

# A cooperative conjugate gradient method for linear systems permitting multithread implementation of low complexity

Amit Bhaya, Pierre-Alexandre Bliman, Guilherme Niedu and Fernando Pazos <sup>\*†</sup>

June 18, 2018

## Abstract

This paper proposes a generalization of the conjugate gradient (CG) method used to solve the equation  $Ax = b$  for a symmetric positive definite matrix  $A$  of large size  $n$ . The generalization consists of permitting the scalar control parameters (= stepsizes in gradient and conjugate gradient directions) to be replaced by matrices, so that multiple descent and conjugate gradient directions are updated simultaneously. Implementation involves the use of multiple agents or threads and is referred to as cooperative CG (cCG), in which the cooperation between agents resides in the fact that the calculation of each entry of the control parameter matrix now involves information that comes from the other agents. For a sufficiently large dimension  $n$ , the use of an optimal number of cores gives the result that the multithread implementation has worst case complexity  $O(n^{2+1/3})$  in exact arithmetic. Numerical experiments, that illustrate the interest of theoretical results, are carried out on a multicore computer.

## 1 Introduction

The paradigm of cooperation between agents in order to achieve some common objective has now become quite common in many areas such as control and distributed computation, while representing a multitude of different situations and related mathematical questions [7, 11, 10].

In the field of computation, the emphasis has been mainly on the paradigm of parallel computing in which some computational task is subdivided into as many subtasks as there are available processors. The subdivision naturally induces a communication structure (or graph), connecting processors and the challenge is to achieve a subdivision that maximizes concurrency of tasks (hence minimizing total computational time), while simultaneously minimizing communication overhead. This paradigm arose as a consequence of the usual architecture of most early multiprocessor machines, in which interprocessor communication is a much slower operation than a mathematical operation carried out in the same processor. Disadvantages of this approach arise from the difficulty of effectively decomposing a large task into minimally connected subtasks, difficulties of analysis and the need for synchronization barriers at which all processors wait for the slowest one, in order to exchange information with the correct time stamps (i.e., without asymmetric delays).

More recently, in the area of control, interest has been focused on multiagent systems, in which a number of agents cooperate amongst themselves, in a distributed manner and also subject to a communication graph that describes possible or allowable channels between agents, in order to achieve some (computational) task. Similarly, in the area of computation, multicore processors have now become common – in these processors, each core accommodates a thread which is executed independently of the threads in the other cores. Thus, in the context of this paper, which is focused on solution of the linear system of equations  $Ax = b$  for a symmetric positive definite matrix  $A$  of large size  $n$ , we will assume that each agent carries out a task that is represented by one thread that executes on one core, so that, in this sense, the words agent and thread can be assumed to represent the same thing. In what follows,

---

<sup>\*</sup>ABs work was supported by grants BPP/CNPq and, additionally, CNE from FAPERJ and Universal/CNPq. GN and FP were supported by DS and PNPD fellowships, respectively, from CNPq.

<sup>†</sup>A. Bhaya, G. Niedu, F. Pazos are with the Department of Electrical Engineering, Federal University of Rio de Janeiro, Rio de Janeiro, Brazil, [amit@nacad.ufrj.br](mailto:amit@nacad.ufrj.br). P.-A. Bliman is with Inria, Rocquencourt BP105, 78153 Le Chesnay cedex, France, [pierre-alexandre.bliman@inria.fr](mailto:pierre-alexandre.bliman@inria.fr)

unless we are specifically talking about numerical implementation, we will give preference to the word agent.

With the advent of ever larger on-chip memory and multicore processors that allow multithread programming, it is now possible to propose a new paradigm in which each thread, with access to a common memory, computes its own estimate of the solution to the whole problem (i.e., decomposition of the problem into subproblems is avoided) and the threads exchange information amongst themselves, this being the cooperative step. The design of a cooperative algorithm has the objective of ensuring that exchanged information is used by the threads in such a way as to reduce overall convergence time.

The idea of information exchange between two iterative processes was introduced into numerical linear algebra long before the advent of multicore processors by Brezinski [3] under the name of *hybrid procedures*, defined as (we quote) “a combination of two arbitrary approximate solutions with coefficients summing up to one...(so that) the combination only depends on one parameter whose value is chosen in order to minimize the Euclidean norm of the residual vector obtained by the hybrid procedure... The two approximate solutions which are combined in a hybrid procedure are usually obtained by two iterative methods.” The objective of minimizing the residue is to accelerate convergence of the overall hybrid procedure. This idea was generalized and discussed in the context of distributed asynchronous computation in [1].

More specifically, this paper explores the idea of cooperation between  $p$  agents (or threads) in the context of the conjugate gradient (CG) algorithm applied to an  $n$ -dimensional linear system  $Ax = b$ , for a symmetric positive definite matrix  $A$  of large size  $n$ . Throughout the paper it is assumed that  $p < n$ , and even that  $p \ll n$ : the number of agents may be “large”, but it is usually “much smaller” than the “huge” size of matrix  $A$ . The famous CG algorithm, proposed in [6], has several interesting properties, both as an algorithm in exact arithmetic and as one in finite precision arithmetic [9, 4]. However, it is well known that, due to its structure, it cannot be parallelized in the conventional sense. In this paper, we revisit the CG algorithm from a multithread perspective, which can be seen as a direct generalization of the control approach to the CG algorithm proposed in [2, pp.77-82], in which the scalar control parameters (stepsizes in gradient and conjugate gradient directions) are replaced by matrices (i.e., multivariable control). The cooperation between agents resides in the fact that the calculation of each entry of the control matrix now involves information that comes from the other agents. The method can also be seen as a generalization of the traditional CG algorithm in which multiple descent and conjugate directions are updated simultaneously.

The paper is organized as follows. Section 2 briefly recalls the construction, as well as the main convergence results, of Conjugate Gradient method. Section 3 then presents the new algorithm, called *cooperative Conjugate Gradient (cCG) method*. In order to simplify this presentation of the new algorithm, the case of  $p = 2$  agents is first introduced in Section 3.1. The general case  $p \geq 2$  is then stated in full generality in Section 3.2, together with analysis results. Complexity issues are broached in Section 4. The results stated therein concerns execution of cCG algorithm in exact arithmetic. Section 5 is then devoted to numerical experiments with the multi-thread implementation. Section 6 provides conclusions and directions for future work.

**Notation** For the fixed symmetric definite positive matrix  $A \in \mathbb{R}^{n \times n}$ , we define  $A$ -norm in  $\mathbb{R}^n$  by

$$\|x\|_A \doteq (x^T A x)^{1/2}, \quad x \in \mathbb{R}^n \quad (1)$$

and define  $A$ -orthogonality (or *conjugacy*) of vectors by:

$$x \perp_A y \Leftrightarrow x^T A y = 0, \quad x, y \in \mathbb{R}^n. \quad (2)$$

We will also have to consider matrices whose columns are vectors of interest. Accordingly, we will say that  $X, Y \in \mathbb{R}^{n \times p}$  are orthogonal (resp.  $A$ -orthogonal) whenever each column of  $X$  is orthogonal (resp.  $A$ -orthogonal) to each column of  $Y$ , that is when

$$X^T Y = 0 \quad (\text{resp. } X^T A Y = 0). \quad (3)$$

For any set of vectors  $r_i \in \mathbb{R}^n$ ,  $i = 0, 1, \dots, k$ , we denote respectively  $\{r_i\}_0^k$  and  $[r_i]_0^k$  the set of these vectors, and the matrix obtained by their concatenation:  $[r_i]_0^k = [r_0 \ r_1 \ \dots \ r_k] \in \mathbb{R}^{n \times (k+1)}$ . The

notation  $\text{span } [r_i]_0^k$  will denote the subspace of linear combinations of the columns of the matrix  $[r_i]_0^k$ . When  $\mathbb{R}^n$  is the ambient vector space, we thus have

$$\text{span } [r_i]_0^k \doteq \left\{ v \in \mathbb{R}^n : \exists \gamma \in \mathbb{R}^{k+1}, v = \sum_{i=0}^k \gamma_i r_i = [r_i]_0^k \gamma \right\}. \quad (4)$$

Similarly, for matrices  $R_i \in \mathbb{R}^{n \times p}$ ,  $i = 0, \dots, k$ , the notation  $\{R_i\}_0^k$  (respectively,  $[R_i]_0^k \in \mathbb{R}^{n \times (k+1)p}$ ) is used for the set of these matrices (respectively, the matrix obtained as concatenation of the matrices  $R_0, R_1, \dots, R_k$ , i.e.,  $[R_i]_0^k = [R_0 \ R_1 \ \dots \ R_k] \in \mathbb{R}^{n \times (k+1)p}$ .) Also, we write  $\text{span } [R_i]_0^k$  for the subspace of linear combinations of the columns of  $[R_i]_0^k$ :

$$\text{span } [R_i]_0^k \doteq \left\{ v \in \mathbb{R}^n : \exists \gamma \in \mathbb{R}^{(k+1)p}, v = [R_i]_0^k \gamma \right\}. \quad (5)$$

Notice that this notation generalizes the definition provided earlier for vectors, and that  $\text{span } [R]$  is already meaningful for a single matrix  $R \in \mathbb{R}^{n \times p}$ . As an example,  $\dim \text{span } [R] = \text{rank } R$ .

Last, for any matrix  $R \in \mathbb{R}^{n \times p}$  and for any set  $J$  of indices in  $\{1, \dots, p\}$ , we will denote

$$R|_{j \in J} \doteq (R_{ij})_{1 \leq i \leq n, j \in J}.$$

## 2 The Conjugate Gradient Method

One approach to solving the equation

$$Ax = b, \quad (6)$$

with  $A$  symmetric positive definite and of large dimension, is to minimize instead the convex quadratic function

$$f(x) = \frac{1}{2} x^T A x - b^T x, \quad (7)$$

since the unique optimal point is  $x^* = A^{-1}b$ . Several algorithms are based on the standard idea of generating a sequence of points, starting from an arbitrary initial guess, and proceeding in the descent direction (negative gradient of  $f(x)$ ), with an adequate choice of the step size. In mathematical terms:

$$x_{k+1} = x_k - \alpha_k r_k, \quad r_k = \nabla f(x_k) = Ax_k - b, \quad (8)$$

where  $\alpha_k$  is the step size. The vector  $r_k$  represents both the *gradient* of the cost function  $f$  at the current point  $x_k$ , and the current *residue* in the process of solving (6).

Amongst the possible choices for  $\alpha_k$ , a most natural one consists in minimizing the value of the function  $f$  at  $x_{k+1}$ , that is in taking

$$\alpha_k = \arg \min_{\alpha \in \mathbb{R}} f(x_k - \alpha r_k) \quad (9)$$

The algorithm obtained using this principle is the *Steepest Descent Method*, and one shows easily that the optimal value is given by the Rayleigh quotient

$$\alpha_k = \frac{r_k^T r_k}{r_k^T A r_k} \quad (10)$$

Algorithm (8)-(10) is convergent, but in general one cannot expect better convergence speed than the one provided by

$$\|x_k - x^*\|_A \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^k \|x_0 - x^*\|_A \quad (11)$$

where  $\kappa$  is the condition number

$$\kappa \doteq \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}. \quad (12)$$

The main weakness of Steepest Descent is the fact that steps taken in the same directions as earlier steps are likely to occur. The *Conjugate Direction Methods* avoid this drawback. Based on a set of  $n$

mutually conjugate nonzero vectors  $\{d_i\}_0^{n-1}$  (that is  $d_i \perp_A d_j$  for any  $i \neq j$ ,  $i, j = 0, \dots, n-1$ ), the family of conjugate direction methods use the sequence generated according to

$$x_{k+1} = x_k + \alpha_k x_k, \quad \alpha_k = -\frac{r_k^\top d_k}{d_k^\top A d_k}. \quad (13)$$

It is a classical result that the residue  $r_{k+1}$  is orthogonal to the directions  $d_i$ ,  $i = 0, \dots, k$ ; and that  $x_{k+1}$  indeed minimizes  $f$  on the affine subspace  $x_0 + \text{span} [d_i]_0^k$  [8]. As a consequence of this last property, conjugate direction methods lead to finite time convergence (in exact arithmetic).

The *Conjugate Gradient method*, developed by Hestenes and Stiefel [6], is the particular method of conjugate directions obtained when constructing the conjugate directions by Gram-Schmidt orthogonalization, achieved at step  $k+1$  on the set of the gradients  $\{r_i\}_0^k$ . A key point here is that this construction can be carried out iteratively. The iterative equations of the Conjugate Gradient method are given in the pseudocode instructions of Algorithm 1. Instructions 6–7 constitute the optimal descent process in the direction  $d_k$ ; while instructions 9–10 achieve iteratively the orthogonalization of the subspaces  $\text{span} [r_i]_0^k$ .

---

**Algorithm 1** Conjugate Gradient (CG) algorithm

---

```

1: choose  $x_0 \in \mathbb{R}^n$ 
2:  $r_0 := Ax_0 - b$ 
3:  $d_0 := r_0$ 
4:  $k := 0$ 
5: while  $d_k \neq 0$  do
6:    $\alpha_k := -r_k^\top d_k (d_k^\top A d_k)^{-1}$ 
7:    $x_{k+1} := x_k + \alpha_k d_k$ 
8:    $r_{k+1} := Ax_{k+1} - b$ 
9:    $\beta_k := -r_{k+1}^\top A d_k (d_k^\top A d_k)^{-1}$ 
10:   $d_{k+1} := r_{k+1} + \beta_k d_k$ 
11:   $k \leftarrow k + 1$ 
12: end while

```

---

We recall the main properties of this algorithm, in an adapted form, to allow for easier comparison with the results to be stated later.

**Theorem 1** (Properties of CG). *As long as the vector  $d_k$  is not zero*

- the vectors  $\{r_i\}_0^k$  are mutually orthogonal, the vectors  $\{d_i\}_0^k$  are mutually  $A$ -orthogonal, and the subspaces  $\text{span} [r_i]_0^k$ ,  $\text{span} [d_i]_0^k$  and  $\text{span} [A^i r_0]_0^k$  are equal and have dimension  $(k+1)$ ;
- the point  $x_{k+1}$  is the minimizer of  $f$  on the affine subspace  $x_0 + \text{span} [d_i]_0^k$ .

When the residue vector is zero, the optimum has been attained, showing that CG terminates in finite time. Apart from the finite time convergence property, the following formula indicates net improvement with respect to Steepest Descent:

$$\|x_k - x^*\|_A \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x_0 - x^*\|_A \quad (14)$$

which represents substantial improvement with respect to (11).

For a proof of this theorem as well as further details on the contents of this section, see [8, 5].

## 3 Statement and Analysis of the Cooperative Conjugate Gradient Method

### 3.1 The two-agent case

In this subsection, in order to aid comprehension and ease notation, the case of two agents (the case  $p = 2$ ) is considered: their estimates at step  $k$  are written as  $x_k, x'_k$  respectively, the residues as  $r_k, r'_k$ ,

and the two descent directions as  $d_k, d'_k$ . The gradients at each one of the current estimates are given as

$$\begin{pmatrix} r_k & r'_k \end{pmatrix} = A \begin{pmatrix} x_k & x'_k \end{pmatrix} - \mathbf{1}^\top b, \quad (15)$$

with  $\mathbf{1} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ . As for CG method, we distinguish two steps.

• **Descent step.** Given the current residues  $r_k, r'_k$  and two descent directions  $d_k, d'_k$ , this step determines the upgraded value of the estimates  $x_{k+1}, x'_{k+1}$  and therefore of the residues  $r_{k+1}, r'_{k+1}$ .

One allows the use of the two descent directions  $d_k, d'_k$ , thus looking for updates of the form

$$\begin{pmatrix} x_{k+1} & x'_{k+1} \end{pmatrix} = \begin{pmatrix} x_k & x'_k \end{pmatrix} + \begin{pmatrix} d_k & d'_k \end{pmatrix} \alpha_k^\top. \quad (16)$$

The matrix  $\alpha_k \in \mathbb{R}^{2 \times 2}$  has to be chosen. In the same spirit as for CG, this choice is made in such a way as to minimize  $f(x_{k+1})$  and  $f(x'_{k+1})$ . This yields in fact two independent minimization problems. Denoting

$$\alpha_j \doteq (\alpha_{j1} \quad \alpha_{j2}), \quad j = 1, 2, \quad (17)$$

the two optimality conditions are given by

$$0 = \begin{pmatrix} d_k & d'_k \end{pmatrix}^\top A \begin{pmatrix} x_k + \begin{pmatrix} d_k & d'_k \end{pmatrix} \alpha_1^\top \end{pmatrix} - \begin{pmatrix} d_k & d'_k \end{pmatrix}^\top b = \begin{pmatrix} d_k & d'_k \end{pmatrix}^\top A \begin{pmatrix} x'_k + \begin{pmatrix} d_k & d'_k \end{pmatrix} \alpha_2^\top \end{pmatrix} - \begin{pmatrix} d_k & d'_k \end{pmatrix}^\top b. \quad (18)$$

This shows that the minimum is uniquely defined, and attained when

$$\alpha_k = \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} = - \begin{pmatrix} r_k & r'_k \end{pmatrix}^\top \begin{pmatrix} d_k & d'_k \end{pmatrix} \left( \begin{pmatrix} d_k & d'_k \end{pmatrix}^\top A \begin{pmatrix} d_k & d'_k \end{pmatrix} \right)^{-1}. \quad (19)$$

Notice that the two descent directions  $d_k, d'_k$  have to be linearly independent for the matrix in (19) to be invertible. Similarly to CG algorithm, we have the following *four* useful properties

$$r_{k+1}, r'_{k+1} \perp d_k, d'_k. \quad (20)$$

• **Orthogonalization step.** The second step consists, given the residues  $r_{k+1}, r'_{k+1}$  and the current descent directions  $d_k, d'_k$ , in determining the next descent directions  $d_{k+1}, d'_{k+1}$ . The latter should be  $A$ -orthogonal to all the previous descent directions. In fact, it will be sufficient to ensure  $A$ -orthogonality to  $d_k, d'_k$ , as for CG. One takes

$$\begin{pmatrix} d_{k+1} & d'_{k+1} \end{pmatrix} = \begin{pmatrix} r_{k+1} & r'_{k+1} \end{pmatrix} + \begin{pmatrix} d_k & d'_k \end{pmatrix} \beta_k^\top. \quad (21)$$

The matrix  $\beta_k \in \mathbb{R}^{2 \times 2}$  is chosen to ensure the *four* conditions

$$d_{k+1}, d'_{k+1} \perp_A d_k, d'_k.$$

This also leads to two independent problems for the two vectors  $d_{k+1}, d'_{k+1}$ : writing now

$$\beta_j \doteq (\beta_{j1} \quad \beta_{j2}), \quad j = 1, 2, \quad (22)$$

the previous orthogonality conditions can be written as:

$$0 = \begin{pmatrix} r_k + \begin{pmatrix} d_k & d'_k \end{pmatrix} \beta_1^\top \end{pmatrix}^\top A \begin{pmatrix} d_k & d'_k \end{pmatrix} = \begin{pmatrix} r'_k + \begin{pmatrix} d_k & d'_k \end{pmatrix} \beta_2^\top \end{pmatrix}^\top A \begin{pmatrix} d_k & d'_k \end{pmatrix} \quad (23)$$

which yields the unique solution

$$\beta_k = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = - \begin{pmatrix} r_k & r'_k \end{pmatrix}^\top A \begin{pmatrix} d_k & d'_k \end{pmatrix} \left( \begin{pmatrix} d_k & d'_k \end{pmatrix}^\top A \begin{pmatrix} d_k & d'_k \end{pmatrix} \right)^{-1} \quad (24)$$

• **Summary of cCG in the case  $p = 2$ .** Putting together the previous findings, we summarize the cCG algorithm in the case  $p = 2$  as

$$(r_k \ r'_k) = A (x_k \ x'_k) - \mathbf{1}^\top b \quad (25a)$$

$$(x_{k+1} \ x'_{k+1}) = (x_k \ x'_k) + (r_k \ r'_k) \alpha_k^\top \quad (25b)$$

$$\alpha_k = - (r_k \ r'_k)^\top (d_k \ d'_k) ((d_k \ d'_k)^\top A (d_k \ d'_k))^{-1} \quad (25c)$$

$$(d_{k+1} \ d'_{k+1}) = (r_{k+1} \ r'_{k+1}) + (d_k \ d'_k) \beta_k^\top \quad (25d)$$

$$\beta_k = - (r_{k+1} \ r'_{k+1})^\top A (d_k \ d'_k) ((d_k \ d'_k)^\top A (d_k \ d'_k))^{-1} \quad (25e)$$

### 3.2 Cooperative CG algorithm: the general case

We now provide generalization to the case of  $p \geq 2$  agents. The extension is indeed straightforward from (25). The matrices whose  $j$ -th column represents respectively the solution estimate, the residue and the descent direction of agent  $j$ ,  $j = 1, \dots, p$ , for iteration  $k$  are denoted  $X_k \in \mathbb{R}^{n \times p}$ ,  $R_k \in \mathbb{R}^{n \times p}$ ,  $D_k \in \mathbb{R}^{n \times p}$ . In other words,  $X_k, R_k, D_k$  stand for the matrices written down  $(x_k \ x'_k), (r_k \ r'_k), (d_k \ d'_k)$  in Section 3.1.

The algorithm cCG in full generality is given as the list of instructions in Algorithm 2. Algorithm cCG is a generalization of CG, which is the case  $p = 1$ . In all algorithms in this paper, comments appear to the right of the symbol  $\triangleright$ .

---

#### Algorithm 2 cooperative Conjugate Gradient (cCG) algorithm

---

```

1: choose  $X_0 \in \mathbb{R}^{n \times p}$ 
2:  $R_0 := AX_0 - \mathbf{1}_p^\top b$ 
3:  $D_0 := R_0$   $\triangleright$  Generically, rank  $D_0 = p$ 
4:  $k := 0$ 
5: while  $D_k$  is full rank do
6:    $\alpha_k := -R_k^\top D_k (D_k^\top A D_k)^{-1}$   $\triangleright \alpha_k \in \mathbb{R}^{p \times p}$ 
7:    $X_{k+1} := X_k + D_k \alpha_k^\top$   $\triangleright X_{k+1} \in \mathbb{R}^{p \times p}$ 
8:    $R_{k+1} := AX_{k+1} - \mathbf{1}_{pk}^\top b$   $\triangleright R_{k+1} \in \mathbb{R}^{p \times p}$ 
9:    $\beta_k := -R_{k+1}^\top A D_k (D_k^\top A D_k)^{-1}$   $\triangleright \beta_k \in \mathbb{R}^{p \times p}$ 
10:   $D_{k+1} := R_{k+1} + D_k \beta_k^\top$   $\triangleright D_{k+1} \in \mathbb{R}^{p \times p}$ 
11:   $k \leftarrow k + 1$ 
12: end while

```

---

**Theorem 2** (Properties of cCG). *As long as the matrix  $D_k$  is full rank (that is rank  $D_k = p$ )*

- the matrices  $\{R_i\}_0^k$  are mutually orthogonal, the matrices  $\{D_i\}_0^k$  are mutually  $A$ -orthogonal, and the subspaces  $\text{span } [R_i]_0^k, \text{span } [D_i]_0^k$  and  $\text{span } [A^i R_0]_0^k$  are equal and have dimension  $(k+1)p$ ;
- for any vector  $e_j$  of the canonical basis of  $\mathbb{R}^p$ , the vector  $X_{k+1} e_j \in \mathbb{R}^n$  (which constitutes the  $j$ -th column of  $X_{k+1}$ ) is the minimizer of  $f$  on the affine subspace  $X_0 e_j + \text{span } [D_i]_0^k$ .

Theorem 2 indicates that, as long as the residue vector  $R_k$  is full rank, the algorithm cCG behaves essentially as does CG, providing  $p$  different estimates at iteration  $k$ , each of them being optimal in an affine subspace constructed from one of the  $p$  initial conditions and the common vector space obtained from the columns of the direction matrices  $D_i$ ,  $i = 0, \dots, k-1$ . This vector space,  $\text{span } [D_i]_0^k$ , has dimension  $(k+1)p$ : each iteration involves the cancellation of  $p$  directions. Notice that different columns of the matrices  $D_k$  are not necessarily  $A$ -orthogonal (in other words,  $D_k^\top A D_k$  is not necessarily diagonal), but, when  $R_k$  is full rank, they constitute a set of  $p$  independent vectors. The statement as well as the proof of this theorem are inspired by the corresponding ones for the conventional CG algorithm given in [8, p. 270ff] and [5, p. 390-391].

*Proof of Theorem 2.*

- We first show that for any  $k$ ,

$$\text{span } [R_i]_0^k = \text{span } [D_i]_0^k. \quad (26)$$

• We show the first point by induction. Clearly,  $D_j \perp_A D_k$  for any  $j < k$ , and  $\text{span } [R_i]_0^k = \text{span } [D_i]_0^k = \text{span } [A^i R_0]_0^k$  when  $k = 0$ , while nothing has to be verified for the orthogonality conditions. Assume that, for some  $k$ , they are verified for any  $i \leq k$ , let us show them for  $k + 1$ , assuming that  $R_{k+1}$  is full rank.

From lines 7–8 of Algorithm 2,  $R_{k+1} = AX_{k+1} - \mathbf{1}_p^\top b = R_k + AD_k \alpha_k^\top$ . By induction, the columns of both the matrices  $R_k$  and  $AD_k$  are located in  $\text{span } [A^i R_0]_0^{k+1}$ . Thus,

$$R_{k+1} \in \text{span } [A^i R_0]_0^{k+1}$$

and consequently

$$\text{span } [R_i]_0^{k+1} \subset \text{span } [A^i R_0]_0^{k+1} .$$

On the other hand, for any vector  $e_j$  of the canonical basis of  $\mathbb{R}^p$ ,

$$R_{k+1} e_j \notin \text{span } [D_i]_0^k \tag{27}$$

because each residue is orthogonal to the previous descent directions, so that  $R_{k+1} e_j \in \text{span } [D_i]_0^k$  for some  $e_j$  would imply  $R_{k+1} e = 0$ , which contradicts the assumption of full rankness of  $R_{k+1}$ . Indeed, for the same reason, one can also show that, for any  $v \in \mathbb{R}^p \setminus \{0\}$ ,  $R_{k+1} v \notin \text{span } [D_i]_0^k$ . Using again the fact that  $\text{span } [R_{k+1}] = p$ , one sees that

$$\text{span } [R_i]_0^{k+1} \supset \text{span } [A^i R_0]_0^{k+1}$$

and indeed

$$\text{span } [R_i]_0^{k+1} = \text{span } [A^i R_0]_0^{k+1} .$$

One shows similarly from lines (9)–(10) of Algorithm 2 that

$$\text{span } [D_i]_0^{k+1} \subset \text{span } [A^i R_0]_0^{k+1} ,$$

and the equality is obtained using the same rank argument.

The dimension of these sets is  $(k + 2)p$ , as they contain the  $p$  independent vectors in  $\text{span } [D_{k+1}]$  orthogonal to  $\text{span } [D_i]_0^k$ .

From line 10 of Algorithm 2, one gets

$$D_i^\top AD_{k+1} = D_i^\top A (R_{k+1} + D_k \beta_k^\top) \tag{28}$$

For  $i < k$ , the first term is zero because  $AD_i \in \text{span } [D_j]_0^{i+1}$  and the gradients constituting the columns of  $R_{k+1}$  are orthogonal to any vector in  $\text{span } [D_j]_0^{i+1}$ ; while the second term is also zero due to the induction hypothesis. For  $i = k$ , the right-hand side of (28) is zero because  $\beta_k$  is precisely chosen to ensure this property. Thus the  $\{D_i\}_0^{k+1}$  are mutually  $A$ -orthogonal.

Orthogonality of  $R_{k+1}$  follows from lines 6 and 8 of Algorithm 2. The induction hypothesis has been proved for  $k + 1$  concluding the proof of the first part of Theorem 2.

• [Optimality]. Arguing as in [5, p.390-391], we can write

$$X_{k+1} e_j = X_0 e_j + \sum_{i=0}^k D_i \gamma_i^\top \tag{29}$$

for some  $\gamma_i \in \mathbb{R}^{1 \times p}$ . Optimality implies that

$$D_i^\top (AX_{k+1} e_j - b) = 0, \tag{30}$$

Substituting (29) in (30) and rewriting in terms of the matrices  $R_i$  and  $\gamma_i$  yields

$$\gamma_i = -e_j^\top R_0^\top D_i (D_i^\top AD_i)^{-1} \tag{31}$$

On the other hand,

$$\nabla f(X_{k+1} e_j) = A(X_0 e_j + \sum_{i=0}^k D_i \gamma_i^\top) - b = R_0 e_j + \sum_{i=0}^k AD_i \gamma_i^\top \tag{32}$$

Thus

$$D_{k+1}^\top \nabla f(X_{k+1}e_j) = D_{k+1}^\top R_0 e_j \quad (33)$$

Substituting (33) in (31), we get:

$$\gamma_{k+1} = -\nabla f(X_{k+1}e_j)^\top (D_{k+1}^\top A D_{k+1})^{-1} = e_j^\top \alpha_{k+1} \quad (34)$$

□

A natural question is now to study the cases where at some point of the execution of the algorithm cCG one gets rank  $R_k < p$ . In the best case, this occurs because one of the columns of  $R_k$  is null, say the  $j$ -th one, meaning that  $\nabla f(X_k e_j) = 0$ , and thus that the  $k$ -th estimate of the  $j$ -th agent is equal to the optimum  $x^* = A^{-1}b$ . But, of course, rank  $R_k$  can be smaller than  $p$  without any column of  $R_k$  being null.

First of all, the following result ensures that this rank degeneracy is, in general, avoided during algorithm execution.

**Theorem 3** (Genericity of the full rank condition of cCG residues matrix). *For an open dense set of initial conditions  $X_0$  in  $\mathbb{R}^{n \times p}$ , one has during any cCG run*

$$\forall 0 \leq k \leq k^* \doteq \lfloor \frac{n}{p} \rfloor, \quad \text{rank } R_k = p. \quad (35)$$

Moreover

$$\dim \text{span } [D_i]_0^{k^*} = p \lfloor \frac{n}{p} \rfloor. \quad (36)$$

Otherwise said: generically, algorithm cCG can be run during  $k^*$  steps, and

- any of the columns of  $X_{k^*}$  minimizes  $f$  on an affine subspace of  $\mathbb{R}^n$  of codimension  $p \lfloor \frac{n}{p} \rfloor$ ;
- application of CG departing from any of the columns of  $X_{k^*}$  yields convergence in at most  $n - p \lfloor \frac{n}{p} \rfloor \leq p - 1$  steps.

The second part of Theorem 3 has to be interpreted as follows. When the size  $n$  of the matrix  $A$  is a multiple of the number  $p$  of agents, then cCG generically ends up in  $\frac{n}{p}$  steps. When this is not the case, the estimates  $X_{k^*}$  obtained for  $k = k^*$  minimize the function  $f$  on an affine subspace whose underlying vector subspace is  $\text{span } [D_i]_0^{k^*}$  (see Theorem 2). The interest of (36) is to show that this subspace is quite large: its codimension is  $p \lfloor \frac{n}{p} \rfloor$ , which is at most equal to  $p - 1$ .

*Proof of Theorem 3.* The main point consists in showing that generically,  $k^*$  iterations of the algorithm cCG can be conducted without occurrence of the rank deficiency condition. As a matter of fact, the other results of the statement are direct consequences of this fact.

To show the latter, use is made of Theorem 2. From the properties stated therein, one sees that, for any  $0 \leq k \leq k^*$ , the rank of  $X_k$  is deficient if and only if a linear combination of the  $p$  column vectors of  $X_0$  pertains to the  $kp$ -dimensional subspace  $\text{span } [D_i]_0^{k-1}$ . In a vector space of dimension  $n > kp$ , this occurs only in the complement of an open dense set.

Now, if the column vectors of  $X_k$  are linearly independent, the same is true for  $R_k$ , see line 8 of Algorithm 2, and as well for  $D_k$ , see line 10. This completes the proof of Theorem 3. □

We now study what can be done in case of rank degeneracy. When  $p_k \doteq \text{rank } D_k$  is such that  $0 < p_k < p$ , this means that trajectories initially independent have come to a point where the estimates in  $X_k$  will converge along directions which are now linearly dependent. The natural solution is then to choose any full-rank subset of trajectories. We thus propose the modified algorithm 3.

**Theorem 4** (Convergence of mcCG algorithm). *For any nonzero initial condition  $X_0$ , algorithm mcCG ends up in  $k^{**}$  iterations for some  $k^{**} \leq n$ . Moreover*

- the sequence  $(p_k)_{0 \leq k \leq k^{**}}$  is nonincreasing;
- any of the columns of  $X_{k^{**}}$  minimizes  $f$  on an affine subspace of  $\mathbb{R}^n$  of codimension  $p_{k^{**}} \lfloor \frac{n}{p_{k^{**}}} \rfloor$ ;
- application of CG departing from any of the columns of  $X_{k^{**}}$  yields convergence in at most  $n - p_{k^{**}} \lfloor \frac{n}{p_{k^{**}}} \rfloor \leq p_{k^{**}} - 1$  steps.

The proof is straightforward and omitted for brevity.

---

**Algorithm 3** modified cooperative Conjugate Gradient (mCCG) algorithm
 

---

```

1: choose  $X_0 \in \mathbb{R}^{n \times p}$ 
2:  $R_0 := AX_0 - \mathbf{1}_p^\top b$ 
3:  $D_0 := R_0$ 
4:  $p_0 := \text{rank } D_0$  ▷  $p_0$  is the initial value of the rank
5:  $k := 0$ 
6: if  $p_0 = p$  then
7:   go to 14
8: else
9:   choose  $J \in \{1, \dots, p\}$  such that  $\text{rank } R_0|_{j \in J} = p_0$ 
10:   $X_0 \leftarrow X_0|_{j \in J}$  ▷  $X_0 \in \mathbb{R}^{p_0 \times p_0}$ 
11:   $R_0 \leftarrow R_0|_{j \in J}$  ▷  $R_0 \in \mathbb{R}^{p_0 \times p_0}$  is full rank
12:   $D_0 \leftarrow D_0|_{j \in J}$  ▷  $D_0 \in \mathbb{R}^{p_0 \times p_0}$  is full rank
13: end if
14: while  $p_k > 0$  do
15:   $\alpha_k := -R_k^\top D_k (D_k^\top A D_k)^{-1}$  ▷  $\alpha_k \in \mathbb{R}^{p_k \times p_k}$ 
16:   $X_{k+1} := X_k + D_k \alpha_k^\top$  ▷  $X_{k+1} \in \mathbb{R}^{p_k \times p_k}$ 
17:   $R_{k+1} := AX_{k+1} - \mathbf{1}_{p_k}^\top b$  ▷  $R_{k+1} \in \mathbb{R}^{p_k \times p_k}$ 
18:   $\beta_k := -R_{k+1}^\top A D_k (D_k^\top A D_k)^{-1}$  ▷  $\beta_k \in \mathbb{R}^{p_k \times p_k}$ 
19:   $D_{k+1} := R_{k+1} + D_k \beta_k^\top$  ▷  $D_{k+1} \in \mathbb{R}^{p_k \times p_k}$ 
20:   $p_{k+1} := \text{rank } D_{k+1}$ 
21:  if  $p_{k+1} = p_k$  then ▷ If  $p_{k+1} = p_k$ , cCG goes on normally
22:    go to 29
23:  else ▷ If  $p_{k+1} < p_k$ ,  $p_k - p_{k+1}$  agents are suppressed
24:    choose  $J \in \{1, \dots, p_k\}$  such that  $\text{rank } R_{k+1}|_{j \in J} = p_{k+1}$ 
25:     $X_{k+1} \leftarrow X_{k+1}|_{j \in J}$  ▷  $X_{k+1} \in \mathbb{R}^{p_{k+1} \times p_{k+1}}$ 
26:     $R_{k+1} \leftarrow R_{k+1}|_{j \in J}$  ▷  $R_{k+1} \in \mathbb{R}^{p_{k+1} \times p_{k+1}}$  is full rank
27:     $D_k \leftarrow D_k|_{j \in J}$  ▷  $D_k \in \mathbb{R}^{p_{k+1} \times p_{k+1}}$  is full rank
28:  end if
29:   $k \leftarrow k + 1$ 
30: end while

```

---

## 4 Computational complexity

This section is concerned with the evaluation of the gain in computation time of the numerical solution of equation (6), when using cCG algorithm with  $p$  agent, i.e., the gain which is expected is due to the parallelism induced by a multithread implementation. We evaluate this issue here *assuming computations in exact arithmetic*. Moreover, thanks to Theorem 3, we adopt the generic assumption that the rank of the residue matrices remains constant (and full), and that the computations are then carried out for  $\lfloor \frac{n}{p} \rfloor$  iterations. Disregarding as marginal the supplementary CG steps (see the statement of Theorem 3), we thus consider it to be realistic to quantify the worst case complexity by evaluating *the numbers of multiplications involved by  $\frac{n}{p}$  iterations of cCG*. Recall that the case  $p = 1$  corresponds to the usual CG algorithm.

We propose the multithread implementation detailed in Table 1.

Table 1: Number of scalar multiplications in  $k$ -th iteration

Operation carried out by $i$ -th processor	Composite result	Dimension of the result	Number of scalar multiplications carried out by the $i$ th processor
$AD_{k,i}$	$AD_k$	$n \times p$	$n^2$
$D_{k,i}^\top AD_k$	$D_k^\top AD_k$	$p \times p$	$np$
$R_{k,i}^\top D_k$	$R_k^\top D_k$	$p \times p$	$np$
$\alpha_{k,i}$ s.t. $\alpha_{k,i}(D_k^\top AD_k) = R_{k,i}^\top D_k$	$\alpha_k = R_k^\top D_k (D_k^\top AD_k)^{-1}$	$p \times p$	$\frac{p(p+1)(2p+1)}{6}$
$R_{k+1,i} = R_{k,i} - AD_k \alpha_{k,i}^\top$	$R_{k+1} = R_k - AD_k \alpha_k^\top$	$n \times p$	$np$
$X_{k+1,i} = X_{k,i} - D_k \alpha_{k,i}^\top$	$X_{k+1} = X_k - D_k \alpha_k^\top$	$n \times p$	$np$
$R_{k+1,i}^\top AD_k$	$R_{k+1}^\top AD_k$	$p \times p$	$np$
$\beta_{k,i}$ s.t. $\beta_{k,i}(D_k^\top AD_k) = -R_{k+1,i}^\top AD_k$	$\beta_k = -R_{k+1}^\top AD_k (D_k^\top AD_k)^{-1}$	$p \times p$	$\frac{p(p+1)(2p+1)}{6}$
$D_{k+1,i} = R_{k+1,i} + D_k \beta_{k,i}^\top$	$D_{k+1} = R_{k+1} + D_k \beta_k^\top$	$n \times p$	$np$
Total number of scalar multiplications per processor and per iteration			$n^2 + 6np + \frac{p(p+1)(2p+1)}{3}$

In Table 1, the first column indicates the task carried out at each stage by every processor, and the last column the corresponding number of multiplications carried out by a processor. The double lines, separating the first row from the second and the second from the third, indicate the necessity of a phase of information exchange: every processor at that stage needs to know results from other processors, also called a synchronization barrier in computing terminology. The second column, labelled composite result, contains the information that is available by pooling the partial results from each processor and the third column gives the dimension of this composite result. The fourth and final column contains the number of multiplications carried out by the  $i$ th processor. The number  $\frac{p(p+1)(2p+1)}{6}$  of scalar multiplications is needed to realize Gaussian elimination realized through LU factorization [13, p. 15].

As indicated by the last line of Table 1, a total of  $n^2 + 6np + \frac{p(p+1)(2p+1)}{3}$  multiplications per processor is needed to complete an iteration. Since, generically speaking, the algorithm ends in at most  $\frac{n}{p}$  iterations (see Theorem 3), an estimate of the worst-case multithread execution time is given by the following result.

**Theorem 5** (Worst-case multithread execution time in exact arithmetic). *Generically, multithread execution of cCG using  $p$  agents for a linear system (6) of size  $n$  requires*

$$N(p) = \frac{n^3}{p} + 6n^2 + n \frac{(p+1)(2p+1)}{3} \quad (37)$$

*multiplications performed synchronously in parallel by each processor.*

This result has straightforward consequences.

**Corollary 6** (Multithread gain). *For problems of size  $n$  at least equal to 5, it is always beneficial to use  $p \leq n$  processors rather than a single one. In other words, when  $n \geq 5$ ,*

$$\forall 1 \leq p \leq n, \quad N(1) \geq N(p). \quad (38)$$

*Proof.* One has

$$N(1) - N(p) = n^3 + \frac{5}{3}n - \left(2n^2 + \frac{2}{3}n^3\right) = \frac{1}{3}(n^3 - 6n^2 + 5n) = \frac{1}{3}n(n-1)(n-5) .$$

Moreover,

$$\frac{dN(1)}{dp} = n \left( \frac{7}{3} - n^2 \right)$$

which is negative for  $n \geq 2$ , while

$$\frac{dN(n)}{dp} = n \left( \frac{4}{3}n^2 \right) \geq 0 .$$

The convexity of  $N$  then yields the conclusion that  $N(p) \leq N(1)$  for any  $1 \leq p \leq n$ .  $\square$

**Corollary 7** (Optimal multithread gain). *For any size  $n$  of the problem, there exists a unique optimal number  $p^*$  of processors minimizing  $N(p)$ . Moreover, when  $n \rightarrow +\infty$ ,*

$$p^* \approx \left( \frac{3}{4} \right)^{\frac{1}{3}} n^{\frac{2}{3}} \tag{39a}$$

$$N(p^*) \approx \left( \left( \frac{4}{3} \right)^{\frac{1}{3}} + \frac{2}{3} \left( \frac{3}{4} \right)^{\frac{2}{3}} \right) n^{2+\frac{1}{3}} \approx 1.651n^{2+\frac{1}{3}} \tag{39b}$$

*Proof.* One has

$$\frac{dN(p)}{dp} = -\frac{n^3}{p^2} + \frac{4}{3}np + n .$$

There exists a unique  $p^*$  canceling this expression. For this value, one has  $n^2 = p^2(\frac{4}{3}p + 1)$ , which yields the asymptotic behavior given in (39a). The value in (39b) is directly deduced.  $\square$

The conclusion of Corollary 7 is quite important. It shows that solution of  $Ax = b$  is possible by the method proposed here with a cost of  $O(n^{2+\frac{1}{3}})$  multiplications. This is to be compared with the classical results [14].

## 5 Numerical experiments with discussion of multithread implementation

This section reports on a suite of numerical experiments carried out on a set of random symmetric matrices of dimensions varying from 1000 to 25000, the latter being the largest dimension that could be accommodated in the fast access RAM memory of the multicore processor. The random symmetric matrices were generated by choosing random diagonal matrices  $\Lambda$ , with positive diagonal entries uniformly distributed between 1 and a prespecified condition number, which were then pre-multiplied (resp. post-multiplied) by a random orthogonal matrix  $U$  (resp. its transpose  $U^T$ ). The random orthogonal matrices  $U$  were generated using a C translation of Shilon's MATLAB code [12], which produces a matrix distribution uniform over the manifold of orthogonal matrices with respect to the induced  $\mathbb{R}^{n^2}$  Lebesgue measure. The right hand sides and initial conditions were also randomly generated, with all entries uniformly distributed on the interval  $[-10, 10]$ . In this preliminary work, the matrices used were dense and the use of preconditioners was not investigated.

In order to evaluate the performance of the algorithm proposed in this paper, a program was written in language C. The compiler used was the GNU Compiler Collection (GCC), running under Linux Ubuntu 10.0.4. For the Linear Algebra calculations, we used the Linear Algebra Package (LAPACK) and the Basic Linear Algebra Subprograms (BLAS). Finally, to parallelize the program, we used the Open Multi Processing (OMP) API. The processor used was an Intel Core2Quad CPU Q8200 running at 2.33 MHz with four cores.

The pseudo-code in Algorithm 4 gives details of the implementation for three ( $p = 3$ ) agents.

---

**Algorithm 4** Implementation of cooperative Conjugate Gradient (cCG) algorithm

---

```
1: choose  $X_0, Y_0, Z_0 \in \mathbb{R}^n$     ▷ All initialized randomly with numbers between  $-10$  and  $10$ 
2:  $r_{0,x} := A \cdot x_0 - b$ 
3:  $r_{0,y} := A \cdot y_0 - b$ 
4:  $r_{0,z} := A \cdot z_0 - b$ 
5:  $d_{0,x} := r_{0,x}$ 
6:  $d_{0,y} := r_{0,y}$ 
7:  $d_{0,z} := r_{0,z}$ 
8:  $k := 0$ 
9:  $minres := \min(\text{norm}(r_{0,x}), \text{norm}(r_{0,y}), \text{norm}(r_{0,z}))$ 
10: while  $minres > tolerance$  do
11:                                     ▷ Compute matrix-vector products  $A \cdot d_{k,i}$ 
12:   agent 1: compute  $A \cdot d_{k,x}$ 
13:   agent 2: compute  $A \cdot d_{k,y}$ 
14:   agent 3: compute  $A \cdot d_{k,z}$ 
15:   Barrier ▷ Synchronizes all 3 agents, before proceeding to the next computations
16:                                     ▷ Compute  $m_{ij}$ 
17:   agent 1:  $m_{11} := d_{k,x}^T \cdot A \cdot d_{k,x}; m_{12} := d_{k,x}^T \cdot A \cdot d_{k,y}$ 
18:   agent 2:  $m_{13} := d_{k,x}^T \cdot A \cdot d_{k,z}; m_{22} := d_{k,y}^T \cdot A \cdot d_{k,y}$ 
19:   agent 3:  $m_{23} := d_{k,y}^T \cdot A \cdot d_{k,z}; m_{33} := d_{k,z}^T \cdot A \cdot d_{k,z}$ 
20:   Barrier ▷ Synchronizes all 3 agents, before proceeding to the next computations
21:   Initialize  $M := \{m_{ij}\}$     ▷ Symmetric matrix needed to compute alpha,  $m_{ij} = m_{ji}$ 
22:                                     ▷ Right-hand sides needed to compute alpha
23:   agent 1:  $n_1 := [r_{k,x}^T \cdot d_{k,x}; r_{k,x}^T \cdot d_{k,y}; r_{k,x}^T \cdot d_{k,z}]$ 
24:   agent 2:  $n_2 := [r_{k,y}^T \cdot d_{k,x}; r_{k,y}^T \cdot d_{k,y}; r_{k,y}^T \cdot d_{k,z}]$ 
25:   agent 3:  $n_3 := [r_{k,z}^T \cdot d_{k,x}; r_{k,z}^T \cdot d_{k,y}; r_{k,z}^T \cdot d_{k,z}]$ 
26:                                     ▷ Computation of alpha
27:   agent 1: Solve  $M \cdot \alpha_1 = n_1$ 
28:   agent 2: Solve  $M \cdot \alpha_2 = n_2$ 
29:   agent 3: Solve  $M \cdot \alpha_3 = n_3$ 
30:                                     ▷ Update estimates of each agent
31:   agent 1:  $x_k \leftarrow x_k + \alpha_{1,1} \cdot d_{k,x} + \alpha_{1,2} \cdot d_{k,y} + \alpha_{1,3} \cdot d_{k,z}$ 
32:   agent 2:  $y_k \leftarrow y_k + \alpha_{2,1} \cdot d_{k,x} + \alpha_{2,2} \cdot d_{k,y} + \alpha_{2,3} \cdot d_{k,z}$ 
33:   agent 3:  $z_k \leftarrow z_k + \alpha_{3,1} \cdot d_{k,x} + \alpha_{3,2} \cdot d_{k,y} + \alpha_{3,3} \cdot d_{k,z}$ 
34:                                     ▷ Update residues of each agent
35:   agent 1:  $r_{k,x} := A \cdot x_k - b$ 
36:   agent 2:  $r_{k,y} := A \cdot y_k - b$ 
37:   agent 3:  $r_{k,z} := A \cdot z_k - b$ 
38:                                     ▷ Right-hand sides needed to compute beta
39:   agent 1:  $n_1 := [r_{k,x}^T \cdot A \cdot d_{k,x}; r_{k,x}^T \cdot A \cdot d_{k,y}; r_{k,x}^T \cdot A \cdot d_{k,z}]$ 
40:   agent 2:  $n_2 := [r_{k,y}^T \cdot A \cdot d_{k,x}; r_{k,y}^T \cdot A \cdot d_{k,y}; r_{k,y}^T \cdot A \cdot d_{k,z}]$ 
41:   agent 3:  $n_3 := [r_{k,z}^T \cdot A \cdot d_{k,x}; r_{k,z}^T \cdot A \cdot d_{k,y}; r_{k,z}^T \cdot A \cdot d_{k,z}]$ 
42:                                     ▷ Computation of beta
43:   agent 1: Solve  $M \cdot \beta_1 = n_1$ 
44:   agent 2: Solve  $M \cdot \beta_2 = n_2$ 
45:   agent 3: Solve  $M \cdot \beta_3 = n_3$ 
46:                                     ▷ Update of directions
47:   agent 1:  $d_{k,x} \leftarrow r_{k,x} + \beta_{1,1} \cdot d_{k,x} + \beta_{1,2} \cdot d_{k,y} + \beta_{1,3} \cdot d_{k,z}$ 
48:   agent 2:  $d_{k,y} \leftarrow r_{k,y} + \beta_{2,1} \cdot d_{k,x} + \beta_{2,2} \cdot d_{k,y} + \beta_{2,3} \cdot d_{k,z}$ 
49:   agent 3:  $d_{k,z} \leftarrow r_{k,z} + \beta_{3,1} \cdot d_{k,x} + \beta_{3,2} \cdot d_{k,y} + \beta_{3,3} \cdot d_{k,z}$ 
50:                                     ▷ Calculate of residual norms
51:   agent 1:  $norm_{r_x} = \text{norm}(r_{k,x})$ 
52:   agent 2:  $norm_{r_y} = \text{norm}(r_{k,y})$ 
53:   agent 3:  $norm_{r_z} = \text{norm}(r_{k,z})$ 
54:    $minres := \min(\text{norm}_{r_x}, \text{norm}_{r_y}, \text{norm}_{r_z})$ 
55:    $k \leftarrow k + 1$ 
56: end while
```

---

## 5.1 Evaluating speedup

The results of the Cooperative 3 agent cCG, in comparison with classic CG, with a tolerance of  $10^{-3}$ , and matrices with different sizes, but all with the same condition number of  $10^6$ , are shown in Figure 1. Multiple tests were performed, using different randomly generated initial conditions (20 different initial conditions for the small matrices and 10 for the bigger ones). Figure 1 shows the mean values computed for these tests.

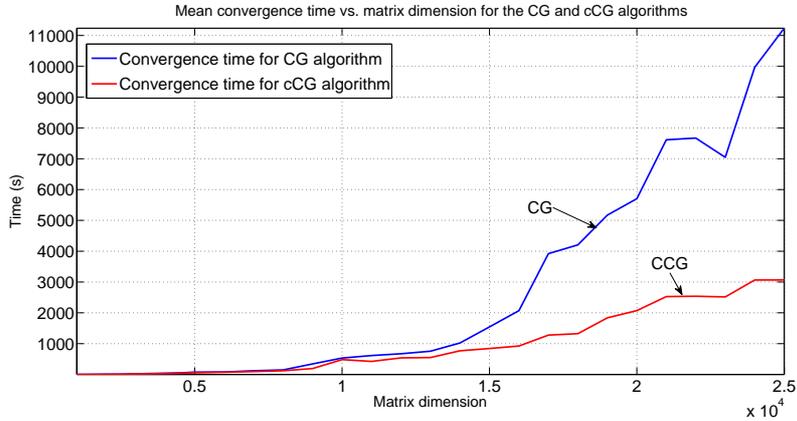


Figure 1: Mean time to convergence for random test matrices of dimensions varying from 1000 to 25000, for 3 agent cCG and standard CG algorithms.

The *iteration speedup* of cCG in comparison with CG is defined as the number of iterations that CG took to converge divided by the number of iterations cCG took to converge and the experimental results are shown in Figure 2, which also shows the classical speed-up, which is the ratio of the time to convergence, i.e., the time taken to run the main loop until convergence, for CG versus cCG.

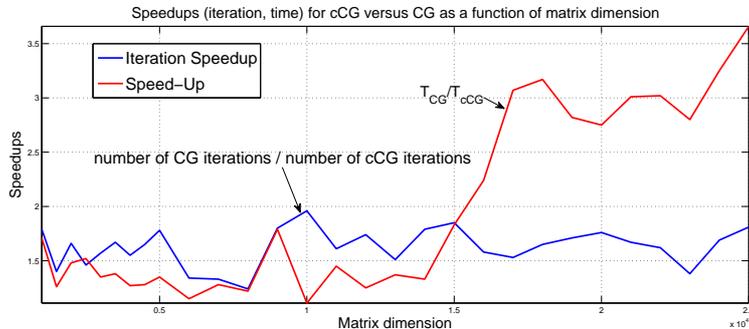


Figure 2: Average speedups of Cooperative 3 agent cCG over classic CG for random test matrices of dimensions varying from 1000 to 25000.

The speedups seem to be roughly equal up to a certain size of matrix ( $n = 16000$ ); however, above this dimension, there is an increasing trend for both speedups.

The numerical results obtained show that cCG, using 3 agents, leads to an improvement in comparison with the usual CG algorithm. The average iteration speedup and the classical speedup of cCG are respectively, 1.62 and 1.94, indicating that cCG converges almost twice as fast as CG for dense matrices with reasonably well-separated eigenvalues.

## 5.2 Verifying the complexity estimates

Figure 3 shows the mean time spent per iteration in seconds (points plotted as squares), versus matrix dimension, as well as the parabola fitted to this data, using least squares. Using the result from the last

row of table 1 and multiplying it by the mean time per scalar multiplication, we obtain the parabola (dash-dotted line in Figure 3) expected in theory. In order to estimate the time per scalar multiplication, we divided the experimentally obtained mean total time spent on each iteration and divided it by the number of scalar multiplications performed in each iteration. This was done for each matrix dimension. Since the same multicore processor is being used for all experiments, each of these divisions should generate the same value of time taken to carry out each scalar multiplication, regardless of matrix dimension. It was observed that these divisions produced a data set which has a mean value of 8.10 nanoseconds per scalar multiplication, with a standard deviation of 1.01 nanoseconds, showing that the estimate is reasonable.

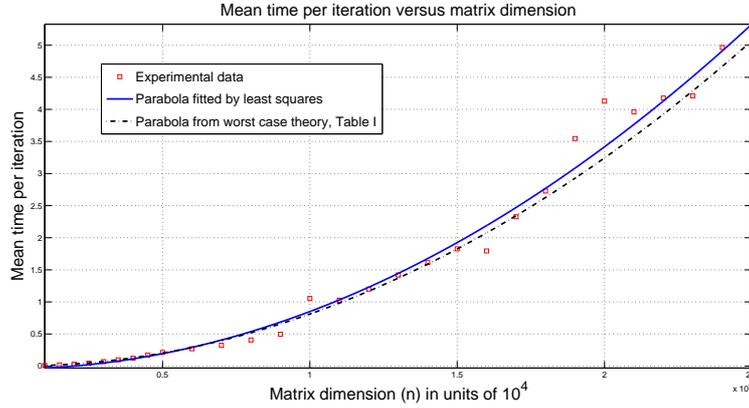


Figure 3: Mean time per iteration versus problem dimension

Now, from equation (37), substituting  $p = 3$ , neglecting small order terms, and multiplying it by the estimated mean time per scalar multiplication (8.10 nanoseconds), the number of matrix multiplications per iteration,  $N(p), p = 3$ , is a cubic polynomial in  $n$ . Thus, the logarithm of the dimension ( $n$ ) of the problem versus the logarithm of time needed to convergence is expected to be a straight line of slope 3. Figure 4 shows this straight line, fitted to the data (squares) by least squares.

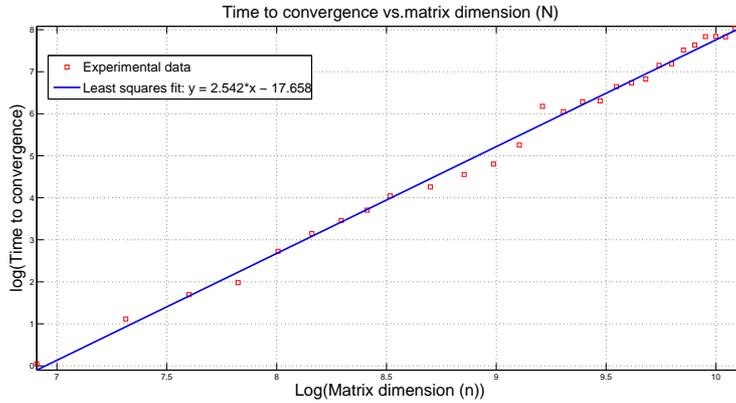


Figure 4: Log-log plot of mean time to convergence versus problem dimension

Its slope (2.542) is fairly close to 3, and data seems to follow a linear trend. The deviation of the slope from the ideal value has several probable causes, the first one being that the exact exponent of 3 is a result of a worst case analysis of CG in exact arithmetic. It is known that CG usually converges, to a reasonable tolerance, in much less than  $n$  iterations, where  $n$  is the matrix dimension [9].

Similarly, the logarithm of the number of iterations needed to convergence versus the logarithm of the dimension of the problem should also follow a linear trend. Since the number of iterations is expected to be  $n/3$ , the slope of this line should be 1. This log-log plot is shown in figure 5, in which the straight

line was fitted by least squares to the original data (red squares). The slope (0.501) of the fitted line is smaller than 1, but is seen to fit the data well (small residuals). The fact that both slopes are smaller than their expected values indicates that the **cCG** algorithm is converging faster than the worst case estimate. Another reason is that a fairly coarse tolerance of  $10^{-3}$  is used, and experiments reported show that decreasing the tolerance favors the **cCG** algorithm even more. Specifically, for a randomly generated matrix of dimension 8000 and condition number  $10^6$ , Table 2 shows the mean number of iterations and time to convergence, calculated for 10 different initial conditions, for the **CG** and **cCG** algorithms, as the tolerance is varied from  $10^{-3}$  to  $10^{-9}$

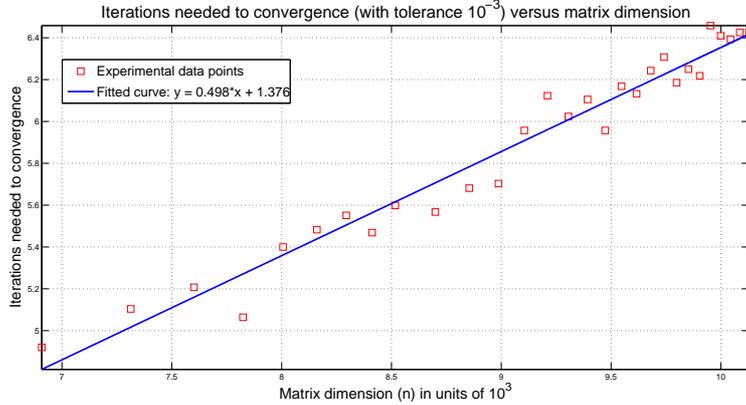


Figure 5: Iterations needed to convergence versus problem dimension

Table 2: Mean number of iterations and mean time to convergence for the **CG** and the **cCG** algorithms, as a function of tolerance used in the stopping criterion.

Tolerance	CG		cCG	
	Time (s)	Iterations	Time (s)	Iterations
$10^{-3}$	148.20	372.20	121.80	299.70
$10^{-4}$	165.20	414.00	125.70	319.40
$10^{-5}$	177.40	444.90	132.00	335.00
$10^{-6}$	192.70	476.80	134.00	352.70
$10^{-7}$	206.50	510.70	135.20	336.20
$10^{-8}$	235.10	538.40	137.10	373.50
$10^{-9}$	275.00	559.70	141.90	381.20

The data used to generate all the graphs in the figures above is shown in tables 3 and 4.

## 6 Concluding Remarks

This paper proposed a new cooperative conjugate gradient (**cCG**) method for linear systems with symmetric positive definite coefficient matrices. This **cCG** method permits efficient implementation on a multicore computer and experimental results bear out the main theoretical properties, namely, that speedups close to the theoretical value of  $p$ , when a  $p$ -core computer is used, are possible, when the matrix dimension is suitably large. The experimental results of the current study were limited to dense randomly generated matrices and only 3 cores of a 4 core computer with a relatively small on-chip shared memory were used. Future work will include the study of the method on matrices that come from real applications and are typically sparse and sometimes ill-conditioned (which will necessitate the use of preconditioners) on larger multi-core machines. The use of larger machines should also permit exploration of the notable theoretical result (corollary 7) that, in the asymptotic limit, as  $n$  becomes large, implying that  $p$  also increases according to (39a), solution of  $Ax = b$  is possible by the method proposed here with a cost of  $O(n^{2+\frac{1}{3}})$  multiplications.

Matrix Dimension	Number of iterations		Time (s)	
	CG	CCG	CG	CCG
1000	245.50	137.05	1.80	1.05
1500	230.15	146.65	3.85	3.05
2000	303.35	182.60	8.05	5.45
2500	231.60	158.20	11.05	7.25
3000	347.00	221.35	20.65	15.30
3500	402.15	240.50	32.15	23.35
4000	399.85	257.45	40.30	31.85
4500	391.85	237.10	51.95	40.70
5000	481.60	270.05	77.00	57.05
6000	351.90	261.70	81.10	70.60
7000	390.90	293.30	121.70	95.00
8000	372.20	299.70	148.20	121.80
9000	659.90	386.50	343.50	191.50
10000	894.70	456.20	532.60	480.80
11000	667.00	413.20	614.40	413.20
12000	780.00	448.20	673.00	537.40
13000	582.80	386.50	753.90	548.60
14000	853.20	477.30	1022.70	769.00
15000	852.40	460.60	1543.00	841.70
16000	813.60	514.40	2070.00	922.70
17000	842.20	548.90	3921.60	1277.20
18000	802.50	485.70	4204.70	1325.20
19000	884.50	518.00	5171.50	1836.70
20000	882.10	501.90	5703.70	2072.60
21000	1064.30	638.00	7614.70	2526.40
22000	7671.80	2537.60	985.10	607.60
23000	7045.60	2516.40	826.30	597.70
24000	9969.20	3065.20	1040.90	617.40
25000	1114.70	617.20	11237.40	3067.50

Table 3: Average results for multiple test matrices of dimensions varying from 1000 to 25000, for Cooperative 3 agent cCGand for classic CG.

## References

- [1] A. Bhaya, P.-A. Bliman, and F. Pazos. Cooperative parallel asynchronous computation of the solution of symmetric linear systems. In *Proc. of the 49th IEEE Conference on Decision and Control*, Atlanta, USA, December 2010.
- [2] A. Bhaya and E. Kaszkurewicz. *Control perspectives on numerical algorithms and matrix problems*. Advances in Control. SIAM, Philadelphia, 2006.
- [3] C. Brezinski and M. Redivo-Zaglia. Hybrid procedures for solving linear systems. *Numerische Mathematik*, 67(1):1–19, 1994.
- [4] A. Greenbaum. *Iterative methods for solving linear systems*. SIAM, Philadelphia, 1997.
- [5] O. Güler. *Foundations of Optimization*. Springer, New York, 2010.
- [6] M. R. Hestenes and E. Stiefel. Methods of conjugate gradients for solving linear systems. *Journal of Research of the National Bureau of Standards*, 49:409–436, 1952.
- [7] V. Kumar, N. Leonard, and A. S. Morse, editors. *2003 Block Island Workshop on Cooperative Control*, Lecture Notes in Control and Information Sciences, vol. 309. Springer, 2005.

Matrix Dimension	Iteration Ratio	Speed-Up
1000	1.79	1.71
1500	1.40	1.26
2000	1.66	1.48
2500	1.46	1.52
3000	1.57	1.35
3500	1.67	1.38
4000	1.55	1.27
4500	1.65	1.28
5000	1.78	1.35
6000	1.34	1.15
7000	1.33	1.28
8000	1.24	1.22
9000	1.80	1.79
10000	1.96	1.11
11000	1.61	1.45
12000	1.74	1.25
13000	1.51	1.37
14000	1.79	1.33
15000	1.85	1.83
16000	1.58	2.24
17000	1.53	3.07
18000	1.65	3.17
19000	1.71	2.82
20000	1.76	2.75
21000	1.67	3.01
22000	1.62	3.02
23000	1.38	2.80
24000	1.69	3.25
25000	1.81	3.66

Table 4: Average gains of Cooperative 3 agent cCG over classic CG for matrices of dimensions varying from 1000 to 25000.

- [8] D. G. Luenberger and Y. Ye. *Linear and nonlinear programming*. Springer, New York, third edition, 2008.
- [9] G. Meurant and Z. Strakoš. The Lanczos and conjugate gradient algorithms in finite precision arithmetic. *Acta Numerica*, 15:471–542, 2006.
- [10] R. M. Murray. Recent research in cooperative control of multi-vehicle systems. *J. Guidance, Control and Dynamics*, 129(5):571–583, September 2007.
- [11] A. Nedic and A. Ozdaglar. *Convex Optimization in Signal Processing and Communications*, chapter Cooperative Distributed Multi-agent Optimization, pages 340–386. Cambridge University Press, 2010.
- [12] O. Shilon. RandOrthMat.m: MATLAB code to generate a random  $n \times n$  orthogonal real matrix, 2006. <http://www.mathworks.com/matlabcentral/fileexchange/authors/23951>.
- [13] G. Strang. *Linear Algebra and its Applications*. Harcourt Brace Jovanovich, Inc., New York, 1988. 3rd ed.
- [14] J. F. Traub and H. Woźniakowski. On the optimal solution of large linear systems. *Journal of the Association for Computing Machinery*, 31(3):545–559, 1984.