# Hiding global synchronization latency in the preconditioned Conjugate Gradient algorithm

P. Ghysels [a,b,*], W. Vanroose [a]

[a] University of Antwerp, Department of Mathematics and Computer Science, Middelheimlaan 1, B-2020 Antwerp, Belgium
[b] Intel ExaScience Lab, Kapeldreef 75, B-3001 Leuven, Belgium

## ARTICLE INFO

## ABSTRACT

Scalability of Krylov subspace methods suffers from costly global synchronization steps that arise in dot-products and norm calculations on parallel machines. In this work, a modified preconditioned Conjugate Gradient (CG) method is presented that removes the costly global synchronization steps from the standard CG algorithm by only performing a single non-blocking reduction per iteration. This global communication phase can be overlapped by the matrix–vector product, which typically only requires local communication. The resulting algorithm will be referred to as *pipelined CG*. An alternative pipelined method, mathematically equivalent to the Conjugate Residual (CR) method that makes different trade-offs with regard to scalability and serial runtime is also considered. These methods are compared to a recently proposed asynchronous CG algorithm by Gropp. Extensive numerical experiments demonstrate the numerical stability of the methods. Moreover, it is shown that hiding the global synchronization step improves scalability on distributed memory machines using the message passing paradigm and leads to significant speedups compared to standard preconditioned CG.

## 1. Introduction

Many high-performance computing applications rely on Krylov subspace methods for their linear algebra. These Krylov methods exploit the sparsity of the matrices that typically appear in scientific applications simulating a problem modeled by a partial differential equation (PDE). The linear systems to be solved are derived by a discretization technique such as finite-differences, finite-volumes or finite-elements in which neighboring variables are related through a stencil. This results in matrices with only a few non-zero elements per row.

The building blocks for Krylov subspace methods are the sparse matrix–vector product (spmv), vector–vector additions and dot-products. Each building block has a different communication pattern on distributed memory machines, resulting in different scaling properties.

The spmv often requires only local communication. The matrix derived from a PDE can, possibly after a permutation of rows and columns, be mapped to the underlying machine architecture in such a way that applying a matrix–vector product only requires communication between neighboring nodes, i.e. nodes that are separated by a small number of hops. In the state-of-the-art literature, many examples can be found of stencil based codes that scale nearly optimal to very large parallel

* Corresponding author at: University of Antwerp, Department of Mathematics and Computer Science, Middelheimlaan 1, B-2020 Antwerp, Belgium.
  E-mail addresses: pieter.ghysels@ua.ac.be (P. Ghysels), wim.vanroose@ua.ac.be (W. Vanroose).

machines [35]. In this work, the focus is on sparse and well-structured matrices or matrix-free linear operators such as stencil applications.

Vector operations such as an AXPY ($y \leftarrow \alpha x + y$) can be calculated locally and do not require communication between nodes. However, on a multicore shared memory system even a simple AXPY operation typically does not scale well over the cores due to memory bandwidth congestion. But this on-chip communication bottleneck is not the focus of the current article.

In contrast, we focus on the global communication. A dot-product of two vectors $v$ and $w$ involves such global communication and requires participation of all processes. Since the dot-product is a single scalar value $\alpha = \sum_{i=1}^{n} v_i^T w_i$, this operation needs the result of each local dot-product and involves a synchronization of all the processes as well. Even if the global communication is extremely fast, for instance when a special reduction network is available, the explicit synchronization of the involved processes makes a dot-product a costly operation. For large systems the cost of the global reduction grows as $\mathcal{O}(\log(P_n))$, the height of the reduction tree, where $P_n$ is the number of nodes. This cost model ignores variability due to OS jitter, core speed variability or load imbalance that start to play an important role for larger systems [25].

Research in reducing the number of global reductions in Krylov methods goes back to the first implementations of Krylov methods on parallel computers [4]. Variations of the Conjugate Gradient (CG) method with only a single global synchronization point have been presented by Saad [32], Meurant [31], D'Azevedo and Romine [14] and Eijkhout [13] as well as by Chronopoulos and Gear [11]. Also, a CG variation exists that is based on two three-term recurrences instead of three two-term recurrences and that only requires one reduction [34]. Yang et al. [47,45,44,46] propose the so called *improved* versions of the QMR, BiCGStab, BiCG and CGS algorithms respectively, which reduce the number of synchronizations to just one per iteration for these methods. The total number of global reductions is reduced even further by $s$-step Krylov methods [11,27,2,29,10,12,28,39,8], where $s$ matrix–vector products are combined and the resulting Krylov basis is orthogonalized simultaneously with a single global reduction, leading to a communication reduction of a factor $s$. Apart from reducing global communication, these methods are also aimed at reducing communication with the slow memory. However, with increasing $s$, the stability of the $s$-step Krylov basis deteriorates. In [30], hierarchical or nested Krylov methods are used that reduce the number of global reductions on the whole machine, but freely allow global communication on a smaller subset of nodes, where communication is cheaper.

Instead of, or in addition to trying to reduce global synchronizations, several authors have come up with ways to overlap expensive communication phases with computations. In [16] it is suggested to overlap the dot-products in CG with the update of the solution vector or with application of a factored preconditioner. In [15], block Gram-Schmidt is used in the Generalized Minimal Residual (GMRES) method where communication for one block can be overlapped with computation on a different block. In the Arnoldi algorithm for eigenvalue computations, Hernandez et al. [22] overlap global communication for reorthogonalization of the Krylov basis with the computations for orthogonalization of the next basis vector. Recently, we have proposed a pipelined GMRES solver [18] where the global reductions are overlapped with the calculation of multiple matrix–vector products. The results of the dot-products are only used with a delay of a few iterations in the algorithm. In [1] a model is developed for the performance of the pipelined GMRES algorithm on large clusters. An asynchronous version of CG has been proposed recently by Gropp [21] where one reduction can be overlapped with the matrix–vector product and the other with the preconditioner. Overlapping global communication with local work has recently become easier and more attractive due to the inclusion of non-blocking collectives in the MPI-3 standard [5,26].

The *aim* of this work is to hide the latency of global communication due to dot-products in the preconditioned CG method. The trade-off between improved scalability and extra floating point operations should result in speedups on medium to large parallel machines without significantly sacrificing numerical robustness.

Our main *contributions* are the development of a preconditioned pipelined CG method that only has a single non-blocking reduction per iteration. This non-blocking reduction can be overlapped with the matrix–vector product and with application of the preconditioner. Also, a preconditioned pipelined Conjugate Residual (CR) method is presented with one non-blocking reduction that can be overlapped with the matrix–vector product. We show that both methods have much improved scalability and runtime compared to standard CG.

The paper is outlined as follows. Section 2 reviews the standard preconditioned CG method. A modified CG method due to Chronopoulos and Gear that will be used to derive pipelined CG is discussed in Section 2.2. Section 3 presents algorithms that can overlap the global reduction with the matrix–vector product, with in Section 3.1 first the unpreconditioned pipelined CG method. Next, in Section 3.2 the preconditioned pipelined CG method is derived. In Section 3.3 a preconditioned pipelined CR method is derived from pipelined CG. Section 3.4 gives the asynchronous CG algorithm due to Gropp. We report on extensive numerical tests that show the stability of the pipelined CG and CR methods in Section 4. Also in Section 4, we hint at the use of a residual replacement strategy to improve the maximum attainable accuracy of the pipelined methods. Finally, the results of benchmarks on a parallel distributed memory computer are presented in Section 5.

## 2. The Conjugate Gradient algorithm

First, the standard preconditioned CG algorithm is reviewed together with the definition of the Krylov subspace. Then a variation of standard CG due to Chronopoulos and Gear is presented that only requires a single global reduction per iteration.

### 2.1. Preconditioned Conjugate Gradients

The mother of all Krylov methods is CG, dating back to a paper from 1952 by Hestenes and Stiefel [23]. Algorithm 1 shows the preconditioned CG iteration [34], which iteratively solves $M^{-1}Ax = M^{-1}b$ where both $A$ and $M$ are symmetric and positive definite square matrices of size $N \times N$.

---

**Algorithm 1.** Preconditioned CG

1: $r_0 := b - Ax_0$; $u_0 := M^{-1}r_0$; $p_0 := u_0$
2: **for** $i = 0, \ldots$ **do**
3:    $s := Ap_i$
4:    $\alpha := (r_i, u_i)/(s, p_i)$
5:    $x_{i+1} := x_i + \alpha p_i$
6:    $r_{i+1} := r_i - \alpha s$
7:    $u_{i+1} := M^{-1}r_{i+1}$
8:    $\beta := (r_{i+1}, u_{i+1})/(r_i, u_i)$
9:    $p_{i+1} := u_{i+1} + \beta p_i$
10:   **end for**

---

The residual vector of the original system is $r_i = b - Ax_i$, while $u_i = M^{-1}r_i$ is the residual of the preconditioned system and $p_i$ is called the search direction. If the exact solution is $\hat{x} = A^{-1}b$, then the error is defined as $e_i = \hat{x} - x_i$. All subsequent approximations $x_i$ lie in a so-called Krylov subspace $\mathcal{K}_i(M^{-1}A, M^{-1}r_0)$, which is defined as

$$\mathcal{K}_i(M^{-1}A, M^{-1}r_0) = \text{span}\{u_0, M^{-1}Au_0, \ldots, (M^{-1}A)^{i-1}u_0\}. \tag{1}$$

It is well known that the preconditioned CG iteration generates a sequence of iterates $x_i \in x_0 + \mathcal{K}_i(M^{-1}A, u_0)$, with the property that at step $i$, $\|e_i\|_A = \sqrt{e_i^T A e_i}$ is minimized.

In terms of communication, the important steps in Algorithm 1 are: application of the sparse matrix–vector product (SPMV) $Ap_i$ in line 3, application of the preconditioner $M^{-1}r_{i+1}$ in line 7, and the two dot-products $(s, p_i)$ and $(r_{i+1}, u_{i+1})$ in lines 4 and 8. All other steps are vector updates that do not require communication between nodes. We assume that the matrix–vector product and application of the preconditioner only require communication among neighboring nodes, which can be implemented in a scalable way. The two dot-products, causing two global synchronization points per iteration become the bottleneck with increasing parallelism.

Note that even when $A$ and $M$ are both symmetric and positive definite, the left or right preconditioned systems $M^{-1}Ax = M^{-1}b$ and $AM^{-1}u = b$ with $x = M^{-1}u$ respectively are no longer symmetric in general. When the preconditioner is available in the form $M = LL^T$, then one way to preserve symmetry is to use split preconditioning $L^{-1}AL^{-T}u = L^{-1}b$. In [16], a split preconditioned CG method is presented where the global reductions can be overlapped with local work. The communication for one reduction can be overlapped with application of $L^{-T}$ and, since the iteration does not depend on $x_i$ (it is only computed to output the final $x_i$ at the end of the iteration), the update for $x_i$ can be postponed by one iteration where it can be overlapped with the second reduction.

Another way to preserve symmetry is based on the observation that $M^{-1}A$ is self-adjoint with respect to the $M$ inner-product $(x, y)_M = (x, My) = (Mx, y)$:

$$\left(M^{-1}Ax, y\right)_M = (Ax, y) = (x, Ay) = (x, M(M^{-1}A)y) = (x, M^{-1}Ay)_M. \tag{2}$$

Replacing the usual Euclidean inner-product in CG with this $M$ inner-product leads to Algorithm 1. Similarly, for right preconditioning, the $M^{-1}$ inner-product also leads to Algorithm 1. Furthermore, the iterates generated by split preconditioned CG are identical to those of Algorithm 1 [34].

### 2.2. Chronopoulos/Gear CG

Several authors have suggested alternatives to CG to reduce the number of global synchronization points to just one. One such method was presented by Saad [32], and was later improved for stability by Meurant [31]. This latter method performs an additional dot-product compared to the standard CG implementation. But, these three dot-products, compared to two for standard CG, can be combined in a single reduction. Similarly, Chronopoulos and Gear [11] presented a CG variation with a single global synchronization point that requires one additional AXPY compared to CG. A slight variation of this method was published by D'Azevedo and Romine [14] and D'Azevedo and Eijkhout [13]. Also, the three term recurrence version of CG [34] can be implemented with a single global reduction. All these CG variations have slightly different properties in terms of memory requirements, total number of flops and stability. The remainder of this work builds on the method by Chronopoulos and Gear [11] as shown in Algorithm 2.

---

**Algorithm 2.** Preconditioned Chronopoulos/Gear CG

1: $r_0 := b - Ax_0$; $u_0 := M^{-1}r_0$; $w_0 := Au_0$
2: $\alpha_0 := (r_0, u_0)/(w_0, u_0)$; $\beta_0 := 0$; $\gamma_0 := (r_0, u_0)$
3: **for** $i = 0, \ldots$ **do**
4:   $p_i := u_i + \beta_i p_{i-1}$
5:   $s_i := w_i + \beta_i s_{i-1}$
6:   $x_{i+1} := x_i + \alpha_i p_i$
7:   $r_{i+1} := r_i - \alpha_i s_i$
8:   $u_{i+1} := M^{-1}r_{i+1}$
9:   $w_{i+1} := Au_{i+1}$
10:   $\gamma_{i+1} := (r_{i+1}, u_{i+1})$
11:   $\delta := (w_{i+1}, u_{i+1})$
12:   $\beta_{i+1} := \gamma_{i+1}/\gamma_i$
13:   $\alpha_{i+1} := \gamma_{i+1}/(\delta - \beta_{i+1}\gamma_{i+1}/\alpha_i)$
14: **end for**

---

Compared to standard CG, Algorithm 2 performs an additional vector update for the vector $s_i = Ap_i$, found by multiplying the recurrence for $p_i$ by $A$

$$Ap_i = Au_i + \beta_i Ap_{i-1}, \quad s_i = Au_i + \beta_i s_{i-1} \tag{3}$$

or since in the Chronopoulos/Gear method $w_i = Au_i$, this becomes $s_i = w_i + \beta_i s_{i-1}$. Algorithm 2 was derived from its unpreconditioned version by replacing the Euclidean inner-product by the $M$ inner-product and is mathematically equivalent to standard preconditioned CG, Algorithm 1.

The communication phase for both dot-products from Algorithm 2, lines 10 and 11, can be combined in a single global reduction. The update for $x_i$ can be postponed and used to overlap the global reduction. However, even for small parallel machines the runtime of a single vector update will not be enough to fully cover the latency of the global communication.

An implementation of Algorithm 2 could be optimized by loading each vector from main memory only once per iteration, allowing for more efficient use of available memory bandwidth. This can be achieved by fusing all AXPYS together, which is not possible in standard CG. This was already suggested by Chronopoulos and Gear in their original paper [11], where they also generalize Algorithm 2 to an $s$-step CG method that only needs one memory sweep per $s$ steps. Algorithm 2 corresponds to the $s = 1$ case.

## 3. Hiding global communication

In this section, modified and reordered versions of Algorithm 2 are presented where the global synchronization can be overlapped by the sparse matrix–vector product. Section 3.1 starts with the unpreconditioned version of this pipelined CG algorithm for simplicity. Section 3.2 adds preconditioning to it and Section 3.3 uses a different inner-product to preserve symmetry for the preconditioned case, which leads to a pipelined CR (pipe-CR) method. Section 3.4 discusses an asynchronous CG method recently proposed by Gropp.

### 3.1. Pipelined Conjugate Gradients

Recall from Section 2 that $r_i = b - Ax_i$ and $s_i = Ap_i$ with $p_i$ the search direction. Lets first consider the unpreconditioned iteration, i.e., $u_i = M^{-1}r_i = r_i$, then $w_i = Au_i = Ar_i$.

Instead of computing $w_{i+1} = Ar_{i+1}$ directly using the SPMV, it also satisfies the recurrence relation

$$Ar_{i+1} = Ar_i - \alpha_i As_i, \quad w_{i+1} = w_i - \alpha_i As_i. \tag{4}$$

This however requires $z_i = As_i \equiv A^2 p_i$, which can be computed via the SPMV as $As_i$, or for which also a recurrence relation can be found

$$As_i = Aw_i + \beta_i As_{i-1}, \quad z_i = Aw_i + \beta_i z_{i-1}. \tag{5}$$

This depends on $q_i = Aw_i \equiv A^2 r_i$, which can be computed from the SPMV. Starting from Algorithm 2, adding the recurrences for $w_{i+1}$ (4) and $z_i$ (5), replacing the SPMV by $q_i = Aw_i$ and reordering the steps leads to Algorithm 3.

---

**Algorithm 3.** Pipelined Chronopoulos/Gear CG

---

1: $r_0 := b - Ax_0$; $w_0 := Ar_0$
2: **for** $i = 0, \ldots$ **do**
3:    $\gamma_i := (r_i, r_i)$
4:    $\delta := (w_i, r_i)$
5:    $q_i := Aw_i$
6:    **if** $i > 0$ **then**
7:      $\beta_i := \gamma_i / \gamma_{i-1}$; $\alpha_i := \gamma_i / (\delta - \beta_i \gamma_i / \alpha_{i-1})$
8:    **else**
9:      $\beta_i := 0$; $\alpha_i := \gamma_i / \delta$
10:   **end if**
11:   $z_i := q_i + \beta_i z_{i-1}$
12:   $s_i := w_i + \beta_i s_{i-1}$
13:   $p_i := r_i + \beta_i p_{i-1}$
14:   $x_{i+1} := x_i + \alpha_i p_i$
15:   $r_{i+1} := r_i - \alpha_i s_i$
16:   $w_{i+1} := w_i - \alpha_i z_i$
17: **end for**

---

Mathematically, Algorithm 3 is still equivalent with standard CG. However, as in the method by Chronopoulos/Gear, the two dot-products in lines 3 and 4 can be combined in a single reduction. Furthermore, since the result of this reduction is only needed in line 6, this global reduction can be overlapped with the spmv. We shall refer to Algorithm 3 as unpreconditioned pipelined CG.

In finite precision arithmetic, pipelined CG will behave differently than standard CG since rounding errors are propagated differently. Numerical stability of this algorithm, and its preconditioned variants presented in the next paragraphs, will be studied using a wide range of matrices in Section 4.

### 3.2. Preconditioned pipelined CG

We shall consider the left preconditioned system $M^{-1}Ax = M^{-1}b$ with both $M$ and $A$ symmetric and positive definite. As for standard CG, pipelined CG, Algorithm 3 cannot be applied to the preconditioned system without modification since $M^{-1}A$ is in general not symmetric positive definite. For preconditioned pipelined CG the same strategy as for standard preconditioned CG will be followed: apply Algorithm 3 to the preconditioned system and replace the classic Euclidean inner-product by the $M$ inner-product $(\cdot, \cdot)_M$, since $M^{-1}A$ is self-adjoint with respect to the $M$ inner-product. The dot-products from pipelined CG (lines 3 and 4) now use the residual of the preconditioned system, $u_i = M^{-1}r_i$, and the $M$ inner-product instead of the Euclidean inner-product

$$\gamma_i = (u_i, u_i)_M = (Mu_i, u_i) = (r_i, u_i) \tag{6}$$

$$\delta_i = \left(M^{-1}Au_i, u_i\right)_M = (Au_i, u_i). \tag{7}$$

Now, let $w_i = Au_i$, $s_i = Ap_i$ and also introduce $q_i = M^{-1}s_i$. A recurrence relation for the preconditioned residual $u_i$ is found by multiplying the recurrence for the original residual on the left by $M^{-1}$

$$M^{-1}r_{i+1} = M^{-1}r_i - \alpha_i M^{-1}s_i, \quad u_{i+1} = u_i - \alpha_i q_i, \tag{8}$$

with a similar recurrence for $q_i = M^{-1}s_i$

$$M^{-1}s_i = M^{-1}w_i + \beta_i M^{-1}s_i, \quad q_i = M^{-1}w_i + \beta_i q_{i-1}. \tag{9}$$

If defined that $m_i = M^{-1}w_i \equiv M^{-1}Au_i \equiv M^{-1}AM^{-1}r_i$, then since $w_{i+1} = Au_{i+1}$,

$$Au_{i+1} = Au_i - \alpha_i Aq_i, \quad w_{i+1} = w_i - \alpha_i Aq_i. \tag{10}$$

Contrary to the unpreconditioned case, now define $z_i = Aq_i$ and the final recurrence relation for $z_i$ becomes

$$Aq_i = AM^{-1}w_i + \beta_i Aq_i, \quad z_i = Am_i + \beta_i z_{i-1} \tag{11}$$

and let $n_i = Am_i \equiv AM^{-1}w_i$.

---

**Algorithm 4.** Preconditioned pipelined CG

1: $r_0 := b - Ax_0$; $u_0 := M^{-1}r_0$; $w_0 := Au_0$
2: **for** $i = 0, \ldots$ **do**
3:     $\gamma_i := (r_i, u_i)$
4:     $\delta := (w_i, u_i)$
5:     $m_i := M^{-1}w_i$
6:     $n_i := Am_i$
7:     **if** $i > 0$ **then**
8:         $\beta_i := \gamma_i/\gamma_{i-1}$; $\alpha_i := \gamma_i/(\delta - \beta_i\gamma_i/\alpha_{i-1})$
9:     **else**
10:        $\beta_i := 0$; $\alpha_i := \gamma_i/\delta$
11:    **end if**
12:    $z_i := n_i + \beta_i z_{i-1}$
13:    $q_i := m_i + \beta_i q_{i-1}$
14:    $s_i := w_i + \beta_i s_{i-1}$
15:    $p_i := u_i + \beta_i p_{i-1}$
16:    $x_{i+1} := x_i + \alpha_i p_i$
17:    $r_{i+1} := r_i - \alpha_i s_i$
18:    $u_{i+1} := u_i - \alpha_i q_i$
19:    $w_{i+1} := w_i - \alpha_i z_i$
20: **end for**

Combining the updates for $r_{i+1}$, $x_{i+1}$, $s_i$ and $p_i$, and Eqs. (6)–(11) with the matrix–vector product $n_i = Am_i$ and application of the preconditioner $m_i = M^{-1}w_i$, directly leads to preconditioned pipelined CG, Algorithm 4. In this algorithm, the reduction for the dot-products (lines 3 and 4) can be overlapped with application of both the preconditioner (line 5) and the matrix–vector product (line 6). However, the number of AXPYS has increased from just three in standard CG or four in Chronopoulos/Gear CG to eight in Algorithm 4. Again, these operations can be fused together (also including the dot-products from the start of the next iteration) such that only a single memory sweep is necessary.

The preconditioned pipelined CG method, Algorithm 4, is based on the $M$ inner-product just like standard preconditioned CG, Algorithm 1, and hence we know that the two methods are mathematically equivalent. They both minimize the $A$ norm of the error, $\|e_k\|_A$ over the Krylov space $\mathcal{K}_k(M^{-1}A, M^{-1}r_0)$.

### 3.3. Preconditioned pipelined Conjugate Residuals

For Algorithm 4, the Euclidean inner-product in Algorithm 3 was replaced by the $M$ inner-product. However, observe that $M^{-1}A$ is also self-adjoint with respect to the $A$ inner-product

$$\left(M^{-1}Ax, y\right)_A = (AM^{-1}Ax, y) = (M^{-1}Ax, Ay) = (M^{-1}Ax, y)_A. \tag{12}$$

Using this inner-product instead leads to

$$\gamma_i = (u_i, u_i)_A = (Au_i, u_i) = (w_i, u_i) \tag{13}$$

$$\delta_i = \left(M^{-1}Au_i, u_i\right)_A = \left(M^{-1}Au_i, Au_i\right) = (m_i, w_i). \tag{14}$$

Using Eqs. (13) and (14) instead of (6) and (7) leads to Algorithm 5. However, since Algorithm 5 is based on a different inner-product, it is no longer equivalent with standard CG but as will become clear it can be seen as a pipelined version of the CR method [34].

Since (14) depends on $m_i = M^{-1}w_i$, the preconditioner now has to be applied before (line 3) the dot-products can be started. This means that the global reduction can only be overlapped with the matrix–vector product (line 6). However, the iteration does not depend on the original unpreconditioned residual $r_i$ and on $s_i = Ap_i$ anymore. These two recurrences can safely be dropped, saving some floating point operations and memory. When the recurrence for $r_i$ is dropped, a stopping criterion can still be based on the preconditioned residual $u_i$.

---

**Algorithm 5.** Preconditioned pipelined CR

---

1: $r_0 := b - Ax_0$; $u_0 := M^{-1}r_0$; $w_0 := Au_0$
2: **for** $i = 0, \ldots$ **do**
3:     $m_i := M^{-1}w_i$
4:     $\gamma_i := (w_i, u_i)$
5:     $\delta := (m_i, w_i)$
6:     $n_i := Am_i$
7:     **if** $i > 0$ **then**
8:       $\beta_i := \gamma_i / \gamma_{i-1}$; $\alpha_i := \gamma_i / (\delta - \beta_i \gamma_i / \alpha_{i-1})$
9:     **else**
10:      $\beta_i := 0$; $\alpha_i := \gamma_i / \delta$
11:    **end if**
12:    $z_i := n_i + \beta_i z_{i-1}$
13:    $q_i := m_i + \beta_i q_{i-1}$
14:    $p_i := u_i + \beta_i p_{i-1}$
15:    $x_{i+1} := x_i + \alpha_i p_i$
16:    $u_{i+1} := u_i - \alpha_i q_i$
17:    $w_{i+1} := w_i - \alpha_i z_i$
18: **end for**

---

Another optimization is possible: note that, apart from lines 14 and 15, the iterates in Algorithm 5 do not depend on $x$ or $p$. Hence the recurrences for $x$ and $p$ can be postponed one iteration and used in the overlap with the reduction. However, for the $p$ recurrence, this requires storage for one extra vector. For long reductions, this would shorten the critical path of the method by two vector updates.

**Theorem 3.1.** *Algorithm 5, when applied to a linear system $M^{-1}Ax = M^{-1}b$, with exact solution $\hat{x}$ and with both $M$ and $A$ symmetric positive definite, minimizes $\|e_k\|_{AM^{-1}A}$, where $e_k = \hat{x} - x_k$ is the error and $x_k \in x_0 + \mathcal{K}_k(M^{-1}A, u_0)$. Without preconditioner, i.e., $M^{-1} = I$, the iteration minimizes the 2-norm of the residual $\|r_k\|_2 = \|e_k\|_{A^2}$.*

**Proof.** We first show the following orthogonality properties

$$(u_k, u_j)_A = u_k^T A u_j = u_k^T w_j = 0, \quad j < k, \tag{15}$$

$$(p_k, M^{-1}Ap_j)_A = p_k^T AM^{-1}Ap_j = p_k^T A q_j = 0, \quad j < k, \tag{16}$$

by induction on $k$. First check that indeed for $k = 1$

$$u_1^T A u_0 = (u_0 - \alpha_0 q_0)^T w_0 = u_0^T w_0 - \alpha_0 (m_0 + \beta_0)^T w_0 \tag{17}$$

is zero if $\beta_0 = 0$ and

$$\alpha_0 = \frac{(u_0, w_0)}{(m_0, w_0)} = \frac{\gamma_0}{\delta_0}, \tag{18}$$

which corresponds to the definitions of $\beta_0$ and $\alpha_0$ in Algorithm 5. Likewise, for (16),

$$p_1^T A q_0 = p_1^T A m_0 = p_1^T n_0 = p_1^T z_0 = (u_1 + \beta_1 p_0)^T \left( \frac{w_0 - w_1}{\alpha_0} \right) \tag{19}$$

$$= \frac{1}{\alpha_0} \left( u_1^T w_0 - u_1^T w_1 + \beta_1 p_0^T w_0 - \beta_1 p_0^T w_1 \right) \tag{20}$$

$$= \frac{1}{\alpha_0} \left( -\gamma_1 + \beta_1 u_0^T w_0 - \beta_1 u_0^T w_1 \right) = \frac{1}{\alpha_0} (-\gamma_1 + \beta_1 \gamma_0) \tag{21}$$

is zero when $\beta_1 = \gamma_1 / \gamma_0$. For $k > 1$, multiply the recurrence relation for $u_k^T$ on the right with $Au_j$ to get

$$u_k^T A u_j = u_{k-1}^T A u_j - \alpha_{k-1} q_{k-1}^T A u_j. \tag{22}$$

For $j < k - 1$, both terms in the right-hand side are zero by induction and because the $u$ and $p$ vectors span the same Krylov space. For $j = k - 1$ the right-hand side of (22) is zero if

$$\alpha_{k-1} = \frac{u_{k-1}^T A u_{k-1}}{q_{k-1}^T A u_{k-1}} = \frac{(u_{k-1}, w_{k-1})}{(q_{k-1}, w_{k-1})} = \frac{\gamma_{k-1}}{(m_{k-1} + \beta_{k-1} q_{k-2}, w_{k-1})} = \frac{\gamma_{k-1}}{(m_{k-1}, w_{k-1}) + \beta_{k-1}(q_{k-2}, w_{k-1})}$$

$$= \frac{\gamma_{k-1}}{\delta_{k-1} + \beta_{k-1}\left(\frac{u_{k-2} - u_{k-1}}{\alpha_{k-2}}, w_{k-1}\right)} = \frac{\gamma_{k-1}}{\delta_{k-1} + \frac{\beta_{k-1}}{\alpha_{k-2}}[(u_{k-2}, w_{k-1}) - (u_{k-1}, w_{k-1})]} = \frac{\gamma_{k-1}}{\delta_{k-1} - \frac{\beta_{k-1}}{\alpha_{k-2}}(u_{k-1}, w_{k-1})} = \frac{\gamma_{k-1}}{\delta_{k-1} - \frac{\beta_{k-1}\gamma_{k-1}}{\alpha_{k-2}}}, \tag{23}$$

this corresponds with the choice for $\alpha$ in Algorithm 5, line 8. Likewise for (16), the right-hand side of

$$p_k^T Aq_j = u_k^T Aq_j + \beta_k p_{k-1}^T Aq_j \tag{24}$$

is zero for $j < k - 1$ by induction and is zero for $j = k - 1$ if

$$\beta_k = -\frac{u_k^T Aq_{k-1}}{p_{k-1}^T Aq_{k-1}} = -\frac{(u_k, q_{k-1})_A}{(p_{k-1}, q_{k-1})_A} = -\frac{\left(u_k, \frac{u_k - u_{k-1}}{-\alpha_{k-1}}\right)_A}{\left(u_{k-1} + \beta_{k-1} p_{k-2}, \frac{u_k - u_{k-1}}{-\alpha_{k-1}}\right)_A}$$

$$= -\frac{(u_k, u_k)_A - (u_k, u_{k-1})_A}{(u_{k-1}, u_k)_A - (u_{k-1}, u_{k-1}) + \beta_{k-1}\left[(p_{k-2}, u_k)_A - (p_{k-2}, u_{k-1})_A\right]} = \frac{(u_k, u_k)_A}{(u_{k-1}, u_{k-1})} = \frac{\gamma_k}{\gamma_{k-1}}, \tag{25}$$

which corresponds to the definition of $\beta$ in Algorithm 5, line 8. This concludes the proof of the orthogonality conditions (15) and (16).

To show the minimization property, first note that $AM^{-1}A$ is symmetric and positive definite and hence $\|\cdot\|_{AM^{-1}A}$ defines a valid norm. To show that $x_k \in x_0 + \mathcal{K}_k(M^{-1}A, u_0)$ minimizes $\|e\|_{AM^{-1}A}$, we consider an arbitrary point $x = x_k - \Delta x \in x_0 + \mathcal{K}_k(M^{-1}A, M^{-1}b)$ with error $e = \hat{x} - x = e_k + \Delta x$. Then

$$\|e\|_{AM^{-1}A}^2 = \|e_k + \Delta x\|_{AM^{-1}A}^2 = (e_k + \Delta x)^T AM^{-1}A(e_k + \Delta x) = e_k^T AM^{-1}Ae_k + \Delta x^T AM^{-1}A\Delta x + 2e_k^T AM^{-1}A\Delta x$$

$$= \|e_k\|_{AM^{-1}A}^2 + \|\Delta x\|_{AM^{-1}A}^2 + 2u_k^T A\Delta x. \tag{26}$$

The last term is zero by orthogonality property (15), and the second term is always positive and only zero for $\Delta x = 0$. Hence, $x_k = x$ is the unique point that minimizes $\|e\|_{AM^{-1}A}$.

For $M = I$, $\|e_k\|_{AM^{-1}A} = \|e_k\|_{A^2} = \|Ae_k\|_2 = \|r_k\|_2$ is minimized.  □

From Theorem 3.1 and the $A$-orthogonality of the residual vectors, relation (15), it is clear that Algorithm 5 is indeed a minimal residual method and is equivalent to CR.

For completeness we also show the original preconditioned CR method, see Algorithm 6.

---

**Algorithm 6.** Preconditioned CR

---

1: $u_0 := M^{-1}(b - Ax_0); \ p_0 := u_0; \ w_0 := Au_0; \ s_0 := w_0$
2: **for** $i = 0, \ldots$ **do**
3:    $q := M^{-1}s_i$
4:    $\alpha := (u_i, w_i)/(s_i, q)$
5:    $x_{i+1} := x_i + \alpha p_i$
6:    $u_{i+1} := u_i - \alpha q$
7:    $w_{i+1} := Au_{i+1}$
8:    $\beta := (u_{i+1}, w_{i+1})/(u_i, w_i)$
9:    $p_{i+1} := u_{i+1} + \beta p_i$
10:   $s_{i+1} := w_{i+1} + \beta s_i$
11: **end for**

---

### 3.4. Method due to Gropp

---

**Algorithm 7.** Gropp's asynchronous CG

---

1: $r_0 := b - Ax_0; \ u_0 := M^{-1}r_0; \ p_0 := u_0; \ s_0 := Ap_0; \ \gamma_0 := (r_0, u_0)$
2: **for** $i = 0, \ldots$
3:    $\delta := (p_i, s_i)$
4:    $q_i := M^{-1}s_i$
5:    $\alpha_i := \gamma_i/\delta$
6:    $x_{i+1} := x_i + \alpha_i p_i$
7:    $r_{i+1} := r_i - \alpha_i s_i$
8:    $u_{i+1} := u_i - \alpha_i q_i$
9:    $\gamma_{i+1} := (r_{i+1}, u_{i+1})$
10:   $w_{i+1} := Au_{i+1}$
11:   $\beta_{i+1} := \gamma_{i+1}/\gamma_i$
12:   $p_{i+1} := u_{i+1} + \beta_{i+1}p_i$
13:   $s_{i+1} := w_{i+1} + \beta_{i+1}s_i$
14: **end for**

---

Recently, an asynchronous CG method has been presented by Gropp [21]. This method, shown as Algorithm 7, still has two global synchronization points per iteration. One reduction (line 9) can be overlapped with the matrix–vector product (line 10) and the other (line 3) with application of the preconditioner (line 4). Compared to standard CG, only two additional AXPYS are required. Algorithm 7 uses the same notation as above: $r_i = b - Ax_i$, $u_i = M^{-1}r_i$, $s_i = Ap_i$, $q_i = M^{-1}s_i$ and $w_i = Au_i$. However, the scalar $\delta$ is defined differently.

Table 1 gives a brief overview of the methods presented in this section with their parallel properties. The five methods listed are standard CG, the Chronopoulos/Gear variant of CG, preconditioned pipelined CG (pipe-CG), pipelined CR (pipe-CR) and the method by Gropp. In Table 1, the column *flops* lists the number of floating point operations for AXPYS and dot-products per iteration ($\times$ the vector length for the total). *Memory* lists how many vectors need to be stored, not counting $x$ and $b$. The last column counts the number of global reductions. The time spent in global communication and application of the sparse matrix–vector product (SPMV) and the preconditioner (PC) is given in the *time* column. Here, G is the time for a global all-reduce (a reduction followed by a broadcast), which is mostly latency bound and can be overlapped with the SPMV or with the preconditioner or with both. The pipelined CG method offers the most potential overlap with the reduction.

## 4. Numerical results

Numerical results are presented with the different CG methods, using different preconditioners, for a wide range of matrices from applications. A possible strategy to improve the maximum attainable accuracy of the methods is given in Section 4.2.

### 4.1. Problems from matrix market

The CG methods presented above have been implemented using PETSc (the Portable, Extensible Toolkit for Scientific Computation[1]). However, PETSc provides a CG method with a single global reduction due to D'Azevedo et al. [13] and D'Azevedo and Romine [14] that is based on the method by Chronopoulos and Gear, Algorithm 2, but differs slightly in the way the scalar $\alpha$ is computed. This single reduction method in PETSc can be used with the command line options `-ksp_type cg` and `-ksp_cg_single_reduction`. We shall compare with this implementation instead of Algorithm 2.

The methods presented above have been tested on a wide range of linear systems. Table 2 lists all square, real and symmetric positive definite matrices from Matrix Market,[2] which covers a wide range of condition numbers, listed in column two of Table 2. The bcsstk and bcsstm matrices are the $K$ and $M$ matrices respectively from a generalized eigenvalue problem $Kx = \lambda Mx$. Columns 3 and 4 give the total number of rows and number of nonzeros for each of the matrices. A linear system for each of these matrices with exact solution $\hat{x}_i = 1/\sqrt{N}$, with $N$ the number of rows, such that $\|\hat{x}\|_2 = 1$ and right-hand side $b = A\hat{x}$ is solved with all of the presented methods. The initial guess was always $x_0 = 0$ and the default PETSc stopping criterion

$$\|u_i\|_2 < \max(10^{-5}\|b\|_2, 10^{-50}) \tag{27}$$

was used with $u_i = M^{-1}(Ax_i - b)$ the preconditioned residual. The rest of the table lists the required number of iterations for the 6 different CG/CR methods

A '–' entry in the table denotes failure to meet the stopping criterion within 10,000 iterations. In all cases, the single reduction CG method behaves similar to standard CG. In most cases, pipe-CG and Gropp-CG need about as many iterations as standard CG, with a few glitches for pipe-CG. The deviation of the number of iterations required for CG1, Gropp-CG and pipe-CG compared to the number of iterations required for standard CG, averaged over all matrices, is given at the bottom of the table. Note that, mainly due to the nos6 matrix, on average, CG1 and Gropp-CG need less iterations than CG. For pipe-CG the average increase in number of iterations compared to CG is only about 3%. For pipe-CR compared to CR, the deviation is about 4%.

Fig. 1 illustrates the convergence history for a few randomly selected matrices from Table 2. Fig. 1 (top-left) shows the relative residual $\|b - Ax_i\|_2/\|b\|_2$ for the bcsstk15 matrix with Jacobi preconditioning. Fig. 1 (top-right) shows the s1rmt3m1 matrix with ICC preconditioner while in (bottom-left) and (bottom-right) the gr-30-30 and the nos2 matrices are shown respectively without preconditioner. For these experiments, the relative tolerance was set to rtol = $10^{-20}$. Observe from Fig. 1 that the convergence for CG, single reduction CG and Gropp-CG is nearly indistinguishable. The pipe-CG method converges as standard CG but levels off sooner, leading to a less accurate solution. The CR methods behaves differently, typically with faster initial convergence, but the same asymptotic behavior. Also, the pipe-CR method has worse final attainable accuracy.

These results are as expected, since all CG variants are mathematically equivalent to standard CG, i.e., they minimize the $A$ norm of the error, $\|\hat{x} - x_i\|_A$, at iteration $i$. The CR and pipe-CR methods minimize the $AM^{-1}A$ norm of the error or, without preconditioner, pipe-CR minimizes $\|e_i\|_{A^2} = \|r_i\|$ in every iteration, which is confirmed by the monotonic convergence of the residual in Fig. 1 (bottom).

---

**Table 1**
Overview of the different preconditioned CG and CR variations. Column *flops* lists the number of flops (×N) for AXPYS and dot-products. The *time* column has the time spent in global all-reduce communication (G), in the matrix–vector product (SPMV) and the preconditioner (PC). Column *#global synchronizations* has the number of global communication phases per iteration. The *memory* column counts the number of vectors that need to be kept in memory (excluding x and b).

| | Flops | Time (excl, AXPYS, DOTS) | #Glob syncs | Memory |
|---|---|---|---|---|
| CG | 10 | 2G + SpMV + PC | 2 | 4 |
| Chron/Gear-CG | 12 | G + SpMV + PC | 1 | 5 |
| CR | 12 | 2G + SpMV + PC | 2 | 5 |
| Pipe-CG | 20 | max(G, SpMV + PC) | 1 | 9 |
| Pipe-CR | 16 | max(G, SpMV) + PC | 1 | 7 |
| Gropp-CG | 14 | max(G, SpMV) + max(G, PC) | 2 | 6 |

**Table 2**
All real, square and symmetric positive definite matrices from Matrix Market, listed with their condition number, number of columns/rows and the total number of nonzeros. A linear system is solved with each of these matrices with the 6 different methods presented. The number of iterations required to reach $\|u_i\|_2 < \max(10^{-5}\|b\|_2, 10^{-50})$ with $u_i$ the preconditioned residual, is given. The deviation of the required number of iterations for CG1, Gropp-CG and p-CG compared to standard CG is averaged over all matrices and listed at the bottom of the table. The same is done for p-CR compared to CR.

| Matrix | Cond (A) | N | #nnz | Iterations | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | CG | CG1 | Gropp | p-CG | CR | p-CR |
| bcsstk14 | 1.3e+10 | 1806 | 63,454 | 1268 | 1369 | 1268 | 1374 | 262 | 270 |
| bcsstk15 | 8e+09 | 3948 | 117,816 | 2325 | 2419 | 2261 | 2474 | 465 | 478 |
| bcsstk16 | 65 | 4884 | 290,378 | 189 | 190 | 190 | 190 | 174 | 176 |
| bcsstk17 | 65 | 10,974 | 428,650 | 4280 | 4378 | 4185 | 4670 | 1803 | 1965 |
| bcsstk18 | 65 | 11,948 | 149,090 | 2858 | 2718 | 2850 | 2991 | 549 | 569 |
| bcsstk27 | 7.7e+04 | 1224 | 56,126 | 428 | 432 | 428 | 438 | 338 | 345 |
| bcsstm19 | 2.3e+05 | 817 | 817 | 166 | 175 | 165 | 187 | 146 | 171 |
| bcsstm20 | 2.6e+05 | 485 | 485 | 127 | 131 | 129 | 136 | 108 | 127 |
| bcsstm21 | 24 | 3600 | 3600 | 3 | 3 | 3 | 3 | 3 | 3 |
| bcsstm22 | 9.4e+02 | 138 | 138 | 42 | 43 | 42 | 42 | 42 | 42 |
| bcsstm23 | 9.5e+08 | 3134 | 3134 | 596 | 603 | 596 | 616 | 321 | 336 |
| bcsstm24 | 1.8e+13 | 3562 | 3562 | 366 | 375 | 365 | 381 | 296 | 317 |
| bcsstm25 | 6.1e+09 | 15,439 | 15,439 | 420 | 369 | 419 | 453 | 252 | 268 |
| bcsstm26 | 2.6e+05 | 1922 | 1922 | 476 | 436 | 451 | 462 | 292 | 301 |
| gr_30_30 | 3.8e+02 | 900 | 7744 | 33 | 33 | 33 | 33 | 33 | 33 |
| nos1 | 2.5e+07 | 237 | 1017 | 953 | 1005 | 950 | 1196 | 372 | 453 |
| nos2 | 6.3e+09 | 957 | 4137 | 643 | 642 | 647 | 963 | 344 | 369 |
| nos3 | 7.3e+04 | 960 | 15,844 | 219 | 219 | 219 | 220 | 215 | 215 |
| nos4 | 2.7e+03 | 100 | 594 | 72 | 72 | 71 | 72 | 71 | 71 |
| nos5 | 2.9e+04 | 468 | 5172 | 227 | 228 | 227 | 228 | 186 | 188 |
| nos6 | 8e+06 | 675 | 3255 | 198 | 165 | 164 | 187 | 121 | 121 |
| nos7 | 4.1e+09 | 729 | 4617 | 3300 | 3286 | 3025 | 664 | 2956 | – |
| s1rmq4m1 | 1.8e+06 | 5489 | 262,411 | 2352 | 2398 | 2419 | 2362 | 1203 | 1253 |
| s1rmt3m1 | 2.5e+06 | 5489 | 217,651 | 2801 | 2849 | 2830 | 3017 | 1455 | 1512 |
| s2rmq4m1 | 1.8e+08 | 5489 | 263,351 | 2189 | 2124 | 2188 | 2301 | 843 | 847 |
| s2rmt3m1 | 2.5e+08 | 5489 | 217,681 | 3281 | 3147 | 3286 | 3161 | 1161 | 1166 |
| s3dkq4m2 | 1.9e+11 | 90,449 | 2,455,670 | 3712 | 3712 | 3712 | 3712 | 1053 | 1053 |
| s3dkt3m2 | 3.6e+11 | 90,449 | 1,921,955 | 4939 | 4939 | 4939 | 4940 | 1455 | 1455 |
| s3rmq4m1 | 1.8e+10 | 5489 | 262,943 | 1539 | 1565 | 1569 | 1688 | 477 | 478 |
| s3rmt3m1 | 2.5e+10 | 5489 | 217,669 | 2104 | 2165 | 2139 | 2328 | 603 | 604 |
| s3rmt3m3 | 2.4e+10 | 5357 | 207,123 | 2976 | 3096 | 3006 | 3204 | 854 | 857 |
| Average deviation wrt CG | | | | | −0.039% | −0.89% | 3.0% | | |
| Average deviation wrt CR | | | | | | | | | 3.9% |

## 4.2. Improving accuracy

In order to improve the maximum attainable accuracy for the pipelined CG methods, we suggest to use a so-called residual replacement strategy, see for example [41,7,36]. In many Krylov iterations, the solution $x_i$ and residual vector $r_i$ are updated as

$$x_{i+1} = x_i + \alpha p_i, \quad r_{i+1} = r_i - \alpha A p_i. \tag{28}$$

Both solution and residual will be affected differently by rounding errors. Any error made in the update for $x$ is not reflected in $r$, since $x$ is typically not used in the rest of the iteration. This leads to the well known problem that the *updated* residual $r_i$ and the *true* residual $b - Ax_i$ start to deviate. A simple remedy is to periodically replace the updated $r_i$ by $r_i = b - Ax_i$. However, if this is done too frequently, the superlinear convergence often observed in CG can be lost. On the other hand, when the
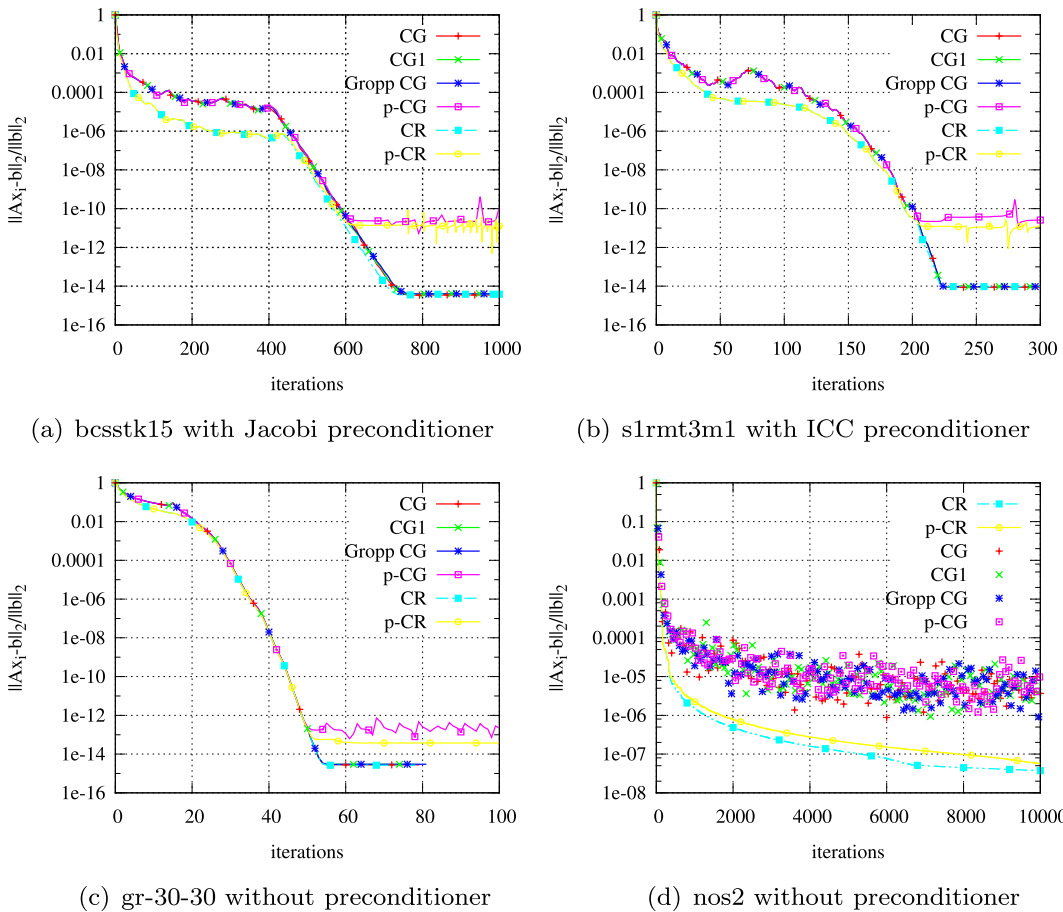
(a) bcsstk15 with Jacobi preconditioner
(b) s1rmt3m1 with ICC preconditioner
(c) gr-30-30 without preconditioner
(d) nos2 without preconditioner

**Fig. 1.** Convergence history for the different CG/CR methods applied to four different symmetric positive definite test matrices from applications (see also Table 2). Convergence of CG, single reduction CG (CG1) and Gropp CG is nearly indistinguishable. The pipelined methods level off somewhat sooner.
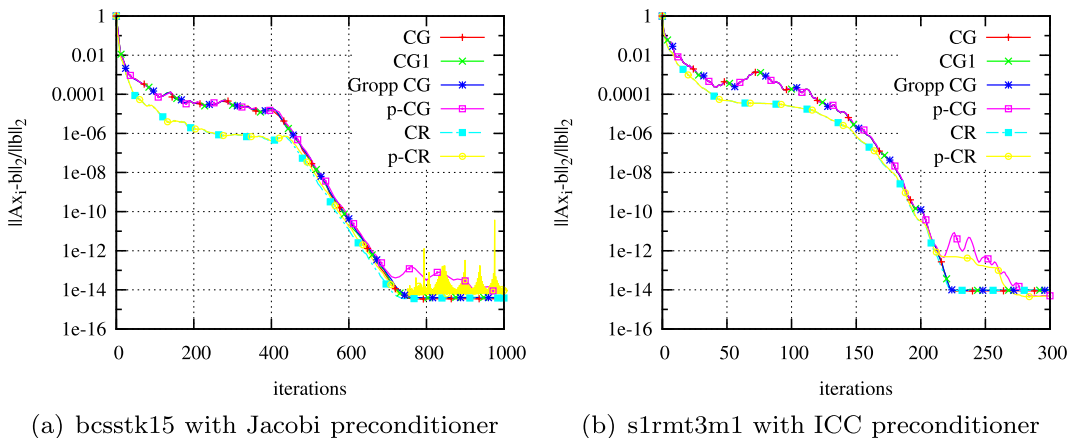


(a) bcsstk15 with Jacobi preconditioner
(b) s1rmt3m1 with ICC preconditioner

**Fig. 2.** Convergence history for the different CG methods applied to two of the matrices also shown in Fig. 1. However, for these tests the updated vectors $r_i, u_i$ and $w_i$ are replaced by $r_i = b - Ax_i$, $u_i = M^{-1}r_i$ and $w_i = Au_i$ every 50th iteration. This yields an improvement in maximum attainable accuracy of several orders of magnitude.

difference between updated and true residual gets too big, convergence stagnates too soon, resulting in a less accurate final solution.

We illustrate the potential of such a residual replacement strategy. In the pipe-CG and pipe-CR methods, the updated residual $r_i$, and the updated preconditioned residual $u_i$, will be replaced every 50th iteration by the true residual

$r_i = b - Ax_i$ and by the true preconditioned residual $u_i = M^{-1}r_i$ respectively. Furthermore, the vector $w_i = Au_i$ is also recomputed. For pipe-CR, the original residual $r_i$ does not need to be stored, but needs to be computed to evaluate $u_i$ anyway. Fig. 2 shows the convergence for two matrices also shown in Fig. 1, now repeated with the residual replacement strategy. For pipe-CG applied to bcsstk15 with Jacobi preconditioner the error after 1000 iterations drops from $\|e_{1000}\|_2 = 8.01e{-}8$ without residual replacement to $4.78e{-}11$ with residual replacement. Likewise, for s1rmt3m1 with ICC preconditioner, the error drops from $1.01e{-}8$ to $3.82e{-}12$. For pipe-CR, bcsstk15, Jacobi preconditioner, the error goes from $2.58e{-}8$ to $1.03e{-}9$ and for s1rmt3m1 with ICC preconditioner, from $2.59e{-}9$ to $1.65e{-}11$. This is always an improvement of at least one order of magnitude.

Of course, the value 50, also used in [36], is arbitrary and does not yield best possible results for all linear systems. Complicated strategies have been presented in the literature to determine good restart point. A strategy proposed by Neumaier [37] was to replace the updated residual by the true residual whenever the residual norm becomes smaller than any previously attained value. At that point, also the solution vector $x$ is updated with the combined contributions since the last restart point, a so-called group update. However, this strategy was proposed with the irregular convergence behavior of methods like BiCG and CGS in mind, where, as is well known, iterations with large residual can lead to early stagnation of the final residual. Van der Vorst and Ye [41] derive an estimate for the deviation between updated and true residual and only do the residual replacement and group update when this bound exceeds some tolerance $\epsilon$, typically the square root of the machine precision. Recently, Carson and Demmel [7] extended the bound from [41] to communication-avoiding (CA) or $s$-step Krylov methods and present an implementation for CA-CG and CA-BiCG. The residual replacement strategy has proven to be an effective strategy to improve the final accuracy of short term Krylov methods. Each replacement incurs additional work, but since replacements occur infrequently the performance impact is limited. However, a general replacement strategy applicable to the pipelined methods presented here remains future work.

Stability is also negatively impacted in the pipelined methods by the extra multiplication with the matrix $A$. A possible improvement might be to add a shift in the matrix–vector product, similar to what is also done in $s$-step Krylov methods [2,27] and in [18] for pipelined GMRES. The CG iteration can provide information on the spectrum of $A$, which can be used to determine good shifts.

## 5. Parallel performance

Different variations of standard CG have been presented in order to overcome the bottleneck of global communication on large parallel machines. In Section 5.1, the parallel performance of the different methods is studied using a test problem on a medium sized parallel machine. Section 5.2 lists several candidate scenarios which could benefit from the pipelined CG/CR solvers.

### 5.1. Benchmark application: ice flow simulation

The parallel experiments are performed on a small cluster with 20 compute nodes. All nodes are interconnected by $4 \times$ QDR InfiniBand technology, providing 32 Gb/s of point-to-point bandwidth for high-performance message passing and I/O. The nodes have two 6-core Intel Xeon X5660 Nehalem 2.80 GHz processors (twelve cores/node) for a total of 240 cores.

The pipe-CG, pipe-CR and Gropp-CG algorithms have been implemented in PETSc [3], and are available from version 3.4.1. PETSc provides a construct for asynchronous dot-products:

– VecDotBegin (…,&dot);
– PetscCommSplitReductionBegin (comm);
– // …other work, like an spMV or application of the preconditioner
– VecDotEnd (…,&dot);

where PetscCommSplitReductionBegin starts a non-blocking reduction by a call to `MPI_Iallreduce`. This returns an `MPI_Request` which is passed as an argument to `MPI_Wait` in VecDotEnd. Note that non-blocking collectives, including `MPI_Iallreduce`, only became available in MPI with the MPI-3 standard. The MPI library used for the experiments is MPICH-3.0.4,[3] with the environment variables `MPICH_ASYNC_PROGRESS=1` and `MPICH_MAX_THREAD_SAFETY=multiple` set.

The test problem we use to asses the parallel performance of the presented methods is a finite element code solving the hydrostatic equations for ice sheet flow, described in [6] and available in the PETSc distribution as example 48 in the nonlinear solvers folder. The simulation domain is discretized using $100 \times 100 \times 50$ $Q1$ finite elements. A nonlinear equation is solved with a line search Newton solver, which repeatedly calls a Krylov solver for the linearized problem. As a preconditioner block Jacobi with ICC (0) inside the blocks is used with one block per core. For the nonlinear solver the relative and absolute tolerances used are $10^{-8}$ and $10^{-15}$ respectively and for the linear solver the tolerances are $10^{-5}$ and $10^{-50}$.
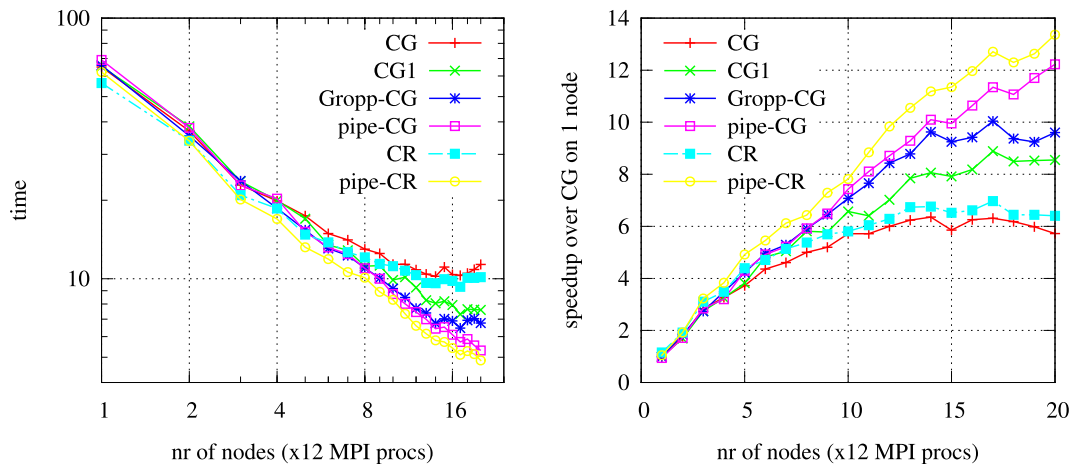
---

[3] http://www.mpich.org/

**Fig. 3.** Left: Time to solution for the 3*D* hydrostatic ice sheet flow simulation using $100 \times 100 \times 50$ *Q*1 finite elements. Right: Speedup as function of number of nodes over standard CG solver on a single node.

Fig. 3 (left) shows the time to solution for the simulation as a function of the number of nodes, with 12 cores and 12 MPI processes per node. For all CG variants the residual norm used in the stopping criterion was the natural norm, i.e., $\sqrt{r_i^T u_i}$. This norm was chosen because the current implementation of CG in PETSc performs one additional global reduction when the preconditioned residual norm is used in the stopping criterion. For a fair comparison this extra reduction was avoided. For the CR methods however, the preconditioned norm $\|u_i\|_2$ was used because pipe-CR does not compute the original residual and hence cannot compute the natural residual norm.

For pipe-CG the maximum speedup compared to CG was $2.14\times$ and $1.43\times$ compared to CG1 and this was achieved on 20 nodes. On 20 nodes, pipe-CG required 7 nonlinear iterations with a total of 1328 linear iterations, exactly the same as for standard CG and CG1. For pipe-CR the maximum speedup compared to CR, $2.09\times$, was also attained on 20 nodes. On 20 nodes, pipe-CR required 7 nonlinear and 1166 linear iterations, the same as CR.

### 5.2. Possible use cases and further optimizations

The aim of this work is to develop a generally applicable highly scalable linear solver. However, the resulting pipelined CG algorithm is a trade-off between more floating point operations, more memory usage, and improved scalability. Therefore, not every application will benefit from this trade-off. We identify a set of scenarios where the improved scalability will likely pay off.

As was already mentioned, to get fast convergence, preconditioning is essential. In a solver with a near optimal preconditioner, most of the time is probably spent in application of the preconditioner and a pipelined outer CG solver will not pay off. However, problems where no good preconditioner is available would still require many CG iterations and could benefit. For instance augmented Lagrangian systems from contact and interior points are examples of nearly-singular symmetric systems for which a multigrid preconditioner is hard to construct.

Also when repeatedly solving linear system, e.g., in a nonlinear solver, in optimization or in implicit time integration, many linear solver iterations are required which could benefit from the pipelined CG solver. In some special cases, a relatively small linear system is solved on a parallel machine, which makes the linear solver mostly latency bound. Hiding the latency as in pipelined CG directly pays off in such cases. Examples are: the solver for the coarsest grid in a geometric multigrid *U*-cycle [43], the field solver in coupled physics methods such as the Particle-In-Cell [24] method where the size of the field is typically small compared to the number of particles used in the simulation, and domain decomposition and Schur complement solvers.

When the reductions are completely overlapped, the pipelined CG methods will scale as the spmv, and the spmv communication might become the bottleneck. In pipe-CG, Algorithm 4, application of the preconditioner (line 5) is followed immediately by the spmv (line 6). A possible optimization would be to fuse the two operations. Consider for instance a polynomial preconditioned linear system $p^s(A)Ax = p^s(A)b$ with $p^s(A)$ a degree *s* polynomial in *A*. Applying both the preconditioner and the matrix–vector product could be done in a single latency by the matrix-powers kernel [17,19]. For structured grid problems, this is relatively easy to implement by communicating a wider ghost region, which is also called communication aggregation, since the same amount of data is transferred, only in less messages [42]. The polynomial preconditioner can then also be implemented efficiently using techniques as cache-oblivious, wavefront or time-skewed stencil loops [9]. Information about the spectrum of *A* can be computed as a side product of the CG iteration and can help to determine an optimal polynomial [33]. However, in a similar way, steps in a multigrid fine-grid smoother can be fused with the matrix–vector product from CG.

## 6. Conclusions

We presented pipelined variations of CG where the global communication can be overlapped with local work, such as the matrix–vector product, application of the preconditioner and local computations. Two different ways to include preconditioning are discussed, leading to two algorithms with different trade-offs regarding scalability and total number of operations. For a strong scaling experiment on a medium sized cluster, we show improved scalability and also faster runtime compared to standard CG. Numerical tests on matrices from Matrix Market with a wide range of condition numbers show that the convergence of the new methods is in line with standard CG. We also compare with a method recently presented by Gropp [21] that has somewhat better numerical properties but offers less overlap with the global communication and thus leads to a less scalable algorithm.

For large machines where a global reduction is expensive compared to a matrix–vector product, it can happen that the latency of the global reduction is not completely overlapped in the pipelined methods. If this turns out to become a bottleneck (perhaps on extremely large or on future systems), the pipelined strategy can be generalized to longer pipelining depths. In this case, the global reduction can be overlapped with multiple CG iterations (matrix–vector products), as was already presented for GMRES in [18]. However, longer pipelines also increase the recurrence length and the required number of operations and memory.

The presented pipelined methods can be seen as an alternative to $s$-step CG [7,39]. Historically, $s$-step CG has had a "bad reputation" [40,33] due to the numerical problems for increasing $s$. By recent advances such as different Krylov bases [2,27] (like Newton or Chebyshev) and a residual replacement strategy [7], these methods have been applied successfully with quite large values of $s$. However, a problem with $s$-step methods remains the combination with preconditioning. Although some progress has been made in this regard recently [20], finding communication avoiding preconditioners that combine well with $s$-step methods remains a challenge. The pipelined methods presented in this work can be used with any preconditioner.

A finite precision round-off analysis for the newly presented methods, similar to the work in [38], is left as future work.

## References

[1] T. Ashby, P. Ghysels, W. Heirman, W. Vanroose, The impact of global communication latency at extreme scales on Krylov methods, Algorithms and Architectures for Parallel Processing (2012) 428–442.

[2] Z. Bai, D. Hu, L. Reichel, A Newton basis GMRES implementation, IMA Journal of Numerical Analysis 14 (4) (1994) 563–581.

[3] S. Balay, J. Brown, K. Buschelman, W.D. Gropp, D. Kaushik, M.G. Knepley, L. Curfman McInnes, B.F. Smith, H. Zhang, PETSc Web page, 2013, <http://www.mcs.anl.gov/petsc>.

[4] R. Barrett, M. Berry, T.F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, H. Van der Vorst, Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods, second ed., SIAM, Philadelphia, PA, 1994.

[5] J. Brown, User-defined nonblocking collectives must make progress, in: IEEE Technical Committee on Scalable Computing (TCSC), 2012.

[6] Jed Brown, Barry F. Smith, Aron Ahmadia, Achieving textbook multigrid efficiency for hydrostatic ice flow, SIAM Journal on Scientific Computing 35 (2) (2013) 359–375. Also, preprint ANL/MCS-P743-1298.

[7] E. Carson, J. Demmel, A residual replacement strategy for improving the maximum attainable accuracy of s-step Krylov subspace methods, Technical Report UCB/EECS-2012-44, University of California, Berkeley, CA, USA, 2012.

[8] E. Carson, N. Knight, J. Demmel, Avoiding communication in two-sided Krylov subspace methods, Technical report, University of California, Berkeley, CA, USA, 2011.

[9] M. Christen, O. Schenk, H. Burkhart, Patus: A code generation and autotuning framework for parallel iterative stencil computations on modern microarchitectures, in: Parallel & Distributed Processing Symposium (IPDPS), 2011 IEEE International, IEEE, 2011, pp. 676–687.

[10] A.T. Chronopoulos, s-Step iterative methods for (non) symmetric (in) definite linear systems, SIAM Journal on Numerical Analysis 28 (6) (1991) 1776–1789.

[11] A.T. Chronopoulos, C.W. Gear, s-Step iterative methods for symmetric linear systems, Journal of Computational and Applied Mathematics 25 (2) (1989) 153–168.

[12] A.T. Chronopoulos, C.D. Swanson, Parallel iterative s-step methods for unsymmetric linear systems, Parallel Computing 22 (5) (1996) 623–641.

[13] E.F. D'Azevedo, V.L. Eijkhout, C.H. Romine, Lapack working Note 56 conjugate gradient algorithms with reduced synchronization overhead on distributed memory multiprocessors, 1999.

[14] E.F. D'Azevedo, C.H. Romine, Reducing communication costs in the conjugate gradient algorithm on distributed memory multiprocessors, Technical report, Oak Ridge National Lab, TN, 1992.

[15] E. De Sturler, H.A. Van der Vorst, Reducing the effect of global communication in GMRES(m) and CG on parallel distributed memory computers, Applied Numerical Mathematics 18 (4) (1995) 441–459.

[16] J. Demmel, M.T. Heath, H.A. Van Der Vorst, Parallel numerical linear algebra, Acta Numerica 2 (-1) (1993) 111–197.

[17] J. Demmel, M. Hoemmen, M. Mohiyuddin, K. Yelick, Avoiding communication in sparse matrix computations, in: 2008 IEEE International Symposium on Parallel and Distributed Processing, 2008, pp. 1–12.

[18] P. Ghysels, T.J. Ashby, K. Meerbergen, W. Vanroose, Hiding global communication latency in the GMRES algorithm on massively parallel machines, SIAM Journal on Scientific Computing 35 (1) (2013) C48–C71.

[19] P. Ghysels, P. Kłosiewicz, W. Vanroose, Improving the arithmetic intensity of multigrid with the help of polynomial smoothers, Numerical linear algebra with applications 19 (2) (2012) 253–267 (Special issue Copper Mountain 2011).

[20] L. Grigori, S. Moufawad, Communication avoiding ILU(0) preconditioner, Rapport de recherche RR-8266, INRIA, March 2013.

[21] W. Gropp, Update on libraries for blue waters, Bordeaux, France, 2010. http://jointlab.ncsa.illinois.edu/events/workshop3/pdf/presentations/Gropp-Update-on-Libraries.pdf>.

[22] V. Hernandez, J.E. Roman, A. Tomas, Parallel Arnoldi eigensolvers with enhanced scalability via global communications rearrangement, Parallel Computing 33 (7–8) (2007) 521–540.

[23] M.R. Hestenes, E. Stiefel, Methods of conjugate gradients for solving linear systems, Journal of Research of the National Bureau of Standards 49 (6) (1952).

[24] R.W. Hockney, J.W. Eastwood, Computer Simulation Using Particles, Adam Hilger, 1988.

[25] T. Hoefler, T. Schneider, A. Lumsdaine, LogGOPSim – simulating large-scale applications in the LogGOPS model, in: Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing, ACM, 2010, pp. 597–604.

[26] T. Hoefler, J. Squyres, G. Bosilca, G. Fagg, A. Lumsdaine, W. Rehm, Non-blocking collective operations for MPI-2, Open Systems Lab, Indiana University, Tech. Rep, 8, 2006.

[27] M. Hoemmen, Communication-avoiding Krylov subspace methods, Ph.D. Thesis, University of California, 2010.

[28] W.D. Joubert, G.F. Carey, Parallelizable restarted iterative methods for nonsymmetric linear systems. Part I: Theory, International Journal of Computer Mathematics 44 (1–4) (1992) 243–267.

[29] S.K. Kim, A.T. Chronopoulos, An efficient parallel algorithm for extreme eigenvalues of sparse nonsymmetric matrices, International Journal of High Performance Computing Applications 6 (4) (1992) 407–420.

[30] L.C. McInnes, B. Smith, H. Zhang, R. Tran Mills. Hierarchical and nested Krylov methods for extreme-scale computing, Technical Report ANL/MCS-P2097-0612, Argonne National Laboratory, 2012.

[31] G. Meurant, Multitasking the conjugate gradient method on the CRAY X-MP/48, Parallel Computing 5 (3) (1987) 267–280.

[32] Y. Saad, Practical use of some Krylov subspace methods for solving indefinite and nonsymmetric linear systems, SIAM Journal on Scientific and Statistical Computing 5 (1) (1984) 203–228.

[33] Y. Saad, Krylov subspace methods on supercomputers, SIAM Journal on Scientific and Statistical Computing 10 (1989) 1200.

[34] Y. Saad, Iterative Methods for Sparse Linear Systems, Society for Industrial Mathematics, 2003.

[35] A. Schäfer, D. Fey, LibGeoDecomp: a grid-enabled library for geometric decomposition codes, in: Proceedings of the 15th European PVM/MPI Users' Group Meeting on Recent Advances in Parallel Virtual Machine and Message Passing Interface, Springer-Verlag, 2008, pp. 285–294.

[36] J.R. Shewchuk, An introduction to the conjugate gradient method without the agonizing pain, 1994.

[37] G.L.G. Sleijpen, H.A. van der Vorst, Reliable updated residuals in hybrid Bi-CG methods, Computing 56 (2) (1996) 141–163.

[38] Z. Strakoš, P. Tichỳ, On error estimation in the conjugate gradient method and why it works in finite precision computations, Electronic Transactions on Numerical Analysis 13 (2002) 56–80.

[39] S.A. Toledo, Quantitative performance modeling of scientific computations and creating locality in numerical algorithms, Ph.D. Thesis, Massachusetts Institute of Technology, 1995.

[40] H.A. Van der Vorst, Iterative Krylov Methods for Large Linear Systems, vol. 13, Cambridge University Press, 2003.

[41] H.A. Van Der Vorst, Q. Ye, Residual replacement strategies for Krylov subspace iterative methods for the convergence of true residuals, SIAM Journal on Scientific Computing 22 (3) (1999) 835–852.

[42] S. Williams, D.D. Kalamkar, A. Singh, A.M. Deshpande, B. Van Straalen, M. Smelyanskiy, A. Almgren, P. Dubey, J. Shalf, L. Oliker, Optimization of geometric multigrid for emerging multi-and manycore processors, in: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, IEEE Computer Society Press, 2012, p. 96.

[43] D. Xie, L.R. Scott, An analysis of parallel U-cycle multigrid method, 2010.

[44] L.T. Yang, The improved CGS method for large and sparse linear systems on bulk synchronous parallel architectures, in: Proceedings of the Fifth International Conference on Algorithms and Architectures for Parallel Processing, 2002, IEEE, 2002, pp. 232–237.

[45] L.T. Yang, R.P. Brent, The improved BiCGStab method for large and sparse unsymmetric linear systems on parallel distributed memory architectures, in: Proceedings of the Fifth International Conference on Algorithms and Architectures for Parallel Processing, 2002, 2002, pp. 324–328.

[46] L.T. Yang, R.P. Brent, The improved parallel BiCG method for large and sparse unsymmetric linear systems on distributed memory architectures, in: Proceedings of the 16th International Parallel and Distributed Processing Symposium, IPDPS 2002, IEEE, 2003.

[47] T.R. Yang, H.X. Lin, The improved quasi-minimal residual method on massively distributed memory computers, in: High-Performance Computing and Networking, Springer, 1997, pp. 389–399.