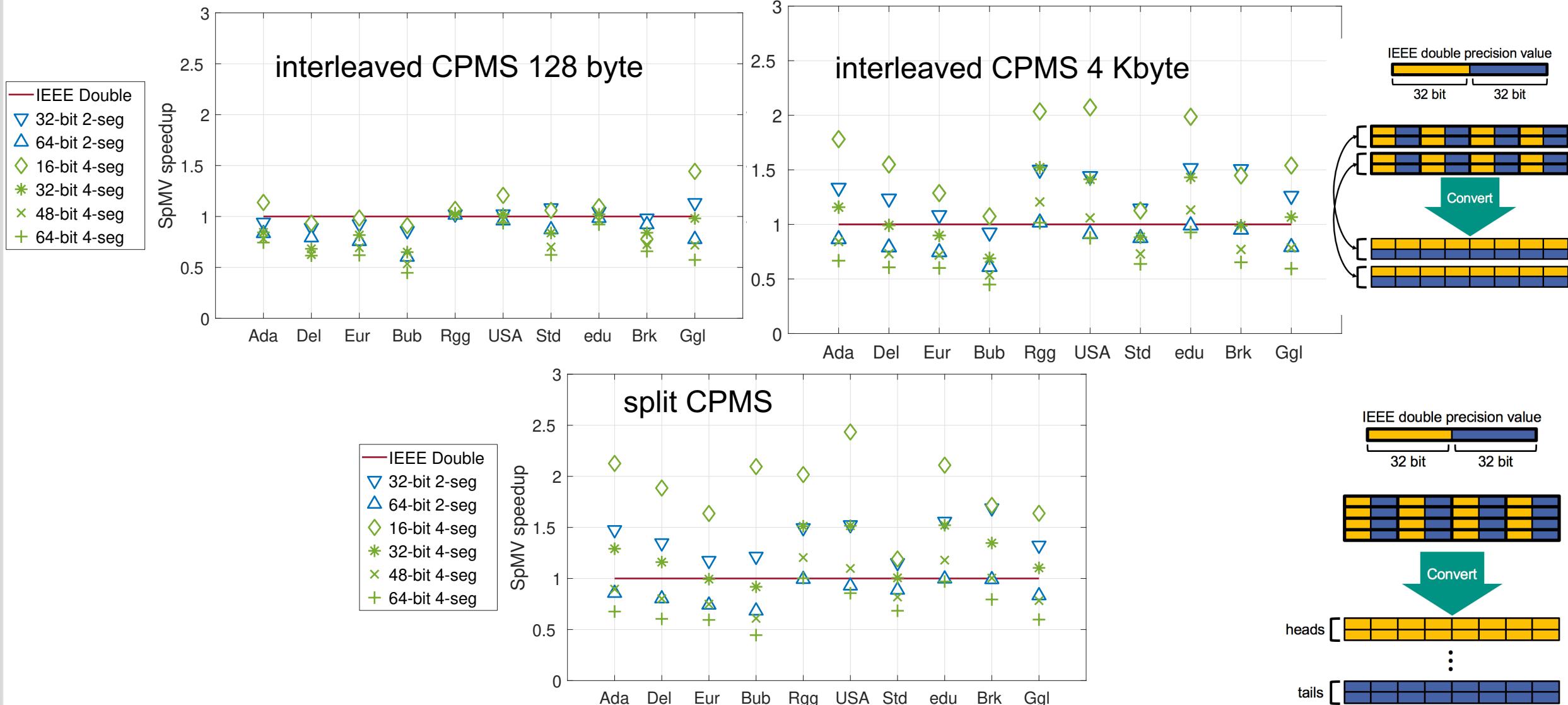
Customized Precision SpMV on CPUs



Customized Precision SpMV on CPUs

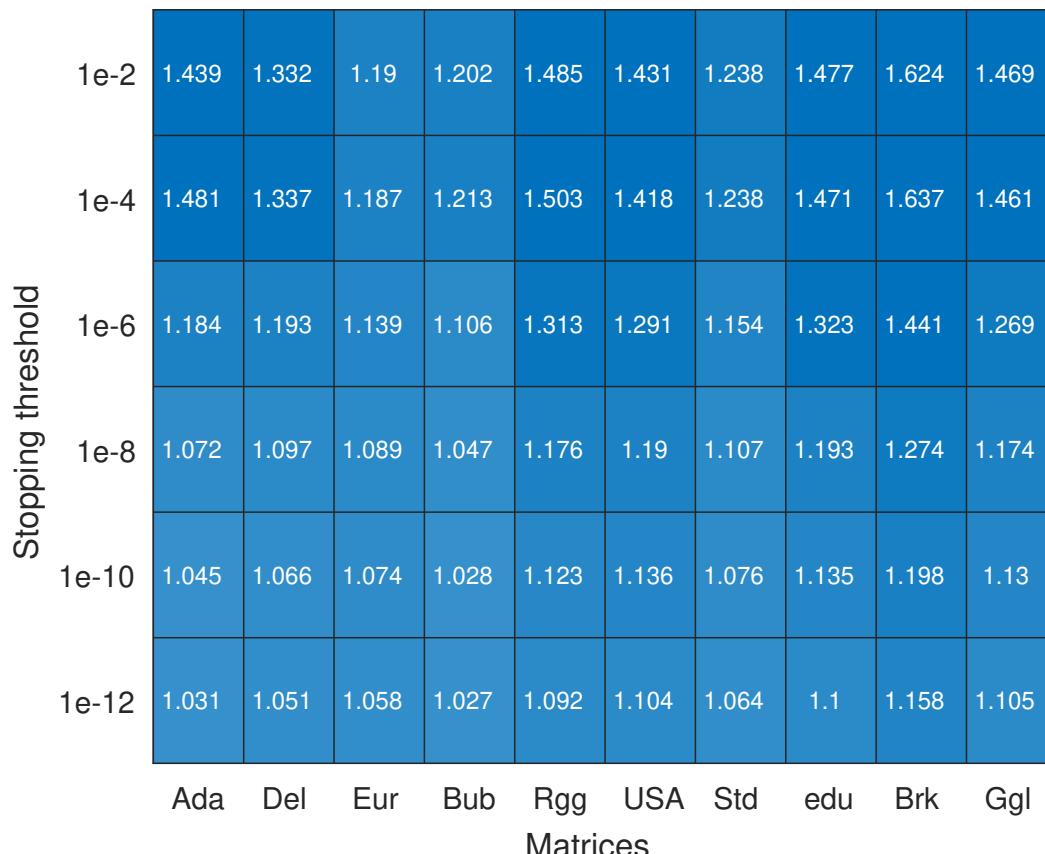


Customized Precision SpMV on CPUs



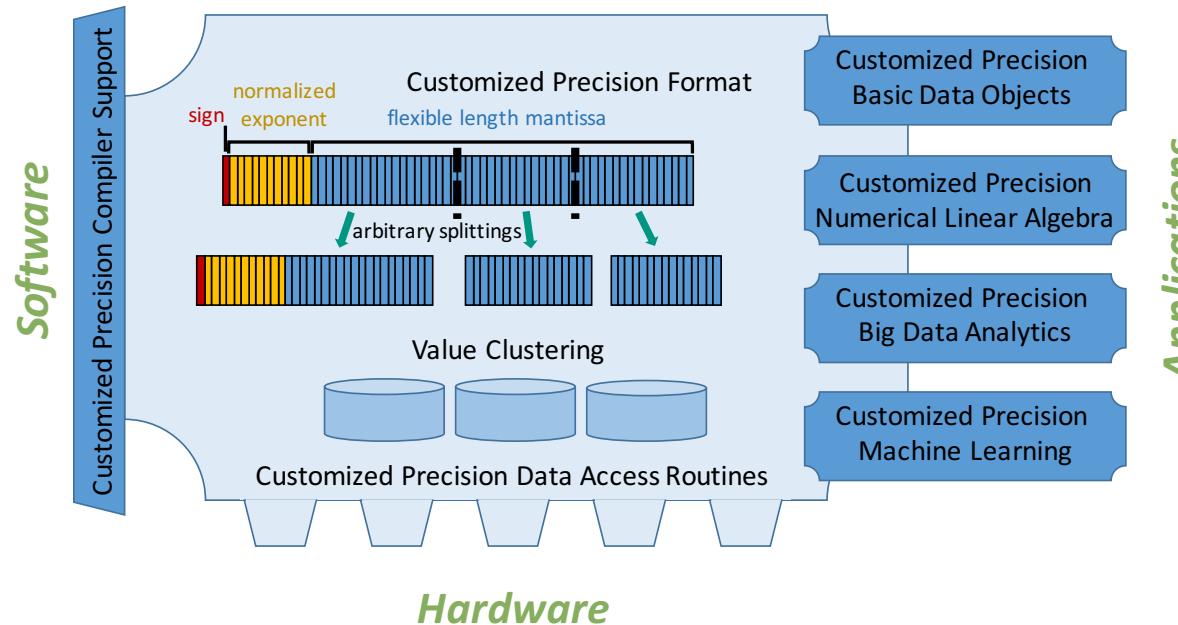
Customized Precision PageRank on CPUs

PageRank speedup over reference implementation (2x Intel(R) Xeon(R) Platinum 8168 "Skylake"). CSR format.



Summary and next steps

- Decouple arithmetic precision from memory precision.
- Using customized precisions for memory operations.
- Dynamically adapt the accuracy to the algorithm properties.
- Speedup of up to 1.3x for Jacobi iterations and PageRank algorithm.
- Speedup of up to 1.3x for adaptive precision block-Jacobi preconditioning.
- Creating a Modular Precision Ecosystem inside  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.