

ROUND OFF ERROR ANALYSIS OF THE CHOLESKYQR2 ALGORITHM*

YUSAKU YAMAMOTO[†], YUJI NAKATSUKASA[‡], YUKA YANAGISAWA[§], AND TAKESHI FUKAYA[¶]

Abstract. We consider the QR decomposition of an $m \times n$ matrix X with full column rank, where $m \geq n$. Among the many algorithms available, the Cholesky QR algorithm is ideal from the viewpoint of high performance computing since it consists entirely of standard level 3 BLAS operations with large matrix sizes, and requires only one reduce and broadcast in parallel environments. Unfortunately, it is well-known that the algorithm is not numerically stable and the deviation from orthogonality of the computed Q factor is of order $O((\kappa_2(X))^2 \mathbf{u})$, where $\kappa_2(X)$ is the 2-norm condition number of X and \mathbf{u} is the unit roundoff. In this paper, we show that if the condition number of X is not too large, we can greatly improve the stability by iterating the Cholesky QR algorithm twice. More specifically, if $\kappa_2(X)$ is at most $O(\mathbf{u}^{-\frac{1}{2}})$, both the residual and deviation from orthogonality are shown to be of order $O(\mathbf{u})$. Numerical results support our theoretical analysis.

Key words. QR decomposition, Cholesky QR, communication-avoiding algorithms, roundoff error analysis.

AMS subject classifications. 15A23, 65F25, 65G50.

1. Introduction. Let $X \in \mathbb{R}^{m \times n}$ be an m by n matrix with $m \geq n$ of full column rank. We consider the computation of its QR decomposition, $X = QR$, where $Q \in \mathbb{R}^{m \times n}$ has orthonormal columns and $R \in \mathbb{R}^{n \times n}$ is upper triangular. This is one of the most fundamental matrix decompositions and is used in various scientific computations. Examples include linear least squares, preprocessing for the singular value decomposition of a rectangular matrix [10], and orthogonalization of vectors arising in block Krylov methods [1, 17] or electronic structure calculations [3, 22]. Frequently in applications, the matrix size is very large, so an algorithm suited for modern high performance computers is desired.

One important feature of modern high performance architectures is that communication is much slower than arithmetic. Here, communication refers to both data transfer between processors or nodes, and data movement between memory hierarchies. Thus, it is essential for higher performance to minimize the frequency and duration of these communications [2]. To minimize interprocessor communications, the algorithm must have a large grain parallelism. To minimize data movement between memory hierarchies, it is effective to reorganize the algorithm to use level 3 BLAS operations as much as possible [10]. Of course, the benefit of using level 3 BLAS operations increases as the size of matrices becomes larger.

Conventionally, three major algorithms have been used to compute the QR decomposition: the Householder QR algorithm, the classical Gram-Schmidt (CGS) algorithm, and the modified Gram-Schmidt (MGS) algorithm. The Householder QR algorithm is widely used due to its excellent numerical stability [11]. MGS, which is less stable, is often preferred when the Q factor is needed explicitly, because it requires only half as much work as the Householder QR in that case. When the matrix A is well conditioned, CGS is also sometimes used since it provides more parallelism. Note that for matrices with 2-norm condition number $\kappa_2(X)$ at

*Received July 18, 2014. Accepted April 10, 2015. Published online on June 14, 2015. Recommended by F. Dopico. The research of Y. Yamamoto was supported by the Ministry of Education, Science, Sports and Culture, Grant-in-Aid for Scientific Research (B)(No. 26286087), Core Research for Evolutional Science and Technology (CREST) Program “Highly Productive, High Performance Application Frameworks for Post Petascale Computing” of Japan Science and Technology Agency (JST). The research of Y. Nakatsukasa was supported by JSPS Scientific Research Grant No. 26870149.

[†]The University of Electro-Communications, Tokyo, Japan / JST CREST, Tokyo, Japan (yusaku.yamamoto@uec.ac.jp).

[‡]University of Tokyo, Tokyo 113-8656, Japan (nakatsukasa@mist.i.u-tokyo.ac.jp).

[§]Waseda University, Tokyo, Japan (yuuka@ruri.waseda.jp).

[¶]RIKEN Advanced Institute for Computational Science, Kobe, Japan / Hokkaido University, Hokkaido, Japan / JST CREST, Tokyo, Japan (fukaya@iic.hokudai.ac.jp).

most $O(\mathbf{u}^{-1})$, where \mathbf{u} is the unit roundoff, repeating CGS or MGS twice leads to algorithms that are as stable as Householder QR [9]. They are known as CGS2 and MGS2, respectively.

For each of these algorithms, variants that can better exploit modern high performance architectures have been developed. There are block versions and recursive versions of Householder QR [6, 18], MGS [13], and CGS [12] that can perform most of the computations in the form of level 3 BLAS. There is also a variant of Householder QR called the tall-and-skinny QR (TSQR) [5], which has large grain parallelism and requires only one reduce and broadcast in a distributed environment.

While these variants have been quite successful, they are not completely satisfactory from the viewpoint of high performance computing. In the block and recursive versions mentioned above, the sizes of matrices appearing in the level 3 BLAS are generally smaller than that of X and become even smaller as the level goes down in the case of recursive algorithms. For the TSQR algorithm, though only one reduce is required throughout the algorithm, the reduction operation is a non-standard one, which corresponds to computing the QR decomposition of a $2n \times n$ matrix formed by concatenating two upper triangular matrices [5]. Thus each reduction step requires $O(n^3)$ work and this tends to become a bottleneck in parallel environments [7]. In addition, the TSQR algorithm requires non-standard level 3 BLAS operations such as multiplication of two triangular matrices [5], for which no optimized routines are available on most machines.

There is another algorithm for the QR decomposition, namely the Cholesky QR algorithm. In this algorithm, one first forms the Gram matrix $A = X^T X$, computes its Cholesky factorization $A = R^T R$, and then finds the Q factor by $Q = X R^{-1}$. This algorithm is ideal from the viewpoint of high performance computing because (1) its computational cost is $2mn^2$ (in the case where $m \gg n$), which is equivalent to the cost of CGS and MGS and half that of Householder QR, (2) it consists entirely of standard level 3 BLAS operations, (3) the first and third steps are highly parallel large size level 3 BLAS operations in which two of the three matrices are of size $m \times n$, (4) the second step, which is the only sequential part, requires only $O(n^3)$ work as opposed to the $O(mn^2)$ work in the first and the third steps, and (5) it requires only one reduce and one broadcast if X is partitioned horizontally. Unfortunately, it is well-known that Cholesky QR is not stable. In fact, deviation from orthogonality of the Q factor computed by Cholesky QR is proportional to $\kappa_2(X)^2$ [19]. Accordingly, standard textbooks like [21] describe the method as “quite unstable and is to be avoided unless we know a priori that R is well conditioned”.

In this paper, we show that the Cholesky QR algorithm can be applied to matrices with a large condition number to give a stable QR factorization if it is repeated *twice*. More specifically, we show that if $\kappa_2(X)$ is at most $O(\mathbf{u}^{-\frac{1}{2}})$, then the Q and R factors obtained by applying Cholesky QR twice satisfy $\|Q^T Q - I\|_F = O(\mathbf{u})$ and $\|X - QR\|_F = O(\mathbf{u})$. Furthermore, we give the coefficients of \mathbf{u} in these bounds explicitly as simple low-degree polynomials in m and n . In the following, we call this method *CholeskyQR2*. Of course, the arithmetic cost of CholeskyQR2 is twice that of Cholesky QR, CGS and MGS, but it is equivalent to the cost of Householder QR, CGS2, and MGS2. Given the advantages stated above, the increase in the computational work might be more than compensated in some cases. Hence, for matrices with $\kappa_2(X) \sim O(\mathbf{u}^{-\frac{1}{2}})$, CholeskyQR2 can be the method of choice in terms of both numerical stability and efficiency on high performance architectures.

Examining the numerical stability of CholeskyQR2 is important because some experimental results suggest it can have very attractive parallel performance. Demmel et al. [5] report the performance of various QR decomposition algorithms, including Cholesky QR (not iterated twice), TSQR, CGS, and conventional Householder QR, on a Pentium III cluster with up to 64 processors and an IBM BlueGene/L with up to 256 processors. They report that on the former

machine, Cholesky QR was more than 6 times faster than TSQR for 1 to 64 processors, while on the latter, Cholesky QR was more than 3 times faster than TSQR for up to 256 processors. Other algorithms were consistently slower than these two. This suggests that CholeskyQR2 would be 1.5 to 3 times faster than TSQR on these machines. In addition, a recent report by the authors [8] compared the performance of CholeskyQR2, TSQR, and conventional Householder QR on the K computer using up to 16384 nodes. In that experiment, the speedup achieved by CholeskyQR2 over TSQR grew with the number of nodes p . Specifically, CholeskyQR2 was about 2 times faster than TSQR when $p = 1024$ and 3 times faster when $p = 16384$. Detailed performance analysis of both algorithms based on performance models is also given in [8].

The idea of performing the QR decomposition twice to get better stability is not new. In his textbook [15], Parlett analyzes Gram-Schmidt orthogonalization of two vectors and introduces the principle of “twice is enough”, which he attributes to Kahan. There is also a classical paper by Daniel, Gragg, Kaufman, and Stewart [4], which deals with the effect of reorthogonalization on the update of the Gram-Schmidt QR decomposition. More recently, Giraud et al. perform a detailed error analysis of CGS2 and MGS2 and show that they give numerically orthogonal Q factor and small residual for matrices with $\kappa_2(X) \sim O(\mathbf{u}^{-1})$ [9]. Stathopoulos et al. experimentally show that the Cholesky QR algorithm can be applied to matrices with a large condition number, if it is applied twice (or more) [19]. Rozložník et al. analyze the CholeskyQR2 algorithm in a more general setting of orthogonalization under indefinite inner product and derive bounds on both the residual and the deviation from orthogonality [16]. However, their bounds are expressed in terms of the computed Q and R factors along with the matrix B that defines the inner product, and do not constitute a priori error bounds in contrast to the bounds derived in this paper. Also, the coefficients of \mathbf{u} are not given explicitly.

Even though the underlying idea of repeating an unstable algorithm twice to improve stability is the same, it is worth noting the inherent disadvantage of CholeskyQR2 when compared with CGS2 and MGS2: numerical breakdown. Specifically, if $\kappa_2(X) \gg O(\mathbf{u}^{-\frac{1}{2}})$, then the Cholesky factorization of $X^T X$ can break down, and so does CholeskyQR2. By contrast, Gram-Schmidt type algorithms are free from such breakdowns (except for very obvious breakdowns due to division by zeros in the normalization) and, as shown in [9], give stable QR factorizations for a much wider class of matrices $\kappa_2(X) \sim O(\mathbf{u}^{-1})$ when repeated twice.

The rest of this paper is organized as follows. In Section 2, after giving some definitions and assumptions, we introduce the CholeskyQR2 algorithm. A detailed error analysis of CholeskyQR2 is presented in Section 3. Numerical results that support our analysis is provided in Section 4. Section 5 gives some discussion on our results. Finally, some concluding remarks are given in Section 6.

2. The CholeskyQR2 algorithm.

2.1. Notation and assumptions. In the following, we consider computing the QR decomposition of an m by n real matrix X , where $m \geq n$. Throughout this paper, we assume that computations are performed using IEEE 754 floating point standard and denote the unit roundoff by \mathbf{u} . Let $\sigma_i(X)$ be the i th largest singular value of X and $\kappa_2(X) = \sigma_1(X)/\sigma_n(X)$ be its condition number. We further assume that

$$(2.1) \quad \delta \equiv 8\kappa_2(X) \sqrt{mn\mathbf{u} + n(n+1)\mathbf{u}} \leq 1.$$

This means that the condition number of X is at most $O(\mathbf{u}^{-\frac{1}{2}})$. From this assumption and $\kappa_2(X) \geq 1$, we also have

$$(2.2) \quad mn\mathbf{u} \leq \frac{1}{64}, \quad n(n+1)\mathbf{u} \leq \frac{1}{64}.$$

Following [11], let us define a quantity γ_k for a positive integer k by

$$\gamma_k = \frac{k\mathbf{u}}{1 - k\mathbf{u}}.$$

Then, it is easy to show that under the assumption (2.1)

$$(2.3) \quad \gamma_m = \frac{m\mathbf{u}}{1 - m\mathbf{u}} \leq 1.1m\mathbf{u}, \quad \gamma_{n+1} = \frac{(n+1)\mathbf{u}}{1 - (n+1)\mathbf{u}} \leq 1.1(n+1)\mathbf{u}.$$

Throughout the paper, a vector norm is always the 2-norm.

2.2. The algorithm. In the Cholesky QR algorithm, we compute the QR decomposition of X by the following procedure

$$\begin{aligned} A &= X^\top X, \\ R &= \text{chol}(A), \\ Y &= XR^{-1}, \end{aligned}$$

where $\text{chol}(A)$ is a function that computes the (upper triangular) Cholesky factor of A . Then, $X = YR$ can be regarded as the QR decomposition of X .

In the CholeskyQR2 algorithm, after obtaining Y and R by the above procedure, we further compute the following

$$\begin{aligned} B &= Y^\top Y, \\ S &= \text{chol}(B), \\ Z &= YS^{-1} \quad (= X(SR)^{-1}), \\ U &= SR. \end{aligned}$$

If the columns of Y are exactly orthonormal, B becomes the identity and $Z = Y$. However, in finite precision arithmetic, this does not hold in general and $Z \neq Y$. In the CholeskyQR2 algorithm the QR decomposition of X is given by $X = ZU$.

3. Error analysis of the CholeskyQR2 algorithm. Our objective is to show that under assumption (2.1), the CholeskyQR2 algorithm delivers an orthogonal factor Z and an upper triangular factor U for which both the orthogonality $\|Z^\top Z - I\|_F$ and residual $\|X - ZU\|_F/\|X\|_2$ are of $O(\mathbf{u})$. Here, the constants in $O(\mathbf{u})$ contain lower order terms in m and n , but not in $\kappa_2(X)$.

This section is structured as follows. In Section 3.1, we formulate the CholeskyQR2 algorithm in floating point arithmetic and prepare several bounds that are necessary to evaluate the orthogonality of the computed orthogonal factor. Using these bounds, the bound on the orthogonality is derived in Section 3.2. In Section 3.3, several bounds that are needed to evaluate the residual are provided, and they are used in Section 3.4 to give a bound on the residual.

3.1. Preparation for evaluating the orthogonality. Let us denote the matrices A , R and Y computed using floating point arithmetic by $\hat{A} = fl(X^\top X)$, $\hat{R} = fl(\text{chol}(\hat{A}))$ and $\hat{Y} = fl(X\hat{R}^{-1})$, respectively. Taking rounding errors into account, the computed quantities satisfy

$$(3.1) \quad \hat{A} = X^\top X + E_1,$$

$$(3.2) \quad \hat{R}^\top \hat{R} = \hat{A} + E_2 = X^\top X + E_1 + E_2,$$

$$(3.3) \quad \hat{\mathbf{y}}_i^\top = \mathbf{x}_i^\top (\hat{R} + \Delta \hat{R}_i)^{-1} \quad (i = 1, 2, \dots, m).$$

Here, \mathbf{x}_i^\top and $\hat{\mathbf{y}}_i^\top$ are the i th row vectors of X and \hat{Y} , respectively. The forward error of the matrix-matrix multiplication $X^\top X$ is denoted by E_1 , while E_2 is the backward error of the Cholesky decomposition of \hat{A} . The matrix $\Delta\hat{R}_i$ denotes the backward error arising from solving the linear simultaneous equation $\mathbf{y}_i^\top \hat{R} = \mathbf{x}_i^\top$ by forward substitution. It would be easier if we could express the backward error of the forward substitution as

$$\hat{Y} = X(\hat{R} + \Delta\hat{R})^{-1},$$

but we have to use the row-wise expression (3.3) instead, because the backward error $\Delta\hat{R}$ depends on the right-hand side vector \mathbf{x}_i^\top .

In the following, we evaluate each of E_1 , E_2 and $\Delta\hat{R}_i$. We also give bounds on the 2-norms of \hat{R}^{-1} and $X\hat{R}^{-1}$ for later use. Furthermore, we derive an alternative form of equation (3.3):

$$\hat{\mathbf{y}}_i^\top = (\mathbf{x}_i^\top + \Delta\mathbf{x}_i^\top)\hat{R}^{-1},$$

in which the backward error enters in the right-hand side vector instead of the coefficient matrix. Equivalently, $\Delta\mathbf{x}_i^\top$ is the residual of the linear system $\mathbf{y}_i^\top \hat{R} = \mathbf{x}_i^\top$. Then, by letting $\Delta X = (\Delta\mathbf{x}_1, \Delta\mathbf{x}_2, \dots, \Delta\mathbf{x}_m)^\top$, we can rewrite (3.3) as

$$\hat{Y} = (X + \Delta X)\hat{R}^{-1},$$

which is more convenient to use. We also evaluate the norm of ΔX .

3.1.1. Forward error in the matrix-matrix multiplication $X^\top X$. Let $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$. Then, the componentwise forward error of the matrix-matrix multiplication $C = AB$ can be evaluated as

$$(3.4) \quad |C - \hat{C}| \leq \gamma_n |A| |B|,$$

where $\hat{C} = fl(AB)$, $|A|$ denotes the matrix whose (i, j) element is $|a_{ij}|$, and the inequality sign means componentwise inequality [11]. The 2-norm of the i th column of X , which we denote by $\tilde{\mathbf{x}}_i$, is clearly less than or equal to $\|X\|_2$. Hence,

$$(3.5) \quad |E_1|_{ij} = |A - \hat{A}|_{ij} \leq \gamma_m (|X|^\top |X|)_{ij} = \gamma_m |\tilde{\mathbf{x}}_i|^\top |\tilde{\mathbf{x}}_j| \leq \gamma_m \|\tilde{\mathbf{x}}_i\| \|\tilde{\mathbf{x}}_j\| \leq \gamma_m \|X\|_2^2.$$

Thus we have

$$(3.6) \quad \|E_1\|_2 \leq \|E_1\|_F \leq \gamma_m n \|X\|_2^2.$$

Simplifying this result using (2.3) leads to

$$(3.7) \quad \|E_1\|_2 \leq 1.1 m n \mathbf{u} \|X\|_2^2.$$

3.1.2. Backward error of the Cholesky decomposition of \hat{A} . Let $A \in \mathbb{R}^{n \times n}$ be symmetric positive definite and assume that the Cholesky decomposition of A in floating point arithmetic runs to completion and the upper triangular Cholesky factor \hat{R} is obtained. Then, there exists $\Delta A \in \mathbb{R}^{n \times n}$ satisfying

$$\hat{R}^\top \hat{R} = A + \Delta A, \quad |\Delta A| \leq \gamma_{n+1} |\hat{R}|^\top |\hat{R}|;$$

see Theorem 10.3 of [11] for details. In our case, we take $\hat{A} = A$ in (3.2) to obtain

$$(3.8) \quad |E_2| \leq \gamma_{n+1} |\hat{R}|^\top |\hat{R}|.$$

Hence,

$$(3.9) \quad \|E_2\|_2 \leq \| |E_2| \|_F \leq \gamma_{n+1} \| |\hat{R}|^\top |\hat{R}| \|_F \leq \gamma_{n+1} \| |\hat{R}| \|_F^2 \leq \gamma_{n+1} n \| \hat{R} \|_2^2.$$

On the other hand, we have from equation (3.2),

$$(3.10) \quad \| \hat{R} \|_2^2 = \| \hat{R}^\top \hat{R} \|_2 \leq \| \hat{A} \|_2 + \| E_2 \|_2.$$

Substituting equation (3.10) into the right hand side of equation (3.9) leads to

$$\| E_2 \|_2 \leq \gamma_{n+1} n (\| \hat{A} \|_2 + \| E_2 \|_2),$$

or,

$$(3.11) \quad \| E_2 \|_2 \leq \frac{\gamma_{n+1} n}{1 - \gamma_{n+1} n} \| \hat{A} \|_2.$$

Noting that

$$(3.12) \quad \| \hat{A} \|_2 \leq \| X^\top X \|_2 + \| E_1 \|_2 \leq \| X \|_2^2 + \gamma_m n \| X \|_2^2 = (1 + \gamma_m n) \| X \|_2^2,$$

from equations (3.1) and (3.6) we have

$$\| E_2 \|_2 \leq \frac{\gamma_{n+1} n (1 + \gamma_m n)}{1 - \gamma_{n+1} n} \| X \|_2^2.$$

This result can be simplified using (2.2) and (2.3) as

$$(3.13) \quad \begin{aligned} \| E_2 \|_2 &\leq \frac{1.1(n+1)\mathbf{u} \cdot n \cdot (1 + 1.1mn\mathbf{u})}{1 - 1.1n(n+1)\mathbf{u}} \| X \|_2^2 \\ &\leq \frac{1.1(n+1)\mathbf{u} \cdot n \cdot (1 + \frac{1.1}{64})}{1 - \frac{1.1}{64}} \| X \|_2^2 \\ &= \frac{7161}{6290} n(n+1)\mathbf{u} \| X \|_2^2 \leq 1.2n(n+1)\mathbf{u} \| X \|_2^2. \end{aligned}$$

3.1.3. Backward error of the forward substitution. Let $U \in \mathbb{R}^{n \times n}$ be a nonsingular triangular matrix. Then, the solution $\hat{\mathbf{x}}$ obtained by solving the linear simultaneous equations $U\mathbf{x} = \mathbf{b}$ by substitution in floating point arithmetic satisfies

$$(3.14) \quad (U + \Delta U)\hat{\mathbf{x}} = \mathbf{b}, \quad |\Delta U| \leq \gamma_n |U|;$$

see Theorem 8.5 of [11]. Note that ΔU depends both on U and \mathbf{b} , although the bound in (3.14) does not. In our case, $U = \hat{R}$, so we have for $1 \leq i \leq m$,

$$(3.15) \quad \| \Delta \hat{R}_i \|_2 \leq \| \Delta \hat{R}_i \|_F = \| |\Delta \hat{R}_i| \|_F \leq \gamma_n \| |\hat{R}| \|_F \leq \gamma_n \sqrt{n} \| \hat{R} \|_2.$$

By inserting equation (3.11) into the right-hand side of (3.10), and using (3.12), we have

$$(3.16) \quad \| \hat{R} \|_2^2 \leq \frac{1}{1 - \gamma_{n+1} n} \| \hat{A} \|_2 \leq \frac{1 + \gamma_m n}{1 - \gamma_{n+1} n} \| X \|_2^2.$$

Inserting this into equation (3.15) leads to

$$\| \Delta \hat{R}_i \|_2 \leq \gamma_n \sqrt{\frac{n(1 + \gamma_m n)}{1 - \gamma_{n+1} n}} \| X \|_2.$$

Simplifying the right-hand side in the same way as in equation (3.13), we obtain

$$\begin{aligned}
 \|\Delta \hat{R}_i\|_2 &\leq 1.1 n \mathbf{u} \sqrt{\frac{n(1 + 1.1 m n \mathbf{u})}{1 - 1.1 n(n + 1) \mathbf{u}}} \|X\|_2 \\
 (3.17) \quad &\leq 1.1 n \mathbf{u} \sqrt{\frac{n \cdot (1 + \frac{1.1}{64})}{1 - \frac{1.1}{64}}} \|X\|_2 \leq 1.2 n \sqrt{n \mathbf{u}} \|X\|_2.
 \end{aligned}$$

3.1.4. Bounding the 2-norm of \hat{R}^{-1} . Next we evaluate the 2-norm of \hat{R}^{-1} . Noting that all the matrices appearing in equation (3.2) are symmetric, we can apply the Bauer-Fike theorem (or Weyl's theorem) to obtain

$$(\sigma_n(X))^2 - (\|E_1\|_2 + \|E_2\|_2) \leq (\sigma_n(\hat{R}))^2.$$

Using assumption (2.1), equations (3.7) and (3.13), we have $\|E_1\|_2 + \|E_2\|_2 \leq \frac{1.2}{64} (\sigma_n(X))^2 \leq (1 - \frac{1}{1.1^2}) (\sigma_n(X))^2$. Hence,

$$\frac{1}{1.1^2} (\sigma_n(X))^2 \leq (\sigma_n(\hat{R}))^2,$$

leading to the bound on \hat{R}^{-1} as

$$(3.18) \quad \|\hat{R}^{-1}\|_2 = (\sigma_n(\hat{R}))^{-1} \leq 1.1 (\sigma_n(X))^{-1}.$$

3.1.5. Bounding the 2-norm of $X \hat{R}^{-1}$. From equation (3.2), we have

$$(3.19) \quad \hat{R}^{-\top} X^\top X \hat{R}^{-1} = I - \hat{R}^{-\top} (E_1 + E_2) \hat{R}^{-1}.$$

Thus,

$$\|X \hat{R}^{-1}\|_2^2 \leq 1 + \|\hat{R}^{-1}\|_2^2 (\|E_1\|_2 + \|E_2\|_2).$$

By using $\|E_1\|_2 + \|E_2\|_2 \leq \frac{1.2}{64} (\sigma_n(X))^2$ again and inserting equation (3.18), we obtain

$$(3.20) \quad \|X \hat{R}^{-1}\|_2 \leq 1.1.$$

3.1.6. Evaluation of the backward error ΔX . From equation (3.3), we have

$$\hat{\mathbf{y}}_i^\top = \mathbf{x}_i^\top (\hat{R} + \Delta \hat{R}_i)^{-1} = \mathbf{x}_i^\top (I + \hat{R}^{-1} \Delta \hat{R}_i)^{-1} \hat{R}^{-1}.$$

Now, let

$$(I + \hat{R}^{-1} \Delta \hat{R}_i)^{-1} = I + \check{R}_i.$$

Then, since $\check{R}_i = \sum_{k=1}^{\infty} (-\hat{R}^{-1} \Delta \hat{R}_i)^k$, we obtain the bound on $\|\check{R}_i\|_2$ as

$$\begin{aligned}
 \|\check{R}_i\|_2 &\leq \sum_{k=1}^{\infty} (\|\hat{R}^{-1}\|_2 \|\Delta \hat{R}_i\|_2)^k = \frac{\|\hat{R}^{-1}\|_2 \|\Delta \hat{R}_i\|_2}{1 - \|\hat{R}^{-1}\|_2 \|\Delta \hat{R}_i\|_2} \\
 (3.21) \quad &\leq \frac{1.1 (\sigma_n(X))^{-1} \cdot 1.2 n \sqrt{n \mathbf{u}} \|X\|_2}{1 - 1.1 (\sigma_n(X))^{-1} \cdot 1.2 n \sqrt{n \mathbf{u}} \|X\|_2},
 \end{aligned}$$

where we used equation (3.17) and (3.18) in the last inequality. The denominator of equation (3.21) can be evaluated as

$$\begin{aligned}
 1 - 1.1(\sigma_n(X))^{-1} \cdot 1.2 n\sqrt{n\mathbf{u}}\|X\|_2 &\geq 1 - \frac{1.1 \cdot 1.2 n\sqrt{n\mathbf{u}}}{8\sqrt{mn\mathbf{u} + n(n+1)\mathbf{u}}} \\
 &\geq 1 - \frac{1.32}{8}\sqrt{n\mathbf{u}} \\
 &\geq 1 - \frac{1.32}{8}\sqrt{\frac{1}{11}} \geq 0.95.
 \end{aligned}$$

Inserting this into equation (3.21) and evaluating the numerator using equation (3.17) again, we have

$$\|\check{R}_i\|_2 \leq \frac{1}{0.95} \cdot 1.1 \kappa_2(X) \cdot 1.2 n\sqrt{n\mathbf{u}} \leq 1.4 \kappa_2(X) n\sqrt{n\mathbf{u}}.$$

Now, let

$$\Delta \mathbf{x}_i^\top = \mathbf{x}_i^\top \check{R}_i.$$

Then,

$$(3.22) \quad \hat{\mathbf{y}}_i^\top = (\mathbf{x}_i^\top + \Delta \mathbf{x}_i^\top) \hat{R}^{-1}.$$

By defining the matrix $\Delta X \in \mathbb{R}^{m \times n}$ as $\Delta X = (\Delta \mathbf{x}_1, \Delta \mathbf{x}_2, \dots, \Delta \mathbf{x}_m)^\top$, we can rewrite equation (3.22) as

$$(3.23) \quad \hat{Y} = (X + \Delta X) \hat{R}^{-1}.$$

The bound on $\|\Delta X\|_F$ can be given as

$$\begin{aligned}
 \|\Delta X\|_F &= \sqrt{\sum_{i=1}^m \|\Delta \mathbf{x}_i^\top\|^2} \leq \sqrt{\sum_{i=1}^m \|\mathbf{x}_i^\top\|^2} \|\check{R}_i\|_2 \\
 (3.24) \quad &\leq 1.4 \kappa_2(X) n\sqrt{n\mathbf{u}} \sqrt{\sum_{i=1}^m \|\mathbf{x}_i^\top\|^2} \leq 1.4 \kappa_2(X) \|X\|_2 n^2 \mathbf{u},
 \end{aligned}$$

where the relationship $\sqrt{\sum_{i=1}^m \|\mathbf{x}_i^\top\|^2} = \|X\|_F \leq \sqrt{n} \|X\|_2$ is used to derive the last inequality.

3.2. Orthogonality of \hat{Y} and \hat{Z} . Based on the bounds given in the previous subsection, we evaluate the orthogonality of \hat{Y} and \hat{Z} as computed by the Cholesky QR and CholeskyQR2 algorithms. The following lemma holds.

LEMMA 3.1. *Suppose that $X \in \mathbb{R}^{m \times n}$, with $m \geq n$, satisfies equation (2.1). Then, the matrix \hat{Y} obtained by applying the Cholesky QR algorithm in floating point arithmetic to X satisfies the following inequality, with δ defined as in (2.1),*

$$\|\hat{Y}^\top \hat{Y} - I\|_2 \leq \frac{5}{64} \delta^2.$$

Proof. By expanding $\hat{Y}^\top \hat{Y}$ using equation (3.23), we have

$$\begin{aligned}
 \hat{Y}^\top \hat{Y} &= \hat{R}^{-\top} (X + \Delta X)^\top (X + \Delta X) \hat{R}^{-1} \\
 &= \hat{R}^{-\top} X^\top X \hat{R}^{-1} + \hat{R}^{-\top} X^\top \Delta X \hat{R}^{-1} + \hat{R}^{-\top} \Delta X^\top X \hat{R}^{-1} + \hat{R}^{-\top} \Delta X^\top \Delta X \hat{R}^{-1} \\
 &= I - \hat{R}^{-\top} (E_1 + E_2) \hat{R}^{-1} + (X \hat{R}^{-1})^\top \Delta X \hat{R}^{-1} + \hat{R}^{-\top} \Delta X^\top (X \hat{R}^{-1}) \\
 &\quad + \hat{R}^{-\top} \Delta X^\top \Delta X \hat{R}^{-1}.
 \end{aligned}$$

Here, we used equation (3.19) to derive the last equality. Thus,

$$\begin{aligned}
 \|\hat{Y}^\top \hat{Y} - I\|_2 &\leq \|\hat{R}^{-1}\|_2^2 (\|E_1\|_2 + \|E_2\|_2) + 2\|\hat{R}^{-1}\|_2 \|X \hat{R}^{-1}\|_2 \|\Delta X\|_2 + \|\hat{R}^{-1}\|_2^2 \|\Delta X\|_2^2 \\
 &\leq \|\hat{R}^{-1}\|_2^2 (\|E_1\|_2 + \|E_2\|_2) + 2\|\hat{R}^{-1}\|_2 \|X \hat{R}^{-1}\|_2 \|\Delta X\|_F + \|\hat{R}^{-1}\|_2^2 \|\Delta X\|_F^2 \\
 &\leq (1.1(\sigma_n(X))^{-1})^2 (1.1 m n \mathbf{u} + 1.2 n(n+1) \mathbf{u}) \|X\|_2^2 \\
 &\quad + 2 \cdot 1.1(\sigma_n(X))^{-1} \cdot 1.1 \cdot 1.4 \kappa_2(X) \|X\|_2 n^2 \mathbf{u} \\
 &\quad + (1.1(\sigma_n(X))^{-1} \cdot 1.4 \kappa_2(X) \|X\|_2 n^2 \mathbf{u})^2 \\
 &\leq \frac{1.1^2 \cdot 1.2}{64} \delta^2 + \frac{2 \cdot 1.1^2 \cdot 1.4}{64} \delta^2 + \left(\frac{1.1 \cdot 1.4}{64} \delta^2 \right)^2 \\
 (3.25) \quad &\leq \frac{5}{64} \delta^2.
 \end{aligned}$$

In the fourth inequality, we used equations (3.7), (3.13), (3.18), (3.20), and (3.24). In the last inequality, we simplified the expression using the assumption $\delta \leq 1$. \square

The next corollary follows immediately from Lemma 3.1.

COROLLARY 3.2. *The condition number of \hat{Y} satisfies $\kappa_2(\hat{Y}) \leq 1.1$.*

Proof. By Lemma 3.1, every eigenvalue λ_i of $\hat{Y}^\top \hat{Y}$ satisfies

$$1 - \frac{5}{64} \leq \lambda_i \leq 1 + \frac{5}{64}.$$

Hence, every singular value $\sigma_i(\hat{Y})$ of \hat{Y} satisfies

$$(3.26) \quad \frac{\sqrt{59}}{8} \leq \sigma_i(\hat{Y}) \leq \frac{\sqrt{69}}{8}.$$

Thus it follows that

$$\kappa_2(\hat{Y}) = \frac{\sigma_1(\hat{Y})}{\sigma_n(\hat{Y})} \leq \sqrt{\frac{69}{59}} \leq 1.1. \quad \square$$

In other words, the matrix \hat{Y} obtained by applying the Cholesky QR algorithm once is extremely well-conditioned, though its deviation from orthogonality, $\|\hat{Y}^\top \hat{Y} - I\|_2$, is still of order 0.1.

Combining Lemma 3.1 and Corollary 3.2, we obtain one of the main results of this paper.

THEOREM 3.3. *The matrix \hat{Z} obtained by applying CholeskyQR2 in floating point arithmetic to X satisfies the following inequality.*

$$\|\hat{Z}^\top \hat{Z} - I\|_2 \leq 6(mn\mathbf{u} + n(n+1)\mathbf{u}).$$

Proof. Noting that $\kappa_2(\hat{Y}) \leq \sqrt{\frac{69}{59}}$, from Corollary 3.2 and applying Lemma 3.1 again to \hat{Y} , we have

$$\begin{aligned}
 \|\hat{Z}^\top \hat{Z} - I\|_2 &\leq \frac{5}{64} \delta^2 \leq \frac{5}{64} \cdot \frac{69}{59} \cdot 64(mn\mathbf{u} + n(n+1)\mathbf{u}) \\
 (3.27) \quad &\leq 6(mn\mathbf{u} + n(n+1)\mathbf{u}). \quad \square
 \end{aligned}$$

3.2.1. Orthogonality error in the Frobenius norm. In the previous sections, we derived the bound on the orthogonality error in terms of the 2-norm, because we wanted to give a bound on the 2-norm condition number of \hat{Y} . However, by tracing the derivation of equation (3.25), we can also derive the following bound in the Frobenius norm,

$$\begin{aligned} \|\hat{Y}^\top \hat{Y} - I\|_F &\leq \|\hat{R}^{-1}\|_2^2 (\|E_1\|_F + \|E_2\|_F) \\ &\quad + 2\|\hat{R}^{-1}\|_2 \|X \hat{R}^{-1}\|_2 \|\Delta X\|_F + \|\hat{R}^{-1}\|_2^2 \|\Delta X\|_F^2. \end{aligned}$$

As it is clear from equations (3.6) and (3.9), the upper bounds on $\|E_1\|_2$ and $\|E_2\|_2$ that were used in equation (3.25) are also bounds on $\|E_1\|_F$ and $\|E_2\|_F$. Thus, the same bound given in equation (3.27) holds for the Frobenius norm as well. We summarize this observation as a corollary as follows.

COROLLARY 3.4. *The matrix \hat{Z} obtained by applying CholeskyQR2 in floating point arithmetic to X satisfies the following inequality.*

$$(3.28) \quad \|\hat{Z}^\top \hat{Z} - I\|_F \leq 6(mn\mathbf{u} + n(n+1)\mathbf{u}).$$

3.3. Preparation for evaluating the residual. Let the matrices B , S , Z , and U , computed by floating point arithmetic, be denoted by $\hat{B} = fl(\hat{Y}^\top \hat{Y})$, $\hat{S} = fl(\text{chol}(\hat{B}))$, $\hat{Z} = fl(\hat{Y} \hat{S}^{-1})$, and $\hat{U} = fl(\hat{S} \hat{R})$, respectively. Then we have

$$(3.29) \quad \begin{aligned} \hat{B} &= \hat{Y}^\top \hat{Y} + E_3, \\ \hat{S}^\top \hat{S} &= \hat{B} + E_4 = \hat{Y}^\top \hat{Y} + E_3 + E_4, \\ \hat{\mathbf{z}}_i^\top &= \hat{\mathbf{y}}_i^\top (\hat{S} + \Delta \hat{S}_i)^{-1}, \quad i = 1, 2, \dots, m, \end{aligned}$$

$$(3.30) \quad \hat{U} = \hat{S} \hat{R} + E_5.$$

Here, $\hat{\mathbf{z}}_i^\top$ is the i th row vector of \hat{Z} , E_3 and E_5 are the forward errors of the matrix multiplications $\hat{Y}^\top \hat{Y}$ and $\hat{S} \hat{R}$, respectively, while E_4 is the backward error of the Cholesky decomposition of \hat{B} . The matrix $\Delta \hat{S}_i$ is the backward error introduced in solving the linear simultaneous equation $\mathbf{z}_i^\top \hat{S} = \hat{\mathbf{y}}_i^\top$ by forward substitution.

As a preparation to evaluating the residual, we first give estimates for the norms of \hat{R} , \hat{S} , $\Delta \hat{S}_i$, E_5 and \hat{Z} .

3.3.1. Evaluation of \hat{R} . From equation (3.16), we have

$$(3.31) \quad \frac{\|\hat{R}\|_2}{\|X\|_2} \leq \sqrt{\frac{1 + \gamma_m n}{1 - \gamma_{n+1} n}} \leq \sqrt{\frac{1 + \frac{mn\mathbf{u}}{1-m\mathbf{u}}}{1 - \frac{n(n+1)\mathbf{u}}{1-(n+1)\mathbf{u}}}} \leq \sqrt{\frac{1 + \frac{\frac{1}{64}}{1-\frac{1}{11}}}{1 - \frac{\frac{1}{64}}{1-\frac{1}{11}}}} = \sqrt{\frac{651}{629}} \leq 1.1.$$

3.3.2. Evaluation of \hat{S} . Noticing that $\|\hat{Y}\|_2 \leq \frac{\sqrt{69}}{8}$ from equation (3.26), we can obtain an upper bound on the norm of \hat{S} by multiplying the bound of equation (3.31) by $\frac{\sqrt{69}}{8}$. Thus,

$$\|\hat{S}\|_2 \leq \sqrt{\frac{651}{629}} \cdot \frac{\sqrt{69}}{8} \leq 1.1.$$

3.3.3. Evaluation of $\Delta \hat{S}_i$. Similarly, multiplying the bound of equation (3.17) by $\frac{\sqrt{69}}{8}$ leads to the following bound on $\Delta \hat{S}_i$

$$\|\Delta \hat{S}_i\|_2 \leq \frac{\sqrt{69}}{8} \cdot 1.2 n \sqrt{n\mathbf{u}} \leq 1.3 n \sqrt{n\mathbf{u}}.$$

3.3.4. Evaluation of E_5 . By using the error bound on matrix multiplication given in equation (3.4), we have

$$|E_5| \leq \gamma_n |\hat{S}| |\hat{R}|.$$

Hence,

$$\begin{aligned} \|E_5\|_2 &\leq \| |E_5| \|_F \leq \gamma_n \| |\hat{S}| |\hat{R}| \|_F \leq \gamma_n \| |\hat{S}| \|_F \| |\hat{R}| \|_F = \gamma_n \| \hat{S} \|_F \| \hat{R} \|_F \\ &\leq n \gamma_n \| \hat{S} \|_2 \| \hat{R} \|_2 \leq n \cdot 1.1 n \mathbf{u} \cdot \sqrt{\frac{651}{629}} \cdot \frac{\sqrt{69}}{8} \cdot \sqrt{\frac{651}{629}} \|X\|_2 \leq 1.2 n^2 \mathbf{u} \|X\|_2. \end{aligned}$$

3.3.5. Evaluation of \hat{Z} . From equation (3.26), we have $\|\hat{Y}\|_F \leq \frac{\sqrt{69}}{8} \sqrt{n}$. This is a bound that does not depend on $\|X\|_2$, so it holds also for $\|\hat{Z}\|_F$. Hence,

$$\|\hat{Z}\|_F \leq \frac{\sqrt{69}}{8} \sqrt{n} \leq 1.1 \sqrt{n}.$$

3.4. Bounding the residual. Based on the above results, we evaluate the residual of the pair (\hat{Z}, \hat{U}) . The following theorem holds, which is also one of our main results.

THEOREM 3.5. *Assume that an $m \times n$ real matrix X ($m \geq n$) satisfies equation (2.1). Then the matrices \hat{Z} and \hat{U} obtained by applying the CholeskyQR2 algorithm in floating point arithmetic to X satisfy the following inequality*

$$(3.32) \quad \frac{\|\hat{Z}\hat{U} - X\|_F}{\|X\|_2} \leq 5n^2 \sqrt{n} \mathbf{u}.$$

Proof. Expanding $\hat{\mathbf{z}}_i^\top \hat{U} - \mathbf{x}_i^\top$ by the equations (3.30), (3.29), and (3.3), leads to

$$\begin{aligned} \|\hat{\mathbf{z}}_i^\top \hat{U} - \mathbf{x}_i^\top\| &= \|\hat{\mathbf{z}}_i^\top (\hat{S}\hat{R} + E_5) - \hat{\mathbf{z}}_i^\top (\hat{S} + \Delta\hat{S}_i)(\hat{R} + \Delta\hat{R}_i)\| \\ &= \|\hat{\mathbf{z}}_i^\top \hat{S}\hat{R} + \hat{\mathbf{z}}_i^\top E_5 - \hat{\mathbf{z}}_i^\top \hat{S}\hat{R} - \hat{\mathbf{z}}_i^\top \hat{S}\Delta\hat{R}_i - \hat{\mathbf{z}}_i^\top \Delta\hat{S}_i \hat{R} - \hat{\mathbf{z}}_i^\top \Delta\hat{S}_i \Delta\hat{R}_i\| \\ &\leq \|\hat{\mathbf{z}}_i^\top\| (\|E_5\|_2 + \|\hat{S}\|_2 \|\Delta\hat{R}_i\|_2 + \|\Delta\hat{S}_i\|_2 \|\hat{R}\|_2 + \|\Delta\hat{S}_i\|_2 \|\Delta\hat{R}_i\|_2) \\ &\leq \|\hat{\mathbf{z}}_i^\top\| (1.2 n^2 \mathbf{u} + 1.1 \cdot 1.2 n \sqrt{n} \mathbf{u} + 1.3 n \sqrt{n} \mathbf{u} \cdot 1.1 \\ &\quad + 1.3 n \sqrt{n} \mathbf{u} \cdot 1.2 n \sqrt{n} \mathbf{u}) \|X\|_2 \\ (3.33) \quad &\leq \|\hat{\mathbf{z}}_i^\top\| \|X\|_2 \cdot 4n^2 \mathbf{u}. \end{aligned}$$

Hence,

$$\begin{aligned} \frac{\|\hat{Z}\hat{U} - X\|_F}{\|X\|_2} &= \frac{\sqrt{\sum_{i=1}^n \|\hat{\mathbf{z}}_i^\top \hat{U} - \mathbf{x}_i^\top\|^2}}{\|X\|_2} \leq 4n^2 \mathbf{u} \sqrt{\sum_{i=1}^n \|\hat{\mathbf{z}}_i^\top\|^2} = 4n^2 \mathbf{u} \|\hat{Z}\|_F \\ &\leq 5n^2 \sqrt{n} \mathbf{u}. \quad \square \end{aligned}$$

4. Numerical results. In this section we evaluate the numerical stability of CholeskyQR2 and compare it with the stability of other popular QR decomposition algorithms, namely, Householder QR, classical and modified Gram-Schmidt (CGS and MGS; we also run them twice, shown as CGS2 and MGS2), and Cholesky QR. To this end, we generate test matrices with a specified condition number by $X := U\Sigma V \in \mathbb{R}^{m \times n}$, where U is an $m \times n$ random orthogonal matrix, V is an $n \times n$ random orthogonal matrix, and

$$\Sigma = \text{diag}(1, \sigma^{\frac{1}{n-1}}, \dots, \sigma^{\frac{n-2}{n-1}}, \sigma).$$

Here, $0 < \sigma < 1$ is some constant. Thus $\|X\|_2 = 1$ and the 2-norm condition number of X is $\kappa_2(X) = 1/\sigma$. We vary $\kappa_2(X)$, m , and n , and investigate the dependence of the orthogonality and residual on them. All computations were done on Matlab R2012b, using IEEE standard 754 binary64 (double precision) so that $\mathbf{u} = 2^{-53} \approx 1.11 \times 10^{-16}$, on a computer running Mac OS X version 10.8, equipped with a 2GHz Intel Core i7 Duo processor.

In Figures 4.1 through 4.6 we show the orthogonality and residual measured by the Frobenius norm under various conditions. Figures 4.1 and 4.2 show the orthogonality $\|\hat{Z}^T \hat{Z} - I\|_F$ and residual $\|\hat{Z} \hat{U} - X\|_F$, respectively, for the case $m = 10000$, $n = 100$, and varying $\kappa_2(X)$. In Figures 4.3 and 4.4, $\kappa_2(X) = 10^5$, $n = 100$, and m varies from 1000 to 10000. In Figures 4.5 and 4.6, $\kappa_2(X) = 10^5$, $m = 1000$, and n varies from 100 to 1000.

It is clear from Figures 4.1 and 4.2 that both the orthogonality and residual are independent of $\kappa_2(X)$ and are of order $O(\mathbf{u})$, as long as $\kappa_2(X)$ is at most $O(\mathbf{u}^{-\frac{1}{2}})$. This is in good agreement with the theoretical prediction and is in marked contrast to the results of CGS, MGS, and Cholesky QR, for which the deviation from orthogonality increases in proportion to $\kappa_2(X)$ and $(\kappa_2(X))^2$, respectively. As it can be seen from Figures 4.3 through 4.6, the orthogonality and residual increase only mildly with m and n , which is also in agreement with the theoretical results, although they are inevitably overestimates. Compared with Householder QR, it was observed that CholeskyQR2 generally produces smaller orthogonality and residual. From these results, we can conclude that CholeskyQR2 is stable for matrices with condition number at most $O(\mathbf{u}^{\frac{1}{2}})$. As is well-known, Gram-Schmidt type algorithms perform well when repeated twice.

5. Discussion. We discuss now four topics related to the stability of CholeskyQR2. First, we compare the orthogonality and residual bounds of CholeskyQR2 given in Theorems 3.4 and 3.5, respectively, with known bounds for Householder QR [11] and CGS2 [9]. Second, we consider how to examine the applicability of CholeskyQR2 for a given matrix. Third, we show that CholeskyQR2 is not only norm-wise stable, but also column-wise stable. Finally, we discuss row-wise stability of CholeskyQR2, which cannot be proved, but is nearly always observed in practice.

5.1. Comparison with the error bounds of Householder QR and CGS2.

5.1.1. Orthogonality. For Householder QR, the Q factor is computed by applying n Householder transformations to $I_{1:m,1:n}$, an $m \times n$ matrix consisting of the first n columns of the identity matrix of order m . Hence, from Lemma 19.3 of [11], the computed Q factor satisfies

$$\hat{Q} = P^T (I_{1:m,1:n} + \Delta I),$$

where P is some $m \times m$ exactly known orthogonal matrix, and ΔI is an $m \times n$ matrix whose column vectors have a norm bounded by $n\gamma_{cm}$, where c is a small positive constant. From this, it is easy to derive the bound

$$\|\hat{Q}^T \hat{Q} - I\|_F \leq n\sqrt{n}\gamma_{c'm} \simeq c'mn\sqrt{n}\mathbf{u}.$$

For CGS2, Giraud et al. show the following bound for deviation from orthogonality under the assumption that $\kappa_2(X)m^2n^3\mathbf{u} = O(1)$ [9].

$$\|\hat{Q}^T \hat{Q} - I\|_2 \leq c''mn\sqrt{n}\mathbf{u}.$$

Although this assumption is hard to satisfy for large matrices (notice that $\kappa_2(X)m^2n^3$ is 10^{19} for the largest matrix appearing in Figure 4.3), it has been observed that CGS2 produces near-orthogonal matrices in many cases where this condition is violated [14].

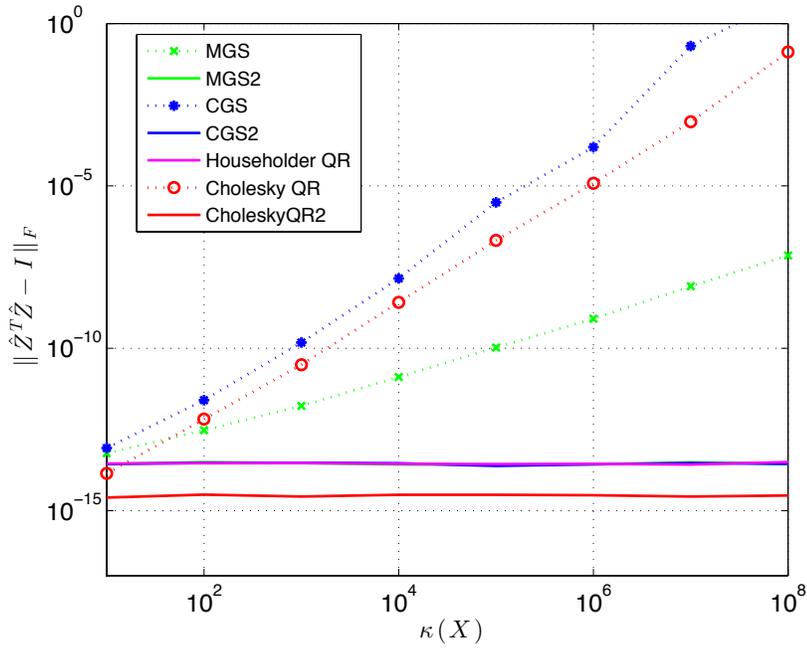


FIG. 4.1. Orthogonality $\|\hat{Z}^T \hat{Z} - I\|_F$ versus $\kappa_2(X)$, for test matrices with $m = 10000$, $n = 100$.

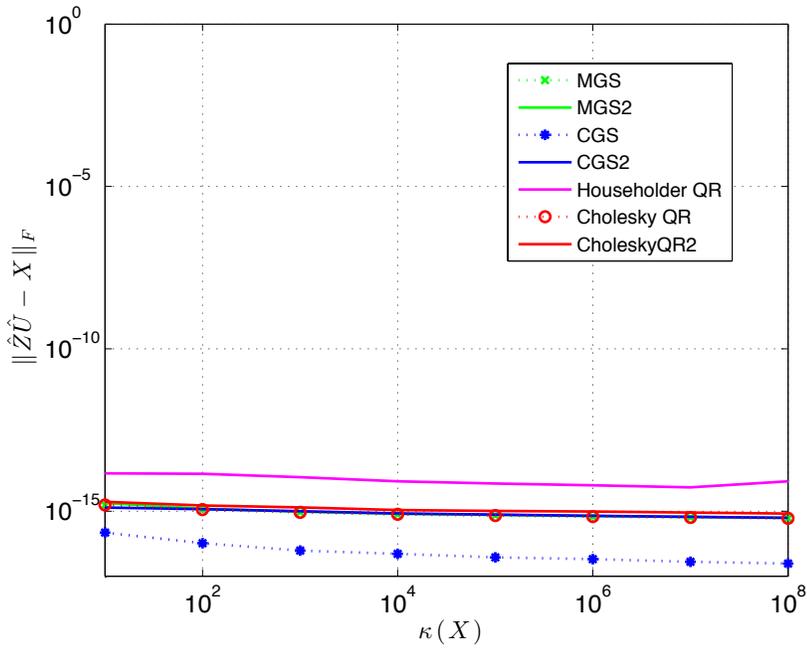


FIG. 4.2. Residual $\|\hat{Z} \hat{U} - X\|_F$ versus $\kappa_2(X)$, for test matrices with $m = 10000$, $n = 100$.

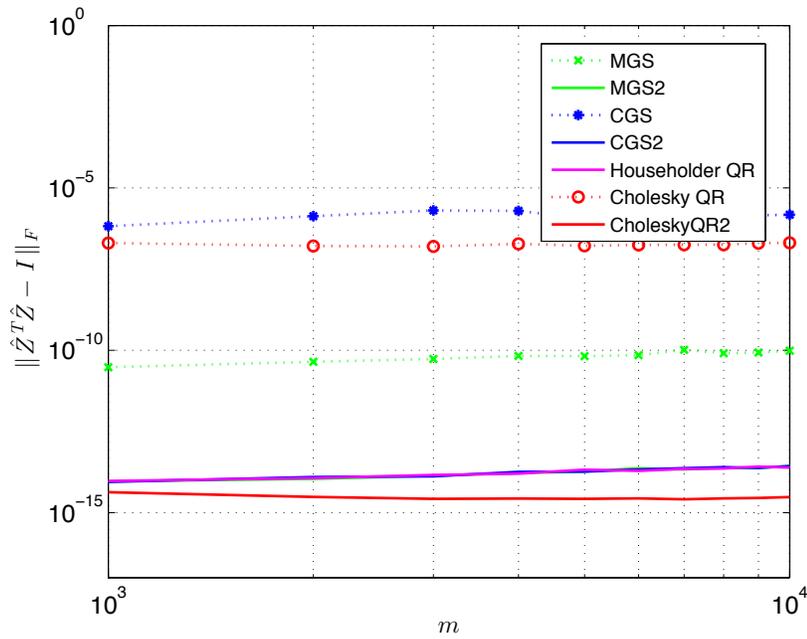


FIG. 4.3. Orthogonality $\|\hat{Z}^T \hat{Z} - I\|_F$ versus m , for test matrices with $\kappa_2(X) = 10^5$, $n = 100$.

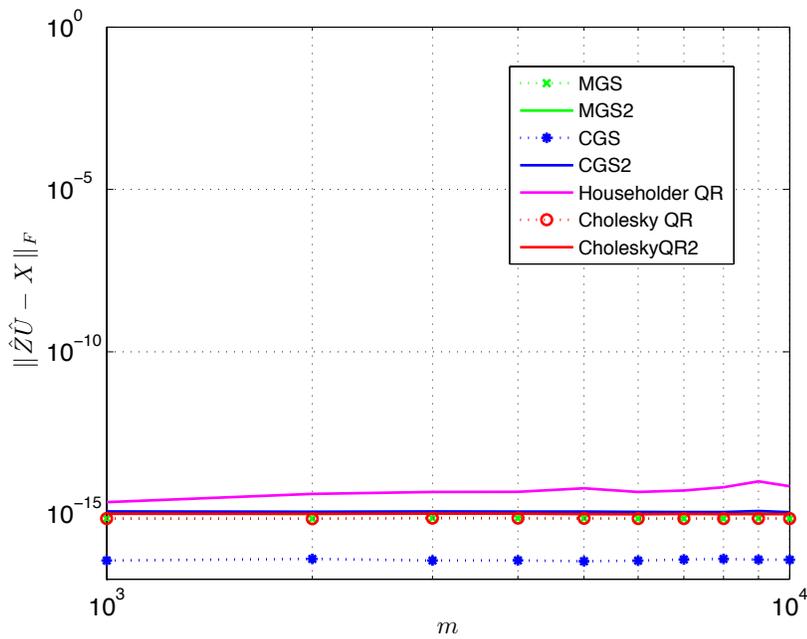


FIG. 4.4. Residual $\|\hat{Z} \hat{U} - X\|_F$ versus m , for test matrices with $\kappa_2(X) = 10^5$, $n = 100$.

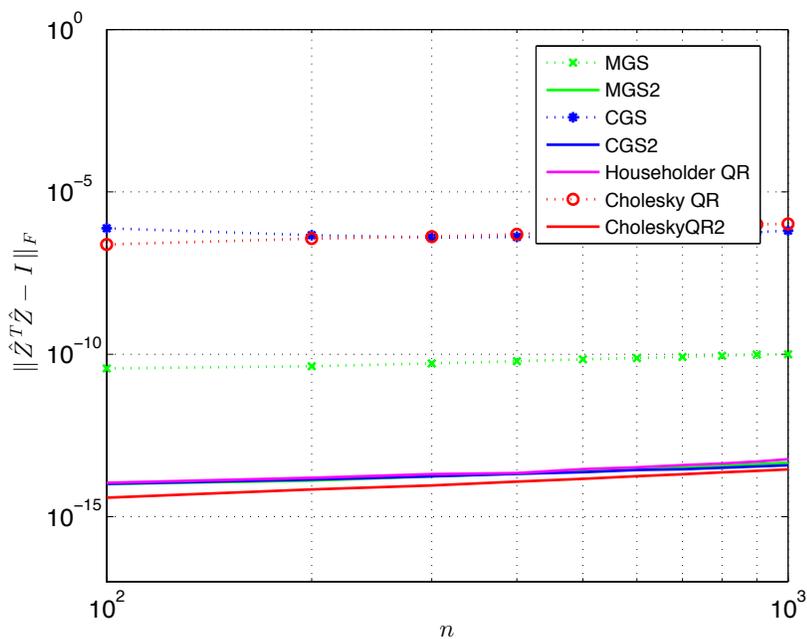


FIG. 4.5. Orthogonality $\|\hat{Z}^T \hat{Z} - I\|_F$ versus n , for test matrices with $\kappa_2(A) = 10^5$, $m = 1000$.

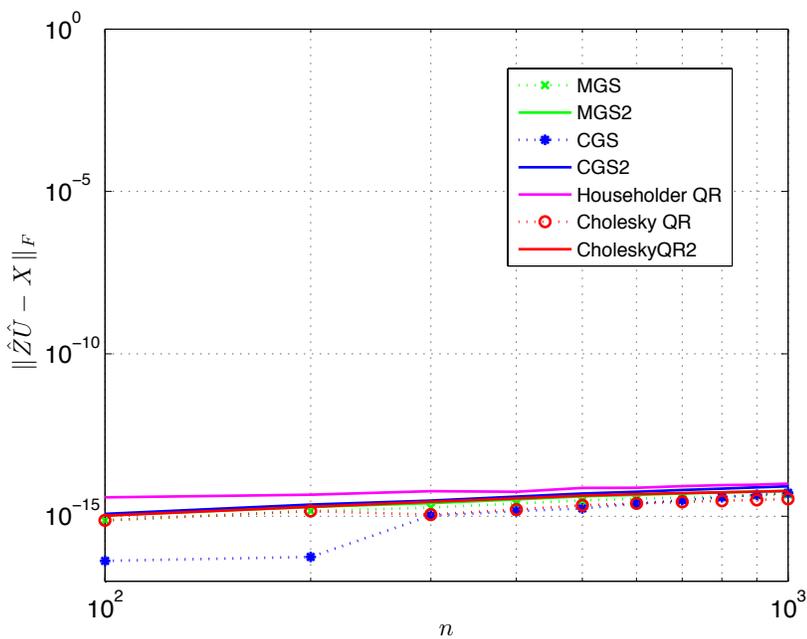


FIG. 4.6. Residual $\|\hat{Z} \hat{U} - X\|_F$ versus n , for test matrices with $\kappa_2(A) = 10^5$, $m = 1000$.

Comparing these bounds with equation (3.28), we observe that the error bound of CholeskyQR2 is smaller by a factor \sqrt{n} . This is in qualitative agreement with the results of the numerical experiments given in the previous section.

Note, however, that this difference should not be overemphasized, because these are merely upper bounds. In fact, the Givens QR algorithm admits an error bound that is smaller than that of Householder QR by a factor n , but no difference in accuracy has been observed in practice [11, p. 368].

5.1.2. Residual. According to [11, Sec. 19.3], the upper bound on the residual of the Householder QR algorithm can be evaluated as $O(mn\sqrt{n}\mathbf{u})$. As for CGS2, it is not difficult to derive a bound of the same order using the results given in [9]. Thus, we can say that the CholeskyQR2 algorithm has a smaller bound also in terms of the residual. This is related to the fact that in the CholeskyQR2 algorithm the computation of Y from X , and Z from Y , is done by row-wise forward substitution. Thus, the backward errors introduced there, or their sum of squares, which is one of the main sources of the residual, do not depend on m when $\|X\|_2$ is fixed. In addition, the forward error in the computation of $\hat{S}\hat{R}$, which is another source of residual, also involves only n . Thus, the residual depends only on n , which is in marked contrast to Householder QR.

A few more comments are in order regarding the bound (3.32). A close examination of equation (3.33) shows that the highest order term in the residual comes from the forward error of the matrix multiplication $\hat{S}\hat{R}$, which we denoted by E_5 . This implies that if we compute this matrix multiplication using extended precision arithmetic, we can reduce the upper bound on the residual to $O(n^2\mathbf{u})$ with virtually no increase in the computational cost (when $m \gg n$). Moreover, in a situation where only the orthogonal factor \hat{Z} is needed, as in the orthogonalization of vectors, we can leave the product $\hat{S}\hat{R}$ uncomputed and say that the triplet $(\hat{Z}, \hat{S}, \hat{R})$ has residual $O(n^2\mathbf{u})$.

5.2. Applicability of CholeskyQR2 for a given matrix. There are some cases in which the condition number of X is known in advance to be moderate. An example is orthogonalization of vectors in first-principles molecular dynamics [3]. In this application, we are interested in the time evolution of an orthogonal matrix $X(t) \in \mathbb{R}^{m \times n}$, whose column vectors are an orthogonal basis for the space of occupied-state wave functions. To obtain $X(t + \Delta t)$, we first compute $\tilde{X} = X(t) - F(X)\Delta t$, where $F(X) \in \mathbb{R}^{m \times n}$ is some nonlinear matrix function of X , and then compute $X(t + \Delta t)$ by orthogonalizing the columns of \tilde{X} . Since $X(t)$ is orthogonal, we can easily evaluate the deviation from orthogonality of \tilde{X} by computing the norm of $F(X)\Delta t$. Usually, the time step Δt is small enough to ensure that $\kappa_2(\tilde{X}) \ll \mathbf{u}^{-\frac{1}{2}}$.

In some cases, however, the condition number of X cannot be estimated in advance and one may want to examine the applicability of CholeskyQR2 from intermediate quantities that are computed in the algorithm. This is possible if \hat{R} has been computed without breakdown in the Cholesky decomposition. Given \hat{R} , one can estimate its largest and smallest singular values using the power method and inverse power method on $R^T R$, respectively. Indeed the MATLAB condition number estimator `condest` first computes the LU factorization of the input matrix, then applies a few iterations of power method to obtain a reliable estimate of the 1-norm condition number. This should not cost too much because \hat{R} is triangular and each step of both methods requires only $O(n^2)$ work. After that, one can evaluate the condition number of X by using the relations (3.1) and (3.2), the bounds (3.7) and (3.13) on $\|E_1\|_2$ and $\|E_2\|_2$, respectively, and the Bauer-Fike theorem.

5.3. Column-wise stability of CholeskyQR2. Thus far, we have investigated the norm-wise residual of CholeskyQR2. Sometimes the columns of X have widely varying norms, and

one may wish to obtain the more stringent column-wise backward stability, which requires

$$\|\tilde{\mathbf{x}}_j - \hat{Q}\hat{\mathbf{r}}_j\|/\|\tilde{\mathbf{x}}_j\| = O(\mathbf{u}), \quad j = 1, \dots, n.$$

Here, $\tilde{\mathbf{x}}_j$ and $\hat{\mathbf{r}}_j$ denote the j th columns of X and \hat{R} , respectively. In this subsection, we prove that CholeskyQR2 is indeed column-wise backward stable.

To see this, we first consider a single Cholesky QR and show that the computed $\|\hat{\mathbf{r}}_j\|$ is of the same order as $\|\tilde{\mathbf{x}}_j\|$. Let us recall equations (3.1) through (3.3). From equation (3.5), we have

$$(5.1) \quad |E_1|_{jj} \leq \gamma_m |\tilde{\mathbf{x}}_j|^\top |\tilde{\mathbf{x}}_j| = \gamma_m \|\tilde{\mathbf{x}}_j\|^2.$$

By considering the (j, j) th element of equation (3.2) and substituting equations (5.1) and (3.8), we obtain

$$\|\hat{\mathbf{r}}_j\|^2 \leq |\hat{A}_{jj}| + \gamma_{n+1} |\hat{\mathbf{r}}_j|^\top |\hat{\mathbf{r}}_j| \leq \|\tilde{\mathbf{x}}_j\|^2 + \gamma_m \|\tilde{\mathbf{x}}_j\|^2 + \gamma_{n+1} \|\hat{\mathbf{r}}_j\|^2.$$

Hence,

$$(5.2) \quad \|\hat{\mathbf{r}}_j\| \leq \sqrt{\frac{1 + \gamma_m}{1 - \gamma_{n+1}}} \|\tilde{\mathbf{x}}_j\| = \|\tilde{\mathbf{x}}_j\| \cdot O(1).$$

Now we demonstrate the column-wise backward stability of a single Cholesky QR. Let the j th column of \hat{Y} be denoted by $\hat{\mathbf{y}}_j$. From equation (3.3), we have

$$X_{ij} = \hat{\mathbf{y}}_i^\top (\hat{\mathbf{r}}_j + \Delta \hat{\mathbf{r}}_j^{(i)}).$$

Here, $\Delta \hat{\mathbf{r}}_j^{(i)}$ is the j th column of $\Delta \hat{R}_i$. Thus,

$$|X_{ij} - \hat{\mathbf{y}}_i^\top \hat{\mathbf{r}}_j| \leq |\hat{\mathbf{y}}_i^\top \Delta \hat{\mathbf{r}}_j^{(i)}| \leq \|\hat{\mathbf{y}}_i\| \|\Delta \hat{\mathbf{r}}_j^{(i)}\| \leq \gamma_n \|\hat{\mathbf{y}}_i\| \|\hat{\mathbf{r}}_j\|.$$

Squaring both sides and summing over i leads to

$$\|\tilde{\mathbf{x}}_j - \hat{Y}\hat{\mathbf{r}}_j\|^2 \leq \gamma_n^2 \|\hat{Y}\|_F^2 \|\hat{\mathbf{r}}_j\|^2.$$

By using $\|\hat{Y}\|_F = O(1)$ (see Lemma 3.1) and equation (5.2), we can establish the column-wise backward stability of Cholesky QR as follows

$$(5.3) \quad \frac{\|\tilde{\mathbf{x}}_j - \hat{Y}\hat{\mathbf{r}}_j\|}{\|\tilde{\mathbf{x}}_j\|} \leq \gamma_n \|\hat{Y}\|_F \cdot \frac{\|\hat{\mathbf{r}}_j\|}{\|\tilde{\mathbf{x}}_j\|} = \gamma_n \cdot O(1) \cdot \sqrt{\frac{1 + \gamma_m}{1 - \gamma_{n+1}}}.$$

To apply the above result to CholeskyQR2, we consider the backward errors in the second QR decomposition $Y = ZS$ and the product of the two upper triangular factors $U = SR$. These backward errors, which we denote by $\Delta \hat{Y}$ and $\Delta \hat{S}$, respectively, satisfy

$$\begin{aligned} \hat{Y} + \Delta \hat{Y} &= \hat{Z}\hat{S}, \\ \hat{\mathbf{u}}_j &= (\hat{S} + \Delta \hat{S})\hat{\mathbf{r}}_j. \end{aligned}$$

Here, $\hat{\mathbf{u}}_j$ is the j th column of \hat{U} . To evaluate $\Delta \hat{Y}$, we note the following inequality, which can be obtained in the same way as equation (5.3)

$$\|\tilde{\mathbf{y}}_j - \hat{Z}\hat{\mathbf{s}}_j\|^2 \leq \gamma_n^2 \|\hat{Z}\|_F^2 \|\hat{\mathbf{s}}_j\|^2.$$

Summing both sides over i and taking the square root gives

$$(5.4) \quad \|\Delta\hat{Y}\|_F = \|\hat{Y} - \hat{Z}\hat{S}\|_F \leq \gamma_n \|\hat{Z}\|_F \|\hat{S}\|_F = \gamma_n \cdot O(1).$$

As for $\Delta\hat{S}$, the standard result on the error analysis of matrix-vector product, combined with $\|\hat{S}\|_F \simeq \|\hat{Y}\|_F = O(1)$, leads to

$$(5.5) \quad \|\Delta\hat{S}\|_F \leq \gamma_n \|\hat{S}\|_F = \gamma_n \cdot O(1).$$

On the other hand,

$$(5.6) \quad \begin{aligned} \tilde{\mathbf{x}}_j - \hat{Z}\hat{\mathbf{u}}_j &= \tilde{\mathbf{x}}_j - \hat{Z}(\hat{S} + \Delta\hat{S})\hat{\mathbf{r}}_j \\ &= \tilde{\mathbf{x}}_j - (\hat{Y} + \Delta\hat{Y} + \hat{Z}\Delta\hat{S})\hat{\mathbf{r}}_j \\ &= (\tilde{\mathbf{x}}_j - \hat{Y}\hat{\mathbf{r}}_j) - (\Delta\hat{Y} + \hat{Z}\Delta\hat{S})\hat{\mathbf{r}}_j. \end{aligned}$$

By substituting equations (5.3), (5.4), and (5.5) into equation (5.6), we finally obtain the column-wise backward stability of CholeskyQR2 as follows.

$$\frac{\|\tilde{\mathbf{x}}_j - \hat{Z}\hat{\mathbf{u}}_j\|}{\|\tilde{\mathbf{x}}_j\|} \leq \frac{\|\tilde{\mathbf{x}}_j - \hat{Y}\hat{\mathbf{r}}_j\|}{\|\tilde{\mathbf{x}}_j\|} + (\|\Delta\hat{Y}\|_F + \|\hat{Z}\|_F \|\Delta\hat{S}\|_F) \cdot \frac{\|\hat{\mathbf{r}}_j\|}{\|\tilde{\mathbf{x}}_j\|} = \gamma_n \cdot O(1).$$

5.4. Row-wise stability of CholeskyQR2. In this subsection, we investigate the row-wise stability of CholeskyQR2, which is defined as

$$(5.7) \quad \|\mathbf{x}_i^\top - \hat{\mathbf{q}}_i^\top \hat{R}\| / \|\mathbf{x}_i^\top\| = O(\mathbf{u}), \quad i = 1, \dots, m.$$

Here \mathbf{x}_i^\top and $\hat{\mathbf{q}}_i^\top$ denote the i th rows of the matrices.

The requirement (5.7) is strictly more stringent than the normwise stability, and indeed the standard Householder QR factorization does not always achieve (5.7). It is known [11, Ch. 19] that when row sorting and column pivoting are used, Householder QR factorization gives row-wise stability. However, pivoting involves an increased communication cost and is best avoided in high-performance computing.

Having established the normwise and column-wise stability of CholeskyQR2, we now examine its row-wise stability. To gain some insight we first run experiments with a semi-randomly generated matrix X , whose row norms vary widely. Specifically, we generate a random $m \times n$ matrix via the MATLAB command $X = \text{randn}(m, n)$, then left-multiply a diagonal matrix $X := DX$, with $D_{jj} = 2^{\frac{j}{2}}$ for $j = 1, \dots, m$. Here we set $m = 100$ and $n = 50$; the matrix thus has rows of exponentially growing norms and $\kappa_2(X) \approx \mathbf{u}^{-\frac{1}{2}}$. Figure 5.1 shows the row-wise residuals of three algorithms: standard Householder QR, Householder QR employing row sorting and column pivoting, and CholeskyQR2.

We make several observations from Figure 5.1. First, we confirm the known fact that the standard Householder QR factorization is not row-wise backward stable, but this can be cured by employing row sorting and column pivoting. Second, CholeskyQR2 gives row-wise stability comparable to Householder QR with row sorting and column pivoting; this is perhaps surprising, considering the fact that CholeskyQR2 employs no pivoting or sorting.

To illustrate the situation, we examine the first step of CholeskyQR2. Recall that $\hat{\mathbf{q}}_i^\top = fl(\mathbf{x}_i^\top \hat{R}^{-1})$. Hence $\|\mathbf{x}_i^\top - \hat{\mathbf{q}}_i^\top \hat{R}\|_2 = \|\mathbf{x}_i^\top - fl(\mathbf{x}_i^\top \hat{R}^{-1})\hat{R}\|_2$, and by standard triangular solve there exist ΔR_i for $i = 1, \dots, m$ such that

$$fl(\mathbf{x}_i^\top \hat{R}^{-1})(\hat{R} + \Delta R_i) = \mathbf{x}_i^\top, \quad \|\Delta R_i\| = O(\mathbf{u})\|\hat{R}\|.$$

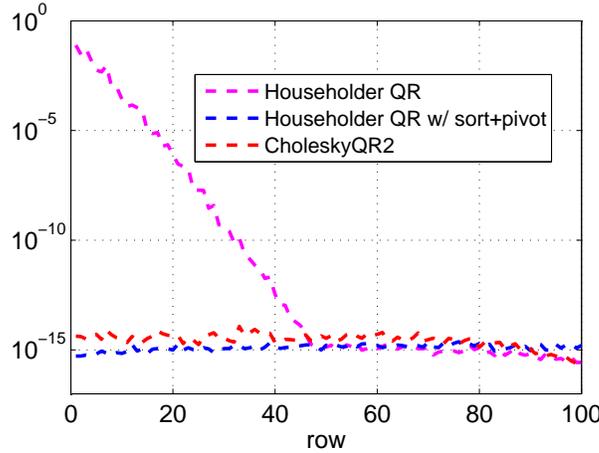


FIG. 5.1. Row-wise residual $\|\mathbf{x}_i^\top - \hat{\mathbf{q}}_i^\top \hat{R}\|_2 / \|\mathbf{x}_i^\top\|_2$.

Hence for row-wise stability we need $\|fl(\mathbf{x}_i^\top \hat{R}^{-1})\Delta R_i\| = O(\mathbf{u})\|\mathbf{x}_i^\top\|$. Since

$$\|fl(\mathbf{x}_i^\top \hat{R}^{-1})\Delta R_i\| \leq O(\mathbf{u}) \|\hat{R}\| \|fl(\mathbf{x}_i^\top \hat{R}^{-1})\|,$$

a sufficient condition is

$$(5.8) \quad \|fl(\mathbf{x}_i^\top \hat{R}^{-1})\| = O(\|\mathbf{x}_i^\top\| / \|\hat{R}\|).$$

Since the general normwise bound for $\|\mathbf{y}^\top R^{-1}\|$ is $\|\mathbf{y}^\top R^{-1}\| \leq \|\mathbf{y}\| / \|R\| \kappa_2(R)$, the condition (5.8) is significantly more stringent when R is ill-conditioned.

Even so, as illustrated in the example above, in all our experiments with random matrices the condition (5.8) was satisfied with $\|\mathbf{x}_i^\top - \hat{\mathbf{q}}_i^\top \hat{R}\| / \|\mathbf{x}_i^\top\| < n\mathbf{u}$ for all i . We suspect that this is due to the observation known to experts that triangular linear systems are usually solved to much higher accuracy than the theoretical bound suggests [11, 20]. However, as with this classical observation, counterexamples do exist in our case: For example, taking R to be the Kahan matrix [11], which is an ill-conditioned triangular matrix known to have special properties, the bound $\|fl(\mathbf{y}^\top R^{-1})\| = O(\|\mathbf{y}^\top\| / \|R\|)$ is typically tight for a randomly generated \mathbf{y}^\top , which means (5.8) is significantly violated. In view of this we form X so that the Cholesky factor R of $X^\top X$ is the Kahan matrix. This can be done by taking $X = QR$ for an $m \times n$ orthogonal matrix Q . To introduce large variation in the row norms of X we construct Q as the orthogonal factor of a matrix as generated in the example above. For every such X with varying size n , (5.8) was still satisfied. Finally, we then appended a row of random elements at the bottom of X , with much smaller norm than the rest, and repeated the experiment. Now the row-wise residual for the last row was significantly larger than $O(\mathbf{u}\|\mathbf{x}_i^\top\|)$, indicating row-wise stability does not always hold. Employing pivoting in the Cholesky factorization did not improve the residual.

A referee has suggested more examples for which CholeskyQR2 fails to have row-wise backward stability. One example is the following: take X to be the off-diagonal parts of the 6×6 Hilbert matrix and set the (3, 3) element to $5e6$.

Experiments suggest nonetheless that cases in which CholeskyQR2 is not row-wise stable are extremely rare.

6. Conclusion. We performed a roundoff error analysis of the CholeskyQR2 algorithm for computing the QR decomposition of an $m \times n$ real matrix X , where $m \geq n$. We showed

that if X satisfies equation (2.1), the computed Q and R factors, which we denote by \hat{Z} and \hat{U} , respectively, satisfy the following error bounds.

$$(6.1) \quad \begin{aligned} \|\hat{Z}^\top \hat{Z} - I\|_F &\leq 6(mn\mathbf{u} + n(n+1)\mathbf{u}), \\ \|\hat{Z}\hat{U} - X\|_F/\|X\|_2 &\leq 5n^2\sqrt{n}\mathbf{u}. \end{aligned}$$

The bounds shown here are of a smaller order than the corresponding bounds for the Householder QR algorithm. Furthermore, it was shown that when only the Q factor is required, the right-hand side of equation (6.1) can be reduced to $O(n^2\mathbf{u})$. Numerical experiments support our theoretical analysis. CholeskyQR2 is also column-wise backward stable, as Householder QR. We also observed that the row-wise stability, which is a more stringent condition than the norm-wise stability shown by equation (6.1), nearly always holds in practice, though it cannot be proved theoretically.

In this paper, we focused on the stability of CholeskyQR2. Performance results of CholeskyQR2 on large scale parallel machines, along with comparison with other QR decomposition algorithms and detailed performance analysis, are given in our recent paper [8].

When the matrix is nearly square, it might be more efficient to partition the matrix into panels and apply the CholeskyQR2 algorithm to each panel successively. Development of such an algorithm remains as future work.

Acknowledgments. We thank Professor Nicholas J. Higham for suggesting the Kahan matrix for a possible counterexample of the row-wise backward stability of CholeskyQR2. We also thank Professor Masaaki Sugihara for valuable comments on the first version of our manuscript. We are grateful to the referees for their suggestions, especially for another counterexample for the row-wise backward stability, and a comment on columnwise stability.

REFERENCES

- [1] Z. BAI, J. DEMMEL, J. DONGARRA, A. RUHE, AND H. VAN DER VORST, eds., *Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide*, SIAM, Philadelphia, 2000.
- [2] G. BALLARD, J. DEMMEL, O. HOLTZ, AND O. SCHWARTZ, *Minimizing communication in numerical linear algebra*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 866–901.
- [3] R. CAR AND M. PARRINELLO, *Unified approach for molecular dynamics and density-functional theory*, Phys. Rev. Lett., 55 (1985), pp. 2471–2474.
- [4] J. W. DANIEL, W. B. GRAGG, L. KAUFMAN, AND G. W. STEWART, *Reorthogonalization and stable algorithms for updating the Gram-Schmidt QR factorization*, Math. Comp., 30 (1976), pp. 772–795.
- [5] J. DEMMEL, L. GRIGORI, M. HOEMMEN, AND J. LANGOU, *Communication-optimal parallel and sequential QR and LU factorizations*, SIAM J. Sci. Comput., 34 (2012), pp. 206–239.
- [6] E. ELMROTH AND F. G. GUSTAVSON, *Applying recursion to serial and parallel QR factorization leads to better performance*, IBM J. Res. Develop., 44 (2000), pp. 605–624.
- [7] T. FUKAYA, T. IMAMURA, AND Y. YAMAMOTO, *Performance analysis of the Householder-type parallel tall-skinny QR factorizations toward automatic algorithm selection*, in High Performance Computing for Computational Science—VECPAR 2014, M. Daydé, O. Marques, and K. Nakajima, eds., Lecture Notes in Computer Science, 8969, Springer, Cham, 2015, pp. 269–283.
- [8] T. FUKAYA, Y. NAKATSUKASA, Y. YANAGISAWA, AND Y. YAMAMOTO, *CholeskyQR2: a simple and communication-avoiding algorithm for computing a tall-skinny QR factorization on a large-scale parallel system*, in Proceedings of the 5th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems, IEEE Computer Society Press, Los Alamitos, 2014, pp. 31–38.
- [9] L. GIRAUD, J. LANGOU, M. ROZLOŽNÍK, AND J. VAN DEN ESHOF, *Rounding error analysis of the classical Gram-Schmidt orthogonalization process*, Numer. Math., 101 (2005), pp. 87–100.
- [10] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 4th ed., Johns Hopkins University Press, Baltimore, 2012.
- [11] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, 2nd ed., SIAM, Philadelphia, 2002.
- [12] J. IWATA, D. TAKAHASHI, A. OSHIYAMA, T. BOKU, K. SHIRAIISHI, S. OKADA, AND K. YABANA, *A massively-parallel electronic-structure calculations based on real-space density functional theory*, J. Comput. Phys., 229 (2010), pp. 2339–2363.

- [13] W. JALBY AND B. PHILIPPE, *Stability analysis and improvement of the block Gram-Schmidt algorithm*, SIAM J. Sci. Statist. Comput., 12 (1991), pp. 1058–1073.
- [14] S. J. LEON, Å. BJÖRCK, AND W. GANDER, *Gram-Schmidt orthogonalization: 100 years and more*, Numer. Linear Algebra Appl., 20 (2013), pp. 492–532.
- [15] B. N. PARLETT, *The Symmetric Eigenvalue Problem*, SIAM, Philadelphia, 1998.
- [16] M. ROZLOŽNÍK, F. OKULICKA-DŁUŽEWSKA, AND A. SMOKTUNOWICZ, *Cholesky-like factorization of symmetric indefinite matrices and orthogonalization with respect to bilinear forms*, Preprint NCMM/2013/31, Nečas Center for Mathematical Modeling, Prague, 2013.
- [17] Y. SAAD, *Numerical Methods for Large Eigenvalue Problems*, Manchester University Press, Manchester, 1992.
- [18] R. SCHREIBER AND C. F. VAN LOAN, *A storage-efficient WY representation for products of Householder transformations*, SIAM J. Sci. Statist. Comp., 10 (1989), pp. 53–57.
- [19] A. STATHOPOULOS AND K. WU, *A block orthogonalization procedure with constant synchronization requirements*, SIAM J. Sci. Comput., 23 (2002), pp. 2165–2182.
- [20] G. W. STEWART: *Introduction to Matrix Computations*, Academic Press, New York, 1973.
- [21] ———, *Matrix Algorithms, Vol. 1: Basic Decompositions*, SIAM, Philadelphia, 1998.
- [22] S. TOLEDO AND E. RABANI, *Very large electronic structure calculations using an out-of-core filter-diagonalization method*, J. Comput. Phys., 180 (2002), pp. 256–269.