# Customized-Precision Block-Jacobi Preconditioning for Krylov Iterative Solvers on Data-Parallel Manycore Processors

GORAN FLEGAR, Departamento de Ingeniería y Ciencia de Computadores, Universidad Jaime I, Spain

HARTWIG ANZT, Karlsruhe Institute of Technology, Germany and University of Tennessee, USA

TERRY COJEAN, Karlsruhe Institute of Technology, Germany

ENRIQUE S. QUINTANA-ORTÍ, Departamento de Informática de Sistemas y Computadores, Universitat Politècnica de València, Spain

Hardware graphics accelerators combined with low precision arithmetic have recently shown promise in the field of machine learning. In this work, we extend these ideas to the conventional problem of solving a sparse linear system. Specifically, we present a first high-performance graphics processing unit (GPU) implementation of the adaptive precision block-Jacobi preconditioner, which selects the precisions used to store the preconditioner data on-the-fly, taking into account the numerical properties of the system matrix. In addition, we explore the use of customized unconventional data formats (i.e., not supported by hardware), decoupling the storage format from the format used for arithmetic computations. Experiments run on the state-of-the-art accelerator hardware show that our implementation offers attractive runtime savings for certain classes of problems compared with the conventional full-precision block-Jacobi preconditioner. The developed implementation is augmented with a simple application programming interface (API), and is available as production code in the open-source Ginkgo sparse linear algebra library.

CCS Concepts: • **Mathematics of computing** → **Mathematical software**; *Arbitrary-precision arithmetic*;

Additional Key Words and Phrases: Sparse linear algebra, adaptive precision, preconditioning, block-Jacobi, Krylov solvers, GPU

## 1 INTRODUCTION

Krylov subspace methods (KSMs) were developed in the second half of the past century, with the derivation of the Lanczos algorithm creating Krylov subspaces [Lanczos 1952] and, in the same year, the publication of the Conjugate Gradient (CG) method for solving linear systems of equations [Hestenes and Stiefel 1952]. When applied to a linear system $Ax = b$ (with sparse coefficient matrix $A \in \mathbb{R}^{n \times n}$, right-hand side $b \in \mathbb{R}^n$, and unknowns $x \in \mathbb{R}^n$) KSMs started with

an initial guess $x_0$ produce a sequence of vectors $x_1, x_2, x_3, \ldots \in \mathbb{R}^n$ that, in general, progressively reduces the norm of the residuals $r_k = b - Ax_k$, eventually yielding an acceptable approximation for the solution of the system $Ax = b$.

The convergence rate of KSMs is largely dictated by the condition number of the system coefficient matrix $A$. Preconditioning schemes aim to accelerate the convergence of this type of solvers by transforming the original problem into the alternative preconditioned system $M^{-1}Ax = M^{-1}b$.

An ideal preconditioner $M^{-1} \in \mathbb{R}^{n \times n}$ is one that yields a transformed coefficient matrix $\hat{A} = M^{-1}A$ with a lower condition number than A, while admitting a software realization such that the preconditioner is relatively cheap to compute and inexpensive to apply (as in $M^{-1}b$).

Block-Jacobi preconditioning schemes are straight-forward extensions of the scalar Jacobi preconditioning [Saad 2003], which exhibit a high degree of parallelism, while offering superior convergence acceleration when applied to matrices exhibiting some inherent block structure. This, for example, is the case for problems arising from a finite element discretization of a partial differential equation (PDE) [Anzt et al. 2017a].

KSMs, enhanced with some form of a simple preconditioner, are memory-bound algorithms, meaning that their performance on current architectures is constrained by the bandwidth between the floating-point units (FPUs) and the memory where the data resides. In case the problem data is too large to fit into the cache memory of the processor(s), the increasing gap between the throughputs of the processor and the main memory (also known as the *memory wall* [Dongarra et al. 2014; Duranton et al. 2015; Lucas et al. 2014]), determines the low performance of this type of algorithms. This is a well-recognized problem, especially in the domain of sparse linear algebra operations, where *communication-avoiding* techniques are particularly appealing; see, e.g., [Cools 2018; Hoemmen 2010] and the references therein. Several works tackle the memory bottleneck in KSMs and related computational kernels by reducing the communication volume and memory footprint. For example, the authors of [Carson and Higham 2018] lower the data movement volume (and arithmetic cost) using the standard IEEE half/single/double precision formats [Commitee 2000] in combination with iterative refinement.[1] The compressed storage block (CSB) format [Buluç et al. 2009] is an orthogonal effort, which aims to decrease the indexing information necessary to maintain the sparse matrix with no numerical side-effects.

In [Anzt et al. 2019], we proposed yet another orthogonal technique that acts on the memory footprint of the problem data in order to address the memory bottleneck by reducing the precision format used to store the preconditioner. The approach was theoretically analyzed for a CG solver equipped with a block-Jacobi preconditioner that operates (that is, performs all arithmetic) in full double precision, while accessing the inverted diagonal blocks of the block-Jacobi preconditioner in a problem-adapted (potentially lower) precision. More precisely, all the problem data is stored in IEEE double precision format, except the blocks of the preconditioner, which are stored in either IEEE half/single/double precision formats, depending on their condition numbers. A type transformation is therefore required every time the preconditioner blocks stored in half or single precision in main memory are moved to the registers (where they are stored in double precision). The data transfer savings were estimated using a theoretical model that takes the floating point format and the convergence impact into account. For a significant portion of the symmetric positive definite matrices available in the SuiteSparse Matrix Collection [Davis and Hu 2011], we observed data transfer savings of up to 70% compared with a solver that handles all (preconditioner) data and arithmetic using double-precision.

---

[1]In the setting of the solution of linear systems, iterative refinement is an old technique, which dates back to the use of the first desk calculators, in the 1940s [Higham 2002].

This paper builds upon our preliminary theoretical analysis by making the following specific contributions:

(1) We improve an existing high performance implementation of the full precision block-Jacobi by leveraging the new cooperative groups CUDA application programming interface (API) to make the implementation compatible with the latest Volta and Turing generation hardware, while preserving performance on older graphics processing units (GPUs).

(2) We abandon the conventional workflow that tightly couples the memory precision format to the arithmetic precision format. Instead, we radically decouple the storage format for the diagonal blocks of the block-Jacobi preconditioner and use a much more compact precision format if the numerical properties allow for it.

(3) We augment the three conventional IEEE formats with three additional specialized floating point formats that can be efficiently converted into standard 64-bit format on current hardware and, compared to the IEEE formats, avoid overflows and underflows by sacrificing some precision. Using all six formats further increases potential memory savings, while using only the three non-IEEE formats that preserve the existing range of the values removes the need of additional preconditioner validity checks, reducing the preconditioner generation time. This contribution implicitly pushes in the direction of *decoupling the storage formats for floating-point numbers from the arithmetic formats* supported in the hardware of current Floating Point Units (FPUs).

(4) We propose an efficient compact layout to store the blocks of the block-Jacobi preconditioner, and evaluate the performance of a production code realizing that scheme in the framework of a high-performance Conjugate Gradient (CG) implementation on a NVIDIA Volta GPU. This experimental evaluation demonstrates the validity of the approach and reveals up to 30% performance improvement for a large range of real-world test problems.

(5) We provide a first high-performance GPU implementation of the adaptive precision block-Jacobi preconditioner and include both the adaptive and the full precision variants in the open-source Ginkgo library.

Our approach shares some of the appealing properties of the prototype in [Anzt et al. 2019]. Concretely, we employ full double precision in the generation and application of the preconditioner, as well as in all other arithmetic computations. Furthermore, we store part of the preconditioner in reduced precision, and convert it into full precision before proceeding with the arithmetic operations in the actual preconditioner application. Thus, our preconditioner still ensures that the preconditioning operator preserves orthogonality in double precision, implying that previously orthogonal Krylov vectors are orthogonal after the preconditioner application. In consequence, there is no need for flexible variants that introduce an additional orthogonalization step to preserve convergence [Golub and Ye 1999].

The rest of the paper is structured as follows. In Section 2 we briefly review the idea of KSMs and block-Jacobi preconditioning. More details about the adaptive precision block-Jacobi and the novel ideas we introduce in this paper are presented in Section 3. We elaborate on the high performance realization of the adaptive precision block-Jacobi in Section 4. We dedicate Section 5 to motivate the need for making novel algorithms and high performance implementations available in sustainable open source software. In Section 6, we present performance results for the block-Jacobi preconditioner generation and application and analyze the effectiveness and efficiency of the adaptive precision block-Jacobi preconditioner. Next, In Section 7 we discuss some central aspects of adaptive precision preconditioning in general and the experimental results in particular, and conclude in Section 8 with a summary of the findings and future research directions.

## 2 COMPUTATIONAL ASPECTS OF BLOCK-JACOBI PRECONDITIONED KSMS

Most instances of KSMs, such as CG, BiCG, GMRES, BiCGStab, etc., are comprised of a sequence of calls to simple computational kernels, such as the dot or inner product (DOT), AXPY-like vector updates and the sparse matrix-vector product (SpMV), inside an iteration loop [Saad 2003]. These kernels are all memory-bound operations, with a ratio between floating-point operations (FLOPs) and memory accesses (MEMOPs) that is $O(1)$, globally yielding a memory-bound solver.

Block-Jacobi preconditioners split the coefficient matrix into $A = L + M + U$, where the preconditioner defined by $M = \text{diag}(D_1, D_2, \ldots, D_m) \in \mathbb{R}^{n \times n}$, with $D_i \in \mathbb{R}^{m_i \times m_i}$ and $\sum_{i=1}^{m} m_i = n$, is a block-diagonal matrix containing the corresponding entries on the diagonal blocks of $A$, while $L, U \in \mathbb{R}^{n \times n}$ contain the elements of the coefficient matrix below and above those of $M$, respectively. (The scalar Jacobi preconditioner is a simple variant of the block counterparts with $m_i = 1$, $i = 1, 2, \ldots, m$, so that $M$ only contains the diagonal of $A$.) The block-Jacobi preconditioner is well defined if the diagonal blocks $D_i$ are all nonsingular. Furthermore, block-Jacobi preconditioning is particularly effective if the system matrix $A$ inherently presents a block structure (which is the case for many problems that arise from a finite element discretization of a PDE [Anzt et al. 2017a]) that is matched by the block structure of the Jacobi preconditioner.

In this work, we integrate a block-Jacobi preconditioner that explicitly computes the block-inverse matrix, $M^{-1} = \text{diag}(D_1^{-1}, D_2^{-1}, \ldots, D_m^{-1}) = \text{diag}(E_1, E_2, \ldots, E_m)$, before the iteration process of the KSM commences. The preconditioner is then applied within the KSM iteration in terms of a dense matrix-vector multiplication (GeMV) per inverse block $E_i$. Thus, the iteration for the preconditioned KSM remains a memory-bound process, as so is the GeMV kernel, independently of the block size $m_i$. In practice, the resulting preconditioner is of a comparable quality to the one computed by the conventional (and numerically more stable) strategy that computes the LU factorization (with partial pivoting) [Golub and Van Loan 1996] of each block ($D_i = L_i U_i$), and then applies the preconditioner using two triangular solves (per factorized block)[Anzt et al. 2018, 2017]. In exchange for a higher cost, the block-Jacobi preconditioner with explicit computation of the inverses presents the appealing property of yielding an application based on a highly parallel kernel (GeMV), compared with the constrained parallelism of the triangular systems that are necessary in the application of the LU-based preconditioning counterpart [Anzt et al. 2017b].

## 3 ADAPTIVE BLOCK-JACOBI PRECONDITIONING

### 3.1 Standard IEEE precision formats

In [Anzt et al. 2019], we proposed an adaptive block-Jacobi preconditioner that individually tunes the storage format of each block $D_i$ depending on its condition number. The scheme adopted in that work relies on three precision formats: 16-bit (fp16), 32-bit (fp32) and 64-bit (fp64), which correspond to the standard IEEE half, single and double precision formats [Commitee 2000], respectively. In detail, the adaptive block-Jacobi preconditioner proceeds as follows:

(1) Before the iteration commences, we explicitly compute the inverse of each block using fp64: $D_i \rightarrow E_i$.
(2) At the same stage (i.e., before the iterative solver is started), we compute $\kappa_1(D_i) = \kappa_1(E_i) = \|D_i\|_1 \|D_i^{-1}\|_1 = \|D_i\|_1 \|E_i\|_1$. As $E_i$ is explicitly available, computing $\kappa_1(D_i)$ is straightforward and inexpensive compared with the inversion of the block [Anzt et al. 2018].
(3) After inverting the diagonal block $D_i$ in fp64, we store the inverted diagonal block $E_i$ in the format determined by its condition number—truncating the entries of the block if necessary.

Precisely, we store $E_i$ in

$$
\begin{cases}
\text{fp16} & \text{if } \tau_h^L < \kappa_1(D_i) \leq \tau_h^U, \\
\text{fp32} & \text{if } \tau_s^L < \kappa_1(D_i) \leq \tau_s^U, \text{ and} \\
\text{fp64} & \text{otherwise,}
\end{cases}
\tag{1}
$$

where the thresholds $\tau$ are set as $\tau_h^L = 0$ and $\tau_h^U = \tau_s^L$.

(4) During the iteration, we recover the block $E_i$, stored in the corresponding format in memory (as determined by (1)), transform its entries to fp64 in the processor's registers, and apply the block in terms of a fp64 GEMV.

Due to the use of the standard formats for half, single and double precision, in the above procedure the truncation can result in either overflows or underflows, whose consequences need to be tackled. Here we only discuss the second case and refer the reader to [Anzt et al. 2019] for the handling of overflows. The risk associated with underflow is that the truncation may turn a non-zero (but close to zero) value in fp64 into a zero which in turn can make $E_i$ an ill-conditioned (or even singular) block, thereby causing numerical difficulties for the convergence of the KSM. In order to avoid this issue, we examine the condition number of the truncated representation of $E_i$, and discard the use of the corresponding reduced precision if it was above a given threshold $\tau_\kappa$.

## 3.2 Unconventional precision formats

In addition to the three floating point formats defined by the IEEE standard, this work augments the set of considered precisions with three additional formats that can be cheaply processed using the instruction set of NVIDIA GPUs. While it would be theoretically possible to employ any combination of exponent and significand bits, the complexity of purely software-based format conversion could prove detrimental to performance. However, conversions for several particular precision configurations can be implemented efficiently.

In particular, if the conversion to lower precision preserves the number of exponent bits and the rounding mode is limited to *round-to-zero*, the conversion to lower precision consists of significand truncation, only. Converting back to full precision then conversely adds zeros as the missing significand bits[Anzt et al. 2019]. Using the notation $fp_{ex,snf}$ (where *ex* and *snf* denotes the number of exponent bits and significand bits, respectively), this procedure can be used on the 64 bit IEEE double precision format ($fp_{11,52}$) with 11 exponent and 52 significand bits to obtain an alternative 32 bit floating point format with 11 exponent and 20 significand bits ($fp_{11,20}$) by dropping the 32 low-order bits of the original format. The range of such a format stays roughly the same as that of IEEE double precision, and the unit roundoff (adjusted for the *round-to-zero* rounding mode) is $u = 9.54e - 7$. A 16-bit format based on IEEE double ($fp_{11,4}$) can also be obtained by dropping the 48 low-order significand bits. The result is a format with 11 exponent and 4 significand bits and unit roundoff $u = 6.25e - 2$. A 16-bit format can also be obtained by basing it on the 32 bit IEEE single precision ($fp_{8,23}$). Such a format ($fp_{8,7}$) has 8 exponent and 7 significand bits, and unit roundoff of $u = 7.81e - 3$.

The additional formats offer a trade-off by providing formats of the same size as their IEEE counterparts, but with larger range and lower precision. They can be used to store a block that is relatively well conditioned (and thus does not require high precision to achieve reasonable accuracy [Anzt et al. 2019]), but the range of values in the block is such that a conventional format would cause catastrophic overflows or underflows. The improved format selection strategy selects the first format from the list $fp_{5,10}, fp_{8,7}, fp_{11,4}, fp_{8,23}, fp_{11,20}, fp_{11,52}$ whose unit roundoff is small enough to deliver the required accuracy, and where the exponent range avoids catastrophic overflows and underflows. The list is sorted by increasing sizes of the formats, which means that
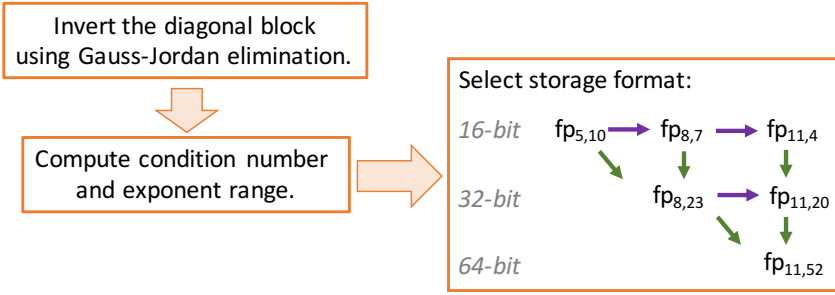
Fig. 1. Workflow for generating the inverse block and selecting a suitable storage format. The horizontal arrows (purple) reflect bitcount-constant traversals addressing overflow and underflow, the green vertical arrows represent significand extensions for increasing the accuracy to the requirements imposed by the condition number.

the procedure selects the smallest format capable of delivering the required accuracy. Within the same format size, the list is sorted so that priority is given to the format that offers more accuracy. In Figure 1 we visualize the process of generating the block-Jacobi preconditioner and selecting a suitable storage format.

## 4 CUDA IMPLEMENTATION

### 4.1 Previous work

As a starting point for the implementation of adaptive block-Jacobi kernels, we use a previous prototype CUDA implementation of full precision block-Jacobi [Anzt et al. 2018]. The implementation includes an optimized kernel for block-Jacobi preconditioner generation which extracts the diagonal blocks from the sparse system matrix stored in Compressed Sparse Row (CSR [Saad 2003]) format, inverts them, and stores the inverses into the GPU main memory. For each block, the entire pipeline is executed using a single warp (a group of 32 GPU cores, roughly equivalent to a 32-wide vector unit) with each core processing a single column of the matrix. The inversion is realized via the highly parallel Gauss-Jordan Elimination (GJE) algorithm, and the explicit inverse diagonal blocks are stored in row-major order to enable coalesced access both when extracting the blocks from the sparse structure, as well as when storing the inverses back into memory. The generation pipeline leverages the extensive register storage available in recent CUDA architectures (up to 32 KB per warp) to keep the entire block in processor registers during the computation and completely avoid expensive data access. This strategy allows to efficiently process double precision blocks of up to 32 rows and columns.

The second component of the prototype is a custom implementation of the preconditioner application procedure. Once again, each warp is responsible for processing a single preconditioner block. First, the section of the input vector corresponding to the block is read into the registers and distributed among the threads of the warp. Then, for each row of the block, the warp collaboratively reads the values in the row, forms a dot product between the input vector (already present in the registers) and the row, and writes the result to the output vector. Processing the blocks stored in row-major in this way ensures contiguous access to the main memory.

A final optimization included in the initial prototype involves the processing of small blocks. If all the preconditioner blocks are smaller than some dimension $k < 32$, a more efficient version of the kernel can be generated by having each thread of the warp use an array large enough to store only $k$ instead of 32 values. This reduces the resource requirements of the warp, allowing the GPU

Sequential block storage

Interleaved block storage

Legend:
- Jacobi block
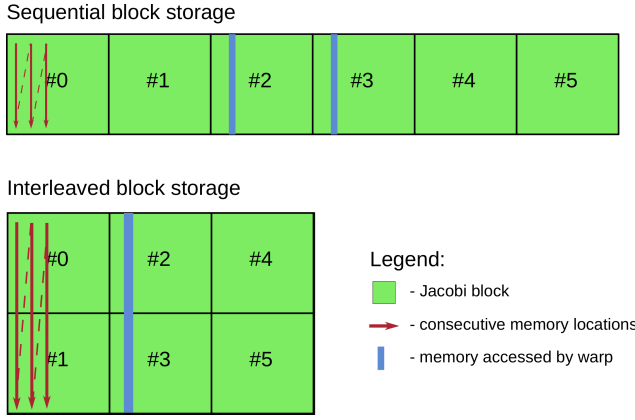- consecutive memory locations
- memory accessed by warp

Fig. 2. Preconditioner storage scheme. Top: sequential storage used by the initial implementation. Bottom: block-interleaved storage used by the new implementation.

to simultaneously schedule more warps per multiprocessor. In addition, for small values of $k$, a warp can be logically split into two (or more) sub-warps; then, instead of using the entire warp to process a single block, each sub-warp can handle the generation of one preconditioner block. Precisely, for a maximum block size $s$, every warp handles $2^{5-\lceil \log_2 s \rceil}$ blocks:

$$
\begin{array}{ll}
32 \geq s > 16 & 1 \text{ block per warp,} \\
16 \geq s > 8 & 2 \text{ blocks per warp,} \\
8 \geq s > 4 & 4 \text{ blocks per warp,} \\
4 \geq s > 2 & 8 \text{ blocks per warp,} \\
2 \geq s > 1 & 16 \text{ blocks per warp,} \\
1 = s & 32 \text{ blocks per warp.}
\end{array}
$$

To enable these optimizations, we generate a kernel optimized for each maximal block size $\hat{k} = 1, 2, \ldots, 32$.

## 4.2 Kernel improvements

Before implementing the adaptive precision version of the block-Jacobi preconditioner, we first incorporate several improvements to the full precision block-Jacobi preconditioner.

Starting with CUDA toolkit version 9.0, NVIDIA updated the warp shuffle and warp vote APIs used for intra-warp communication to support the new Volta architecture that features relaxed warp execution constraints [NVIDIA Corporation 2018]. While the APIs used by the previous implementation of block-Jacobi kernels are still available (albeit deprecated), using them causes the kernel to stall[2] when run on the Volta architecture. In addition to the updated low-level APIs, the CUDA toolkit version 9.0 also includes a new cooperative group APIs which encapsulates the details of the low-level APIs. Instead of using the low level API directly, we decided to modify our code to use this high-level alternative as it provides more flexibility and can potentially enable better compatibility with future CUDA versions.

We also identified several additional performance optimizations concerning the memory layout of the block-Jacobi preconditioner, specifically the question of storing the blocks in row-major

---

[2]Since only a subset of the warp was calling the API in the original implementation.

vs. column-major layout. A detailed analysis of the preconditioner application kernel explained in Section 4.1 revealed that the time needed for intra-warp communication in the collaborative computation of the dot product (necessary in a row-major block storage) is significant compared with the time needed to load the data from memory, so improving that part of the kernel can render performance gains. For this reason, we change the data layout of the preconditioner blocks to use column-major instead of row-major storage. This enables efficient column-wise access of the block – equivalent to a column-major GEMV for each Jacobi block. The downside of this approach is that the block data has to be transposed after the inversion, which results in suboptimal memory accesses during the preconditioner generation step. However, since the preconditioner is generated only once, but applied multiple times (at least once per KSM iteration), we expect this change in storage layout will render performance improvements for most use cases.

The final improvement aims at processing small blocks more efficiently. The original implementation stores consecutive blocks in sequence, as depicted in the top part of Figure 2. With such storage, memory access during preconditioner application is optimal for large blocks. However, as soon as the maximal block size becomes small enough to split the warp into sub-warps, so that several blocks are processed by the same warp, this no longer holds. Since the corresponding columns of consecutive blocks are not consecutive in memory, reading them causes suboptimal strided memory access. To eliminate this problem, we replace the sequential storage scheme with the block-interleaved storage shown in the bottom part of Figure 2. The new scheme groups all blocks processed by a warp together, and interleaves the storage of their columns. Precisely, the scheme initially stores the first columns of all blocks in the group, then proceeds with storing the second columns, etc.; this strategy ensures contiguous memory accesses during preconditioner application.

The last two optimizations are essential to enable performance improvements via low precision storage. Without the former, communication would dominate the cost of preconditioner application, severely limiting the benefit of reduced data transfers. Without the latter, accessing small blocks would incur unnecessary data loads into cache. Since the size of the cache lines is fixed, reducing the size of the individual elements would just increase the amount of memory being wasted, without reducing the total data movement volume.

### 4.3 Adaptive block-Jacobi

Extending the full precision block-Jacobi to the adaptive precision variant requires adding the precision detection logic to the preconditioner generation, storing the blocks in appropriate precision together with metadata specifying which precision is employed for the distinct blocks and, during preconditioner application, restoring the original block on the fly from low precision storage using the metadata.

The precision selection method we employ is that explained in Section 3.1, enhanced with the additional formats introduced in Section 3.2. The condition number of the block is determined by computing the matrix 1-norm of the block before and after inverting it [Anzt et al. 2018]. The condition number is then evaluated against the unit roundoffs to select the optimal format using the format priority list we introduced in Section 3.2. For precisions that require additional protection against catastrophic underflow or overflow (IEEE singe and half), the conditioning of the inverse stored in lower precision is computed by converting each value of the inverse block to lower precision, converting it back to double precision, followed by norm calculation, inversion, and another norm calculation — all in double precision. This way, the condition number is computed with high accuracy. Before reducing the precision, a copy of the full precision inverse is backed up to main GPU memory. This allows to retrieve the full precision inverse afterwards (if necessary). When a group of blocks is processed by a single warp (in case of small blocks), the precision is not

393 decided individually, but jointly for the entire group of blocks, using the first precision in the list
394 from Section 3.2, which is suitable for storing all blocks. This is done for performance reasons, as
395 trying to execute different instructions by threads belonging to the same warp — which would be
396 necessary to read values stored in different precisions — would lead to thread divergence, and the
397 serialization of these instructions, ultimately resulting in a significant slowdown.

398    Since the final precisions are not known before inverting the blocks, a memory workspace large
399 enough to store all blocks in double precision is allocated before launching the preconditioner
400 generation kernel. Once the storage precision is decided, low precision blocks are stored using only
401 the first part of the workspace they are assigned to, while the rest of the workspace remains unused
402 (fragmentation). While it would be possible to post-process the block storage structure to remove
403 unused "gaps" via de-fragmentation, doing so would not reduce the total memory transfer volume
404 during preconditioner application, since the total storage required for the group of blocks is a
405 multiple of the cache line size in any precision, as long as the block size is at least 2. In consequence,
406 the "gaps" will never be transferred from main memory to the cache. Thus, removing gaps is only
407 attractive in case the total memory footprint of the preconditioner is a relevant factor. We refrain
408 from incorporating de-fragmentation in our implementation.

409    A distinct memory block is used to store the information about the precisions used for the
410 inverted blocks. The precision of each block is encoded using 8 bits, which is the smallest amount of
411 data that can be independently stored and loaded from memory. This information is retrieved during
412 the preconditioner application stage to determine the storage locations and precision formats of
413 individual blocks and select the correct conversion procedure.

## 5   USABILITY, REPRODUCIBILITY AND SUSTAINABILITY EFFORTS

416 As not only modern hardware but also the software that can effectively utilize the hardware
417 resources becomes increasingly complex, it can no longer be expected that novel algorithms or high
418 performance implementations presented in scientific publications are adequately explained so that
419 the readers can reproduce an implementation of equivalent quality. Furthermore, domain scientists
420 who can potentially benefit from such work, should not be required to understand low-level
421 optimization techniques needed to produce a high performance implementation. In consequence,
422 it is becoming increasingly important to openly publish high performance implementations and
423 simplify their integration into other software ecosystems.

424    To address these issues, we integrate both the full-precision block-Jacobi preconditioner as well
425 as the adaptive precision variant into the open source Ginkgo linear algebra package[3]. Ginkgo
426 is a C++ library originally designed for the iterative solution of sparse linear systems. It features
427 various matrix formats and solvers with high performance implementations for both GPU and
428 CPU architectures, and allows for the easy integration into existing software stacks. At this point,
429 the (adaptive) block-Jacobi preconditioner is available in "reference mode" (a single threaded
430 straightforward CPU implementation that can be used for correctness checking and evaluating the
431 convergence benefits of the preconditioner) as well as in "CUDA mode", with the latter featuring
432 the high performance GPU implementation described in this work. A high performance CPU
433 implementation based on OpenMP parallelization is planned, but not yet available.

434    Adding the block-Jacobi preconditioner into a larger software effort provides the benefits of
435 reusing existing workflows: Ginkgo's low-level building blocks are utilized to simplify the imple-
436 mentation; its unit testing framework extended to include tests for the block-Jacobi preconditioner
437 which are then automatically run using Ginkgo's continuous integration (CI) system; the bench-
438 marks for block-Jacobi are integrated into Ginkgo's Continuous Benchmarking (CB) framework

---

440 [3]https://ginkgo-project.github.io

```
1     // Read data
2     auto A = share(gko::read<mtx>(std::ifstream("data/A.mtx"), exec));
3     auto b = gko::read<vec>(std::ifstream("data/b.mtx"), exec);
4     auto x = gko::read<vec>(std::ifstream("data/x0.mtx"), exec);
5
6     // Generate solver
7     auto solver_gen =
8         cg::build()
9  +            .with_preconditioner(
10 +                gko::preconditioner::Jacobi<>::build().on(exec))
11             .with_criteria(
12                 gko::stop::Iteration::build().with_max_iters(20u).on(exec),
13                 gko::stop::ResidualNormReduction<>::build()
14                     .with_reduction_factor(1e-20)
15                     .on(exec))
16             .on(exec);
17     auto solver = solver_gen->generate(A);
18
19     // Solve system
20     solver->apply(lend(b), lend(x));
```

```
1     // Read data
2     auto A = share(gko::read<mtx>(std::ifstream("data/A.mtx"), exec));
3     auto b = gko::read<vec>(std::ifstream("data/b.mtx"), exec);
4     auto x = gko::read<vec>(std::ifstream("data/x0.mtx"), exec);
5
6     // Generate solver
7     auto solver_gen =
8         cg::build()
9  +            .with_preconditioner(
10 +                gko::preconditioner::Jacobi<>::build()
11 +                    .with_storage_optimization(
12 +                        gko::precision_reduction::autodetect())
13 +                    .on(exec))
14             .with_criteria(
15                 gko::stop::Iteration::build().with_max_iters(20u).on(exec),
16                 gko::stop::ResidualNormReduction<>::build()
17                     .with_reduction_factor(1e-20)
18                     .on(exec))
19             .on(exec);
20     auto solver = solver_gen->generate(A);
21
22     // Solve system
23     solver->apply(lend(b), lend(x));
```

Fig. 3. Changes needed to enhance Ginkgo's `simple_solver` usage example with the full precision block-Jacobi preconditioner (top) and adaptive precision block-Jacobi preconditioner (bottom).

and can be run separately, or in conjunction with the rest of the benchmarks [Anzt et al. 2019]. From the user's perspective, the integration reduces the amount of software that has to be installed as well as simplifies the installation (as Ginkgo uses the well-established CMake build system) and integration of the software. If the user is already using Ginkgo as part of an application, adding the adaptive block-Jacobi preconditioner requires only a few of additional lines of code, as shown in Figure 3. If the application does not (yet) use Ginkgo, its library interoperability features can be used to wrap existing data structures into Ginkgo objects, which can then be used to construct and apply the preconditioner. Finally, as all the unit tests and benchmarks contributed in this work are distributed as part of the Ginkgo ecosystem, it should be relatively simple for the user to verify the correctness and reproduce the performance results we present in this paper, or even to evaluate the kernels' performance on a different CUDA-supporting architecture.

## 6 EXPERIMENTAL EVALUATION

This section evaluates the numerical properties and the effectiveness, efficiency, and performance of the developed adaptive precision algorithm and the low level kernels on recent CUDA-supporting GPUs. Initially, we assess the performance gains available from the improvements to the full
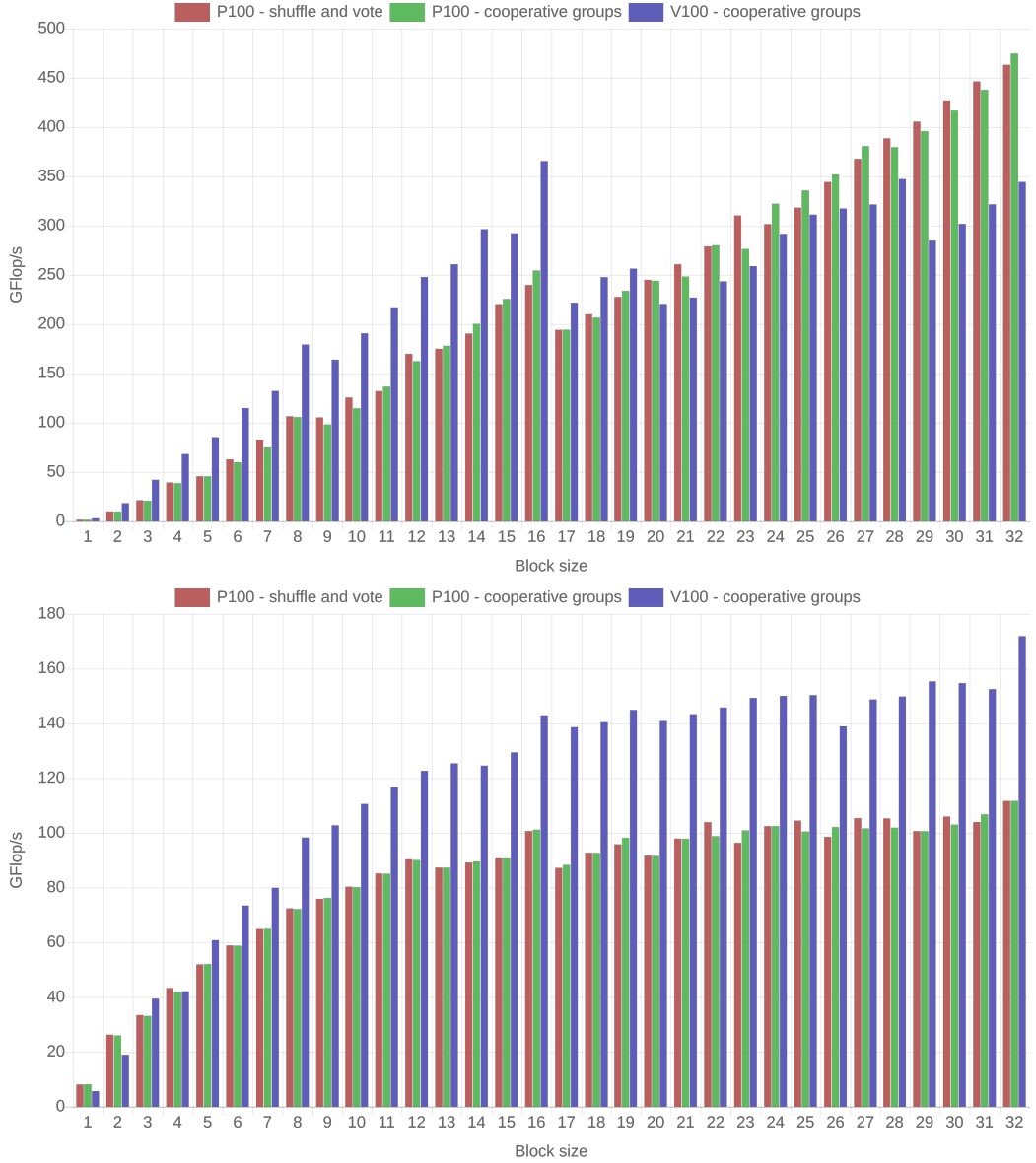
Fig. 4. Effects of using the higher-level cooperative groups API (green and blue) over the low-level warp shuffle and warp vote APIs (red). The results with cooperative groups are shown for the older P100 GPU (green) and newer V100 GPU (blue). The top plot shows the performance of the preconditioner generation (not including block size detection) and the bottom plot the performance of the preconditioner application.

precision block-Jacobi preconditioner. Then, the optimized full precision kernels are compared with the adaptive precision variant.

Two hardware setups are used in the experiments. The first one is a GPU-accelerated node of a compute cluster at the University of Jaume I (UJI). The node is composed of an 8-core Intel Xeon

E5-2620 v4 CPU with 32 GB of RAM and an NVIDIA TESLA P100 (PCI-e form factor) GPU with 16 GB of HBM2 memory. The accelerator achieves a peak double precision performance of 4.7 TFlop/s and a peak memory bandwidth of 732 GB/s.

The second setup is the Summit supercomputer at the Oak Ridge National Laboratory. Our experiments use a single node containing two 22-core IBM POWER9 CPUs with 256 GB of RAM and 6 NVIDIA TESLA V100 (SXM2 form factor) GPUs with 16 GB of HBM2 memory. For our experiments, we use only a single NVIDIA V100 GPU with a peak double precision performance of 7.8 TFlop/s and a peak memory bandwidth of 900 GB/s.

## 6.1 Effects of using the cooperative group APIs and the newer Volta architecture

While the C++ language and its compilers are designed to enable zero-overhead abstractions, there is always a possibility that a particular abstraction is not properly translated by the compiler and does not generate sufficiently optimized code. Thus, we first evaluate the effect of replacing the low-level warp shuffle and vote APIs used in the original code with the higher-level cooperative groups API. Since the initial version using the low-level APIs does not work correctly on the new Volta generation hardware used by the Summit system, the evaluation was realized on the UJI cluster, which features the older Pascal generation P100 GPU. For the new implementation that supports the recent Volta architecture, we also include the results obtained on the Summit system and the V100 GPU to study the effects of switching to newer hardware.

The experiments were performed using synthetically-generated block-diagonal matrices with varying block sizes and a total of 50,000 equally-sized blocks per matrix. The maximal size of preconditioner blocks was set to match the block size of the matrix. The nonzero locations were filled with randomly-generated small floating point numbers uniformly distributed between $-1$ and 1.

The results presented in Figure 4 reveal that there is no significant performance difference when moving to the higher level API. At the same time, the version using the cooperative groups API supports the new Volta architecture. The preconditioner generation stage takes a slight performance hit for large blocks due to the Volta architecture increasing register pressure to support the relaxed warp execution model [Anzt et al. 2018]. However, the more relevant preconditioner application stage does not suffer from this problem, and exhibits performance improvements between 40% and 50% on the V100.

## 6.2 Memory layout improvements

Next, we evaluate the effect of the two preconditioner memory layout optimizations: the block-level optimization that employs column-major instead of row-major storage and the preconditioner-level optimization based on block-interleaved instead of sequential block storage. We run the experiments on Summit's V100 GPU, using the same synthetic benchmarks as in Section 6.1.

Figure 5 demonstrates that changing the storage scheme from row-major to column-major slightly reduces the performance of preconditioner generation as storing the preconditioner data causes non-coalesced memory access. On the other hand, the performance of preconditioner application increases for all block sizes, due to the availability of a more efficient matrix-vector product algorithm. As mentioned earlier, since the preconditioner is generated only once and applied multiple times (and the cost of generating the block-Jacobi preconditioner is negligible compared with the solver runtime [Anzt et al. 2018]), the overall performance of the solver is improved.

The second optimization only has impact on blocks with at most 16 rows and columns. This is expected, as for larger blocks the two storage schemes result in exactly the same data layout. For those cases where the storage layout is different, marginal benefits can be identified in favour
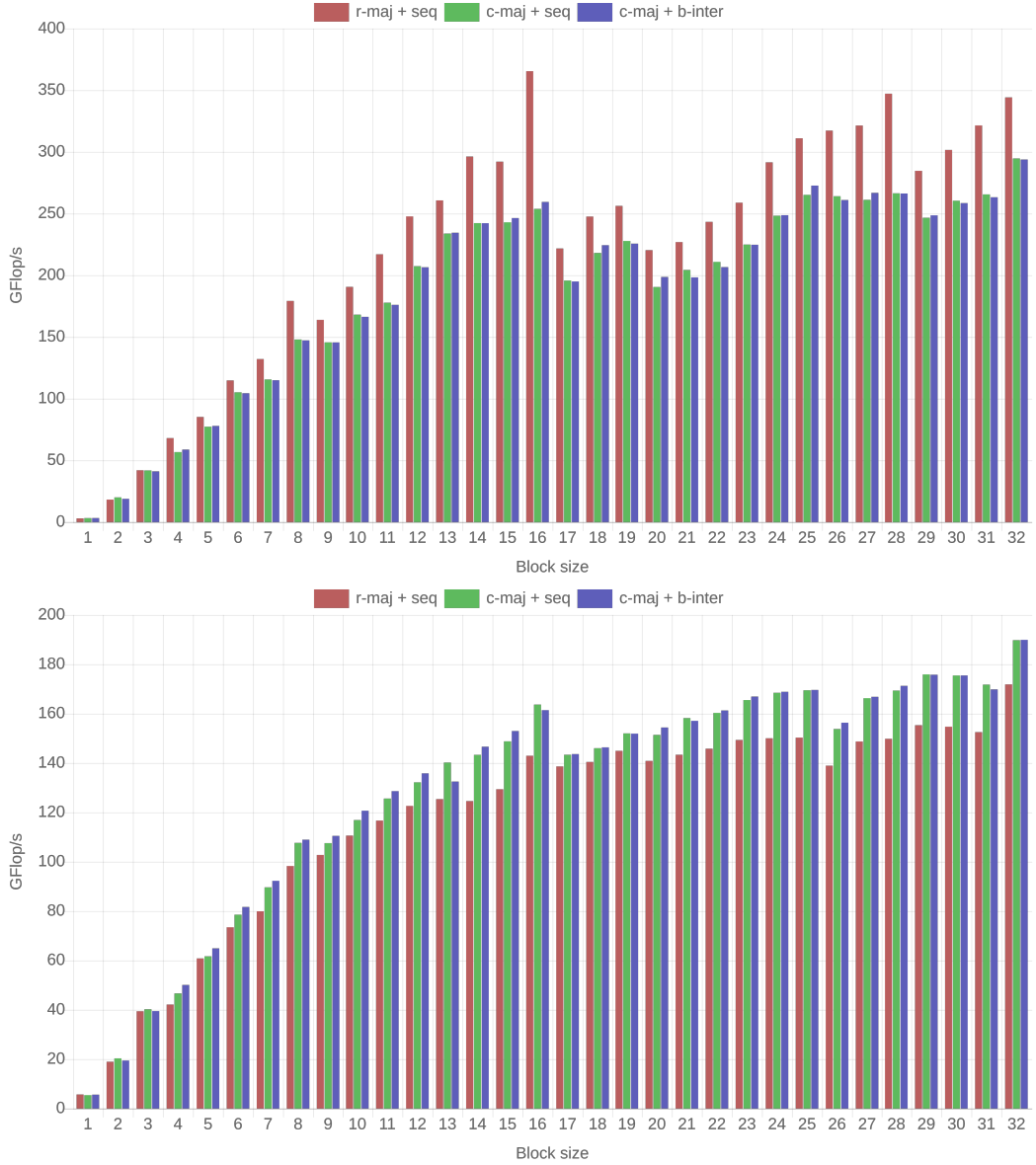
Fig. 5. Performance of memory layout optimizations on the V100 GPU. The original row-major, sequential block storage is shown in red, column-major, sequential block storage in green and column-major, block-interleaved block storage in blue. The top plot shows the performance of the preconditioner generation (without block size detection) and the bottom plot the performance of the preconditioner application.

of the block-interleaved layout. We expect that these benefits become more pronounced in the adaptive block-Jacobi variant, since unfavorable cache access is more detrimental when dealing with smaller data types.

## 6.3 Adaptive precision

Having analyzed the effects of the additional improvements applied to the full precision block-Jacobi, we now turn our attention to the adaptive variant. Once again, we run the experiments on Summit's V100 GPU and use synthetic benchmarks described in Section 6.1. For the full precision version, we evaluate a kernel incorporating all the optimizations described in previous sections: cooperative groups API, column-major storage and block-interleaved block storage. For the adaptive version, a kernel featuring all these optimization steps and additional support for adaptive precision is used. We report results where the autodetection system was disabled, and the precision used for all blocks is fixed beforehand to the same value. This offers an upper bound on the theoretical performance improvement that can be expected if all blocks can be stored in the same precision. On real-world problems (covered in the next section), the actual performance improvement highly depends on the condition number distribution of the diagonal blocks of the system matrix, which is difficult to replicate with synthetic benchmarks. However, since disabling autodetection means that part of the preconditioner generation kernel is skipped, we additionally report performance results that account for the autodetection: a variant that selects only between the three precision formats that do not require additional condition number calculation, and a variant with full autodetection using all six supported precisions formats.

Figure 6 shows the results for the generation and application stages of the adaptive precision block-Jacobi preconditioner and all six supported precisions. While there are some performance improvements available when using lower precision in the generation stage due to the reduction of the total time needed to store low-precision blocks, the improvements in the application stage are of higher importance. These improvements are more pronounced for larger block sizes, where the preconditioner's memory footprint becomes more relevant compared with the footprint of the input and output vectors. In total, low-precision blocks can yield up to 1.7× and 2× speedups for 32-bit and 16-bit storage schemes, respectively. Another interesting observation concerns the effect of automatic precision detection on the preconditioner generation time. When using only the three non-standard formats that preserve the representable range of values, there is virtually no impact on the performance of preconditioner generation. On the other hand, using all six formats can lead to performance degradation of up to a factor of 2×. This implies that, in case the solver is expected to converge in a few iterations (where preconditioner generation represents a high fraction of the total runtime), it may be beneficial to only use the three non-standard precisions, or to maintain the full precision block-Jacobi.

## 6.4 Full solver runtime

Finally, we compare the effectiveness of a solver enhanced with adaptive precision block-Jacobi with that of a solver enhanced with the full precision variant. For this experiment, we use a set of matrices from the SuiteSparse matrix collection[4], arising in real-world applications. Only symmetric positive definite problems that have between $10^6$ and $5 \times 10^8$ nonzeros, and where a block-Jacobi enhanced Conjugate Gradient (CG) solver needs more than 100 iterations are considered, as these problems justify the use of a preconditioned CG solver on GPUs. The CG solver available in the Ginkgo library was employed. For both preconditioners, the double precision standard block-Jacobi as well as the adaptive precision block-Jacobi, we automatically detect natural diagonal blocks using the supervariable amalgamation algorithm. While we recognize that it may be possible to design a more advanced method of block-detection, this is still an area of active research [Goetz and Anzt 2018] that remains outside the scope of this paper. The block size upper limit was set to 32. The solvers were run for at most 10,000 iterations, or until the initial residual norm was
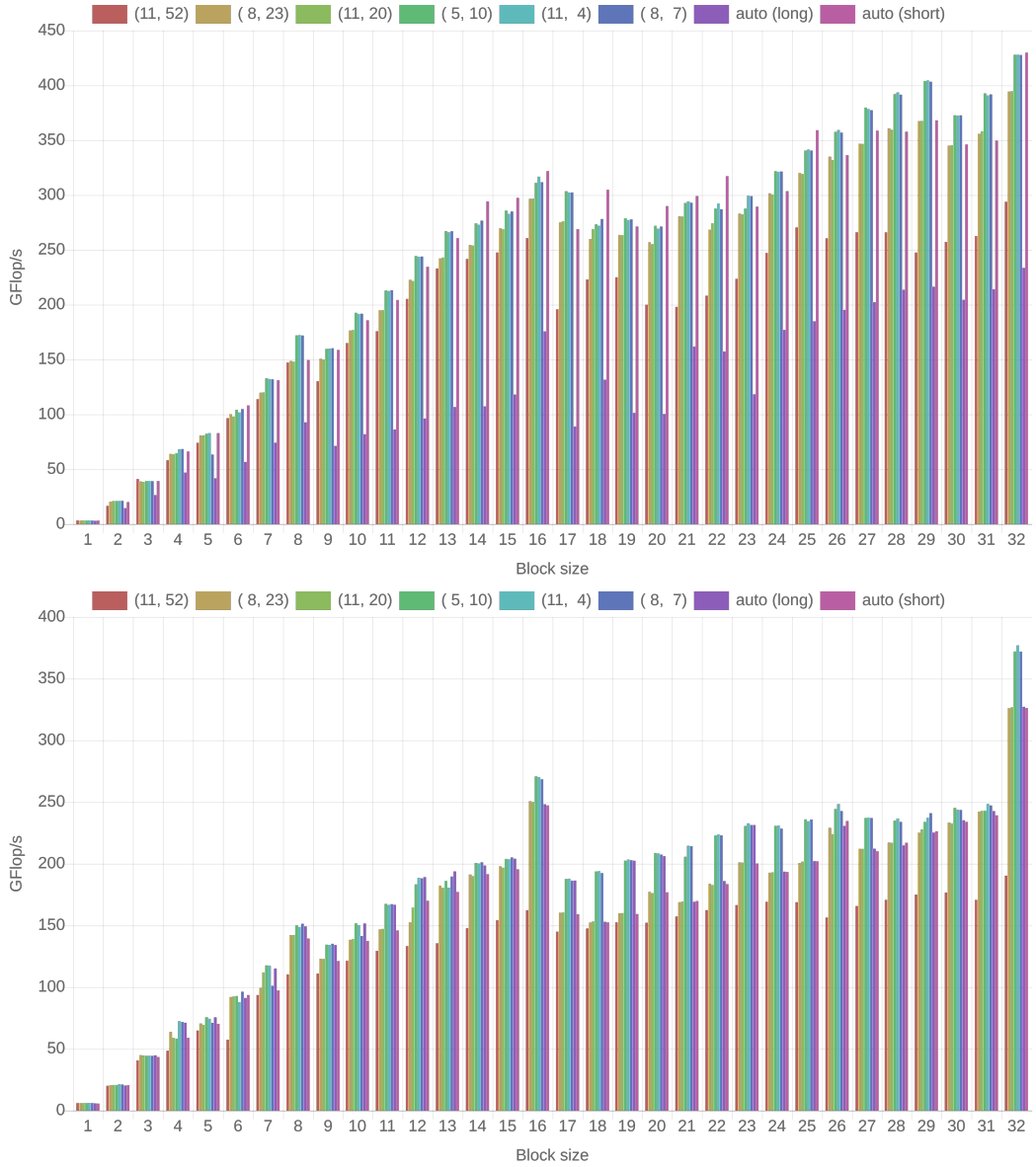
---

[4]https://sparse.tamu.edu/

Fig. 6. Performance of adaptive precision block-Jacobi on the V100 GPU. The storage format is encoded as ($ex$, $snf$). The first bar (red), is the full precision block-Jacobi. The next five bars represent lower storage precision without accounting for the detection of the suitable precision. The last two bars include the autodetection. For the seventh bar (purple), all six precisions were considered (requiring the calculation of two additional condition numbers). In the last bar (pink), only the three precisions that do not require additional condition number calculation are considered. All performance numbers only count the floating-point operations required by the full-precision variant. The top plot shows the performance of preconditioner generation (without block size detection), and the bottom plot the performance of preconditioner application.

reduced by at least 10 orders of magnitude. We used the automatic precision detection method (with all six precisions) described earlier to select the precision of each block in the adaptive precision variant. We run two parameter settings where the automatic precision detection procedure was instructed to assign precisions such that either 1 or 2 decimal digits are preserved when applying the preconditioner. This reflects the assumption that the preconditioner provides 1 and 2 digits of accuracy, respectively. For these two settings the distribution of the distinct precision formats in the preconditioner blocks is shown in Figure 7. In case of preserving 1 decimal digit of the preconditioner (top plot in Figure 7), we observe a that a significant amount of the Jacobi blocks can be be stored in less than double precision. Many blocks are stored in single or half precision, but the non-standard $fp_{8,7}$ format is also employed for a notable fraction of the Jacobi blocks. The alternative non-standard formats, $fp_{11,20}$ and $fp_{11,4}$, are irrelevant. As expected, the situation changes for the setting where we preserve 2 decimal digits of the preconditioner (bottom plot in Figure 7). Since the precision reduction is overall more conservative, no blocks are stored in the $fp_{8,7}$ format.

Figure 8 shows the iteration count and runtime of the CG solver integrated with either the full or the adaptive precision block-Jacobi preconditioner. For the adaptive precision block-Jacobi we again consider two settings where 1 digit (top plot) and 2 digits (bottom plot) of the preconditioner are preserved, respectively. A first observation is that CG enhanced with any of the variants converged for all problems (black and gray dots on top of the plot). Furthermore, the benefit of adaptive precision highly depends on the problem and the parameter setting. If most of the blocks are relatively well conditioned, the majority of the preconditioner can be stored in lower precision, yielding improvements between 10% and 30%. For problems with ill-conditioned blocks, there is no difference between the two variants, since all blocks need to be stored in full precision in order to preserve the quality of the preconditioner. In that case, it is even possible that the adaptive variant becomes slightly slower due to the additional operations needed to read and process the information about the precisions of the blocks. In particular in the setting where only 1 digit of the preconditioner is preserved, there also exist several cases were the adaptive block-Jacobi preconditioning fails to preserve the effectiveness of the preconditioner (i.e. the preconditioner is of higher quality than one digit), which results in an increase in the number of iterations which the adaptive variant needs to converge (top plot in Figure 8). This effect is mitigated if 2 digits of the preconditioner are preserved (bottom plot in Figure 8). Furthermore, we observe that there are only few cases where the preconditioner carries more accuracy than two orders of magnitude. At the same time, the benefits of the adaptive precision block-Jacobi preserving 2 digits over the standard double precision block-Jacobi are only marginally smaller than for the more aggressive setting preserving only one digit.

From this analysis, we may conclude that a setting preserving 2 digits of the preconditioner provides a good default choice, while problem-specific optimization can enable performance advantages.

## 7  DISCUSSION

In this section, we provide a concise discussion of the central numerical aspects coming with the precision adaptation in general, and the setting preserving 2 decimal digits of the preconditioner in particular:

(1) **Do we need a flexible Krylov solver if the preconditioner matrix is stored in lower than working precision?** No, storing the preconditioner in adaptive (lower) precision is independent of the need for a method accepting non-constant preconditioners. As elaborated

Fig. 7. Distribution of floating point formats among the distinct blocks when preserving 1 (top) or 2 (bottom) digits of the preconditioner blocks. Each column represents one of the selected matrices. The fraction of the column filled with a certain color depicts the fraction of blocks stored in the format represented by that color.
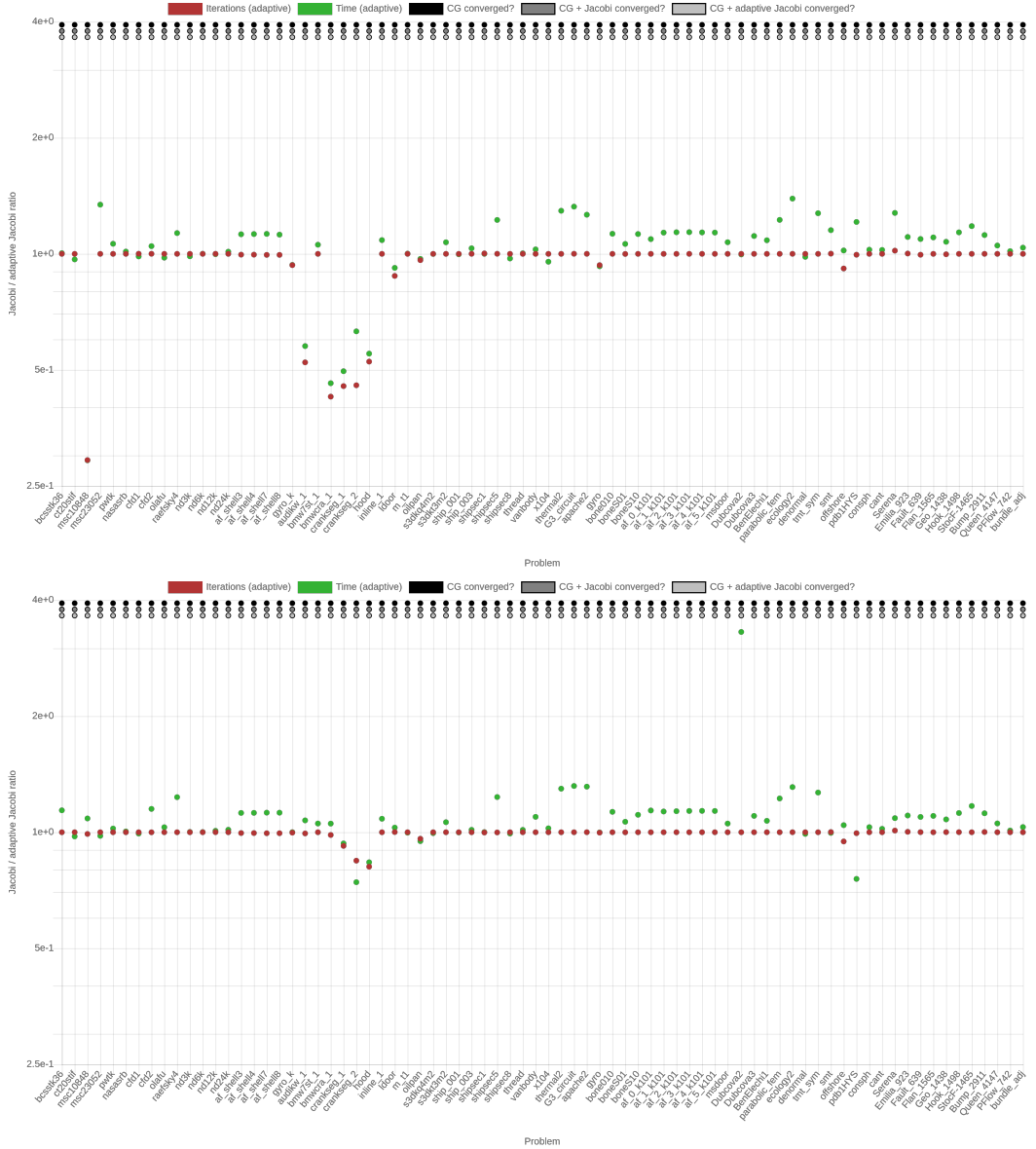
Fig. 8. Iteration count and runtime of the Conjugate Gradient (CG) solver enhanced with the adaptive precision block-Jacobi preconditioner relative to the CG solver with the full precision variant of block-Jacobi. The results include all symmetric positive definite matrices with at least $10^6$ nonzeros from the SuiteSparse matrix collection for which a CG solver needs at least 100 iterations to converge. Black and gray dots on top of the plots represent (from top to bottom) whether CG, CG+(full precision) block-Jacobi and CG+adaptive block-Jacobi converged for that matrix. The absence of a dot means that the method did not converge. The red dots represent the relative number of iterations, while the green dots the relative time of adaptive block-Jacobi compared with the full-precision variant. A value greater than 1 means that the adaptive variant outperforms the full precision block-Jacobi for that specific problem. The adaptive precision preserves 1 digit (top) or two digits (bottom) of the full precision block-Jacobi preconditioner.

in [Anzt et al. 2019], the preconditioner operator is constant as long as all arithmetic operations are handled in working precision.

(2) **Can the adaptive precision block-Jacobi matrix become singular?** No, the automatic precision adaption scheme strictly preserves the regularity of the preconditioner matrix.

(3) **Can the default setting preserving 2 decimal digits of the preconditioner introduce an iteration overhead to the outer solver?** Yes, it is possible that the block-Jacobi preconditioner has higher accuracy than 2 decimal digits. In the extreme case of the system matrix decomposing into independent problems of size smaller than the upper limit for the Jacobi blocks, the preconditioner presents the exact inverse of the system matrix and any format reduction introduces an accuracy loss. However, our analysis suggests that the block-Jacobi preconditioner rarely exceeds 2 decimal digits.

(4) **Are larger runtime savings possible by reducing the memory precision format more aggressively?** Yes, as the results in Figure 8 (top) indicate, preserving only 1 digit of the preconditioner (and therewith reducing the precision format more aggressively) can potentially augment the runtime savings for moderately-accurate block-Jacobi preconditioners. However, preserving only 1 digit in general increases the chance of loosing some preconditioner quality, and therewith increasing the iteration count.

(5) **Is it possible to control how many digits of the preconditioner are preserved?** Yes, the implementation allows to control the number of preserved preconditioner digits via a parameter.

(6) **Is the source code of the adaptive precision block-Jacobi preconditioner publicly available?** Yes, the adaptive precision block-Jacobi preconditioner is part of the Ginkgo open source software package.[5]

(7) **Can the adaptive precision block-Jacobi preconditioner be used inside other Krylov-type solvers?** Yes, the adaptive block-Jacobi preconditioner is independent of the Conjugate Gradient method used in this work and can be employed by any solver that is amenable for preconditioning.

(8) **Can the adaptive precision block-Jacobi preconditioner be used for non-symmetric positive definite problems?** Yes, the adaptive block-Jacobi preconditioner generation is based on Gauss-Jordan elimination enhanced with pivoting [Anzt et al. 2018] and can handle general non-singular problems.

## 8 CONCLUSION AND OUTLOOK

In this work we presented the first practical GPU implementation of the block-Jacobi preconditioner. In addition, we augmented the set of available precisions, allowing larger memory transfer savings than in the original description of adaptive block-Jacobi. The implementation of the full-precision block-Jacobi was improved to work on newer hardware and yield higher performance than the previously presented versions. Finally, both preconditioners, the full precision and adaptive precision block-Jacobi, are made available in production-code as part of the open source Ginkgo library.

Our experiments show that solvers enhanced with adaptive precision block-Jacobi can offer attractive runtime savings of between 10% and 30% compared with those integrating a full precision block-Jacobi preconditioner. The actual savings highly depend on the numerical properties of the problem. Even larger improvements might be available by fine-tuning the preconditioner quality to the specific problem characteristics instead of using the default block-Jacobi preconditioner accuracy of 2 decimal digits.

---

[5]https://ginkgo-project.github.io

We believe that this work offers a nearly-optimal GPU implementation of the block-Jacobi preconditioner. While further improvements might be possible, current technical limitations (e.g. long compile times of the GPU kernels) make further experimentation difficult. Thus, in the future we plan to turn our attention to related topics such as block-diagonal detection algorithms and further exploration of adaptive precision in the context of other iterative solvers and preconditioners.

## ACKNOWLEDGMENTS

## REFERENCES

Hartwig Anzt, Yen-Chen Chen, Terry Cojean, Jack Dongarra, Goran Flegar, Pratik Nayak, Enrique S. Quintana-Ortí, Yuhsiang M. Tsai, and Weichung Wang. 2019. Towards Continuous Benchmarking: An Automated Performance Evaluation Framework for High Performance Software. In *Proceedings of the Platform for Advanced Scientific Computing Conference (PASC '19)*. ACM, New York, NY, USA, Article 9, 11 pages. DOI:http://dx.doi.org/10.1145/3324989.3325719

Hartwig Anzt, Jack Dongarra, Goran Flegar, Nicholas J. Higham, and Enrique S. Quintana-Ortí. 2019. Adaptive precision in block-Jacobi preconditioning for iterative sparse linear system solvers. *Concurrency and Computation: Practice and Experience* 31, 6 (2019), e4460. DOI:http://dx.doi.org/10.1002/cpe.4460

Hartwig Anzt, Jack Dongarra, Goran Flegar, and Enrique S. Quintana-Ortí. 2017a. Batched Gauss-Jordan Elimination for Block-Jacobi Preconditioner Generation on GPUs. In *8th Int. Workshop Programming Models & Appl. for Multicores & Manycores (PMAM)*. 1–10.

Hartwig Anzt, Jack Dongarra, Goran Flegar, and Enrique S. Quintana-Ortí. 2017b. Variable-Size Batched LU for Small Matrices and Its Integration into Block-Jacobi Preconditioning. In *2017 46th International Conference on Parallel Processing (ICPP)*. 91–100.

Hartwig Anzt, Jack Dongarra, Goran Flegar, and Enrique S. Quintana-Ortí. 2018. Variable-size batched Gauss–Jordan elimination for block-Jacobi preconditioning on graphics processors. *Parallel Comput.* (jan 2018). DOI:http://dx.doi.org/10.1016/j.parco.2017.12.006

Hartwig Anzt, Jack Dongarra, Goran Flegar, Enrique S. Quintana-Ortí, and Andrés E. Tomás. 2017. Variable-Size Batched Gauss-Huard for Block-Jacobi Preconditioning. *Procedia Computer Science* 108 (2017), 1783 – 1792. DOI:http://dx.doi.org/https://doi.org/10.1016/j.procs.2017.05.186 International Conference on Computational Science, {ICCS} 2017, 12-14 June 2017, Zurich, Switzerland.

Hartwig Anzt, Jack Dongarra, Goran G. Flegar, and Thomas Grützmacher. 2018. Variable-Size Batched Condition Number Calculation on GPUs. In *2018 30th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*. 132–139. DOI:http://dx.doi.org/10.1109/CAHPC.2018.8645907

Aydin Buluç, Jeremy T. Fineman, Matteo Frigo, John R. Gilbert, and Charles E. Leiserson. 2009. Parallel Sparse Matrix-vector and Matrix-transpose-vector Multiplication Using Compressed Sparse Blocks. In *Proceedings of the Twenty-first Annual Symposium on Parallelism in Algorithms and Architectures (SPAA '09)*. ACM, New York, NY, USA, 233–244. DOI:http://dx.doi.org/10.1145/1583991.1584053

Erin Carson and Nicholas J. Higham. 2018. Accelerating the Solution of Linear Systems by Iterative Refinement in Three Precisions. *SIAM J. Scientific Computing* 40, 2 (2018), A817–A847. DOI:http://dx.doi.org/10.1137/17M1140819

IEEE Standard Commitee. 2000. IEEE Standard for Modeling and Simulation (M Amp;S) High Level Architecture (HLA) - Framework and Rules. *IEEE Std. 1516-2000* (2000), i –22. DOI:http://dx.doi.org/10.1109/IEEESTD.2000.92296

Siegfried Cools. 2018. Numerical stability analysis of the class of communication hiding pipelined Conjugate Gradient methods. *CoRR* abs/1804.02962 (2018). http://arxiv.org/abs/1804.02962

Timothy A. Davis and Yifan Hu. 2011. The University of Florida Sparse Matrix Collection. *ACM Trans. Math. Softw.* 38, 1, Article 1 (Dec. 2011), 25 pages. DOI:http://dx.doi.org/10.1145/2049662.2049663

J. Dongarra and others. 2014. *Applied mathematics research for exascale computing*. Technical Report. U.S. Dept. of Energy, Office of Science, Advanced Scientific Computing Research Program. https://science.energy.gov/~/media/ascr/pdf/

research/am/docs/EMWGreport.pdf.

M. Duranton, K. De Bosschere, A. Cohen, J. Maebe, and H. Munk. 2015. HiPEAC Vision 2015. (2015). https://www.hipeac.org/publications/vision/.

Markus Goetz and Hartwig Anzt. 2018. Machine Learning-Aided Numerical Linear Algebra: Convolutional Neural Networks for the Efficient Preconditioner Generation. In *2018 IEEE/ACM 9th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems (scalA)*. 49–56. DOI:http://dx.doi.org/10.1109/ScalA.2018.00010

G. H. Golub and C. F. Van Loan. 1996. *Matrix Computations* (3rd ed.). The Johns Hopkins University Press, Baltimore.

Gene H. Golub and Qiang Ye. 1999. Inexact Preconditioned Conjugate Gradient Method with Inner-Outer Iteration. *SIAM Journal on Scientific Computing* 21, 4 (1999), 1305–1320. DOI:http://dx.doi.org/10.1137/S1064827597323415

Magnus R. Hestenes and Eduard Stiefel. 1952. Methods of Conjugate Gradients for Solving Linear Systems. *J. Res. Nat. Bur. Standards* 49, 6 (Dec. 1952), 409–436.

Nicholas J. Higham. 2002. *Accuracy and Stability of Numerical Algorithms* (second ed.). Society for Industrial and Applied Mathematics, Philadelphia, PA, USA.

Mark Hoemmen. 2010. *Communication-avoiding Krylov Subspace Methods*. Ph.D. Dissertation. Berkeley, CA, USA. Advisor(s) Demmel, James W. AAI3413388.

C. Lanczos. 1952. Solution of systems of linear equations by minimized iterations. *J. Res. Nat. Bur. Standards* 49, 1 (Dec. 1952), 33–53.

R. Lucas and others. 2014. Top ten Exascale research challenges. (2014). http://science.energy.gov/~/media/ascr/ascac/pdf/meetings/20140210/Top10reportFEB14.pdf.

NVIDIA Corporation 2018. *NVIDIA CUDA Toolkit* (9.0 ed.). NVIDIA Corporation.

Y. Saad. 2003. *Iterative Methods for Sparse Linear Systems* (2nd ed.). SIAM.