

STABILITY ANALYSIS AND IMPROVEMENT OF THE BLOCK GRAM-SCHMIDT ALGORITHM*

W. JALBY† AND B. PHILIPPE†

Abstract. The advent of supercomputers with hierarchical memory systems has imposed the use of block algorithms for the linear algebra algorithms. Although block algorithms may result in impressive improvements in performance, their numerical properties are quite different from their scalar counterpart and deserve an in-depth study. In this paper, the numerical stability of block Gram-Schmidt orthogonalization is studied and a variant is proposed which has numerical properties similar to the modified Gram-Schmidt procedure while retaining most of the performance advantages of the block formulation.

Key words. Gram-Schmidt orthogonalization, stability, block algorithms

AMS(MOS) subject classifications. 65F25, 65G05, 68Q10

1. Introduction. The Gram-Schmidt algorithm is one of the most widely used algorithms to orthogonalize a set of vectors. Let $A \in \mathfrak{R}^{n \times m}$ ($n \geq m$) be the matrix of which columns (a_i) have to be orthogonalized. The Gram-Schmidt method generates the orthogonal matrix $Q = [q_1; \dots; q_m]$, which corresponds to the Q factor of the QR decomposition of A , according to the following scheme:

$$\begin{aligned} q_1 &:= a_1 / \|a_1\|; \\ \text{for } k &= 1, m-1 \\ a'_{k+1} &:= L_k a_{k+1}; \\ q_{k+1} &:= a'_{k+1} / \|a'_{k+1}\|; \\ \text{endfor;} \end{aligned}$$

where L_k stands for the projection onto the orthogonal complement of the subspace spanned by (q_1, \dots, q_k) . Several implementations of the procedure lie under this common presentation (depending upon the representation of L_k).

If V_k denotes the matrix $[q_1; \dots; q_k]$, L_k may be expressed as the matrix $I - V_k V_k^t$; this corresponds to the classical version (CGS). The matrix L_k may also be expressed as the composition of the projections onto the orthogonal complements of the one-dimensional subspaces spanned by q_1, \dots, q_k :

$$(1) \quad L_k = (I - q_k q_k^t) \cdots (I - q_1 q_1^t).$$

The corresponding algorithm is the modified version (MGS).

Although (CGS) and (MGS) are equivalent in exact arithmetic, they behave quite differently in finite arithmetic: (CGS) is considered unreliable, as we will see in the next section. The advent of vector and parallel computers using hierarchical memory systems has emphasized the use of matrix multiplication as an efficient primitive which allows efficient data management and has considerably renewed the interest in block algorithms based on matrix multiplications [3], [4]. For example, in the case of the Gram-Schmidt procedure, a natural block version can be obtained by using blocks of

* Received by the editors January 24, 1990; accepted for publication (in revised form) August 13, 1990. This research was supported primarily by National Science Foundation grant CCR-8717942 and AT&T grant AFFL 67 Sameh, when the authors were visiting the Center for Supercomputing Research and Development, University of Illinois, Urbana, Illinois 61801.

† IRISA, Campus de Beaulieu, 35042 Rennes Cedex, France.

consecutive columns to orthogonalize the blocks with respect to each other, (MGS) being used to orthogonalize inside a block. In such a case, L_k can be expressed as

$$(2) \quad L_k = (I - q_k q_k^t) \cdots (I - q_l q_l^t) (I - Q_\beta Q_\beta^t) \cdots (I - Q_1 Q_1^t)$$

by assuming that q_{k+1} is a column of the block $Q_{\beta+1}$, whose first column is q_l . The corresponding algorithm is:

```

Q1 := MGS(A1);
for β = 1, ν - 1
    Bβ+1 := Aβ+1;
    for α = 1, β
        S := Qαt Bβ+1;
        Bβ+1 := Bβ+1 - Qα * S;
    endfor;
    Qβ+1 := MGS(Bβ+1);
endfor;

```

In the sequel, the block procedure described above will be denoted as (BGS).

The paper addresses the estimation of the numerical reliability of (BGS) which, when quickly examined, exhibits some similarities with (CGS). After a brief description of the loss of orthogonality in the Gram-Schmidt procedure (§ 2), the error introduced by a projection is studied in detail (§ 3). Based on this result, the overall error resulting from (BGS) procedure is analyzed and bounded (§ 4). The estimation of the loss of orthogonality is heavily based on results obtained by Björck, who first characterized the situation for (MGS) [1]. Then, a variant of the (BGS) procedure is introduced and shown to achieve similar numerical quality as (MGS) (§ 5). This variant based on adding extra reorthogonalization steps has also been proposed by Björck [2]. A detailed performance analysis of the variant as well as experimental results are presented in § 6. The major interest of the variant is that numerical properties are greatly enhanced at a moderate price in terms of performance.

2. An illustration of the scene. Björck [1] proved that if (\bar{Q}, \bar{R}) are the factors obtained by the (MGS) procedure, there exist constants $(K_i)_{i=1,3}$, only depending upon n and m , such that¹:

$$\begin{aligned} \|\bar{Q}\bar{R} - QR\|_F &\leq K_1 \|A\|_F \varepsilon, \\ \|\bar{Q}'\bar{Q} - I\|_2 &\leq K_2 \|A\|_F \|R^{-1}\|_2 \varepsilon, \\ \|\bar{U}\bar{R}\|_F &\leq K_3 \|A\|_F \varepsilon \end{aligned}$$

where (Q, R) are the exact factors of A , U is the upper triangular part of the matrix $\bar{Q}'\bar{Q} - I$, and ε is the precision parameter. These results demonstrate the importance of the condition number of A ($\chi_2(A)$), which satisfies

$$\chi_2(A) \leq \|A\|_F \|R^{-1}\|_2 \leq \sqrt{m} \chi_2(A).$$

¹ In this paper, bounds are expressed with constants K , which only depend upon n and m , and the block size p . We do not provide precise estimation of these constants because we feel that they are very pessimistic and would only muddle the presentation. They are interesting because of their existence. For the same reason, we systematically neglect the $O(\varepsilon^2)$ terms.

For instance, let us consider the following matrix:

$$A = \begin{bmatrix} 1 & 1 & 1 \\ \sigma & 0 & 0 \\ 0 & \sigma & 0 \\ 0 & 0 & \sigma \end{bmatrix}.$$

Its condition number is $\chi(A) = \sqrt{3 + \sigma^2} / \sigma$. The loss of orthogonality when applying (CGS) or (MGS) is reported in Table 1 for several values of σ .

TABLE 1
Loss of orthogonality in (MGS) and (CGS). (Experiments performed with MATLAB and $\epsilon = 2^{-52} \approx 2.10^{-16}$.)

σ	$\chi(A)$	$\ \bar{Q}^t \bar{Q} - I\ _2$	
		MGS	CGS
10^{-4}	$\sqrt{3} \times 10^4$	3×10^{-13}	2×10^{-9}
10^{-5}	$\sqrt{3} \times 10^5$	7×10^{-13}	5×10^{-8}
10^{-6}	$\sqrt{3} \times 10^6$	7×10^{-11}	5×10^{-5}
10^{-7}	$\sqrt{3} \times 10^7$	2×10^{-9}	1×10^{-2}

It should be noted that for a block size of 1 or m , (BGS) corresponds exactly to (MGS). The only case where (BGS) corresponds to (CGS) and not to (MGS) arises when $m = 3$ for a block size of 2. Then, the very special example described above can be viewed as an illustration of the potential bad properties of (BGS).

3. Error in a projection. The elementary step of (BGS) is the projection of a block A_β^α onto the orthogonal complement of $r(Q_\alpha)$ where $\alpha < \beta$ and $r(Q_\alpha)$ is the vector space span by the columns of Q_α . This step corresponds to the computation involved in the innermost loop of (BGS). The sequence A_β^α is then defined by

$$(3) \quad \begin{aligned} A_\beta^1 &= A_\beta \quad \text{and} \\ A_\beta^{\alpha+1} &= P_\alpha A_\beta^\alpha \quad \text{where } P_\alpha = I - Q_\alpha Q_\alpha^t \quad \text{for } 1 < \alpha < \beta \leq \nu. \end{aligned}$$

Due to rounding error effects, the computed sequences are not $\{A_\beta^\alpha\}$ and $\{Q_\alpha\}$, but $\{\bar{A}_\beta^\alpha\}$ and $\{\bar{Q}_\alpha\}$ and the recurrence formula (3) becomes

$$(4) \quad \begin{aligned} \bar{A}_\beta^1 &= A_\beta \quad \text{and} \\ \bar{A}_\beta^{\alpha+1} &= \bar{A}_\beta^\alpha - \bar{Q}_\alpha \bar{R}_\beta^\alpha + \Delta_\beta^\alpha, \quad 1 \leq \alpha < \beta \leq \nu, \end{aligned}$$

where $\bar{R}_\beta^\alpha = \bar{Q}_\alpha^t \bar{A}_\beta^\alpha$. Let $\bar{P}_\alpha = I - \bar{Q}_\alpha \bar{Q}_\alpha^t$.

Because our purpose is to look at the loss of orthogonality, we do not assume that $\bar{Q}_\alpha^t \bar{Q}_\alpha = I$. Let $D_\alpha = \bar{Q}_\alpha^t \bar{Q}_\alpha - I$. In such a case, \bar{P}_α is, a priori, no longer a projection. Let us decompose any vector u as the sum of its components over $r(\bar{Q}_\alpha)$ and $r(\bar{Q}_\alpha)^\perp$, i.e.,

$$u = \bar{Q}_\alpha v + w \quad \text{with } \bar{Q}_\alpha^t w = 0.$$

Then, we obtain

$$\bar{P}_\alpha u = -\bar{Q}_\alpha D_\alpha v + w,$$

which proves that \bar{P}_α is the identity over $r(\bar{Q}_\alpha)^\perp$, but is not the null operator over $r(\bar{Q}_\alpha)$, as it would be if \bar{Q}_α was orthogonal.

Below are some estimations which will be useful for the following section.

LEMMA 3.1. *Let \bar{Q}_α and D_α be defined as above; then the following inequalities hold:*

- $\|\bar{Q}_\alpha\|_2 = \sqrt{1 + \|D_\alpha\|_2} \leq 1 + \frac{1}{2}\|D_\alpha\|_2$;
- *Under the assumption that $\|\bar{Q}_\alpha D_\alpha\|_2 \leq 1$:*

$$\|\bar{P}_\alpha\|_2 = 1.$$

Proof. The proof is obvious. \square

Let η be the largest number such that $\eta(1 + \eta/2) \leq 1$ ($\eta \approx 0.78$). We will assume that

$$(5) \quad \text{for all } \alpha, \quad \|D_\alpha\|_2 \leq \eta$$

and therefore $\|\bar{P}_\alpha\|_2 = 1$.

The error Δ_β^α is bounded by the following estimation.

PROPOSITION 3.1. *There exists a sequence of constant numbers K_β^α such that*

$$\|\Delta_\beta^\alpha\|_F \leq K_\beta^\alpha \|\bar{A}_\beta^\alpha\|_F \varepsilon.$$

Proof. The computation of $\bar{A}_\beta^{\alpha+1}$ consists of three steps. Dropping indices α and β for the sake of clarity, these steps are:

$$\begin{aligned} W^{(1)} &= fl(\bar{Q}'\bar{A}) = \bar{Q}'\bar{A} + V^{(1)}, \\ W^{(2)} &= fl(\bar{Q}W^{(1)}) = \bar{Q}\bar{Q}'\bar{A} + \bar{Q}V^{(1)} + V^{(2)}, \\ W^{(3)} &= fl(\bar{A} - W^{(2)}) = \bar{A} - \bar{Q}\bar{Q}'\bar{A} - \bar{Q}V^{(1)} - V^{(2)} + V^{(3)}. \end{aligned}$$

Using the classical results on the error bounds for matrix multiplications [5], we derive the following inequalities:

$$\begin{aligned} \|V^{(1)}\|_F &\leq K(n) \|\bar{Q}\|_2 \|\bar{A}\|_F \varepsilon, \quad \text{and then} \\ \|\bar{Q}V^{(1)}\|_F &\leq K(n) \|\bar{Q}\|_2^2 \|\bar{A}\|_F \varepsilon, \quad \text{and} \\ \|V^{(2)}\|_F &\leq K(p) \|\bar{Q}\|_2^2 \|\bar{A}\|_F \varepsilon + O(\varepsilon^2) \end{aligned}$$

where $K(\cdot)$ is a polynomial of low order and $\bar{Q} \in \mathfrak{R}^{n \times p}$.

The error due to the subtraction may be bounded by

$$\begin{aligned} \|V^{(3)}\|_F &\leq \|\bar{A} - W^{(2)}\|_F \varepsilon \\ &\leq \|\bar{I} - \bar{Q}\bar{Q}'\|_2 \|\bar{A}\|_F \varepsilon + O(\varepsilon^2) \\ &\leq \|\bar{A}\|_F \varepsilon + O(\varepsilon^2) \end{aligned}$$

since we have assumed (5). Putting together all these bounds we obtain

$$\begin{aligned} \|\Delta\|_F &\leq \|\bar{Q}V^{(1)}\|_F + \|V^{(2)}\|_F + \|V^{(3)}\|_F, \\ \|\Delta\|_F &\leq K \|\bar{A}\|_F \varepsilon. \end{aligned} \quad \square$$

LEMMA 3.2. *There exists a constant $\tau = O(\varepsilon)$ such that $\|\bar{A}_\beta^\alpha\|_F \leq (1 + \tau)^{\alpha-1} \|A_\beta\|_F$ for all $\alpha \leq \beta \leq \nu$.*

Proof. The proof is obvious from Lemma 3.1 and Proposition 3.1 with

$$\tau = \max_{\alpha \leq \beta \leq \nu} (K_\beta^\alpha \varepsilon). \quad \square$$

Once \bar{A}_β^β is computed, a Modified Gram-Schmidt process is applied to get \bar{Q}_α .

Let R'_β and \tilde{R}_β be the exact and the computed R -factor of the QR decomposition of \bar{A}_β^β , respectively. From [1], we have the following estimations:

$$\begin{aligned} \|\bar{A}_\beta^\beta - \bar{Q}_\beta \tilde{R}_\beta\|_F &\leq K_1 \|\bar{A}_\beta^\beta\|_F \varepsilon, \\ \|\bar{Q}_\beta^t \bar{Q}_\beta - I\|_2 &\leq K_2 \|\bar{A}_\beta^\beta\|_F \|R'^{-1}_\beta\|_2 \varepsilon, \\ \|(\bar{Q}_\beta^t \bar{Q}_\beta - I) \tilde{R}_\beta\|_F &\leq K_3 \|\bar{A}_\beta^\beta\|_F \varepsilon. \end{aligned}$$

Let $C(\beta) = \|\bar{A}_\beta^\beta\|_F \|R'^{-1}_\beta\|_2$. If orthogonalization was exact, we would have $\|A\|_F \|R^{-1}\|_2 \cong C(\beta)$ and generally, it would be far from equal. We assume that it is still the case with finite arithmetic. The quantity $C = \max_{\beta \leq \nu} C(\beta)$ is of importance in the sequel.

Let \bar{R}_β^β be the matrix $\bar{Q}_\beta^t \bar{A}_\beta^\beta$. From the equality

$$\bar{R}_\beta^\beta - \tilde{R}_\beta = (\bar{Q}_\beta^t \bar{Q}_\beta - I) \tilde{R}_\beta + \bar{Q}_\beta^t (\bar{A}_\beta^\beta - \bar{Q}_\beta \tilde{R}_\beta)$$

we obtain

$$\|\bar{R}_\beta^\beta - \tilde{R}_\beta\|_F \leq (K_3 + \|\bar{Q}_\beta\|_2 K_1) \|\bar{A}_\beta^\beta\|_F \varepsilon,$$

and if

$$(6) \quad \Delta_\beta^\beta = -\bar{A}_\beta^\beta + \bar{Q}_\beta \bar{R}_\beta^\beta$$

we may write

$$\Delta_\beta^\beta = -\bar{A}_\beta^\beta + \bar{Q}_\beta \tilde{R}_\beta + \bar{Q}_\beta (-\tilde{R}_\beta + \bar{R}_\beta^\beta),$$

which implies the following proposition.

PROPOSITION 3.2. *There exists a constant K_β^β such that*

$$\|\Delta_\beta^\beta\|_F \leq K_\beta^\beta \|\bar{A}_\beta^\beta\|_F \varepsilon.$$

4. Loss of orthogonality in the whole process. By adapting Björck's procedure for blocks, we define the following matrices:

$$U_\beta^\alpha = \begin{cases} \bar{Q}_\alpha^t \bar{Q}_\beta & \text{for } 1 \leq \alpha < \beta \leq \nu, \\ 0 & \text{for } 1 \leq \beta \leq \alpha \leq \nu, \end{cases}$$

and

$$D_\alpha = \bar{Q}_\alpha^t \bar{Q}_\alpha - I \quad \text{for } 1 \leq \alpha \leq \nu.$$

Let U and D be the upper block triangular and the block diagonal matrices defined by the blocks $\{U_\beta^\alpha\}$ and $\{D_\alpha\}$, respectively.

From the previous section, we may ensure that

$$(7) \quad \|D\|_2 \leq K_2 C \varepsilon$$

where $C = \max_\alpha C(\alpha)$.

To estimate a bound on $\|U\|_2$, we first estimate a bound of $\|U\bar{R}\|_2$ where \bar{R} is the upper triangular matrix defined by the blocks $\{\bar{R}_\beta^\alpha\}_{1 \leq \alpha \leq \beta \leq \nu}$.

PROPOSITION 4.1. *There exist two constants K_4 and K_5 such that*

$$\|U\bar{R}\|_F \leq (K_4 C + K_5) \|A\|_F \varepsilon.$$

Proof. By writing (4) for $\alpha = \nu + 1, \dots, \beta - 1$, and from the definition of Δ_β^β , we obtain the following equation by accumulating the equalities:

$$\bar{A}_\beta^{\nu+1} = \sum_{\alpha=\nu+1}^{\beta} \bar{Q}_\alpha \bar{R}_\beta^\alpha - \sum_{\alpha=\nu+1}^{\beta} \Delta_\beta^\alpha.$$

The block (μ, β) of $U\bar{R}$ can then be expressed as

$$\begin{aligned}
 [U\bar{R}]_{\mu\beta} &= \sum_{\alpha=\mu+1}^{\beta} U_{\mu\alpha} \bar{R}_{\beta}^{\alpha} \\
 &= \bar{Q}_{\mu}^t \left(\bar{A}_{\beta}^{\mu+1} + \sum_{\alpha=\mu+1}^{\beta} \Delta_{\beta}^{\alpha} \right).
 \end{aligned}$$

But from (4),

$$\begin{aligned}
 \bar{Q}_{\mu}^t \bar{A}_{\beta}^{\mu+1} &= (\bar{Q}_{\mu}^t - (I + \bar{D}_{\mu}) \bar{Q}_{\mu}^t) \bar{A}_{\beta}^{\mu} + \bar{Q}_{\mu}^t \Delta_{\beta}^{\mu} \\
 &= -\bar{D}_{\mu} \bar{Q}_{\mu}^t \bar{A}_{\beta}^{\mu} + \bar{Q}_{\mu}^t \Delta_{\beta}^{\mu},
 \end{aligned}$$

which implies

$$\|\bar{Q}_{\mu}^t \bar{A}_{\beta}^{\mu+1}\|_F \leq \|\bar{D}_{\mu}\|_2 \|\bar{Q}_{\mu}\|_2 \|\bar{A}_{\beta}^{\mu}\|_F + \|\bar{Q}_{\mu}\|_2 \|\Delta_{\beta}^{\mu}\|_F$$

and

$$\|[U\bar{R}]_{\mu\beta}\|_F \leq \|\bar{D}_{\mu}\|_2 \|\bar{Q}_{\mu}\|_2 \|\bar{A}_{\beta}^{\mu}\|_F + \|\bar{Q}_{\mu}\|_2 \sum_{\alpha=\mu}^{\beta} \|\Delta_{\beta}^{\alpha}\|_F.$$

From

$$\|\bar{U}\bar{R}\|_F \leq \sum_{1 \leq \mu \leq \beta \leq \nu} \|[U\bar{R}]_{\mu\beta}\|_F$$

and from Lemma 3.2 and Propositions 3.1 and 3.2, the result of the proposition is obtained. \square

Let $E_1 = \bar{A} - QR$ where $\bar{A} = \bar{Q}\bar{R}$ and (Q, R) is the result of the QR factorization of A . Then the following result can be proved.

PROPOSITION 4.2. *There exists a constant K_6 such that*

$$\|E_1\|_F = \|\bar{Q}\bar{R} - QR\|_F \leq K_6 \|A\|_F \varepsilon$$

and the condition $\|E_1\|_F \|R^{-1}\|_2 < \sqrt{2} - 1$ implies that \bar{A} has full rank.

Proof. By summing equations (3) from $\alpha = 1$ to $\alpha = \beta - 1$, we obtain

$$\bar{A}_{\beta}^{\beta} = A_{\beta} - \sum_{\alpha=1}^{\beta-1} \bar{Q}_{\alpha} \bar{R}_{\beta}^{\alpha} + \sum_{\alpha=1}^{\beta-1} \Delta_{\beta}^{\alpha}.$$

Since $\bar{A}_{\beta}^{\beta} = \bar{Q}_{\beta} \bar{R}_{\beta}^{\beta} - \Delta_{\beta}^{\beta}$ it follows that

$$\|E_1\|_F = \sum_{\beta=1}^{\nu} \sum_{\alpha=1}^{\beta} \|\Delta_{\beta}^{\alpha}\|_F.$$

Propositions 3.1 and 3.2 provide bounds for $\{\Delta_{\beta}^{\alpha}\}$:

$$\|E_1\|_F \leq \left(\sum_{\beta=1}^{\nu} \sum_{\alpha=1}^{\beta} K_{\beta}^{\alpha} \|\bar{A}_{\beta}^{\alpha}\|_F \right) \varepsilon.$$

Let $K_6 = \max_{\alpha \leq \beta} K_{\beta}^{\alpha} (1 + \tau)^{\alpha-1}$, where τ is the $O(\varepsilon)$ quantity that has been introduced in Lemma 3.2. Then

$$\|E_1\|_F \leq K_6 \sum_{\beta=1}^{\nu} \|\bar{A}_{\beta}\|_F \varepsilon \leq K_6 \|\bar{A}\|_F \varepsilon,$$

which corresponds to the first part of the proposition.

For the second part, Björck's proof can be used:

$$(8) \quad \bar{A}'\bar{A} = R'(I + F_1)R$$

where $F_1 = (Q'E_1R^{-1})' + Q'E_1R^{-1} + (E_1R^{-1})'E_1R^{-1}$.

Since

$$\|F_1\|_2 \leq 2\|E_1\|_F\|R^{-1}\|_2 + \|E_1\|_F^2\|R^{-1}\|_2^2,$$

condition $\|E_1\|_F\|R^{-1}\|_2 < \sqrt{2} - 1$ implies $\|F_1\|_2 < 1$ and therefore the nonsingularity of $(I + F_1)$. \square

THEOREM 4.1. *Let $C = \max_\alpha C(\alpha)$. For sufficiently small $(C\varepsilon)$ and $(\|A\|_F\|R^{-1}\|_2\varepsilon)$, there exists a constant K_7 such that*

$$\|I - \bar{Q}'\bar{Q}\|_2 \leq K_7C\|A\|_F\|R^{-1}\|_2\varepsilon.$$

Proof. From the definitions of matrices U and D , we have

$$\|I - \bar{Q}'\bar{Q}\|_2 \leq 2\|U\|_2 + \|D\|_2.$$

Since $\|D\|_2$ has already been bounded in Theorem 4.1, we only have to consider $\|U\|_2$; because $\|U\|_2 \leq \|UR\|_F\|\bar{R}^{-1}\|_2$, we only need to focus on $\|\bar{R}^{-1}\|_2$. Here again, we adapt Björck's proof by using the following formula:

$$(9) \quad \bar{R}'\bar{R} = R'(I + F_2)R$$

where

$$F_2 = F_1 - (R^{-1})'(UR\bar{R}^{-1})'(\bar{R}R^{-1}) - (\bar{R}R^{-1})'(UR\bar{R}^{-1})R^{-1} - (\bar{R}R^{-1})'D(\bar{R}R^{-1})$$

with F_1 being the matrix introduced at the end of the proof of Proposition 4.2. From (9), it follows that

$$\|\bar{R}R^{-1}\|_2^2 \leq 1 + \|F_2\|_2.$$

Let $x = \sqrt{1 + \|F_2\|_2}$ and $\alpha = \|A\|_F\|R^{-1}\|_2\varepsilon$. From the definition of F_1 and from Proposition 4.2, we obtain

$$\begin{aligned} \|F_1\|_2 &\leq 2\|E_1\|_F\|R^{-1}\|_2 + \|E_1\|_F^2\|R^{-1}\|_2^2 \\ &\leq 2K_6a + K_6^2a^2 + O(\varepsilon^2). \end{aligned}$$

Then from the definition of F_2 , from (9), and from Proposition 4.2, we may ensure that

$$x^2 \leq 1 + 2K_6a + K_6^2a^2 + 2a(K_4C + K_5)x + K_2Cx^2\varepsilon.$$

Then x is such that the following relation is true:

$$(1 - K_2C\varepsilon)x^2 - 2a(K_4C + K_5)x - 1 - 2K_6a - K_6^2a^2 \leq 0.$$

If $(C\varepsilon)$ is small enough to keep $(1 - K_2C\varepsilon)$ positive, then the polynomial has two zeros: one is negative and the other is greater than 1. The second root is denoted $1 + \rho$; $\rho > 0$ and $\rho = O(\varepsilon)$. In conclusion, we have $\|F_2\|_2 \leq \rho$.

Let us assume that the quantities $(\|A\|_F\|R^{-1}\|_2\varepsilon)$ and $(C\varepsilon)$ are small enough to ensure that $\rho < 1$. Then, by considering the inverse of both sides of (9), we obtain

$$\|\bar{R}^{-1}\|_2^2 \leq \|R^{-1}\|_2^2\|(I + F_2)^{-1}\|_2 \leq (1 - \rho)^{-1}\|R^{-1}\|_2^2.$$

This last bound ends the proof of the theorem. \square

The worst bound can be obtained when one block has a condition number almost equal to the condition number of the whole system. In that situation, the quantity $(C\|A\|_F\|R^{-1}\|_2)$ is of the order of the square of the condition number of A .

5. (B2GS) A block method as stable as (MGS). Theorem 4.1 shows the importance of the constant C in the error bound. Even when it is small, the error that occurs within the orthogonalization of a block may significantly impact on the precision of the following steps. By reorthonormalizing every block, i.e., by applying twice (MGS) to every block, the constant C would disappear from the estimation, since it would correspond to the QR factorization of the $\{\bar{Q}_\beta\}$ blocks, which have condition numbers close to unity. This remark, leads directly to the algorithm (B2GS), which has exactly the same structure as (BGS) except that the reorthonormalization procedure is applied twice on each block.

To illustrate the stability improvement in (B2GS) compared to (BGS), some numerical results are presented below.

For the first series of experiments, the results obtained by (BGS) and (B2GS) are compared using MATLAB. The matrix under consideration is the Hilbert matrix:

$$A = (a_{ij}) \in \mathfrak{R}^{20 \times 10} \quad \text{where } a_{ij} = 1/(i+j-1).$$

The characteristics of the matrix are

$$\begin{aligned} \|A\|_F &= 1.9, \\ \|R^{-1}\|_2 &= 1.4 \times 10^{11}, \\ \chi(A) &= 2.6 \times 10^{11}. \end{aligned}$$

The value of the residuals, with respect to the block size, are reported in Table 2.

A second series of experiments was performed by implementing the codes corresponding to the three algorithms ((MGS), (BGS), and (B2GS)) in FORTRAN on the CRAY2. The matrix used was a matrix A of size 1024×512 obtained by the multiplication of three matrices such that

$$A = (I + \alpha H_1) M H_2 \quad \text{where}$$

$$M \in \mathfrak{R}^{1024 \times 512} \quad \begin{cases} m_{1j} = 1, & j = 1, \dots, 512, \\ m_{ij} = \varepsilon \delta_{i-1,j}, & i = 2, \dots, 1024 \quad \text{and } j = 1, \dots, 512, \end{cases}$$

$$H_1 \in \mathfrak{R}^{1024 \times 1024} \quad h_{1ij} \text{ randomly selected in } [-1, +1],$$

$$H_2 \in \mathfrak{R}^{512 \times 512} \quad h_{2ij} \text{ randomly selected in } [-1, +1]$$

with $\alpha = 10^{-3}$ and $\varepsilon = 10^{-2}$. The condition number of the resulting matrix A was estimated by a LINPACK procedure and its value was $\chi(A) = 0.505E + 07$.

Table 3 compares the Frobenius norm of the residuals for various block sizes and for (BGS) and (B2GS). These results show clearly the interest of (B2GS) which gives results with an improvement in the order of magnitude of 3 to 4 compared to the standard (BGS). Furthermore, for all the block sizes, (B2GS) achieves results very close to the standard (MGS) procedure (results of which are given in the first row of Table 3).

TABLE 2

Comparison of loss of orthogonality for (BGS) and (B2GS) with reorthogonalization for each block size on a Hilbert matrix (MATLAB).

Block size	1	2	3	4	5
(BGS)	5.2×10^{-6}	2.8×10^{-6}	2.3×10^{-5}	1.1×10^{-4}	5.2×10^{-3}
(B2GS)	5.2×10^{-6}	3.0×10^{-6}	4.3×10^{-6}	3.1×10^{-6}	4.0×10^{-6}

TABLE 3
Comparison of residuals for (MGS), (BGS), and (B2GS).

(MGS) residual: 0.171×10^{-7} .		
Block size	(BGS) residuals	(B2GS) residuals
16	0.253×10^{-4}	0.592×10^{-7}
32	0.703×10^{-4}	0.696×10^{-7}
48	0.103×10^{-3}	0.527×10^{-7}
64	0.126×10^{-3}	0.464×10^{-7}
80	0.180×10^{-3}	0.501×10^{-7}
96	0.186×10^{-3}	0.466×10^{-7}
112	0.258×10^{-3}	0.623×10^{-7}
128	0.304×10^{-3}	0.577×10^{-7}
144	0.277×10^{-3}	0.511×10^{-7}
160	0.242×10^{-3}	0.534×10^{-7}
176	0.268×10^{-3}	0.483×10^{-7}
192	0.327×10^{-3}	0.619×10^{-7}
208	0.350×10^{-3}	0.674×10^{-7}

6. Performance analysis of (B2GS). Because one of the major advantages of block algorithms is to provide a good data locality together with a good potential for vectorization and parallelization, our performance analysis will be carried out along two axes: number and characteristics of the arithmetic operations (i.e., number of operations, vectorization, and parallelization properties) and data locality. First, the methodology for the performance analysis will be briefly described, then the two axes mentioned above will be studied. The behavior of (B2GS) will be systematically compared with that of (BGS) and the choice of the block size will be analyzed. Finally, some experimental results on an ALLIANT FX80 and a CRAY-2 will be presented.

6.1. Performance model. Throughout § 6, we will use the same framework used in [3] and [4] for studying block algorithms. Let us recall briefly its major characteristics. The target machine will be assumed to be a shared memory multi-(vector) processor, using a memory system consisting of a cache (of size CS) with fast access and a large memory with slower access. According to the methodology developed in [3] and [4], the total execution time will be split into two components:

(1) T_a : Arithmetic time. The total computation time assuming that the cache has an infinite size and that initially all the data reside in cache.

(2) T_l : Load time. The extra time spent in the loads to be performed from the memory due to the finite cache size (including the initial data loads as well as the ones occurring because all the data cannot fit in the cache).

The first component (T_a) takes into account all the parallelization and vectorization properties of the algorithm studied, while the second one (T_l) attempts to quantify its data locality characteristics. This second measure is extremely difficult to evaluate accurately because it depends upon many intrinsic details of a particular cache organization; as it was proposed in [3] and [4], instead of considering the problem of optimizing T_l we will consider the problem of optimizing N_l , which is the total number of loads from memory. For determining N_l , we will assume that the data transfers between the cache and the memory are under software control (i.e., we can specify which data will reside in cache). Such a measure is much simpler to compute and allows capture

of most of the trends in the performance behavior of the algorithm. Furthermore, we will evaluate the ratio $\mu = N_i/N_a$, where N_a is the total number of floating point operations. This quantity enables us to get a better appreciation of the relative cost of the loads from memory.

Our analysis of (B2GS) is greatly simplified because the differences between (BGS) and (B2GS) are minor. For example, since (B2GS) is using exactly the same computational primitives as (BGS), all the results relative to the parallelization and vectorization properties can be carried out from (BGS) to (B2GS); the only difference is the different weight associated with each primitive.

6.2. Optimization of the arithmetic time. First, let us evaluate the cost in terms of number of floating point operations of the extra computation involved in the reorthonormalization. The regular version of (BGS) has the same complexity as (MGS) or (CGS). For an $n \times m$ system, this number of operations is given by

$$N_{\text{op}}(\text{BGS}) = N_{\text{op}}(\text{MGS}) \approx 2nm^2.$$

Let us assume that $m = k\omega$ where ω is the block size. Then, an extra reorthonormalization of every block leads to a total of number of operations for (B2GS) given by

$$(10) \quad N_{\text{op}}(\text{B2GS}) \approx 2nm^2 + 2k\omega^2 = N_{\text{op}}(\text{BGS}) + 2nm\omega.$$

This clearly shows the interest of keeping ω relatively small. Additionally, it also indicates that the idea of varying block sizes is not worthwhile from the arithmetic point of view. In fact, if we assume nonconstant block sizes (i.e., a sequence $\{\omega_i\}_{i=1,k}$), the extra cost associated with the reorthogonalization step is $\sum_{i=1}^k n\omega_i^2$. For a fixed number of blocks, this quantity will be minimal if all the block sizes are constant and equal to m/k . In terms of relative costs,

$$\frac{N_{\text{op}}(\text{B2GS})}{N_{\text{op}}(\text{BGS})} = 1 + \frac{\omega}{m} = 1 + \frac{1}{k}.$$

Besides the extra number of operations, (B2GS) differs from (BGS) by the work repartition between the different primitives. Although all the primitives used in (BGS) (and therefore in (B2GS)) offer a good degree of vectorization and parallelization, there will be a sensible discrepancy in performance between the block factorization ((MGS) primitive on a block) and the matrix multiply primitives, because matrix multiply lends itself to a very efficient use of chaining and registers and has a much simpler synchronization graph than the (MGS) primitive.

More formally, if V_{MGS} (respectively, V_{MAT}) denotes the speed in megaflops of the (MGS) primitive (respectively, the matrix multiply primitive), we end up with

$$T_a(\text{BGS}) = \frac{2mn\omega}{V_{\text{MGS}}} + \frac{2mn(n-\omega)}{V_{\text{MAT}}},$$

$$T_a(\text{B2GS}) = \frac{4mn\omega}{V_{\text{MGS}}} + \frac{2mn(n-\omega)}{V_{\text{MAT}}}.$$

The difference between (BGS) and (B2GS) is therefore

$$\frac{T_a(\text{B2GS}) - T_a(\text{BGS})}{T_a(\text{BGS})} = \frac{V_{\text{MAT}}/V_{\text{MGS}}}{(V_{\text{MAT}}/V_{\text{MGS}}) + k - 1}.$$

If we assume that $V_{\text{MAT}}/V_{\text{MGS}} \leq 6$, which is reasonable for many practical machines, this gives

$$\frac{T_a(\text{B2GS}) - T_a(\text{BGS})}{T_a(\text{BGS})} \leq \frac{6}{k+5}.$$

This last estimation insures that as long as the number of blocks (k) is large enough, the discrepancy in performance (due to the extra arithmetic operations) between BGS and B2GS will be low. In practice, for large problems, it is relatively easy to obtain a large number of blocks: for instance, for 512 vectors and a block size of 16, the discrepancy is less than 16 percent.

6.3. Optimization of data locality. Again, because (B2GS) uses the same primitives as (BGS), similar conclusions can be derived: increasing the block size increases data locality. For simplifying the analysis, we can assume without loss of generality that $\omega < CS/n$. This corresponds to the case where the orthogonalization of a block $n \times \omega$ fits in cache. In the case of (B2GS), such a restriction is really not too constraining because the use of block sizes larger than CS/n would generally result in a prohibitive penalty for the number of floating operations.

The good point of (B2GS) is that the two orthogonalizations primitives are applied one after another for each block. This means that no additional loads will be involved due to reorthonormalization:

$$N_i(\text{B2GS}(\omega)) = N_i(\text{BGS}(\omega)).$$

Using the results presented in [4] on (BGS) behavior, we derive that

$$\mu(\text{B2GS}) = \frac{1}{m} \left(1 + k - \frac{2}{k} \right) + \frac{1}{n} \left(1 - \frac{1}{k} \right).$$

As expected, μ is an increasing function with respect to k . This trend is exactly contradictory with the one observed for the optimization of the arithmetic time. A precise determination of the best choice of the number of blocks requires a precise knowledge of the cost of the memory fetch versus a floating point operation in order to find the right trade-off between arithmetic time minimization and memory load minimization. However, it should be noted that a value of $k = 10$, which gives a relatively good arithmetic optimization, would result in a value for $\mu(\text{B2GS})$ given by

$$\mu(\text{B2GS}) \simeq \frac{10.8}{m} + \frac{0.9}{n},$$

which gives reasonably good value for μ under the conditions that m and n are large enough. In fact, if these conditions are not met, as we will see in the next section on experimental results, (B2GS) as well as (BGS) performances are sensibly affected.

6.4. Experimental results. In this section, experimental results obtained on a CRAY-2 (1CPU) and an ALLIANT FX80 (8 processors) are presented. These results support the performance analysis results by exhibiting trends as predicted and showing that (B2GS) offers comparable performance with BGS. These machines were chosen because they had both a hierarchical memory system with three levels: vector registers (eight vector registers of 64 elements for the CRAY-2, eight vector registers of 32 elements for the ALLIANT), an intermediate memory level (a private local memory of 16K words for the CRAY-2, a shared cache of 64K words for the ALLIANT), and finally, the main memory.

6.4.1. Various block sizes. In this experiment, (BGS) and (B2GS) were run on a 1024×512 matrix with various block sizes (machine used: CRAY-2, one CPU). The timings obtained are presented in Fig. 1, where the result for the classical (MGS) is plotted as a reference point. As predicted by the performance analysis, we clearly see two conflicting trends: at the beginning, increasing the block size decreases the total execution time (the benefit from minimizing the load is bigger than the extra arithmetic

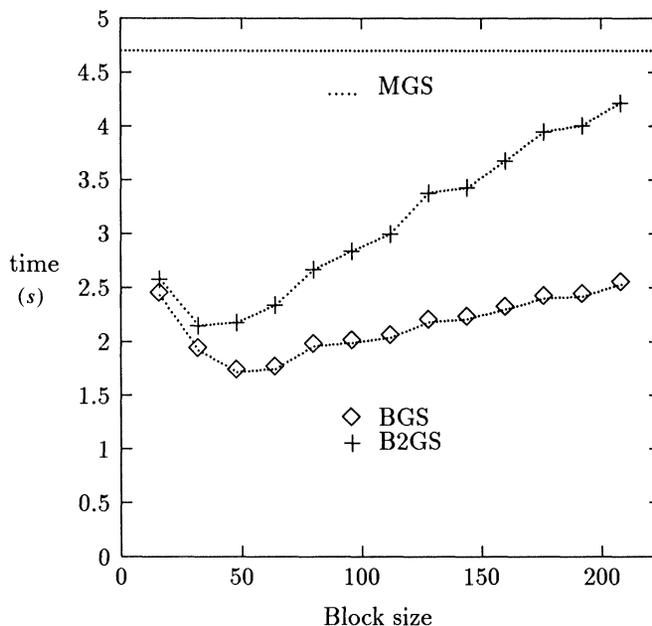


FIG. 1. Running times for (MGS), (BGS), and (B2GS), rectangular matrix (1024×512), (CRAY-2 1CPU).

work), then continuing to increase the block size reverses the situation: the penalty due to the extra operations is overwhelming, resulting in an increase of the execution time.

For (BGS), the minimal execution time of 1.72s is obtained for a block size of 48, while for (B2GS), the minimal execution time of 2.15s corresponds to a block size of 32. The extra cost of (B2GS) is 25 percent. A part of this penalty is due to the fact that our implementation of (B2GS) on the CRAY2 did not keep the block of columns in the local memory between two successive orthonormalizations (requiring a slight modification of the performance analysis), which implies that this implementation (B2GS) does not only increase the number of operations but also the number of loads. Such a characteristic also explains why (B2GS) times are ramping up much faster for large block sizes than (BGS). However, (B2GS) still achieves a speedup of over 2 when compared to the standard (MGS).

6.4.2. Various matrix sizes. In these experiments, (MGS), (BGS), and (B2GS) were run for different matrix shapes and sizes on an ALLIANT FX80. The results are presented in Fig. 2 (square matrix $n \times n$), Fig. 3 (rectangular matrix $1024 \times n$), and Fig. 4 (rectangular matrix $2048 \times n$). In all these figures, the y -axis corresponds to the normalized megaflops, which is obtained by dividing for the three codes the same number of floating point operations (corresponding to (MGS)) by the timing. This allows a fair comparison among the three algorithms: essentially, if the normalized megaflop rate is three times higher for (BGS) than for (MGS), that (BGS) execution time is three times smaller.

The experiments were performed by surrounding the code to be measured by a repetition loop in order to reduce the impact of the clock accuracy. This has an adverse effect for the small matrix sizes, which entirely fit in cache; in such cases, the first iteration will load the matrix in the cache, then all the subsequent operations will be performed with the data entirely resident in cache. This explains the strange shape of the curves for MGS, where the performance first increases then decreases when the

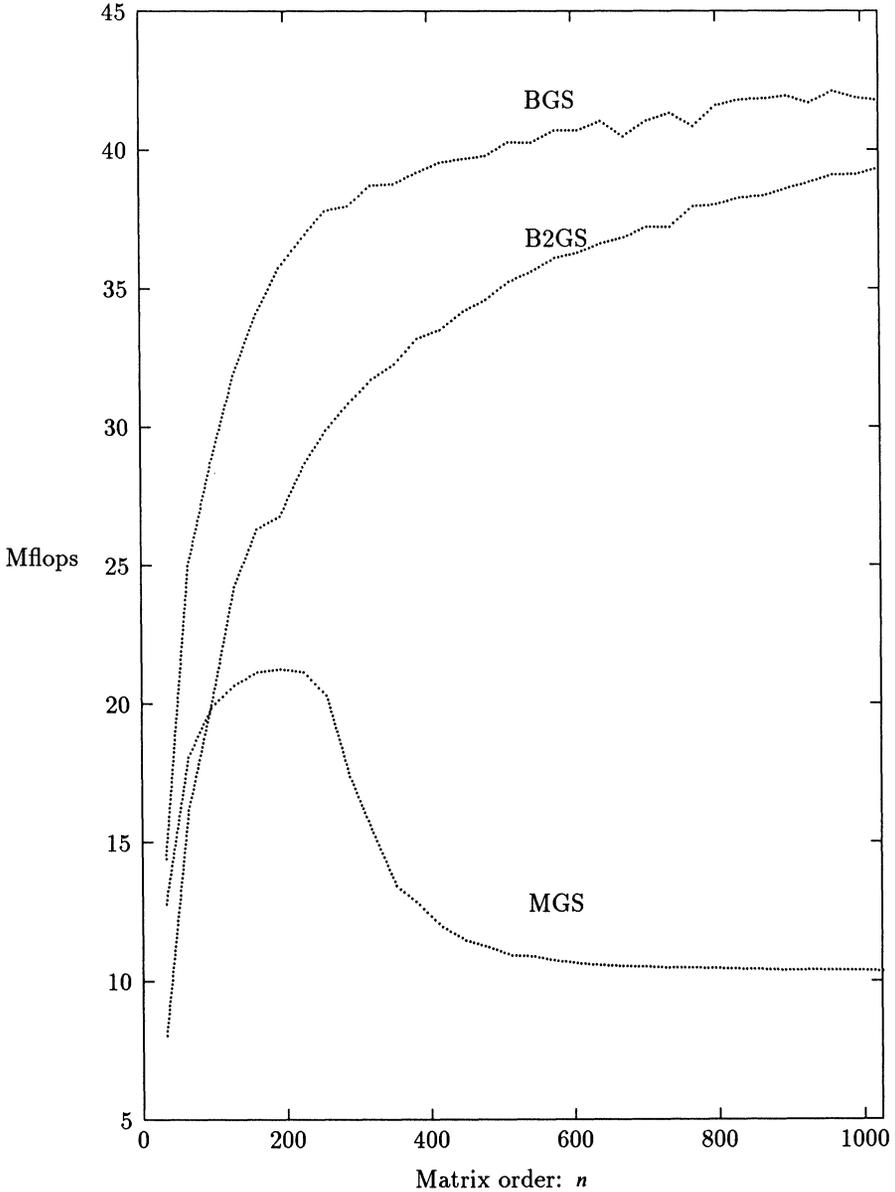


FIG. 2. Normalized megaflops for (MGS), (BGS), and (B2GS), square matrix ($n \times n$) (ALLIANT FX80).

data do not fit any more in cache. It should be noted that (BGS) and (B2GS) are not affected by such a phenomenon because their structure in blocks allows them to keep most of their references to cache.

The block sizes were chosen according to the following rules:

- Square matrix, $n \times n$ (Fig. 2):

BGS: Block size = 16 for $32 \leq n \leq 96$

Block size = 32 for $96 < n \leq 1024$

B2GS: Block size = 16 for $32 \leq n \leq 160$

Block size = 32 for $160 < n \leq 1024$.

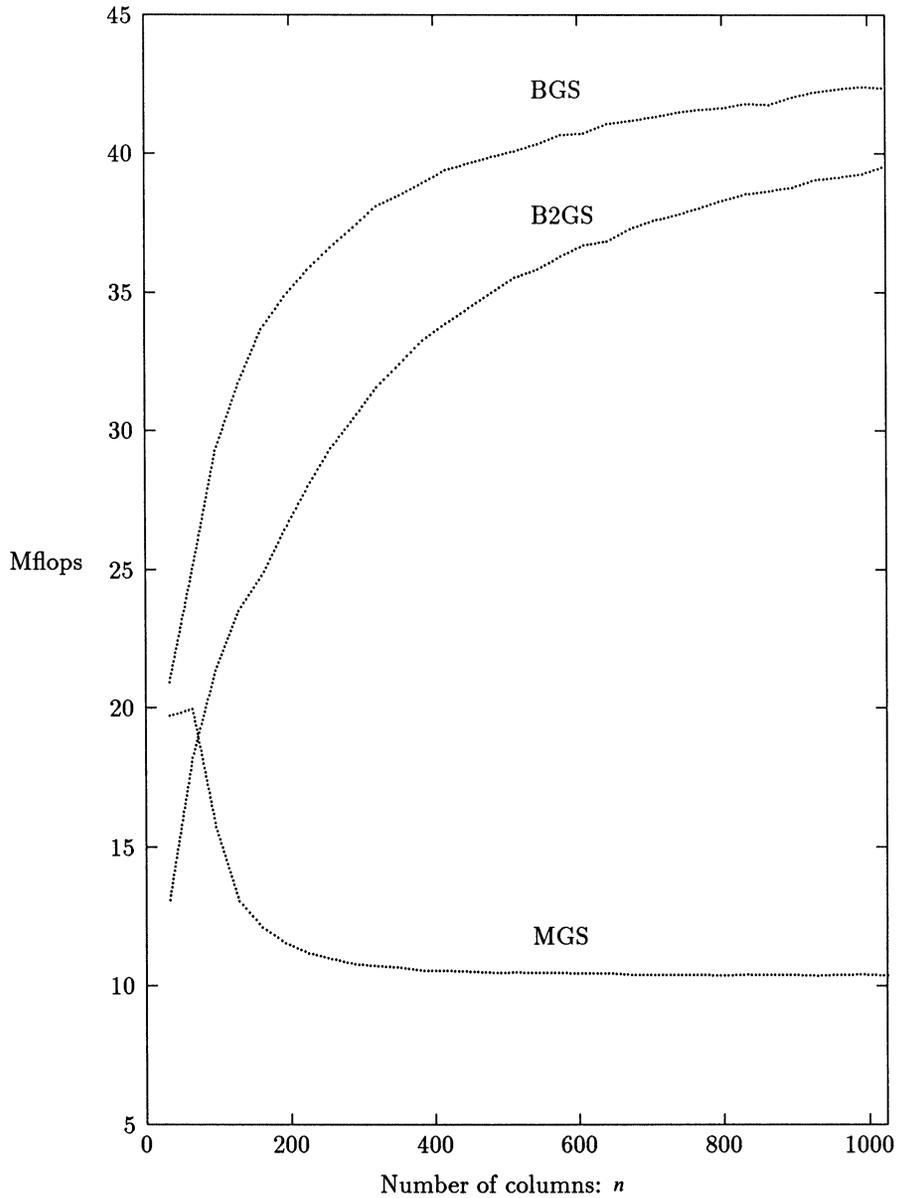


FIG. 3. Normalized megaflops for (MGS), (BGS), and (B2GS), rectangular matrix ($1024 \times n$) (ALLIANT FX80).

• Rectangular matrix, $1024 \times n$ (Fig. 3):

BGS: Block size = 16 for $32 \leq n \leq 64$

Block size = 32 for $64 < n \leq 1024$

B2GS: Block size = 16 for $32 \leq n \leq 128$

Block size = 32 for $128 < n \leq 1024$.

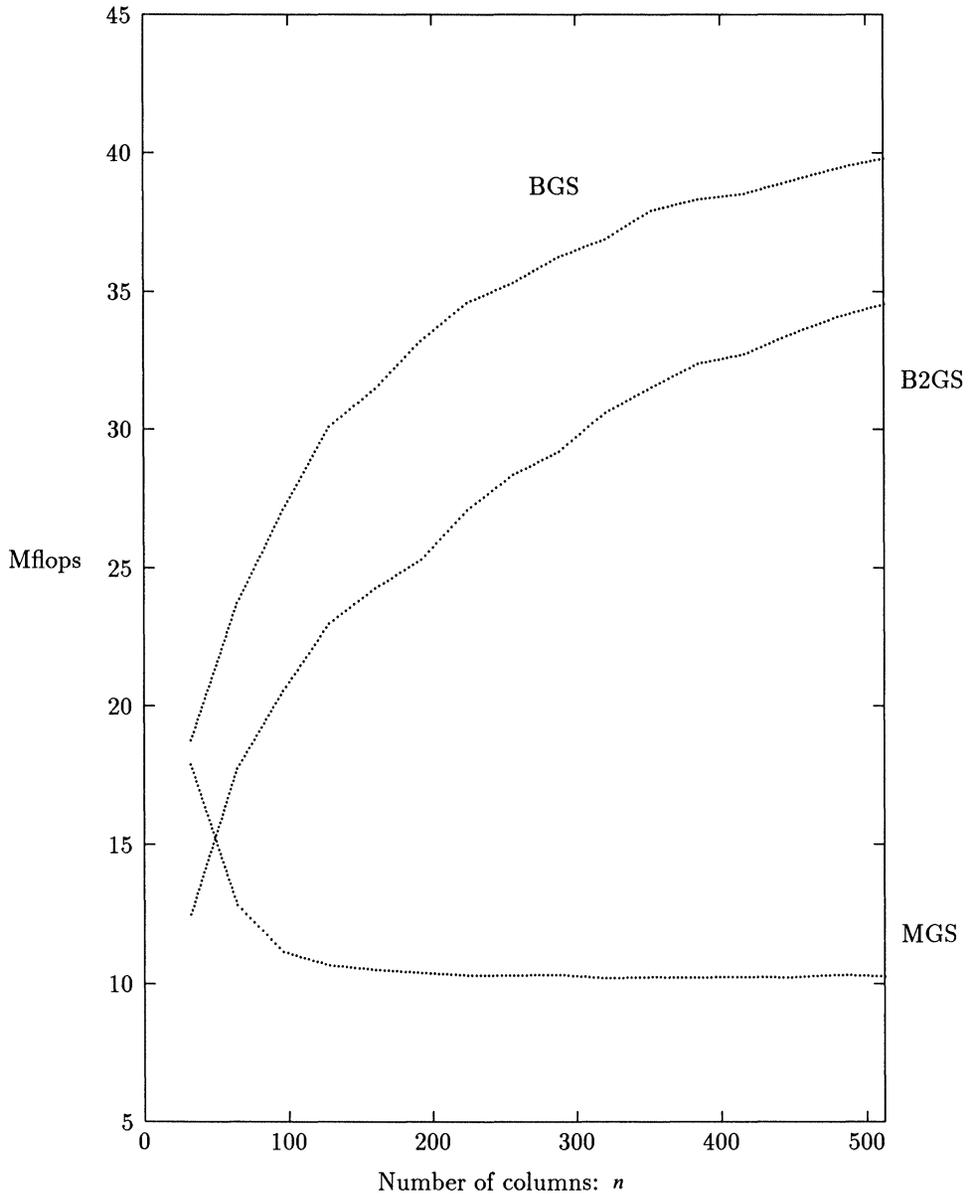


FIG. 4. Normalized megaflops for (MGS), (BGS), and (B2GS), rectangular matrix ($2048 \times n$) (ALLIANT FX80).

- Rectangular matrix, $2048 \times n$ (Fig. 4):

BGS: Block size = 16 for $32 \leq n \leq 64$

Block size = 32 for $64 < n \leq 512$

B2GS: Block size = 16 for $32 \leq n \leq 160$

Block size = 32 for $128 < n \leq 512$.

Such block sizes were chosen according to experimental results.

The main conclusion of the experimental results is that (B2GS) costs between 12 and 25 percent in terms of performance, compared with (BGS); however, the speedups over classical (MGS) are still impressive—over 3 for large matrices. The only negative point is small size problems, where (B2GS) is not competitive, but it should be noted that (BGS) performance is also severely affected in such cases. However, for such small problems, the best choice seems to be classical (MGS), because in such cases the whole data set is close to fit into the cache and therefore does not require specific block algorithms.

7. Conclusion. In this paper it has been proven that by using a reorthogonalization procedure on blocks at an additional cost of low order, the (B2GS) algorithm numerically behaves very closely to the (MGS) algorithm. Another alternative consists of replacing the (MGS) reorthonormalization procedure on each block by another orthogonalization procedure such as the one based on polar decomposition. Since a block \bar{Q}_β is nearly orthonormal, its orthogonalization may be performed in the following way [6]:

$$\begin{aligned}\bar{D} &= \bar{Q}_\beta^t \bar{Q}_\beta - I && np^2 \text{ flops (not considering symmetry),} \\ T &= (I + \bar{D})^{-1/2} && O(p^3) \text{ flops,} \\ \bar{Q}_\beta^c &= \bar{Q}_\beta T && np^2 \text{ flops.}\end{aligned}$$

With such a choice, using reorthonormalization costs twice as much as using (MGS), but only involves matrix multiplication. The use of such primitives (resulting in speedup over 2 compared to Blas1 or Blas2 primitives) may in fact offset the additional price in number of operations, and constitute an interesting alternative.

REFERENCES

- [1] A. BJÖRCK, *Solving linear least squares problems by Gram-Schmidt orthogonalization*, BIT, 7 (1967), pp. 1-21.
- [2] —, *Pivoting and Condition Estimation in Algorithms for Linear Least Squares Problems*, 1990 IMSL User Group Conference, Bologna, Italy, March 26-28, 1990.
- [3] K. GALLIVAN, W. JALBY, AND U. MEIER, *The use of BLAS3 in linear algebra on a parallel processor with a hierarchical memory*, SIAM J. Sci. Statist. Comput., 8 (1987), pp. 1079-1084.
- [4] K. GALLIVAN, W. JALBY, U. MEIER, AND A. SAMEH, *Impact of hierarchical memory systems on linear algebra algorithm design*, Internat. J. Supercomputer Appl., 2 (1988), pp. 12-48.
- [5] G. GOLUB AND C. VAN LOAN, *Matrix Computations*, The Johns Hopkins University Press, Baltimore, MD, 1983.
- [6] B. PHILIPPE, *An algorithm to improve nearly orthonormal sets of vectors on a vector processor*, SIAM J. Algebraic Discrete Methods, 8 (1987), pp. 396-403.