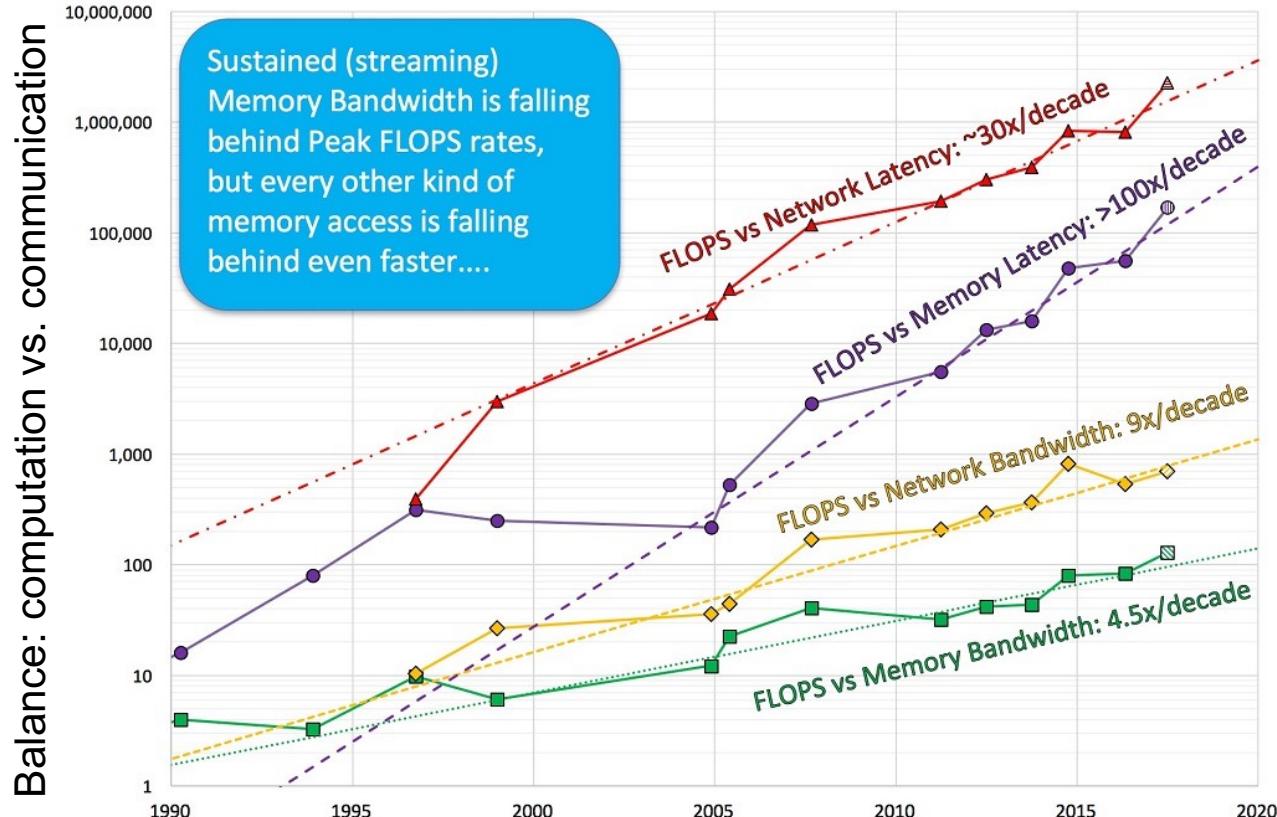


Can Mixed Precision Accelerate Sparse Solvers?

Hartwig Anzt
Technical University of Munich
University of Tennessee



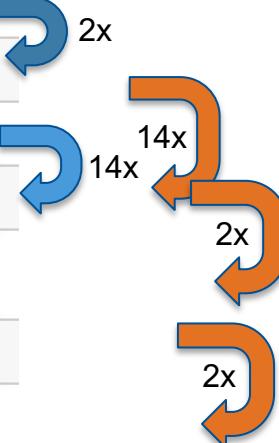
This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration.



Trends in the relative performance of floating-point arithmetic and several classes of data access for select HPC servers over the past 25 years. Source: John McCalpin



Form Factor	H100 SXM
FP64	34 teraFLOPS
FP64 Tensor Core	67 teraFLOPS
FP32	67 teraFLOPS
TF32 Tensor Core	989 teraFLOPS ²
BFLOAT16 Tensor Core	1,979 teraFLOPS ²
FP16 Tensor Core	1,979 teraFLOPS ²
FP8 Tensor Core	3,958 teraFLOPS ²
INT8 Tensor Core	3,958 TOPS ²
GPU memory	80GB
GPU memory bandwidth	3.35TB/s



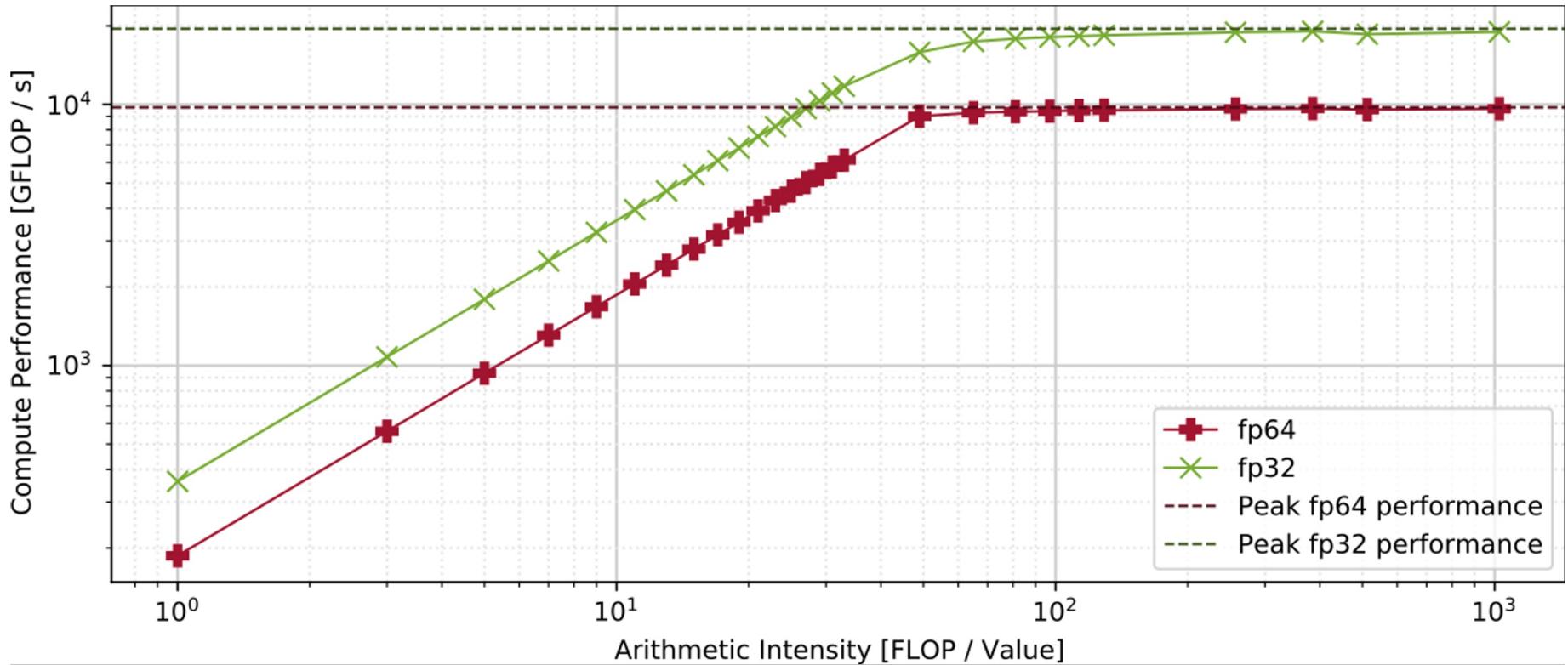
PERFORMANCE

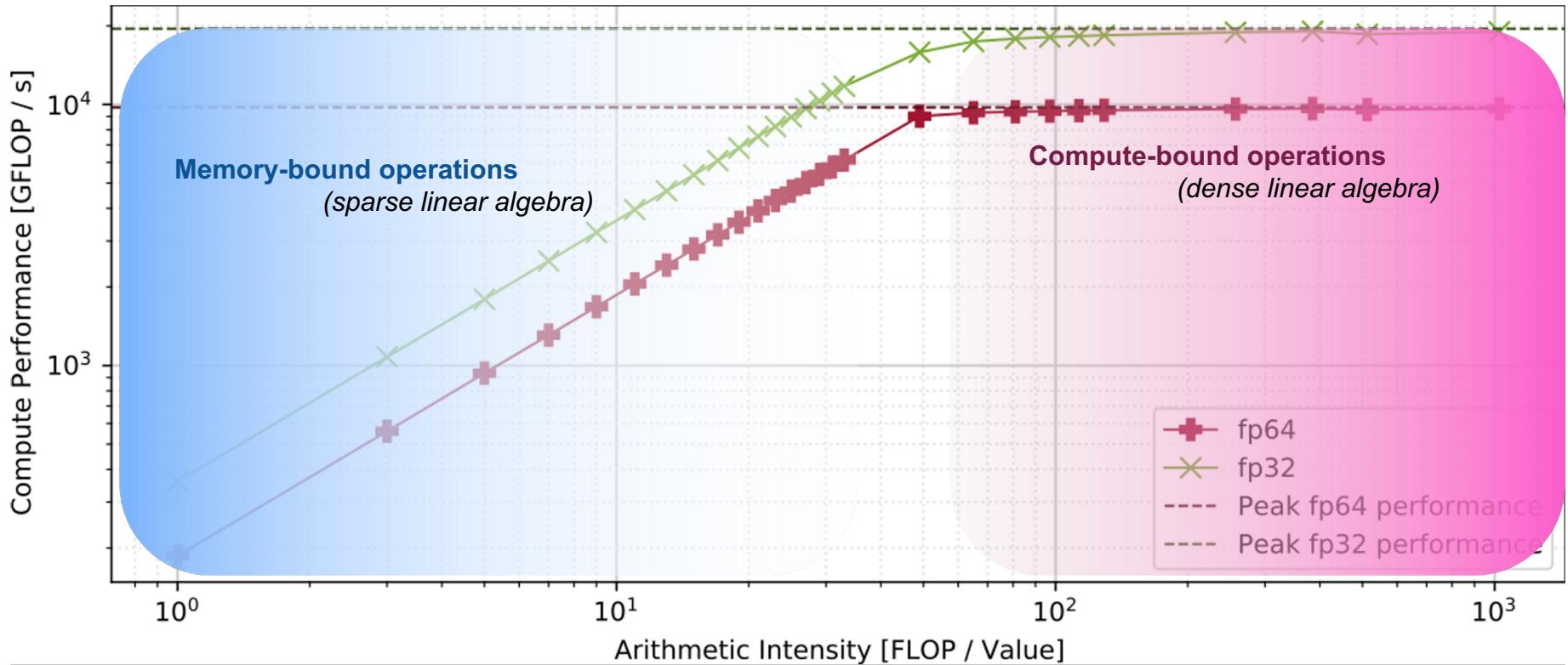
	MI250
Compute Units	208CU
Stream Processors	13,312
Peak FP64/FP32 Vector	45.3 TFLOPS
Peak FP64/FP32 Matrix	90.5 TFLOPS
Peak FP16/BF16	362.1 TFLOPS
Peak INT4/INT8	362.1 TOPS

MEMORY

Memory Size	128GB HBM2e
Memory Interface	8,192 bits
Memory Clock	1.6GHz
Memory Bandwidth	up to 3.2TB/sec ²

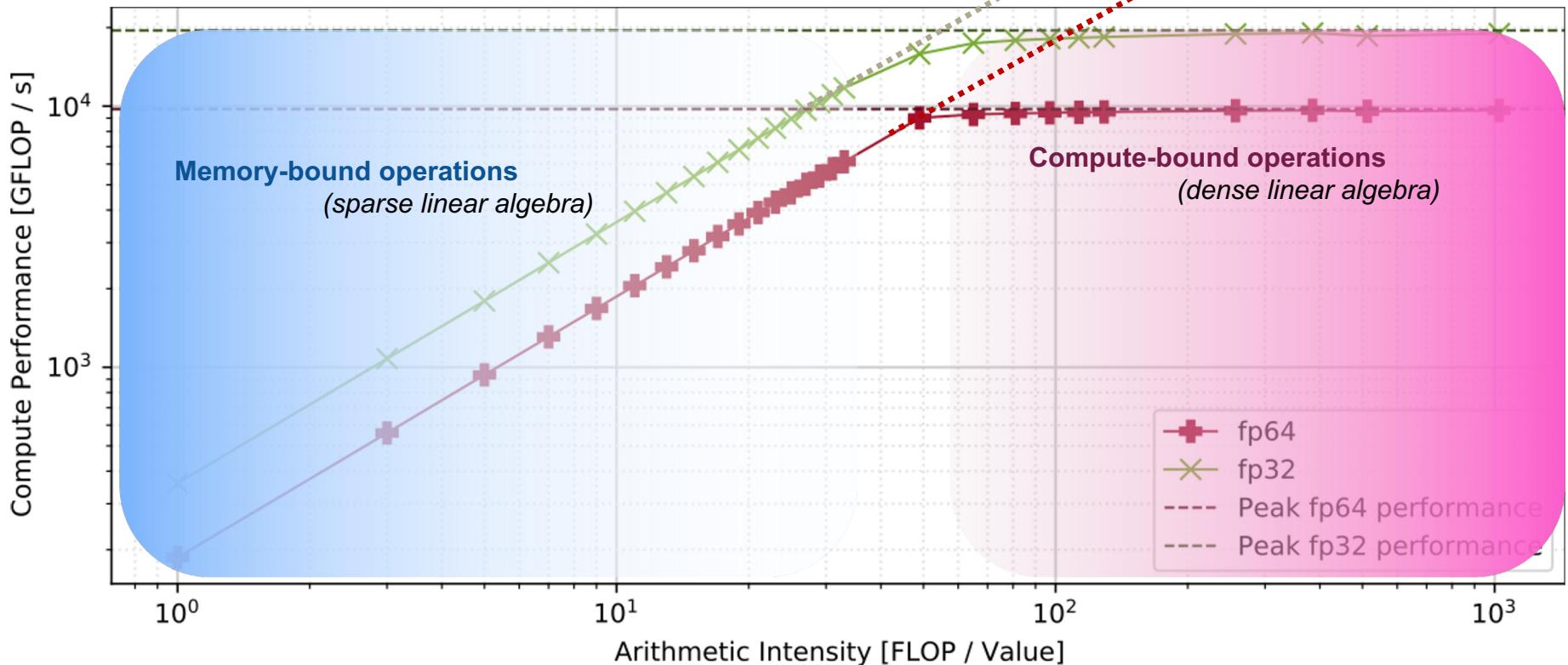
- (Dense) Matrix Performance >> Vector Operation Performance
- Low Precision Performance >> High Precision Performance





Matrix fp32

Matrix fp64

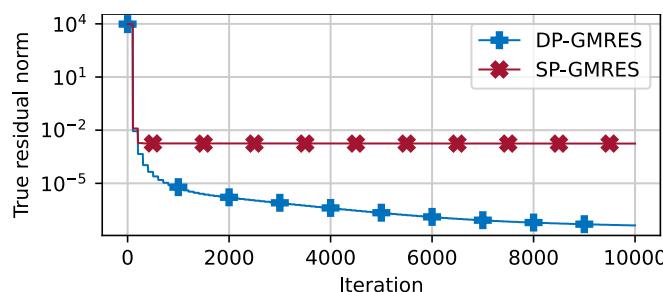


Linear System $Ax=b$ with $\text{cond}(A) \approx 10^7$
 (apache2 from SuiteSparse) NVIDIA V100 GPU

Double precision GMRES
 Initial residual norm $\text{sqrt}(r^T r)$: 9670.36
 Final residual norm $\text{sqrt}(r^T r)$: $9.6639e-09$
 GMRES iteration count: 23271
 GMRES execution time: 43801 ms

Single precision GMRES
 Initial residual norm $\text{sqrt}(r^T r)$: 9670.36
 Final residual norm $\text{sqrt}(r^T r)$: 0.00175464
 GMRES iteration count: 25000
 GMRES execution time: 27376 ms

~2x faster!

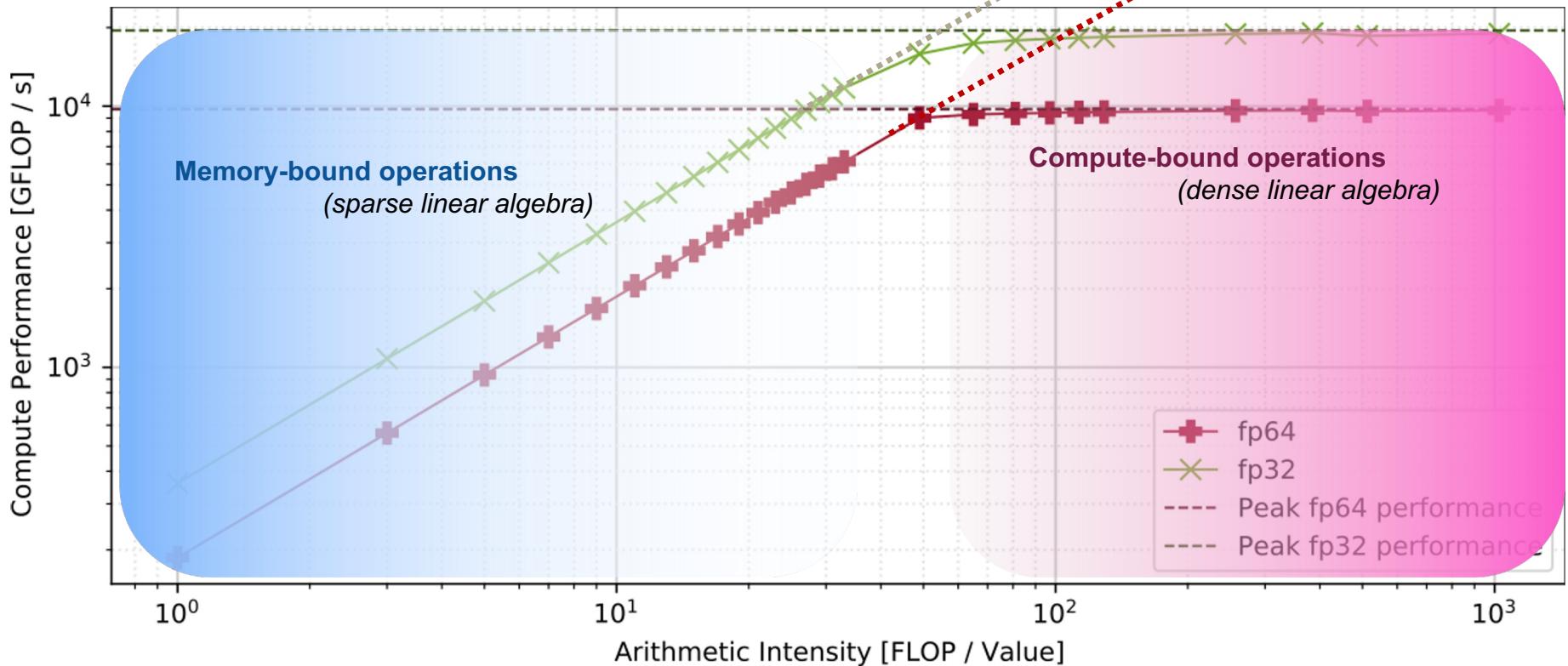


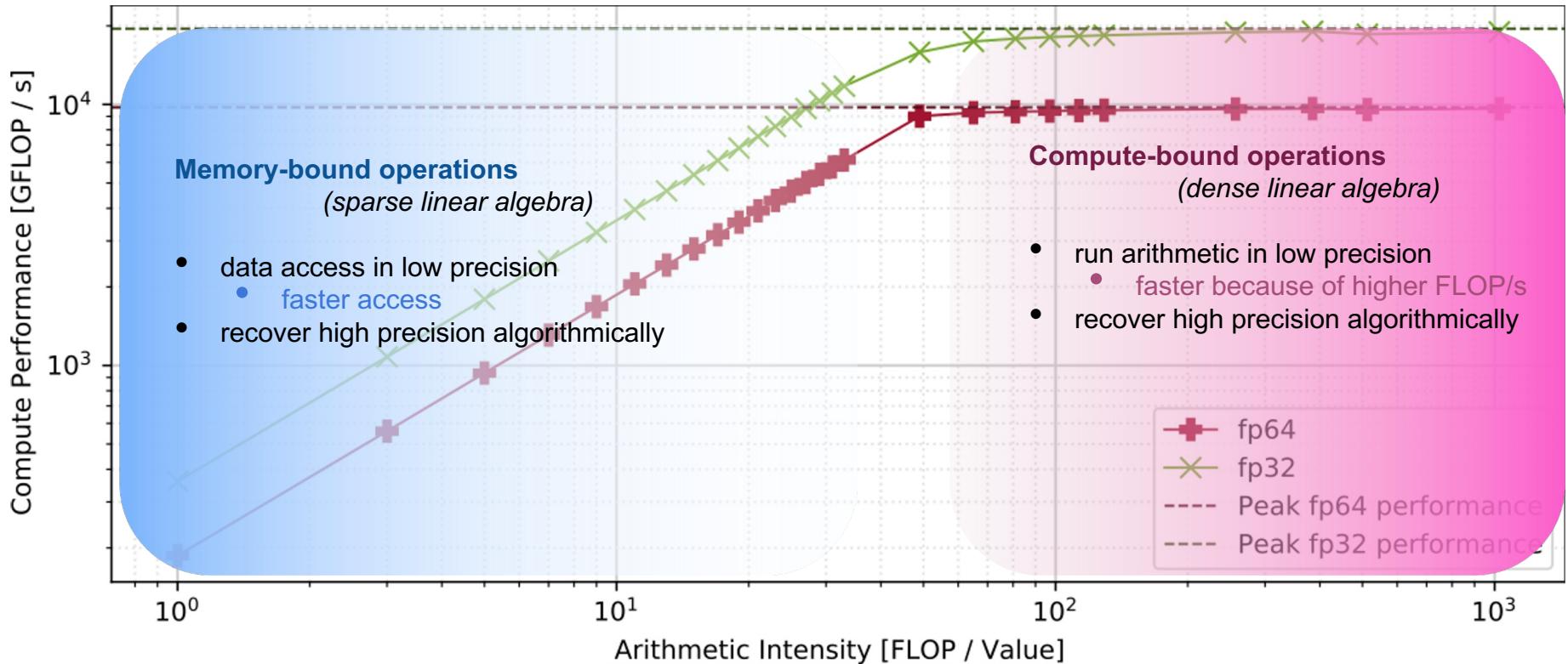
forward error \approx (unit round-off) * (linear system's condition number)

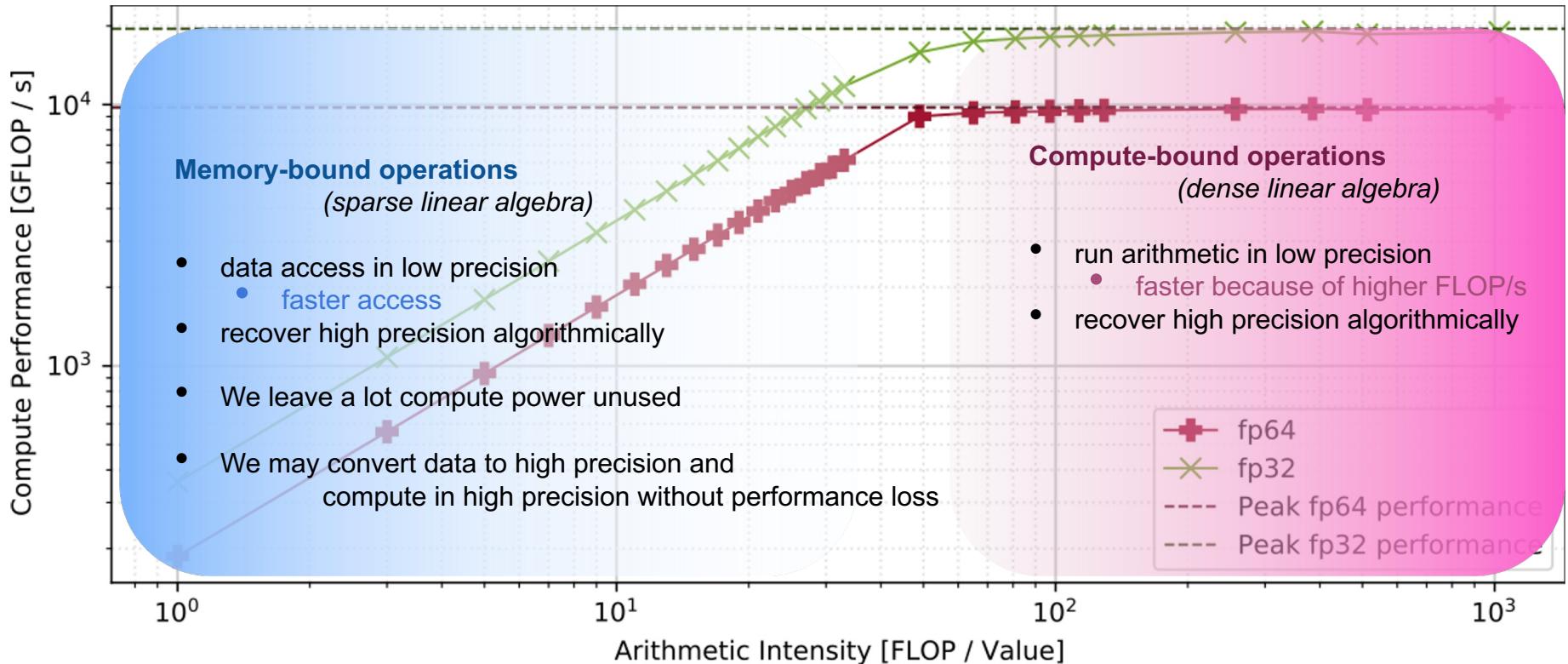
N. Higham: Accuracy and stability of numerical algorithms. SIAM, 2002.

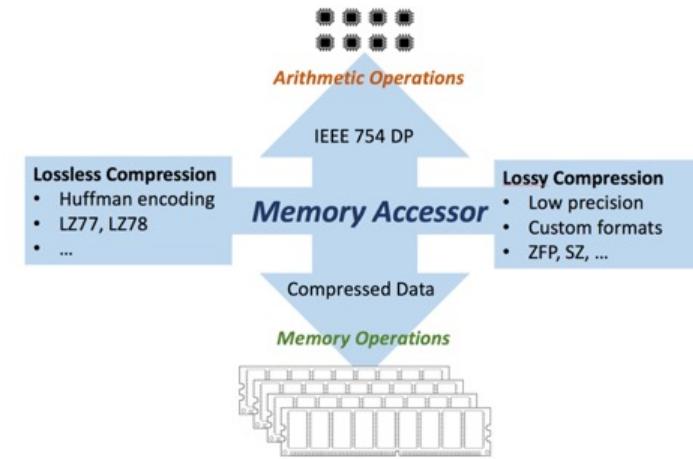
Matrix fp32

Matrix fp64

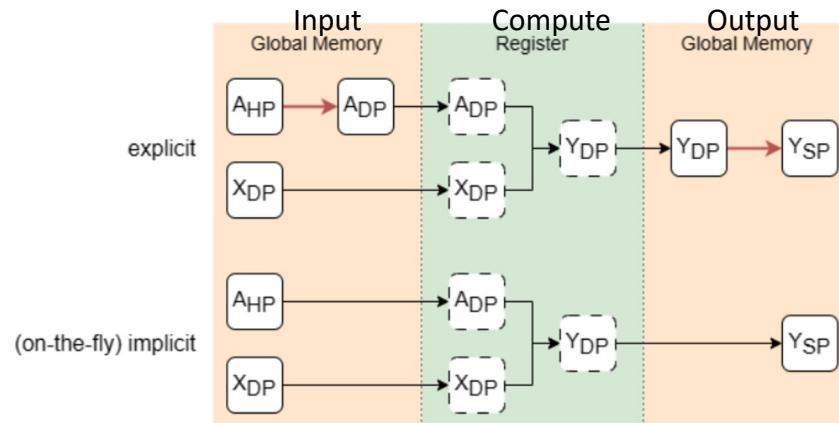




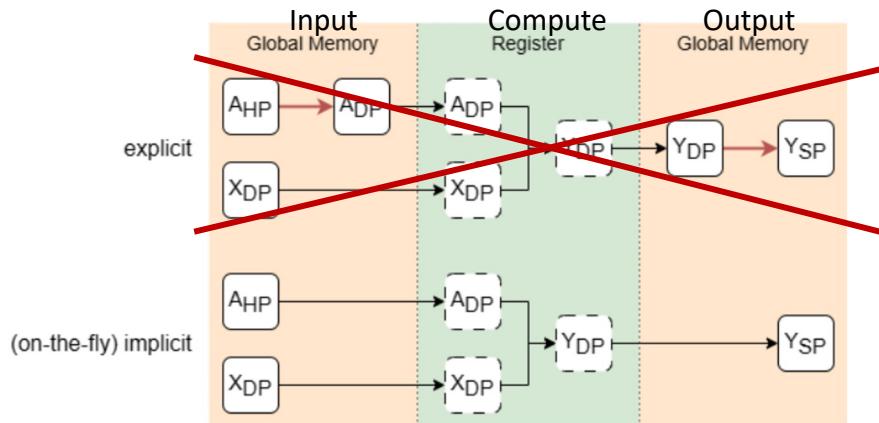
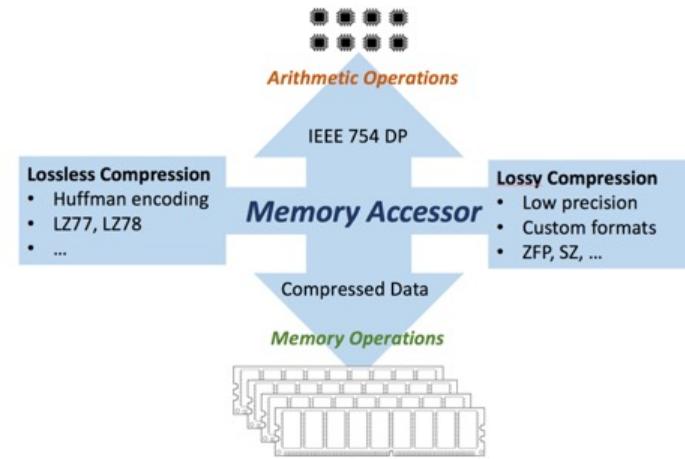


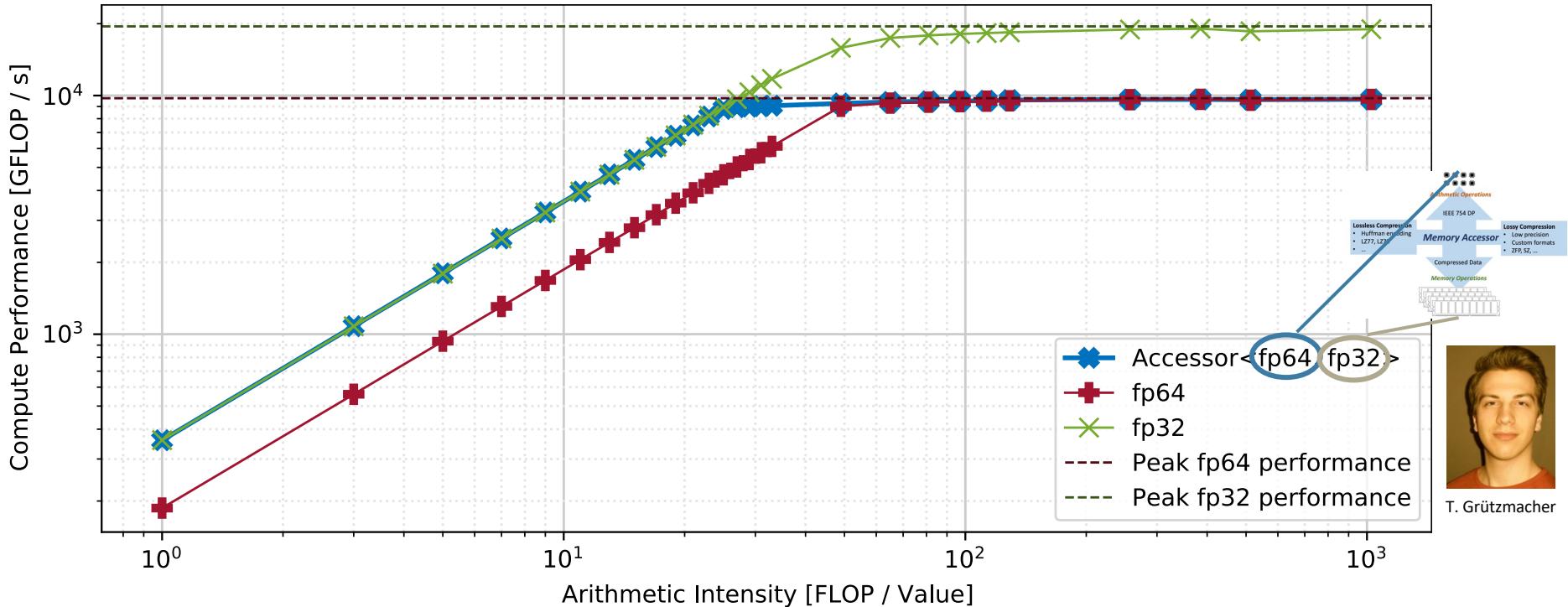


- Traditionally, we use a strong coupling between the precision formats used for **arithmetic operations** and **storing data**.
- *Maybe this is not the right thing?*
- *We should compute in fp64*
- *Maybe we should use the free compute cycles (vector/tensor cores) to compress the data*
- *Compression / Conversion needs to happen on-the-fly*



- Traditionally, we use a strong coupling between the precision formats used for **arithmetic operations** and **storing data**.
- Maybe this is not the right thing?*
- We should compute in fp64*
- Maybe we should use the free compute cycles (vector/tensor cores) to compress the data*
- Compression / Conversion needs to happen on-the-fly*





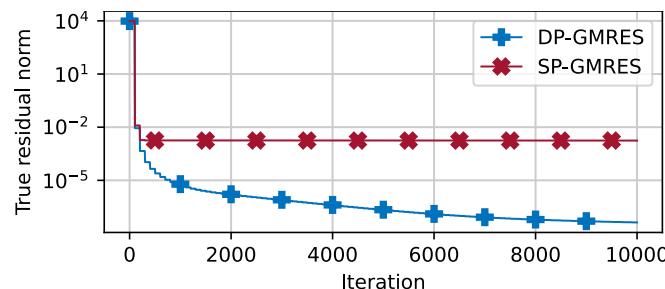
Linear System $Ax=b$ with $\text{cond}(A) \approx 10^7$
 (apache2 from SuiteSparse) NVIDIA V100 GPU

Double precision GMRES

Initial residual norm $\text{sqrt}(r^T r)$: 9670.36
 Final residual norm $\text{sqrt}(r^T r)$: $9.6639e-09$
 GMRES iteration count: 23271
 GMRES execution time: 43801 ms

Single precision GMRES

Initial residual norm $\text{sqrt}(r^T r)$: 9670.36
 Final residual norm $\text{sqrt}(r^T r)$: 0.00175464
 GMRES iteration count: 25000
 GMRES execution time: 27376 ms

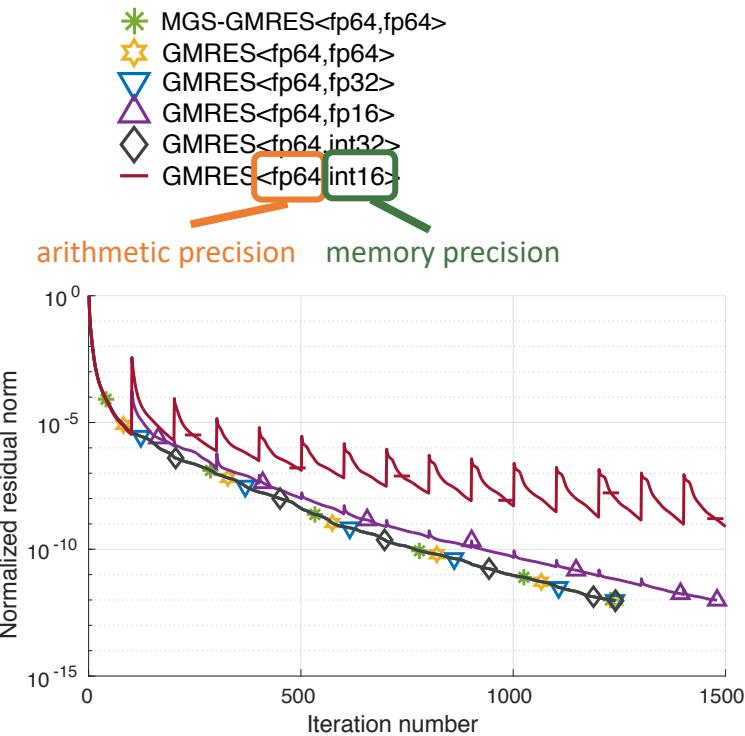


forward error \approx (unit round-off) * (linear system's condition number)

N. Higham: Accuracy and stability of numerical algorithms. SIAM, 2002.

Compressed Basis (CB-) GMRES

- Use double precision in all arithmetic operations;
- Store Krylov basis vectors in lower precision;
 - Search directions are no longer DP-orthogonal;
 - Hessenberg system maps solution to “perturbed” Krylov subspace;
 - Additional iterations may be needed;
 - As long as the loss-of-orthogonality is moderate, we should see moderate convergence degradation;



Linear System $Ax=b$ with $\text{cond}(A) \approx 10^7$
 (apache2 from SuiteSparse) NVIDIA V100 GPU

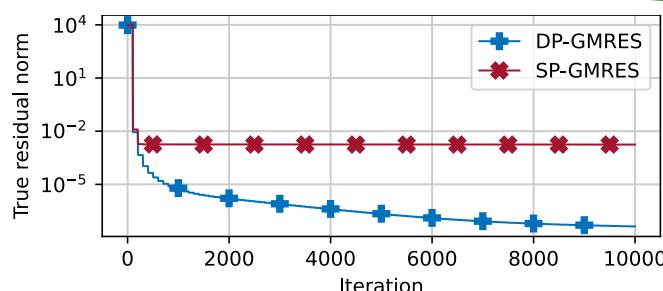
Double precision GMRES

Initial residual norm $\text{sqrt}(r^T r)$: 9670.36
 Final residual norm $\text{sqrt}(r^T r)$: $9.6639e-09$
 GMRES iteration count: 23271
 GMRES execution time: 43801 ms

Single precision GMRES

Initial residual norm $\text{sqrt}(r^T r)$: 9670.36
 Final residual norm $\text{sqrt}(r^T r)$: 0.00175464
 GMRES iteration count: 25000
 GMRES execution time: 27376 ms

~2x faster!



forward error \approx (unit round-off) * (linear system's condition number)

N. Higham: Accuracy and stability of numerical algorithms. SIAM, 2002.

Linear System $Ax=b$ with $\text{cond}(A) \approx 10^7$
(apache2 from SuiteSparse) NVIDIA V100 GPU

Double precision GMRES

Initial residual norm $\sqrt{r^T r}$: 9670.36
Final residual norm $\sqrt{r^T r}$: $9.6639e-09$
GMRES iteration count: 23271
GMRES execution time: 43801 ms

Single precision GMRES

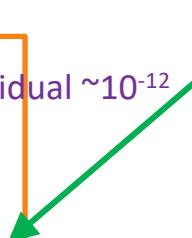
Initial residual norm $\sqrt{r^T r}$: 9670.36
Final residual norm $\sqrt{r^T r}$: 0.00175464
GMRES iteration count: 25000
GMRES execution time: 27376 ms

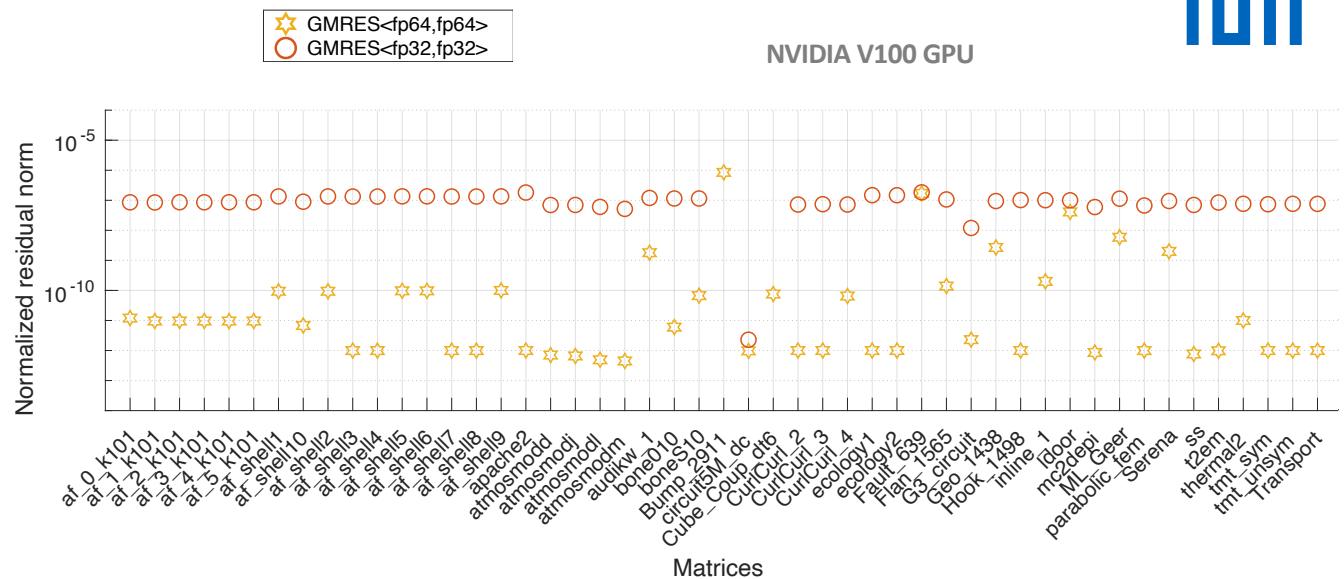
Compressed Basis GMRES

Initial residual norm $\sqrt{r^T r}$: 9670.36
Final residual norm $\sqrt{r^T r}$: $9.6591e-09$
GMRES iteration count: 23271
GMRES execution time: 29369 ms

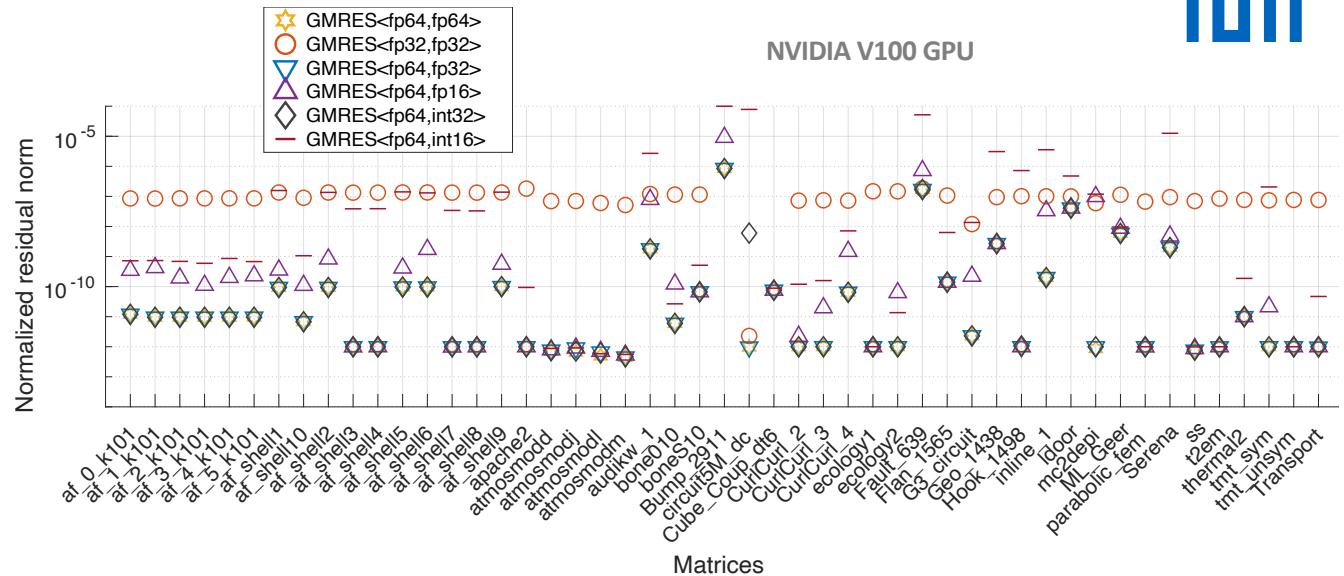
Relative residual $\sim 10^{-12}$

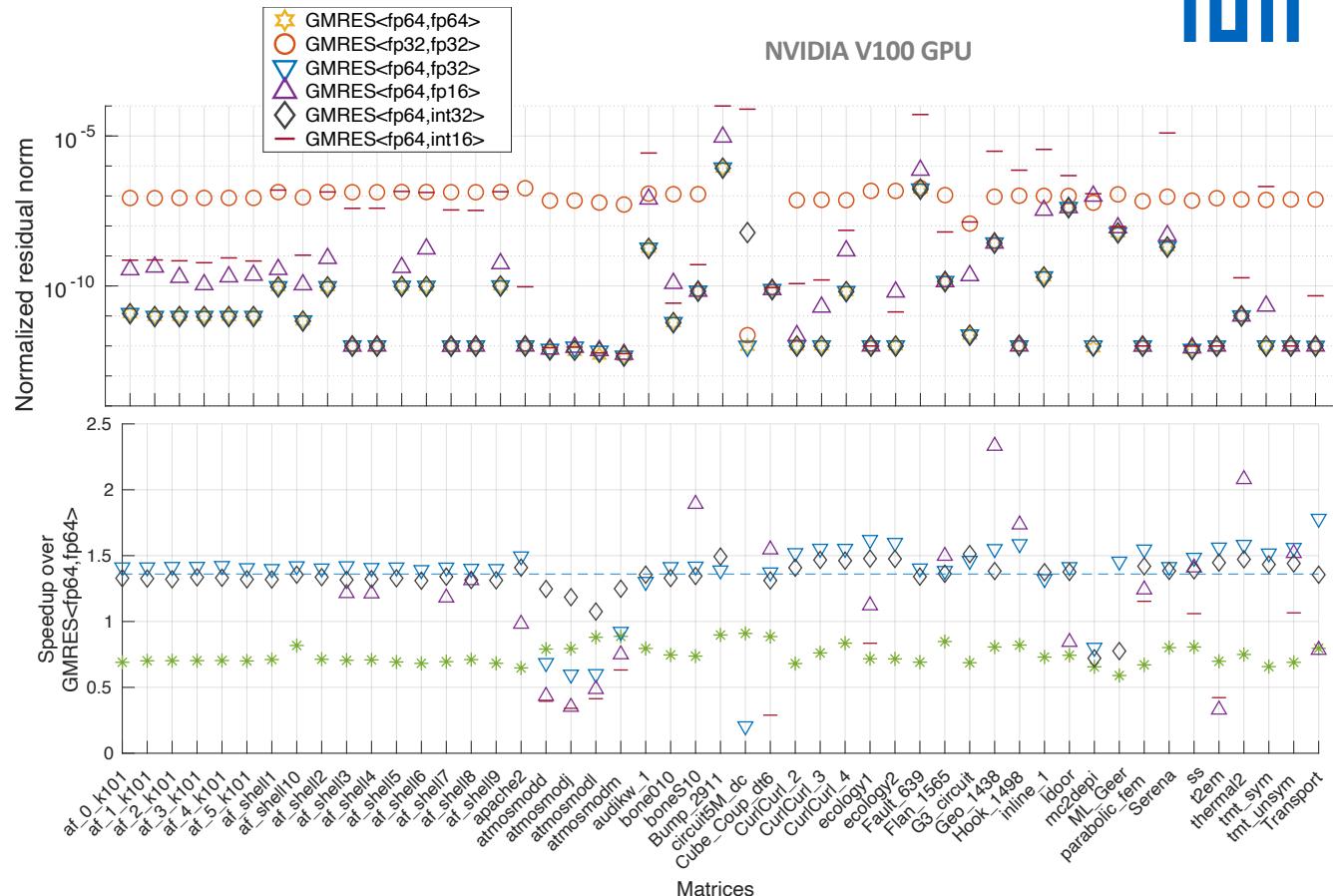
Accuracy of DP GMRES
Performance similar to SP GMRES





- CB-GMRES using 32-bit storage preserves DP accuracy
(SP-GMRES does not)



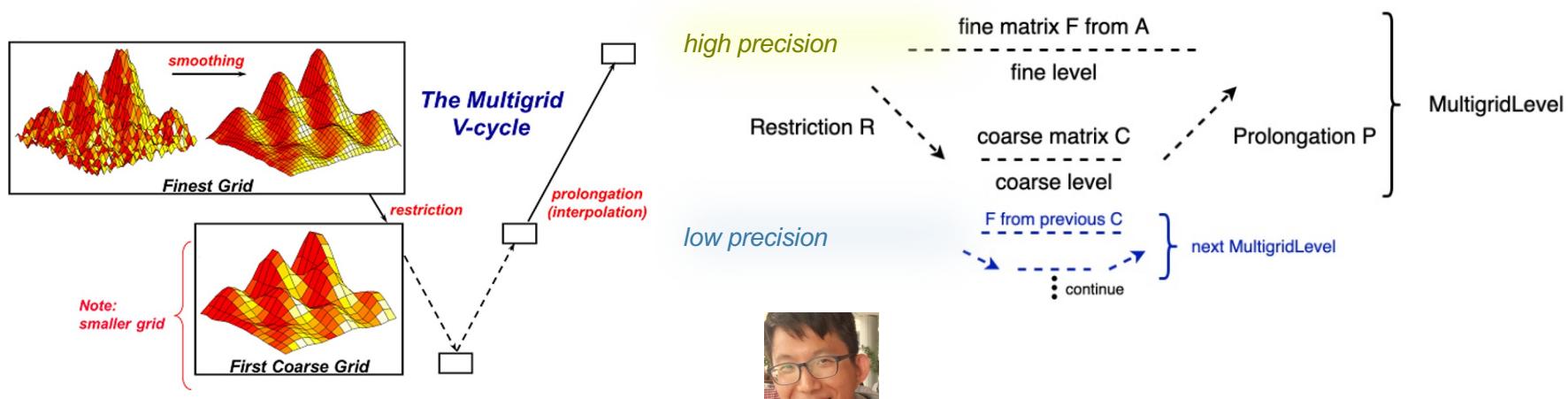


- CB-GMRES using 32-bit storage preserves DP accuracy (SP-GMRES does not)
- Speedups problem-dependent
- Speedup $\varnothing 1.4x$ (for restart 100)
- 16-bit storage mostly inefficient



Aliaga JI, Anzt H, Grützmacher T, Quintana-Ortí ES, Tomás AE. Compressed basis GMRES on high-performance graphics processing units. *The International Journal of High Performance Computing Applications*. 2022;0(0). doi:[10.1177/10943420221115140](https://doi.org/10.1177/10943420221115140)

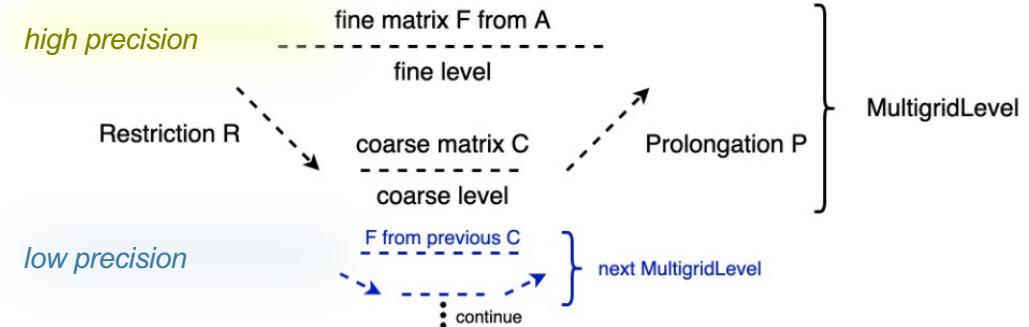
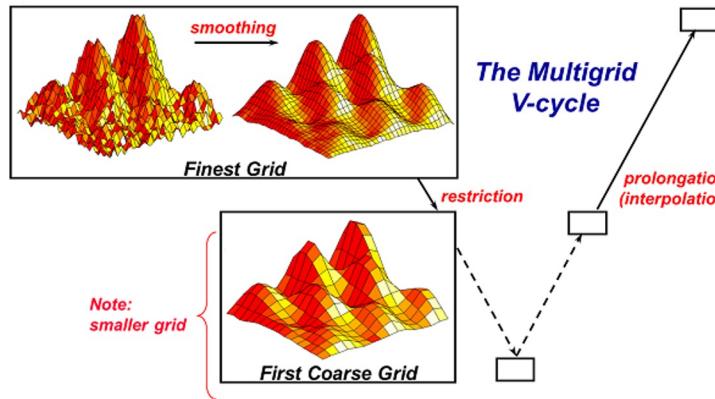
Mixed Precision Multigrid



Mike Tsai

MS3

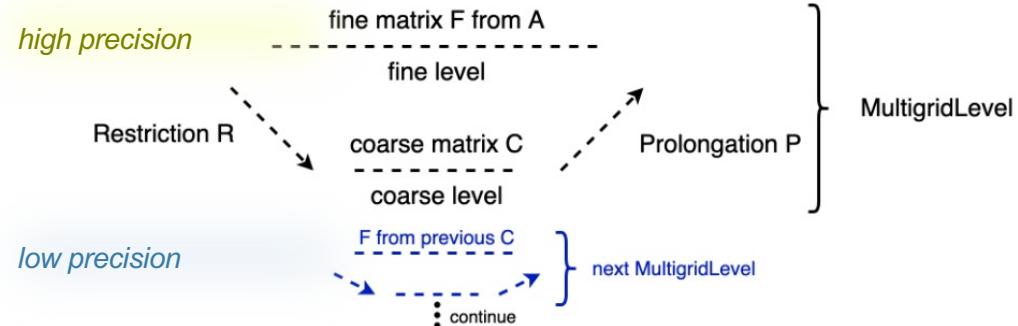
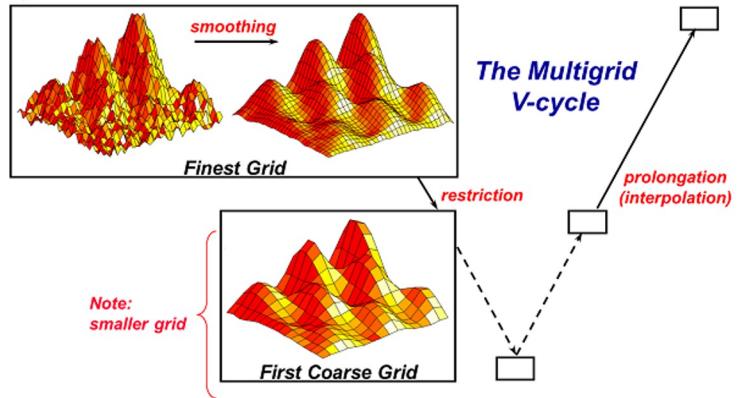
Linear Solvers in Large Distributed Applications - Part I of III



```

1 multigrid::build()
2   .with_max_levels(10u) // equal to NVIDIA/AMGX 11 max levels
3   .with_min_coarse_row(64u)
4   .with_pre_smootheser(sm, sm_f)
5   .with_mg_level(pgm, pgm_f)
6   .with_level_selector(
7     [](const size_type level, const LinOp*) -> size_type {
8       // Only the first level is generated by MultigridLevel(double).
9       // The subsequent levels are generated by MultigridLevel(float)
10      return level >= 1 ? 1 : 0;
11    })
12   .with_coarest_solver(coarest_solver_f)

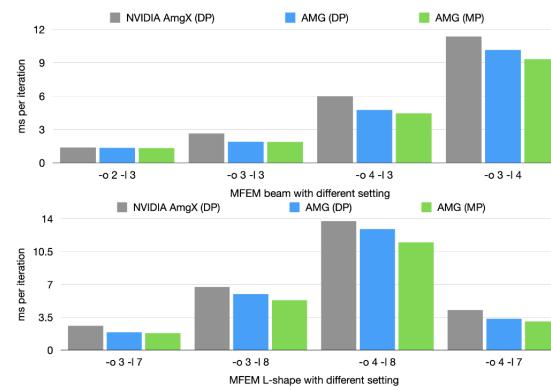
```

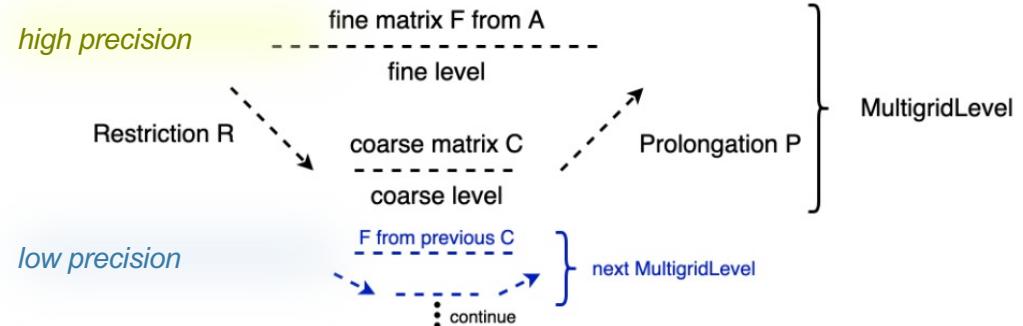
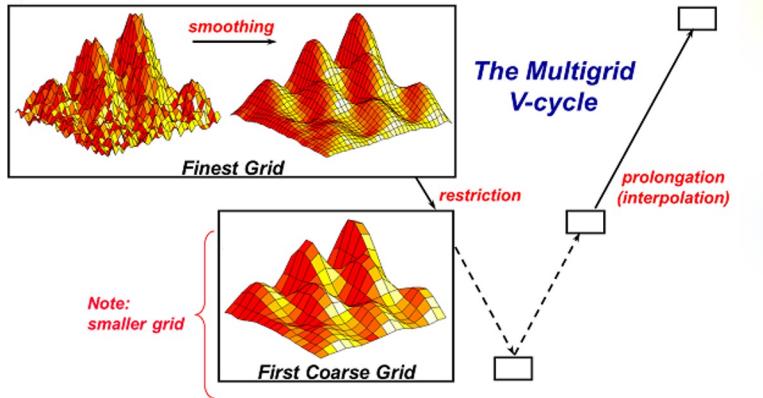


```

1 multigrid::build()
2   .with_max_levels(10u) // equal to NVIDIA/AMGX 11 max levels
3   .with_min_coarse_row(64u)
4   .with_pre_smoothes(sm, sm_f)
5   .with_mg_level(pgm, pgm_f)
6   .with_level_selector(
7     [](const size_type level, const LinOp*) -> size_type {
8       // Only the first level is generated by MultigridLevel(double).
9       // The subsequent levels are generated by MultigridLevel(float)
10      return level >= 1 ? 1 : 0;
11    })
12   .with_coarest_solver(coarest_solver_f)

```

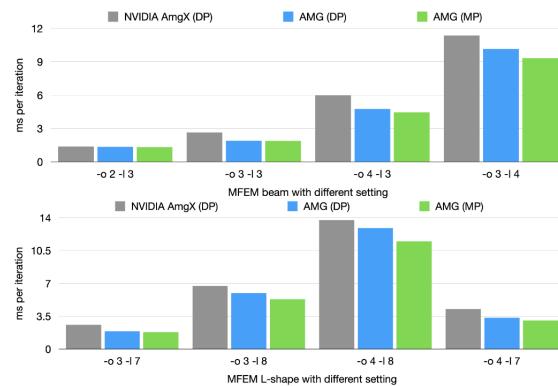




```

1 multigrid::build()
2   .with_max_levels(10u) // equal to NVIDIA/AMGX 11 max levels
3   .with_min_coarse_row(64u)
4   .with_pre_smother(sm, sm_f)
5   .with_mg_level(pgm, pgm_f)
6   .with_level_selector(
7     [](const size_type level, const LinOp*) -> size_type {
8       // Only the first level is generated by MultigridLevel(double).
9       // The subsequent levels are generated by MultigridLevel(float)
10      return level >= 1 ? 1 : 0;
11    })
12   .with_coarest_solver(coarest_solver_f)

```



Mike Tsai



Natalie Beams

YHM Tsai, N Beams, H Anzt
Mixed Precision Algebraic Multigrid on GPUs
International Conference on Parallel Processing and Applied Mathematics, 113-125, 2022

DP-SP-HP

double



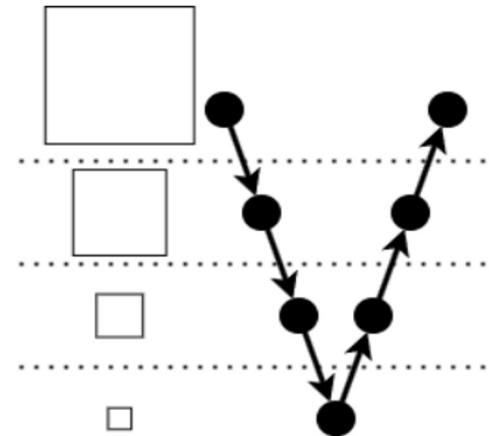
single



half



half



DP-SP-HP

double



single



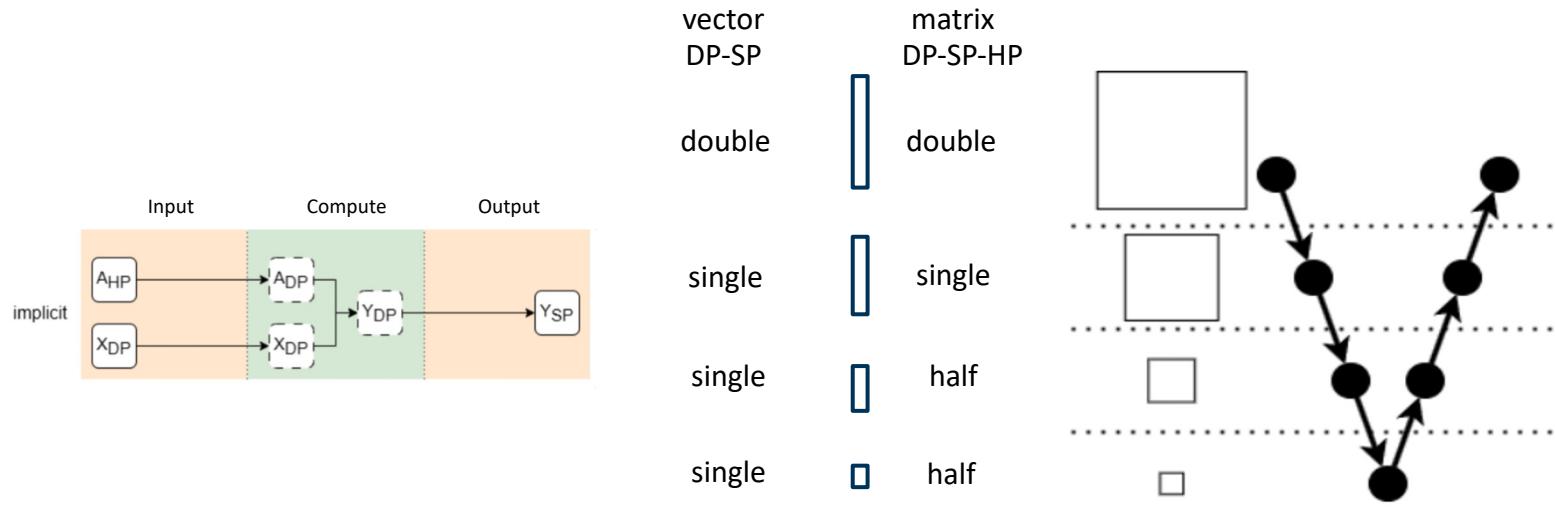
half



half



problem	GINKGO's AMG (DP)			GINKGO's AMG (DP-SP)			GINKGO's AMG (DP-SP-HP)			GINKGO's AMG (DP-HP)		
	res.	norm	#iter	time[ms]	res.	norm	#iter	time[ms]	res.	norm	#iter	time[ms]
2cubes	6.56e-09	20	14.62	6.56e-09	20	12.91	NaN	800	493.44	NaN	800	485.06
cage13	3.68e-10	11	15.49	3.68e-10	11	14.09	4.24e-10	11	14.18	6.15e-10	15	17.79
cage14	3.41e-10	10	32.98	3.41e-10	10	29.80	3.73e-10	10	29.04	7.02e-10	13	35.15
offshore	1203.64	800	770.56	1533.85	800	688.64	NaN	800	706.05	NaN	800	671.65
thermal2	2.18e-06	349	489.08	2.38e-06	425	536.09	NaN	800	952.22	NaN	800	919.08
tmt_sym	6.95e-05	359	371.70	7.41e-05	401	379.21	NaN	800	703.40	NaN	800	684.14
beam(o3l3)	3.05e-15	44	43.65	3.05e-15	44	40.89	5.53e-15	86	80.72	7.27e-15	127	117.62
l-shape(o3l7)	4.64e-14	160	195.31	4.78e-14	171	191.62	3.74e-11	800	869.77	2.19e-07	800	868.83



(a) fp32: Single-precision IEEE Floating Point Format

Range: $\sim 1e^{-38}$ to $\sim 3e^{38}$



(b) fp16: Half-precision IEEE Floating Point Format

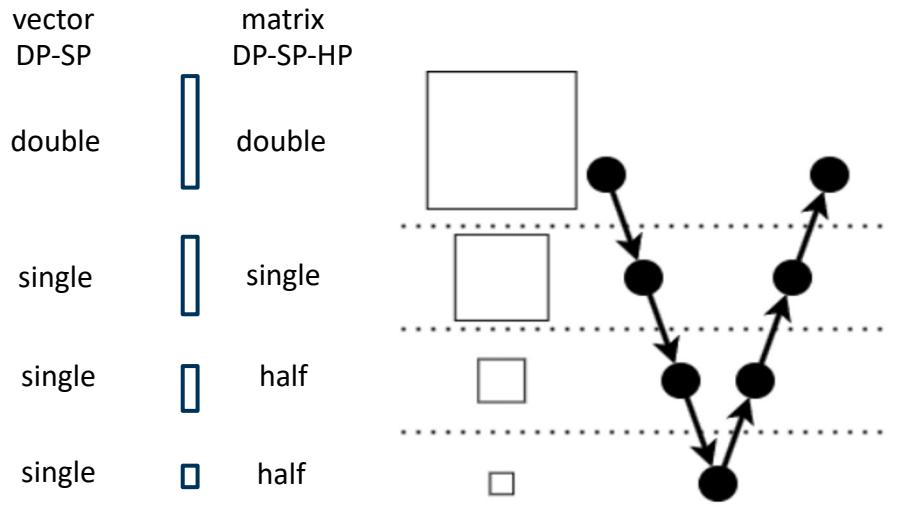
Range: ~5.96e-8 to 65504



(c) bfloat16: Brain Floating Point Format

Range: $\sim 1e^{-38}$ to $\sim 3e^{38}$





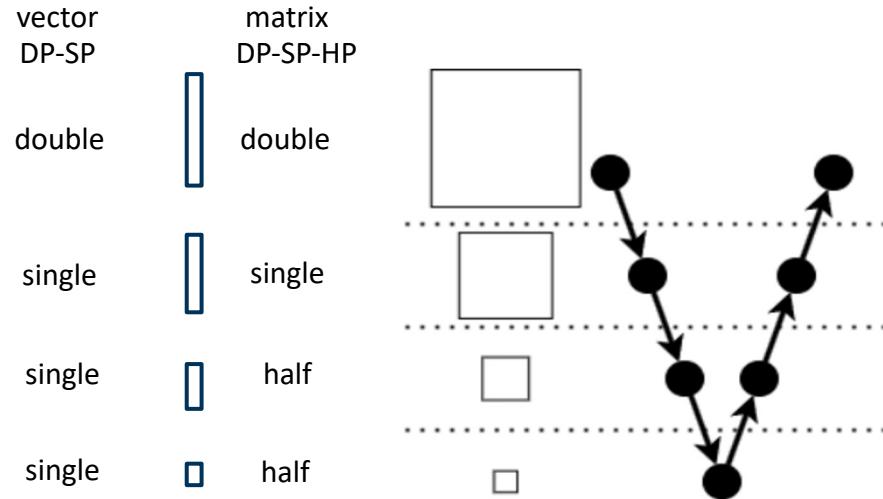
problem	GINKGO's AMG (DP-SP-HP)			GINKGO's AMG (DP-HP)			GINKGO's AMG (DP-SP-BF)			GINKGO's AMG (DP-BF)			
	res.	norm	#iter	time[ms]	res.	norm	#iter	time[ms]	res.	norm	#iter	time[ms]	
2cubes	NaN	800	493.44		NaN	800	485.06	6.70e-09	20	13.53	1.15e-08	20	13.31
cage13	4.24e-10	11	14.18		6.15e-10	15	17.79	3.67e-10	11	13.88	3.71e-10	11	13.76
cage14	3.73e-10	10	29.04		7.02e-10	13	35.15	3.40e-10	10	29.46	3.42e-10	10	27.60
offshore	NaN	800	706.05		NaN	800	671.65	1416.76	800	713.00	2045.25	800	692.75
thermal2	NaN	800	952.22		NaN	800	919.08	2.80e-06	597	730.73	2.92e-06	635	723.24
tmt_sym	NaN	800	703.4		NaN	800	684.14	9.59e-05	649	593.51	9.87e-05	754	665.21
beam(o3l3)	5.53e-15	86	80.72		7.27e-15	127	117.62	3.15e-15	44	42.90	3.29e-15	44	42.99
l-shape(o3l7)	3.74e-11	800	869.77		2.19e-07	800	868.83	5.01e-14	186	203.06	5.10e-14	192	208.69



YHM Tsai, N Beams, H Anzt
Three-precision algebraic multigrid on
GPUs Future Generation Computer
Systems 149, 280-293, 2023.

Mike Tsai

Natalie Beams



problem	GINKGO's AMG (DP-SP-HP)			GINKGO's AMG (DP-HP)			GINKGO's AMG (DP-SP-BF)			GINKGO's AMG (DP-BF)			
	res.	norm	#iter	time[ms]	res.	norm	#iter	time[ms]	res.	norm	#iter	time[ms]	
2cubes	NaN	800	493.44		NaN	800	485.06	6.70e-09	20	13.53	1.15e-08	20	13.31
cage13	4.24e-10	11	14.18		6.15e-10	15	17.79	3.67e-10	11	13.88	3.71e-10	11	13.76
cage14	3.73e-10	10	29.04		7.02e-10	13	35.15	3.40e-10	10	29.46	3.42e-10	10	27.60
offshore	NaN	800	706.05		NaN	800	671.65	1416.76	800	713.00	2045.25	800	692.75
thermal2	NaN	800	952.22		NaN	800	919.08	2.80e-06	597	730.73	2.92e-06	635	723.24
tmt_sym	NaN	800	703.4		NaN	800	684.14	9.59e-05	649	593.51	9.87e-05	754	665.21
beam(o3l3)	5.53e-15	86	80.72		7.27e-15	127	117.62	3.15e-15	44	42.90	3.29e-15	44	42.99
l-shape(o3l7)	3.74e-11	800	869.77		2.19e-07	800	868.83	5.01e-14	186	203.06	5.10e-14	192	208.69

Mixed Precision Preconditioning

- **Preconditioning iterative solvers**

- Idea: Approximate inverse of system matrix to make the system “easier to solve”: $P^{-1} \approx A^{-1}$
and solve $Ax = b \Leftrightarrow P^{-1}Ax = P^{-1}b \Leftrightarrow \tilde{A}x = \tilde{b}$.

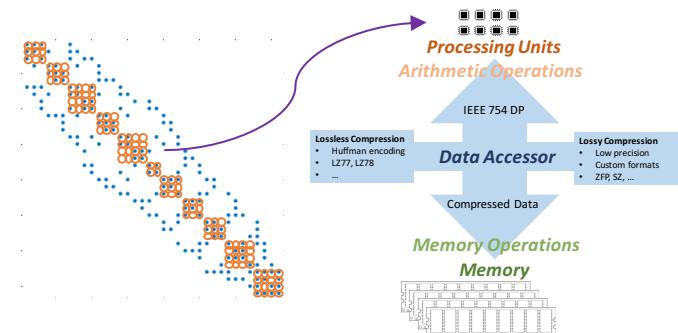
- **Block-Jacobi preconditioner** is based on **block-diagonal scaling**: $P = \text{diag}_B(A)$

- Each block corresponds to one (small) linear system.
 - *Larger* blocks typically improve convergence.
 - *Larger* blocks make block-Jacobi more expensive.

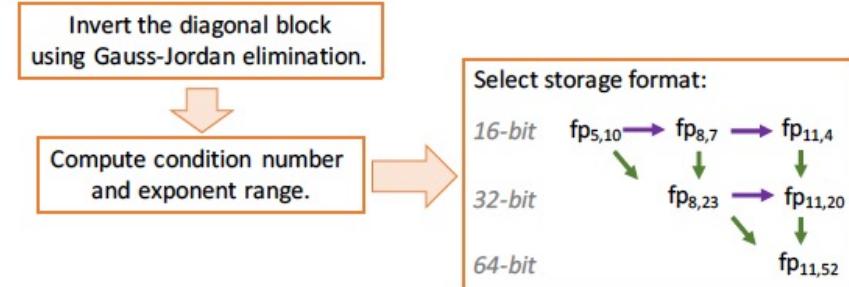
- *Why should we store the preconditioner matrix P^{-1} in full (high) precision?*

- Use the accessor to store the inverted diagonal blocks in lower precision.

- *Be careful to preserve the regularity of each inverted diagonal block!*



- Choose how much accuracy of the preconditioner should be preserved in the selection of the storage format.
- All computations use double precision, but store blocks in lower precision.



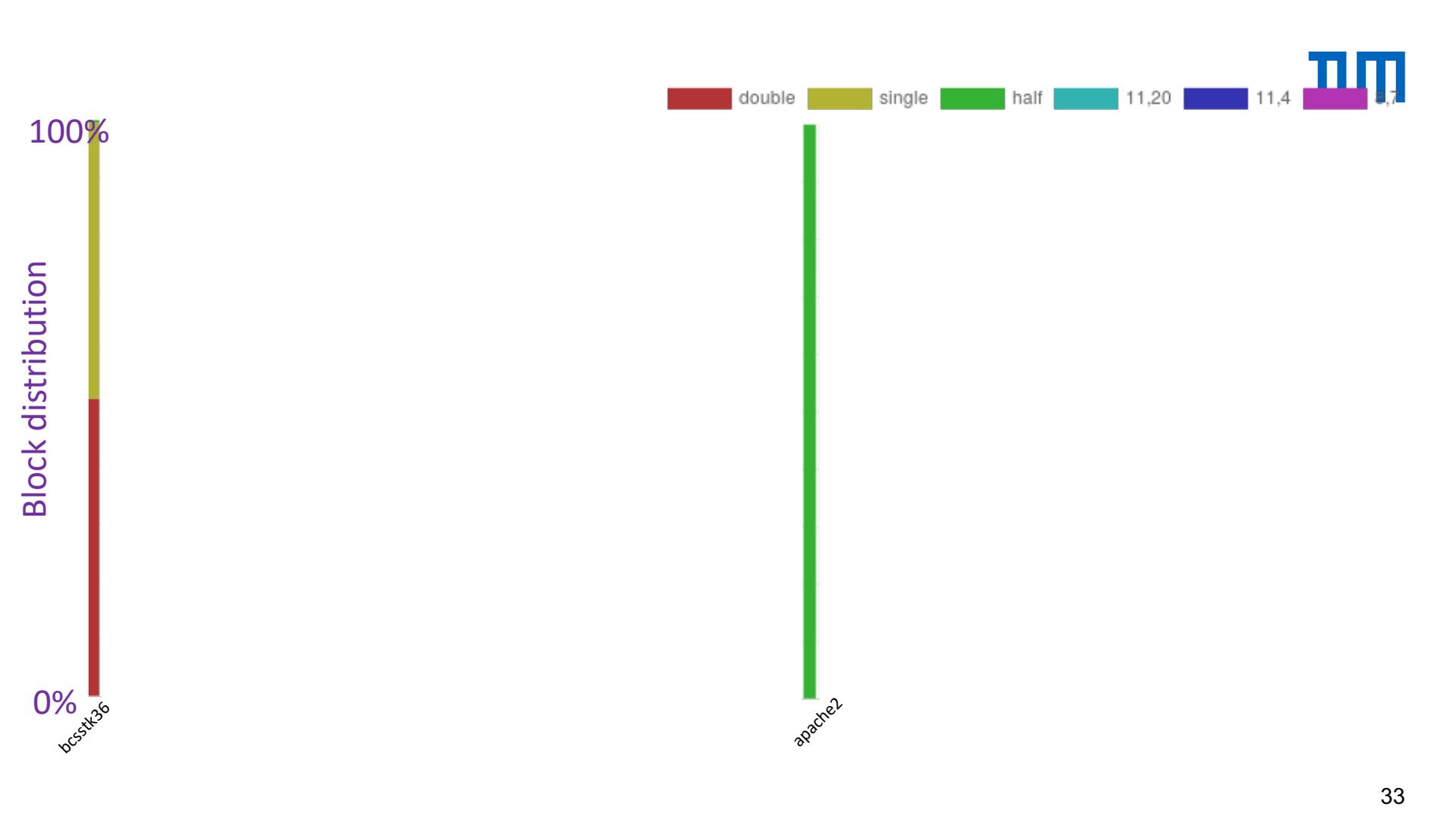
- + Regularity preserved;
- + Flexibility in the accuracy;
- + "Not a low precision preconditioner"
 - + Preconditioner is a constant operator;
 - + No flexible Krylov solver needed ;

- Overhead of the precision detection
(condition number calculation);
- Overhead from storing precision information
(need to additionally store/retrieve flag);
- Speedups / preconditioner quality problem-dependent;

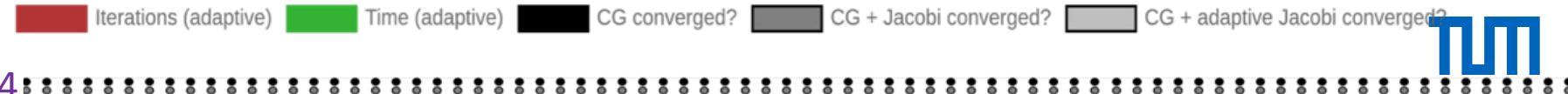
Block distribution

100%

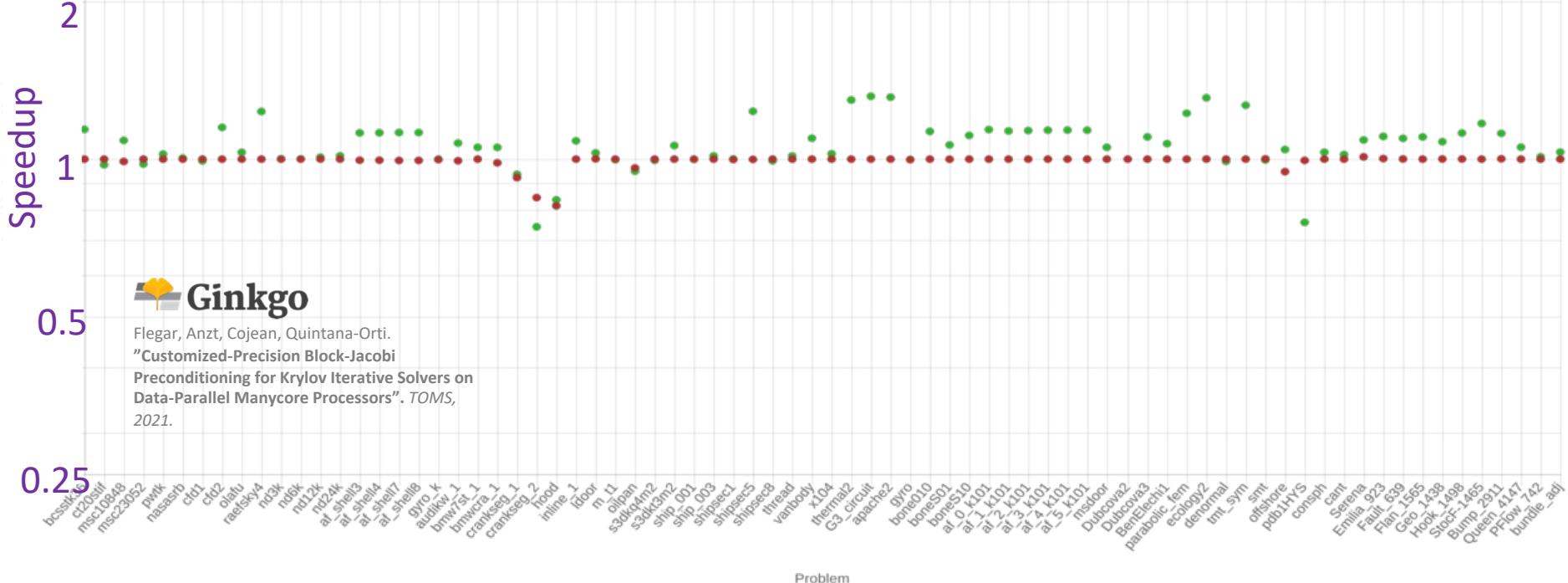




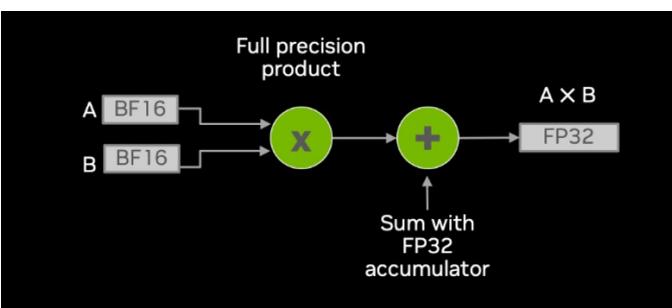
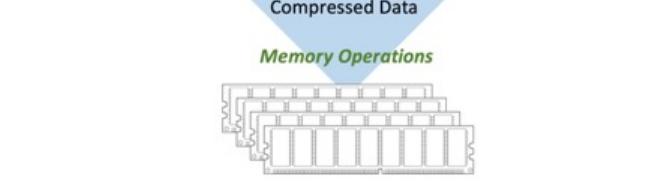
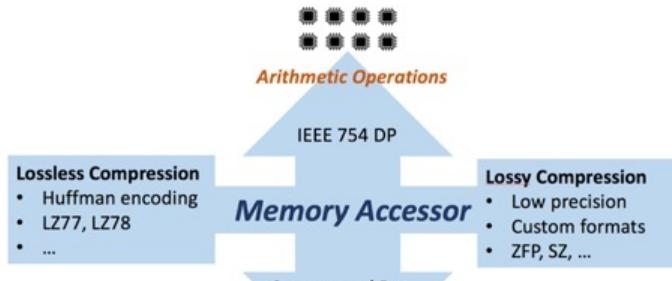




NVIDIA A100 GPU



Can mixed precision accelerate Sparse Solvers? – Yes, but we have to work harder...



NVIDIA Tensor Core operation.

- The performance is limited by the bandwidth
- Hardware increasingly features powerful matrix engines
- We should use these "free" FLOP/s
- One strategy is to use lossy/lossless data compression for memory operations and communication
- We need efficient on-chip data compression

We have no standard for sparse operations

cuSPARSE, rocSPARSE, oneAPI, PETSc, Trilinos...

all have different interfaces & functionality

We try to define an interface that allows for horizontal and vertical compatibility:

- Useful as building blocks for high-level algorithms
- Vendors can wrap their current interface
- Horizontal: bindings for Fortran, C, ...

SparseBLAS workshop

<https://icl.utk.edu/workshops/sparseblas2023/index.html>

Next week:

<https://icl.utk.edu/workshops/sparseblas2024/index.html>

Want to participate in the discussion

– reach out hantz@icl.utk.edu



Hartwig Anzt · Sie

Director of the Innovative Computing Laboratory (ICL) an...

1 Woche ·

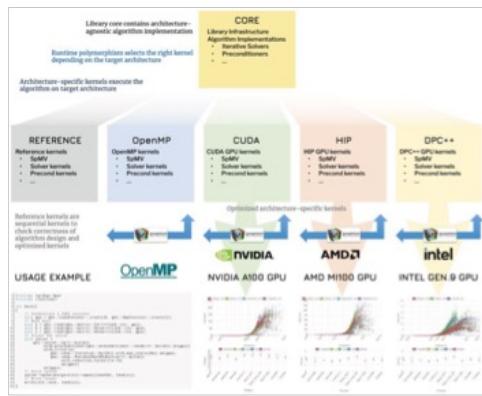
Three days of good talks, discussions, and prototyping are over: Thank you for the sparseBLAS workshop held at the Innovative Computing Laboratory (ICL) of the Tickle College of Engineering at the University of Tennessee. Together with experts from Intel Corporation, NVIDIA, AMD, Arm, MathWorks, University of California, Berkeley, Intel Labs, Computing at ORNL, Karlsruher Institut für Technologie (KIT), RIKEN, Lawrence Livermore National Laboratory, Sandia National Laboratories, and Massachusetts Institute of Technology, we worked on defining a common understanding and interface design for sparse linear algebra functionality.





DESIGN

Ginkgo is a C++ framework for sparse numerical linear algebra. Using a universal linear operator abstraction, Ginkgo provides basic building blocks such as the sparse matrix vector product for a variety of matrix formats, iterative solvers, and preconditioners. Ginkgo targets multi- and many-core systems, and currently features back-ends for AMD GPUs, Intel CPU/CPUs, NVIDIA GPUs, and OpenMP®-supporting architectures. Core functionality is separated from hardware-specific kernels for easy extension to other architectures, with runtime polymorphism selecting the correct kernels.



SUSTAINABILITY

Ginkgo is part of the extreme-scale Scientific Software Stack (EAS) and the extreme-scale Software Development Kit (sSDK), and adopts the sSDK community policies for sustainable software development and high software quality. The source code of the Ginkgo library can be accessed in a public git repository on GitHub. Code development in Ginkgo is realized in a Continuous Integration Continuous Benchmarking framework. GitLab runners are used on private servers and HPC clusters where Docker/Enroot is used to provide different compilation and execution environments. To test the correct execution, each functionality is complemented by unit tests. The unit testing is realized using the Google Test framework.

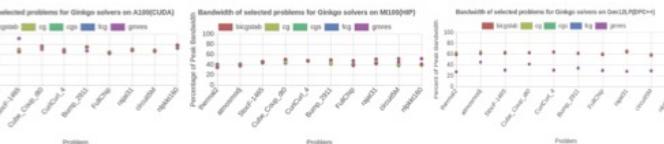


This research was supported by the Esascale Computing Project (10-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration, the Horizon2020 EuroHPC Project MICROCARD, and the Helmholtz Impuls und Vernetzungsfond VH-NI-124.



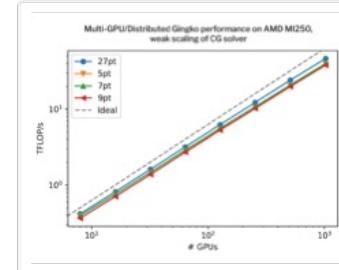
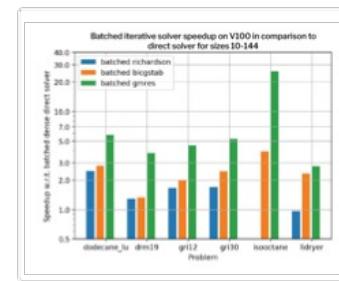
PERFORMANCE

NVIDIA

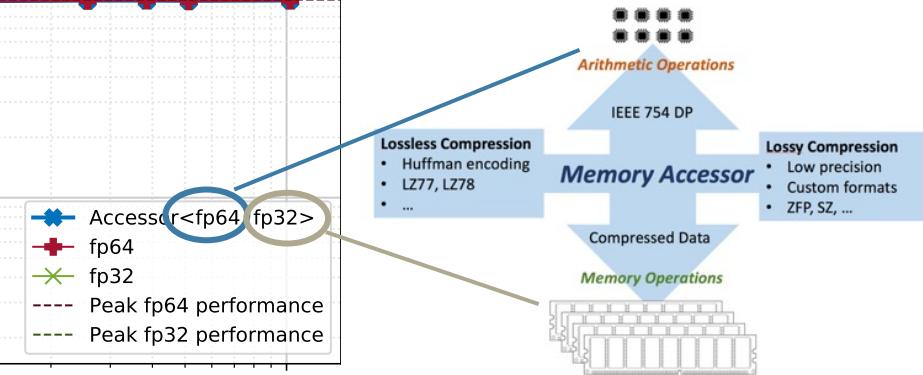
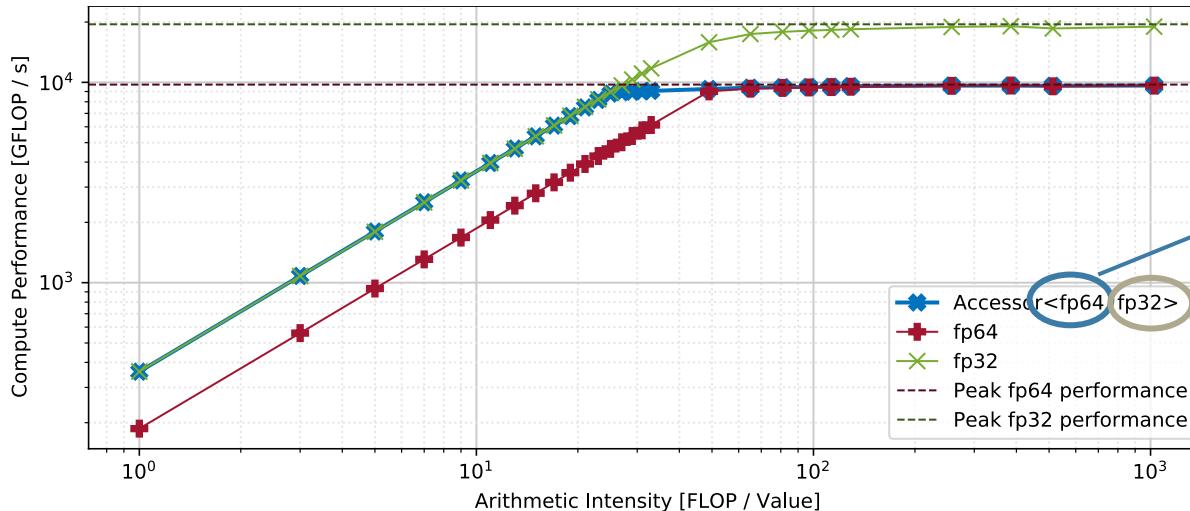


FUNCTIONALITY

Functionality		OMP	CUDA	HIP	DPC++
Basic	SpMV	✗	✗	✗	✗
	SpMM	✗	✗	✗	✗
Krylov solvers	SpGeMM	✗	✗	✗	✗
	BICG	✗	✗	✗	✗
BICGSTAB	BICGSTAB	✗	✗	✗	✗
	CG	✗	✗	✗	✗
CGS	CGS	✗	✗	✗	✗
	GMRES	✗	✗	✗	✗
IDR	IDR	✗	✗	✗	✗
	(Block-)Jacobi	✗	✗	✗	✗
ILU/IC	ILU/IC	✗	✗	✗	✗
	Parallel ILU/IC	✗	✗	✗	✗
Preconditioners	Parallel ILUT/ICT	✗	✗	✗	✗
	Sparse Approximate Inverse	✗	✗	✗	✗
Batched	Batched BICGSTAB	✗	✗	✗	✗
	Batched CG	✗	✗	✗	✗
Batched	Batched GMRES	✗	✗	✗	✗
	Batched ILU	✗	✗	✗	✗
Batched	Batched ISAI	✗	✗	✗	✗
	Batched Jacobi	✗	✗	✗	✗
AMG	AMG preconditioner	✗	✗	✗	✗
	AMG solver	✗	✗	✗	✗
Sparse direct	Parallel Graph Match	✗	✗	✗	✗
	Symbolic Cholesky	✗	✗	✗	✗
AMG	Numeric Cholesky	UNDER DEVELOPMENT			
	Symbolic LU	UNDER DEVELOPMENT			
Utilities	Numeric LU	✗	✗	✗	✗
	Sparse TRSV	✗	✗	✗	✗
Utilities	On-Device Matrix Assembly	✗	✗	✗	✗
	MC64/RCM reordering	✗			
Utilities	Wrapping user data		✓		
	Logging		✓		
Utilities	PAPI counters		✓		



Memory Accessor for NVIDIA A100 GPU



T. Grützmacher