

## EFFICIENT RANDOMIZED ALGORITHMS FOR THE FIXED-PRECISION LOW-RANK MATRIX APPROXIMATION\*

WENJIAN YU<sup>†</sup>, YU GU<sup>‡</sup>, AND YAOHANG LI<sup>§</sup>

**Abstract.** Randomized algorithms for low-rank matrix approximation are investigated, with the emphasis on the fixed-precision problem and computational efficiency for handling large matrices. The algorithms are based on the so-called QB factorization, where  $Q$  is an orthonormal matrix. First, a mechanism for calculating the approximation error in the *Frobenius norm* is proposed, which enables efficient adaptive rank determination for a large and/or sparse matrix. It can be combined with any QB-form factorization algorithm in which  $B$ 's rows are incrementally generated. Based on the blocked randQB algorithm by Martinsson and Voronin, this results in an algorithm called randQB\_EI. Then, we further revise the algorithm to obtain a pass-efficient algorithm, randQB\_FP, which is mathematically equivalent to the existing randQB algorithms and also suitable for the fixed-precision problem. Especially, randQB\_FP can serve as a single-pass algorithm for calculating leading singular values, under a certain condition. With large and/or sparse test matrices, we have empirically validated the merits of the proposed techniques, which exhibit remarkable speedup and memory saving over the blocked randQB algorithm. We have also demonstrated that the single-pass algorithm derived by randQB\_FP is much more accurate than an existing single-pass algorithm. And with data from a scenic image and an information retrieval application, we have shown the advantages of the proposed algorithms over the adaptive range finder algorithm for solving the fixed-precision problem.

**Key words.** adaptive rank determination, randomized algorithm, low-rank matrix approximation, pass-efficient algorithm, fixed-precision problem

**AMS subject classifications.** 15A18, 65F30, 65F15, 68W20, 60B20

**DOI.** 10.1137/17M1141977

**1. Introduction.** Low-rank matrix factorizations, like the partial singular value decomposition (SVD) and the rank-revealing QR factorization, play a crucial role in data analysis and scientific computing. In recent years, techniques based on randomization have been investigated for performing the computation and low-rank factorization of large matrices [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]. They involve the same or fewer floating-point operations (*flops*) than classical algorithms and are more efficient by exploiting modern computing architectures.

A basic idea of the randomized techniques is using random projection to approximate the dominant subspace of a matrix. For an  $m \times n$  matrix  $A$ , suppose the orthogonal basis vectors of this approximate subspace form an  $m \times k$  orthonormal matrix  $Q$ . Then, we have [1, 2]

$$(1) \quad A \approx QB,$$

where  $B$  is a  $k \times n$  matrix, and

$$(2) \quad B = Q^T A.$$

\*Received by the editors August 3, 2017; accepted for publication (in revised form) by P. Drineas July 17, 2018; published electronically August 30, 2018.

<http://www.siam.org/journals/simax/39-3/M114197.html>

**Funding:** This work was partially supported by the NSFC under grants 61872206, 6172811, and NSF under grant 1066471.

<sup>†</sup>BNRist, Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China (yu-wj@tsinghua.edu.cn).

<sup>‡</sup>Department of Computer Science and Technology and Institute for Interdisciplinary Information Sciences, Tsinghua University, Beijing 100084, China (guyu13@mails.tsinghua.edu.cn).

<sup>§</sup>Department of Computer Science, Old Dominion University, Norfolk, VA 23529 (yaochang@cs.odu.edu).

The basic randQB algorithm	The blocked randQB algorithm
<b>Input:</b> $\mathbf{A}$ , $k$ , $s$ . <b>Output:</b> $\mathbf{Q}$ , $\mathbf{B}$ . (1) $\mathbf{\Omega} = \text{randn}(n, k + s)$ (2) $\mathbf{Q} = \text{orth}(\mathbf{A}\mathbf{\Omega})$ (3) $\mathbf{B} = \mathbf{Q}^T \mathbf{A}$	<b>Input:</b> $\mathbf{A}$ , $\varepsilon$ , $b$ . <b>Output:</b> $\mathbf{Q}$ , $\mathbf{B}$ . (1) <b>for</b> $i = 1, 2, 3, \dots$ (2) $\mathbf{\Omega}_i = \text{randn}(n, b)$ (3) $\mathbf{Q}_i = \text{orth}(\mathbf{A}\mathbf{\Omega}_i)$ (4) $\mathbf{Q}_i = \text{orth}(\mathbf{Q}_i - \sum_{j=1}^{i-1} \mathbf{Q}_j \mathbf{Q}_j^T \mathbf{Q}_i)$ (5) $\mathbf{B}_i = \mathbf{Q}_i^T \mathbf{A}$ (6) $\mathbf{A} = \mathbf{A} - \mathbf{Q}_i \mathbf{B}_i$ (7) <b>if</b> $\ \mathbf{A}\  < \varepsilon$ <b>then stop</b> (8) <b>end for</b> (9) $\mathbf{Q} = [\mathbf{Q}_1 \ \dots \ \mathbf{Q}_i]$ ; $\mathbf{B} = [\mathbf{B}_1^T \ \dots \ \mathbf{B}_i^T]^T$ .
(a) The randQB algorithm	(b) The randQB_b algorithm

FIG. 1. The randomized algorithms for the QB factorization.

Standard factorizations, e.g., SVD, can be further performed on the smaller matrix  $\mathbf{B}$ , to obtain the low-rank factorizations of  $\mathbf{A}$ .

The approximation presented by (1) and (2) can also be regarded as a kind of low-rank factorization of  $\mathbf{A}$ , called *QB factorization* or *QB approximation* in this work. In [2], a basic randomized scheme for computing the QB approximation was presented, as shown in Figure 1(a). For producing close to optimal rank- $k$  approximation, the over-sampling scheme using a random Gaussian matrix  $\mathbf{\Omega}$  with  $k + s$  columns is employed, where  $s$  is a small integer. We use `randQB` to denote this algorithm.

Usually, the problem of low-rank matrix approximation falls into two categories:

- the *fixed-rank* problem, where the rank parameter  $k$  is given,
- the *fixed-precision* problem, where we seek  $\mathbf{Q}$  and  $\mathbf{B}$  with as small as possible size such that  $\|\mathbf{A} - \mathbf{QB}\| < \varepsilon$ , where  $\varepsilon$  is a given accuracy tolerance.

A blocked variant of the `randQB` algorithm proposed in [1], i.e., the `randQB_b` algorithm in Figure 1(b), is suitable for the fixed-precision problem. It incrementally builds the factors  $\mathbf{Q}$  and  $\mathbf{B}$  based on the combination of the `randQB` algorithm and the *blocked* Gram-Schmidt scheme and measures the approximation error by explicitly maintaining the residual matrix. However, it is inefficient or even fails for handling a large matrix, because maintaining the residual matrix is costly in runtime and memory usage.

In this work, the randomized algorithms for the fixed-precision problem are investigated considering their adaptability to large and/or sparse matrices. First, a mechanism is proposed for calculating the error of QB approximation in the *Frobenius norm* during the iterative process of building  $\mathbf{Q}$  and  $\mathbf{B}$ . It does not require maintaining the residual matrix (or updating matrix  $\mathbf{A}$ ) and thus avoids fill-in while handling a sparse  $\mathbf{A}$ . This mechanism is also applicable to other iterative computing procedures, e.g., the `trQRCP` algorithm [14]. Second, the algorithm is further revised to largely reduce the number of passes over matrix  $\mathbf{A}$ , in order to adapt the scenarios where the cost of accessing matrix  $\mathbf{A}$  is expensive. These techniques result in two algorithms called `randQB_EI` and `randQB_FP`, which inherit the merits of `randQB`/`randQB_b` algorithms and have extra benefits. Numerical experiments are carried out on a multicore computer to validate the efficiency and accuracy of the proposed algorithms for handling large or sparse matrices in practical scenarios. The results show that our `randQB_EI` and `randQB_FP` algorithms have up to 3X speedup and 3X memory

saving over an implementation of the `randQB_b` algorithm for dense matrices. They also exhibit up to 22X speedup over the implementation of the `randQB_b` algorithm for sparse matrices. Compared with the single-pass algorithm and the adaptive randomized range finder in [2], the proposed algorithms exhibit much better accuracy or avoid the large overestimation of the sizes of  $\mathbf{Q}$  and  $\mathbf{B}$ . For reproducibility, we have shared the MATLAB codes of the proposed algorithms and experimental data on [https://github.com/WenjianYu/randQB\\_auto](https://github.com/WenjianYu/randQB_auto).

**2. Technical preliminaries.** This section summarizes the background we need for presenting the proposed techniques. Throughout the paper, we measure vectors with their Euclidean norm. Two kinds of matrix norm are usually considered: Frobenius norm and spectral norm ( $l_2$ -norm). The spectral norm of a matrix is relatively difficult to calculate, though it is often more informative for noisy data [23]. The Frobenius norm  $\|\mathbf{A}\|_F = (\sum_{i,j} |\mathbf{A}(i,j)|^2)^{1/2}$  is easier for calculation and thus more widely used in data analysis and machine learning applications [10]. We measure matrices with their Frobenius norm by default. We also assume that all matrices are real valued, although the generalization to complex matrices is of no difficulty.

**2.1. Randomized algorithms.** To produce a rank- $k$  factorization for an  $m \times n$  matrix  $\mathbf{A}$ , with the basic `randQB` algorithm in Figure 1(a) one obtains an  $m \times l$  orthonormal matrix  $\mathbf{Q}$  and an  $l \times n$  matrix  $\mathbf{B}$ , where  $l > k$  due to oversampling. With this QB approximation, the standard factorizations can be efficiently computed. For example, the standard SVD algorithm can be performed on  $\mathbf{B}$ , which results in  $\mathbf{B} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^T$ . Then,

$$(3) \quad \mathbf{A} \approx \mathbf{QB} = \mathbf{Q}\tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^T.$$

The first  $k$  columns of matrices  $\mathbf{Q}\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{V}}$  and the  $k \times k$  upper-left submatrix of  $\tilde{\Sigma}$  approximate the rank- $k$  SVD factors of  $\mathbf{A}$ . Similarly, by changing the factorizations made on  $\mathbf{B}$ , one obtains the approximate QR factorization and CUR factorization, etc. [1, 4]. Notice that the accurate truncated SVD provides the optimal low-rank approximation [15]. However, computing accurate SVD of a large matrix is costly, and in many applications the optimality is not necessary. It is thus acceptable to compute the approximate low-rank factorizations for gains in computational efficiency.

The approximation error of the randomized algorithm is a random variable. The authors of [2] have studied the properties of the error in terms of spectral and Frobenius norms and given the bounds on their expectation and variance.

The truncated QR factorization with column pivoting can also be used for low-rank matrix approximation [11, 12]. Classical pivoted QR factorization has the disadvantage that it is hard to be parallelized or to take usage of BLAS-3 operation. Recently, the pivoted QR factorization was largely accelerated through utilizing a randomized technique, which achieves the efficiency comparable to the unpivoted QR decomposition [13, 14]. This makes the truncated pivoted QR factorization competitive for low-rank approximation. Notice that QR factorization can be regarded as a special case of the QB factorization. Therefore, one of the techniques proposed here (cf. section 3) could benefit the solution of the fixed-precision problem based on QR factorizations, as well.

The major computation of the `randQB` algorithm lies at the multiplication of  $\mathbf{A}$  and  $\mathbf{\Omega}$ . If  $\mathbf{A}$  is sparse or a structured matrix (often *implicitly* defined) for which matrix-vector products can be rapidly evaluated, the cost of multiplication can be

largely reduced (even to  $O(m+n)$  flops). The implicitly defined structured matrix often arises from physical problems, such as a discretized integral operator applied via the fast multipole method and similar techniques [24] and is sometimes referred to as an *implicit sparse matrix*.

For the fixed-precision problem, an adaptive randomized range finder was proposed in [2]. It employs the incremental sampling approach with a probabilistic error estimator to determine the size of  $\mathbf{Q}$  and  $\mathbf{B}$ . It is based on the statement that

$$(4) \quad \|\mathbf{A} - \mathbf{Q}\mathbf{Q}^T\mathbf{A}\|_2 \leq 10\sqrt{\frac{2}{\pi}} \max_{i=1,\dots,r} \|(\mathbf{A} - \mathbf{Q}\mathbf{Q}^T\mathbf{A})\boldsymbol{\omega}^{(i)}\|$$

with probability at least  $1 - 10^{-r}$  [2]. Here  $\|\cdot\|_2$  stands for the spectral norm,  $\boldsymbol{\omega}^{(i)}$  is a random vector, and  $r$  is a small integer, e.g.,  $r = 10$ . However, the error estimator often overestimates the approximation error, yielding much larger output matrices than what is necessary.

The `randQB_b` algorithm in [1] is based on the single-vector version of `randQB` algorithm with the Gram–Schmidt procedure, which allows one to construct the QB factorization and to evaluate its error step by step. In order to exploit *blocking* to attain high performance of linear algebraic computation, the algorithm is then converted to the blocked `randQB` algorithm in Figure 1(b), where step (6) is for calculating the residual matrix. Therefore, the `randQB_b` algorithm allows a precise error calculation for the fixed-precision problem. Notice that in step (4) there is a reorthogonalization operation which eases the accumulation of round-off error under floating point arithmetic.

The blocked `randQB` algorithm is able to produce  $\mathbf{Q}$  and  $\mathbf{B}$  in smaller sizes than the adaptive randomized range finder. However, explicitly maintaining the residual matrix brings extra time and memory cost if a large matrix  $\mathbf{A}$  is handled. This disadvantage becomes more serious if  $\mathbf{A}$  is also sparse, because the fill-in phenomena leads to a dense residual matrix.

In order to reveal the difference among the relevant randomized algorithms, we give a brief comparison of them for the fixed-rank problem, presented as Table 1. `svds` denotes the MATLAB built-in command for truncated SVD [21], which is based on a Krylov subspace iterative method. `trQRCP` denotes the randomized pivoted QR factorization [14]. `randQB_FP` is one of the contributions in this paper, which is a pass-efficient algorithm (cf. section 4). To depict the performance of the algorithms in the situation where the cost of accessing matrix entries is expensive (e.g.,  $\mathbf{A}$  are stored in slow memory) [2], we include the number of passes over matrix  $\mathbf{A}$  in Table 1. From the table, we see that the `randQB_FP` algorithm inherits the merits of the `randQB` algorithm and is also suitable for the fixed-precision problem.

TABLE 1

*Properties of the randomized algorithms for the fixed-rank matrix approximation. ( $P$  is a small integer for power scheme, and  $b$  is the block size.)*

	<code>svds</code>	<code>randQB</code>	<code>randQB_b</code>	<code>trQRCP</code>	<code>randQB_FP</code>
Computational efficiency	low	high	high	high	high
Adaptive rank determination	no	no	yes	yes <sup>†</sup>	yes
Number of passes over matrix	$\alpha k^*$	2 or $2+2P$	$4k/b$	$3k/b + 1$	1 or $1+2P^{\ddagger}$

\*  $\alpha$  is a number larger than 1.

† This was not mentioned in [14]. For more detail, please see Remark 3.2.

‡ Suppose matrix  $\mathbf{A}$  is stored in the row-major format.

**3. An efficient Frobenius-norm error indicator and its application.** In this section, we first propose an error indicator for measuring the approximation error in the Frobenius norm and an efficient framework for solving the fixed-precision problem of QB factorization. Then, the `randQB_EI` algorithm is derived based on the blocked `randQB` algorithm. Finally, the accuracy and validity of the error indicator in floating point arithmetic is discussed.

**3.1. An error indicator.** We first give a theorem regarding the Frobenius norm of the error of the QB approximation.

**THEOREM 1.** *Let  $\mathbf{A}$  be an  $m \times n$  matrix.  $\mathbf{Q}$  denotes an  $m \times k$  orthonormal matrix ( $k < m$ ), and  $\mathbf{B} = \mathbf{Q}^T \mathbf{A}$ . Then,*

$$(5) \quad \|\mathbf{A} - \mathbf{QB}\|_F^2 = \|\mathbf{A}\|_F^2 - \|\mathbf{B}\|_F^2.$$

*Proof.* Due to the property of the Frobenius norm, we know for any matrix  $\mathbf{M}$ ,

$$(6) \quad \|\mathbf{M}\|_F^2 = \text{tr}(\mathbf{M}^T \mathbf{M}),$$

where  $\text{tr}(\cdot)$  calculates the trace of a matrix. Because  $\mathbf{Q}$  is orthonormal and  $\mathbf{B} = \mathbf{Q}^T \mathbf{A}$ ,

$$(7) \quad \begin{aligned} (\mathbf{A} - \mathbf{QB})^T (\mathbf{A} - \mathbf{QB}) &= (\mathbf{A} - \mathbf{QQ}^T \mathbf{A})^T (\mathbf{A} - \mathbf{QQ}^T \mathbf{A}) \\ &= \mathbf{A}^T \mathbf{A} - 2\mathbf{A}^T \mathbf{QQ}^T \mathbf{A} + \mathbf{A}^T \mathbf{QQ}^T \mathbf{QQ}^T \mathbf{A} \\ &= \mathbf{A}^T \mathbf{A} - \mathbf{A}^T \mathbf{QQ}^T \mathbf{A} \\ &= \mathbf{A}^T \mathbf{A} - \mathbf{B}^T \mathbf{B}. \end{aligned}$$

Now, applying the trace operation to both sides of (7), and according to (6), we obtain (5).  $\square$

Theorem 1 suggests that, if we have access to  $\mathbf{B}$  and  $\|\mathbf{A}\|_F$  is known a priori, we can calculate the error of QB approximation without referring to the residual matrix. This leads to a framework for solving the fixed-precision problem, presented as Algorithm 1. It suits any algorithm that incrementally generates the rows of  $\mathbf{B}$ , including `randQB_b`, `trQRCP`, and the `randQB_FP` algorithm presented in section 4. Below we prove the correctness of Algorithm 1.

---

**Algorithm 1** A framework for solving the fixed-precision QB factorization problem

---

**Input:** an  $m \times n$  matrix  $\mathbf{A}$ ; desired accuracy tolerance  $\varepsilon$ .

**Output:**  $\mathbf{Q}$ ,  $\mathbf{B}$ , such that  $\|\mathbf{A} - \mathbf{QB}\|_F < \varepsilon$ .

```

1:  $\mathbf{Q} = [ ]$ ;  $\mathbf{B} = [ ]$ ; # empty matrices
2:  $E = \|\mathbf{A}\|_F^2$  # initialization of the error indicator
3: for  $i = 1, 2, 3, \dots$ , do
4:   Generate  $\mathbf{Q}_i$  and  $\mathbf{B}_i$ , s.t.  $[\mathbf{Q}, \mathbf{Q}_i]$  is orthonormal and  $\mathbf{B}_i = \mathbf{Q}_i^T \mathbf{A}$ .
5:    $\mathbf{Q} = [\mathbf{Q}, \mathbf{Q}_i]$ 
6:    $\mathbf{B} = \begin{bmatrix} \mathbf{B} \\ \mathbf{B}_i \end{bmatrix}$ 
7:    $E = E - \|\mathbf{B}_i\|_F^2$  # update the error indicator
8:   if  $E < \varepsilon^2$  then stop
9: end for

```

---

THEOREM 2. After the  $i$ th iteration of loop in Algorithm 1 is executed,

$$(8) \quad E^{(i)} = \|\mathbf{A} - \mathbf{Q}^{(i)}\mathbf{B}^{(i)}\|_{\text{F}}^2,$$

where  $E^{(i)}$ ,  $\mathbf{Q}^{(i)}$ , and  $\mathbf{B}^{(i)}$  denote the values of  $E$ ,  $\mathbf{Q}$ , and  $\mathbf{B}$  after the  $i$ th iteration of the loop is executed, respectively.

*Proof.* Based on step 7 of Algorithm 1 and the property of the Frobenius norm,

$$(9) \quad E^{(i)} = \|\mathbf{A}\|_{\text{F}}^2 - \sum_{j=1}^i \|\mathbf{B}_j\|_{\text{F}}^2 = \|\mathbf{A}\|_{\text{F}}^2 - \|\mathbf{B}^{(i)}\|_{\text{F}}^2.$$

Because  $\mathbf{Q}^{(i)}$  is an orthonormal matrix and  $\mathbf{B}^{(i)} = \mathbf{Q}^{(i)T}\mathbf{A}$ , (8) can be obtained by applying Theorem 1.  $\square$

From Theorem 2, we see that  $E$  in Algorithm 1 equals to the square of the approximation error. It is an error indicator updated through calculating the Frobenius norm of  $\mathbf{B}_i$ . This yields two benefits: we no longer need to maintain the residual  $\mathbf{A} - \mathbf{Q}\mathbf{B}$ , and the approximation error can be calculated with very small cost.

Steps 7 and 8 in Algorithm 1 can be replaced by a row-by-row calculation scheme.

---

```

7a: for  $j = 1, 2, \dots, m_i$  do      #  $m_i$  denotes the number of rows of  $\mathbf{B}_i$ 
7b:    $E = E - \|\mathbf{B}_i(j, :)\|_{\text{F}}^2$ 
8a:   if  $E < \varepsilon^2$  then
8b:     remove  $\mathbf{B}_i(j + 1 : m_i, :)$  from  $\mathbf{B}$ ; remove  $\mathbf{Q}_i(:, j + 1 : m_i)$  from  $\mathbf{Q}$ 
8c:     stop
8d:   end if
8e: end for

```

---

This adds negligible cost but allows us to determine the certain row of  $\mathbf{B}_i$  where the accuracy tolerance is just attained. It makes the column (row) number of outputted  $\mathbf{Q}$  ( $\mathbf{B}$ ) an arbitrary integer, instead of a multiple of block size  $b$  in `randQB_b` algorithm.

*Remark 3.1.* If  $\mathbf{A}$  is an implicit sparse matrix, the proposed framework needs more effort for calculating its Frobenius norm. The columns of  $\mathbf{A}$  can be solved by multiplying  $\mathbf{A}$  with the canonical basis vectors. Therefore, the cost for calculating  $\|\mathbf{A}\|_{\text{F}}$  will be  $O(n(m+n))$  flops, if each matrix-vector product costs  $O(m+n)$  flops. This might be affordable, as the Frobenius norm is a property of the matrix and we need to compute it just once.

**3.2. The `randQB_EI` algorithm.** The combination of Algorithm 1 and the `randQB_b` algorithm results in Algorithm 2 (called `randQB_EI`), whose steps 4–7 replace step 4 in Algorithm 1. Notice that step 7 in Algorithm 2 looks the same as step (5) in the `randQB_b` algorithm but is actually different. And, step (3) in the `randQB_b` algorithm becomes step 5 in Algorithm 2, which is the blocked Gram–Schmidt orthogonalization of  $\mathbf{A}\mathbf{\Omega}_i$  (notice  $\mathbf{B} = \mathbf{Q}^T\mathbf{A}$ ).

Due to Theorem 2 and the orthogonality of  $\mathbf{Q}$ , we have the following proposition.

PROPOSITION 1. The `randQB_EI` algorithm (Algorithm 2) is equivalent to the `randQB_b` algorithm, when executed in exact arithmetic.

Assuming that multiplying two dense matrices of sizes  $m \times n$  and  $n \times l$  costs  $C_{mm}mnl$  flops, and performing an economic QR factorization of an  $m \times n$  dense matrix

---

**Algorithm 2** The randQB\_EI algorithm for the fixed-precision problem

---

**Input:** an  $m \times n$  matrix  $\mathbf{A}$ ; desired accuracy tolerance  $\varepsilon$ ; block size  $b$ .

**Output:**  $\mathbf{Q}$ ,  $\mathbf{B}$ , such that  $\|\mathbf{A} - \mathbf{QB}\|_{\mathbb{F}} < \varepsilon$ .

```

1:  $\mathbf{Q} = []$ ;  $\mathbf{B} = []$ ;
2:  $E = \|\mathbf{A}\|_{\mathbb{F}}^2$ 
3: for  $i = 1, 2, 3, \dots$ , do
4:    $\Omega_i = \text{randn}(n, b)$ 
5:    $\mathbf{Q}_i = \text{orth}(\mathbf{A}\Omega_i - \mathbf{Q}(\mathbf{B}\Omega_i))$ 
6:    $\mathbf{Q}_i = \text{orth}(\mathbf{Q}_i - \mathbf{Q}(\mathbf{Q}^T \mathbf{Q}_i))$            # reorthogonalization
7:    $\mathbf{B}_i = \mathbf{Q}_i^T \mathbf{A}$                                # no need to calculate  $\mathbf{Q}_i^T(\mathbf{A} - \mathbf{QB})$ 
8:    $\mathbf{Q} = [\mathbf{Q}, \mathbf{Q}_i]$ 
9:    $\mathbf{B} = \begin{bmatrix} \mathbf{B} \\ \mathbf{B}_i \end{bmatrix}$ 
10:   $E = E - \|\mathbf{B}_i\|_{\mathbb{F}}^2$ 
11:  if  $E < \varepsilon^2$  then stop
12: end for

```

---

costs  $C_{qr}mn \min(m, n)$  flops, we can analyze the flop counts of the relevant algorithms and compare their performance for handling a *dense*  $\mathbf{A}$ . If  $T_{\text{randQB}}$  and  $T_{\text{randQB}_b}$  denote the runtime of the algorithms randQB and randQB\_b, respectively, [1] then

$$(10) \quad T_{\text{randQB}} \sim 2C_{mm}mnl + C_{qr}ml^2,$$

$$(11) \quad T_{\text{randQB}_b} \sim 3C_{mm}mnl + C_{mm}ml^2 + \frac{2}{t}C_{qr}ml^2,$$

where  $l$  is the number of columns in the resulting matrix  $\mathbf{Q}$ , and  $t$  satisfies  $l = tb$ . Note that  $b$  is practically much smaller than  $l$  (say,  $b = 10$  or  $20$ ), although the optimal choice of block size depends strongly on what hardware is used [1].

For the randQB\_EI algorithm, the runtime can be similarly depicted:

$$(12) \quad T_{\text{randQB}_{EI}} \sim 2C_{mm}mnl + C_{mm}(2m + n)l^2 + \frac{2}{t}C_{qr}ml^2.$$

Because  $l$  is usually much smaller than  $m$  and  $n$ , we see that the flop count of randQB\_EI is about 2/3 of that of randQB\_b and is comparable to that of the basic randQB algorithm. Notice that  $C_{qr}$  is several times larger than  $C_{mm}$ . So, the flop count of randQB\_EI could be smaller than that of randQB in the situation where  $t$  is a large number. If  $\mathbf{A}$  is sparse, this would more possibly happen, because the randQB algorithm loses the benefit brought by the BLAS-3 operation.

The advantage of randQB\_EI over randQB\_b becomes more prominent if  $\mathbf{A}$  is a sparse matrix. With the proposed error indicator, we no longer need the residual matrix. In contrast, it is always a dense matrix in the randQB\_b algorithm and costs much larger memory and induces many more computations.

*Remark 3.2.* Algorithm 1 can also be combined with the trQRCP algorithm [14]. Although trQRCP generates a column permutation matrix as well, it does not affect the Frobenius norm of each partial or the entire matrix of  $\mathbf{B}$ . Therefore, this will produce another efficient algorithm for adaptive low-rank matrix approximation, which also adapts to sparse matrices.

**3.3. Floating point arithmetic.** Below we discuss the accuracy of the error indicator in floating point arithmetic. We use  $\mathcal{E}(\cdot)$  and  $\mathcal{E}_r(\cdot)$  to denote the functions of error and relative error, respectively.

As the error indicator  $E = \|\mathbf{A}\|_{\mathbb{F}}^2 - \|\mathbf{B}\|_{\mathbb{F}}^2$ , its calculated value  $\hat{E}$  cannot be accurate when  $E$  is very small, due to the cancellation in calculation. In floating-point arithmetic, the machine precision  $\epsilon_{mach}$  characterizes the maximum relative error of converting a real number to its floating-point representation, i.e.,

$$(13) \quad \forall x, \quad |\mathcal{E}_r(x)| \leq \epsilon_{mach}.$$

According to the definition of Frobenius norm,  $\|\mathbf{A}\|_{\mathbb{F}}^2$  is the summation of squares of matrix entries. Therefore, the relative error of  $\|\mathbf{A}\|_{\mathbb{F}}^2$  is bounded by  $2\epsilon_{mach}$ . The same thing applies to  $\|\mathbf{B}\|_{\mathbb{F}}^2$ . So,

$$(14) \quad \begin{aligned} |\mathcal{E}(E)| &= |\mathcal{E}(\|\mathbf{A}\|_{\mathbb{F}}^2 - \|\mathbf{B}\|_{\mathbb{F}}^2)| \leq |\mathcal{E}(\|\mathbf{A}\|_{\mathbb{F}}^2)| + |\mathcal{E}(\|\mathbf{B}\|_{\mathbb{F}}^2)| \\ &\leq 2\epsilon_{mach}(\|\mathbf{A}\|_{\mathbb{F}}^2 + \|\mathbf{B}\|_{\mathbb{F}}^2) \\ &< 4\epsilon_{mach}\|\mathbf{A}\|_{\mathbb{F}}^2. \end{aligned}$$

If we want to guarantee that  $E$  has a relative error no more than  $\delta$ , i.e.,  $|\mathcal{E}(E)| \leq \delta E$ , we shall enforce

$$(15) \quad 4\epsilon_{mach}\|\mathbf{A}\|_{\mathbb{F}}^2 \leq \delta E.$$

This means the preset accuracy tolerance  $\varepsilon$ , which is larger than  $\sqrt{E}$  at the termination of the algorithm, should satisfy

$$(16) \quad \varepsilon > \sqrt{E} \geq \sqrt{\frac{4\epsilon_{mach}\|\mathbf{A}\|_{\mathbb{F}}^2}{\delta}} = \sqrt{\frac{4\epsilon_{mach}}{\delta}}\|\mathbf{A}\|_{\mathbb{F}}.$$

So, we obtain the following theorem.

**THEOREM 3.** *Suppose matrix  $\mathbf{A}$  and accuracy tolerance  $\varepsilon$  are the input to the `randQB_EI` algorithm. If  $\varepsilon > \sqrt{\frac{4\epsilon_{mach}}{\delta}}\|\mathbf{A}\|_{\mathbb{F}}$ , the relative error of the calculated error indicator  $E$  must be no more than  $\delta$ , e.g., if  $\varepsilon > 2.1 \times 10^{-7}\|\mathbf{A}\|_{\mathbb{F}}$ , the error of  $E$  is within 1% in the double-precision floating arithmetic, where  $\epsilon_{mach} \approx 1.11 \times 10^{-16}$ .*

Notice that, as an error indicator for the fixed-precision problem,  $E$  should have sufficient accuracy (e.g., with relative error  $\sim 1\%$  or less). Otherwise, the outputted QB factorization would not satisfy the preset accuracy tolerance.

Besides, the orthogonality of  $\mathbf{Q}$  also affects the accuracy of  $E$ . As the number of columns in  $\mathbf{Q}$  increases, its orthogonality gradually degrades. This issue occurs for  $\mathbf{Q}$  produced either by a single run of QR factorization (based on Householder transformation) or by a Gram-Schmidt procedure followed by the reorthogonalization step. We will investigate its effect in the following experiment.

An  $n \times n$  matrix  $\mathbf{A}$  is constructed to have singular values according to a decaying exponential. Two instances are tested, with singular value  $\sigma_j = e^{-j/20}$  and  $\sigma_j = e^{-j/200}$ ,  $j = 1, 2, \dots$ , respectively. The results obtained from executing `randQB_EI` algorithm ( $b = 10$ ) and `randQB` algorithm with different values of rank parameter  $l$  are shown in Figure 2. Note that for some large value of  $l$ , the error indicator can be of negative value, such that it cannot be drawn in the log-scale plot. From the figure, we can validate the correctness of Theorem 3. Providing that the square of



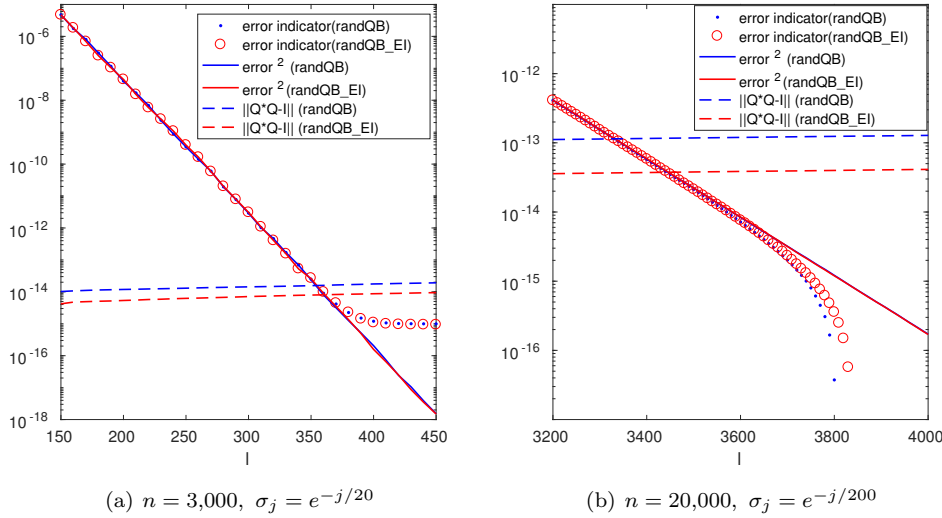


FIG. 2. Normalized error indicator  $E/\|A\|_F^2$ , error square  $\|A - QB\|_F^2/\|A\|_F^2$ , and the loss of orthogonality of  $Q$  vs.  $Q$ 's column number  $l$  in randQB algorithm and randQB\_EI algorithm.

error  $\|A - QB\|_F^2/\|A\|_F^2 > (2.1 \times 10^{-7})^2 = 4.4 \times 10^{-14}$ , the error indicator matches the square of error very well. This holds even when  $l$  is larger than 3,000, which corresponds to the situation with larger accumulated round-off error. In Figure 2, the value of  $\|Q^T Q - I\|_\infty$ <sup>1</sup> is also plotted, which reveals the loss of orthogonality of  $Q$ . The results show that this issue is not severe in both the randQB and randQB\_EI algorithms, although it gradually increases as the columns of  $Q$  are increased.

*Remark 3.3.* Theorem 3 suggests the limitation of the error indicator and the proposed algorithms for the fixed-precision problem. It means that the efficient framework for adaptive rank determination would not work, in double-precision floating arithmetic, for the problem with the accuracy tolerance  $\varepsilon$  less than  $2.1 \times 10^{-7} \|A\|_F$ .

**4. A pass-efficient algorithm for the QB factorization.** The technique in last section efficiently solves the fixed-precision problem measured in the Frobenius norm. However, the randQB\_EI algorithm is not suitable for the scenarios where accessing matrix entries is very expensive (e.g.,  $A$  is too large and has to be stored on hard disk), as it visits  $A$  for considerable time (see Table 1). In this section, we propose a pass-efficient algorithm which largely reduces the number of passes over  $A$ .

**4.1. The version without reorthogonalization.** We first consider the fixed-rank problem where the rank parameter  $k$  is given. A preliminary pass-efficient algorithm (presented as Algorithm 3) can be derived from the randQB\_b or randQB\_EI algorithm. The steps correspond to those of randQB\_EI in Algorithm 2, one by one, except that the applications of  $A$  are moved out of the loop, and step 6 for reorthogonalization is ignored. In Algorithm 3, step 6 is the same as step 4 of randQB\_EI algorithm, and steps 7 and 8 correspond to step 5 of randQB\_EI. In step 9,  $R_i^{-T}$  is the inverse of the transpose of an upper triangular matrix  $R_i$ . This step can be regarded as solving linear equations with the coefficient matrix  $R_i^T$ , which is implemented by the “\” operator in MATLAB.

<sup>1</sup>This measure of loss of orthogonality follows Cleve Moler’s blog with the title “Compare Gram-Schmidt and Householder Orthogonalization Algorithms” posted on October 17, 2016.

---

**Algorithm 3** A preliminary pass-efficient algorithm for fixed-rank QB factorization

---

**Input:**  $A \in \mathbb{R}^{m \times n}$ ,  $k$ ,  $s$ , block size  $b$ .

**Output:**  $Q$ ,  $B$ .

```

1:  $Q = []$ ;  $B = []$ ;  $l = k + s$ ;
2:  $\Omega = \text{randn}(n, l)$ 
3:  $G = A\Omega$ 
4:  $H = A^T G$ 
5: for  $i = 1, 2, 3, \dots, l/b$  do
6:    $\Omega_i = \Omega(:, (i-1)b + 1 : ib)$ 
7:    $Y_i = G(:, (i-1)b + 1 : ib) - Q(B\Omega_i)$ 
8:    $[Q_i, R_i] = \text{qr}(Y_i)$ 
9:    $B_i = R_i^{-T}(H(:, (i-1)b + 1 : ib)^T - \Omega_i^T B^T B)$ 
10:   $Q = [Q, Q_i]$ 
11:   $B = \begin{bmatrix} B \\ B_i \end{bmatrix}$ 
12: end for

```

---

Because  $Y_i = Q_i R_i$ ,

$$(17) \quad Q_i = Y_i R_i^{-1}.$$

Substituting it into step 7 of the randQB\_EI algorithm, we have

$$\begin{aligned}
 B_i &= (Y_i R_i^{-1})^T A \\
 &= R_i^{-T} Y_i^T A \\
 (18) \quad &= R_i^{-T} (G(:, (i-1)b + 1 : ib)^T - \Omega_i^T B^T Q^T) A \\
 &= R_i^{-T} ((A^T G(:, (i-1)b + 1 : ib))^T - \Omega_i^T B^T Q^T A) \\
 &= R_i^{-T} (H(:, (i-1)b + 1 : ib)^T - \Omega_i^T B^T B).
 \end{aligned}$$

This means that step 9 of Algorithm 3 is equivalent to step 7 of randQB\_EI. Therefore, we obtain the following proposition.

**PROPOSITION 2.** *Algorithm 3 is mathematically equivalent to the fixed-rank version of the randQB\_EI algorithm without reorthogonalization.*

**4.2. The version with reorthogonalization.** In reality, the loss of orthogonality among the columns of  $\{Q_1, Q_2, \dots\}$  occurs due to the accumulation of round-off error. This means the reorthogonalization step, i.e., step 6 in Algorithm 2, cannot be ignored. Below we derive the revisions to Algorithm 3 to address this issue.

The reorthogonalization step can be expressed as

$$(19) \quad \tilde{Q}_i \tilde{R}_i = Q_i - Q Q^T Q_i,$$

where  $\tilde{Q}_i \neq Q_i$  and  $\tilde{R}_i \neq I$  due to the loss of orthogonality.  $\tilde{Q}_i$  is better orthogonal to the previously generated  $\{Q_1, Q_2, \dots, Q_{i-1}\}$  than  $Q_i$ . Now, we need to derive a formula for calculating  $B_i$  which does not involve  $A$  explicitly. Based on (17),

$$(20) \quad \tilde{Q}_i = (I - Q Q^T) Y_i R_i^{-1} \tilde{R}_i^{-1}.$$

So, the corresponding formula for  $\mathbf{B}_i$  is

$$\begin{aligned}
 \tilde{\mathbf{B}}_i &= \tilde{\mathbf{Q}}_i^T \mathbf{A} \\
 (21) \quad &= (\tilde{\mathbf{R}}_i \mathbf{R}_i)^{-T} \mathbf{Y}_i^T (\mathbf{I} - \mathbf{Q} \mathbf{Q}^T) \mathbf{A} \\
 &= (\tilde{\mathbf{R}}_i \mathbf{R}_i)^{-T} (\Omega_i^T \mathbf{A}^T - \Omega_i^T \mathbf{B}^T \mathbf{Q}^T) (\mathbf{A} - \mathbf{Q} \mathbf{B}),
 \end{aligned}$$

where the formula of  $\mathbf{Y}_i$  and the equality  $\mathbf{B} = \mathbf{Q}^T \mathbf{A}$  is made use of. The product of the last two brackets can be further simplified.

$$\begin{aligned}
 (22) \quad &(\Omega_i^T \mathbf{A}^T - \Omega_i^T \mathbf{B}^T \mathbf{Q}^T) (\mathbf{A} - \mathbf{Q} \mathbf{B}) \\
 &= \mathbf{H}_i^T - \mathbf{G}_i^T \mathbf{Q} \mathbf{B} - \Omega_i^T \mathbf{B}^T \mathbf{B} + \Omega_i^T \mathbf{B}^T \mathbf{Q}^T \mathbf{Q} \mathbf{B} \\
 &= \mathbf{H}_i^T - \mathbf{Y}_i^T \mathbf{Q} \mathbf{B} - \Omega_i^T \mathbf{B}^T \mathbf{B}.
 \end{aligned}$$

In the deduction,  $\mathbf{G}_i$  and  $\mathbf{H}_i$  denote  $\mathbf{G}(:, (i-1)b+1 : ib)$  and  $\mathbf{H}(:, (i-1)b+1 : ib)$  in Algorithm 3, respectively. Therefore,

$$(23) \quad \tilde{\mathbf{B}}_i = (\tilde{\mathbf{R}}_i \mathbf{R}_i)^{-T} (\mathbf{H}_i^T - \mathbf{Y}_i^T \mathbf{Q} \mathbf{B} - \Omega_i^T \mathbf{B}^T \mathbf{B}).$$

Based on (19) and (23), we can derive the version with reorthogonalization for Algorithm 3. We just need to replace the ninth step with the following steps.

$$\begin{aligned}
 &\underline{9a: [\mathbf{Q}_i, \tilde{\mathbf{R}}_i] = \text{qr}(\mathbf{Q}_i - \mathbf{Q}(\mathbf{Q}^T \mathbf{Q}_i))} \\
 &9b: \mathbf{R}_i = \tilde{\mathbf{R}}_i \mathbf{R}_i \\
 &\underline{9c: \mathbf{B}_i = \mathbf{R}_i^{-T} (\mathbf{H}(:, (i-1)b+1 : ib)^T - \mathbf{Y}_i^T \mathbf{Q} \mathbf{B} - \Omega_i^T \mathbf{B}^T \mathbf{B})}
 \end{aligned}$$

Notice that  $\mathbf{Q}_i$  and  $\mathbf{B}_i$  are overwritten to stand for  $\tilde{\mathbf{Q}}_i$  and  $\tilde{\mathbf{B}}_i$ , respectively. Based on Proposition 2 and the above deduction, we see that the pass-efficient algorithm with reorthogonalization is also mathematically equivalent to the fixed-rank versions of randQB\_EI and randQB\_b algorithms.

This algorithm with fewer passes over  $\mathbf{A}$  is called randQB\_FP. Based on the notation in section 3, its flop count analysis is as follows.

$$(24) \quad T_{\text{randQB\_FP}} \sim 2C_{mm} mnl + 2C_{mm} (m+n)l^2 + \frac{2}{t} C_{qr} ml^2,$$

where  $t$  satisfies  $l = tb$ . Compared with the randQB\_EI algorithm, the randQB\_FP algorithm has a slightly larger flop count. However, while handling a dense  $\mathbf{A}$  its actual runtime may be shorter because it lumps the multiplications with  $\mathbf{A}$ .

*Remark 4.1.* The round-off error may affect the accuracy of  $\mathbf{B}_i$ , and it increases as the number of iterations increases. However, this may not be an issue for practical low-rank approximation problems. In section 5, we will present numerical experiments to validate the effectiveness of the randQB\_FP algorithm, which shows it works very well for many applications with the rank parameter up to several thousands or the relative Frobenius-norm error of approximation as small as  $10^{-7}$ .

*Remark 4.2.* The randQB\_FP algorithm can derive a single-pass algorithm, if matrix  $\mathbf{A}$  is stored in the row-major format or is revealed row(s) by row(s). Suppose  $\mathbf{A}_{i,:}$  denotes the  $i$ th row of  $\mathbf{A}$ . With it we have the  $i$ th row of  $\mathbf{G}$ ,  $\mathbf{G}_{i,:} = \mathbf{A}_{i,:} \mathbf{\Omega}$ . Then, because  $\mathbf{H} = \mathbf{A}^T \mathbf{G} = \sum_i \mathbf{A}_{i,:}^T \mathbf{G}_{i,:}$ , the  $i$ th term in this summation can be obtained. With all rows of  $\mathbf{A}$ , in this way we can accomplish steps 3 and 4 in the randQB\_FP algorithm with only one pass over  $\mathbf{A}$ . It should be pointed out that this algorithm is not a general single-pass algorithm, as it has the restriction of the matrix. For more general single-pass algorithms for low-rank matrix approximation, please refer to the recent work [22].

**4.3. The inclusion of power iteration scheme.** The error of randomized QB factorization could be large for the matrix whose singular value decays slowly [2]. So, the *power iteration scheme* has been proposed to relieve this weakness [1, 2, 7]. Conceptually, the power iteration means replacing  $\mathbf{A}$  with  $(\mathbf{A}\mathbf{A}^T)^P\mathbf{A}$ , where  $P$  is an integer. However, in floating-point computation any singular components smaller than  $\sigma_1\epsilon_{mach}^{1/(2P+1)}$  will be lost. This makes the orthonormalization steps after the applications of  $\mathbf{A}$  and  $\mathbf{A}^T$  necessary, and  $P$  should not be set to a large number. Incorporating the power iteration, we have the randQB\_FP algorithm for the fixed-precision problem presented as Algorithm 4, where the error indicator  $E$  is utilized.

---

**Algorithm 4** The randQB\_FP with power scheme for the fixed-precision problem

---

**Input:**  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , desired accuracy tolerance  $\epsilon$ , block size  $b$ , power parameter  $P$ .

**Output:**  $\mathbf{Q}$ ,  $\mathbf{B}$ , such that  $\|\mathbf{A} - \mathbf{Q}\mathbf{B}\|_F \leq \epsilon$ .

```

1:  $\mathbf{Q} = []$ ;  $\mathbf{B} = []$ ;
2:  $\mathbf{\Omega} = \text{randn}(n, \tilde{l})$ , where  $\tilde{l}$  is a sufficiently large number.
3: for  $i = 1 : P$  do
4:    $\mathbf{G} = \text{orth}(\mathbf{A}\mathbf{\Omega})$ 
5:    $\mathbf{\Omega} = \text{orth}(\mathbf{A}^T\mathbf{G})$ 
6: end for
7:  $\mathbf{G} = \mathbf{A}\mathbf{\Omega}$ 
8:  $\mathbf{H} = \mathbf{A}^T\mathbf{G}$ 
9:  $E = \|\mathbf{A}\|_F^2$ 
10: for  $i = 1, 2, 3, \dots$ , do
11:    $\mathbf{\Omega}_i = \mathbf{\Omega}(:, (i-1)b + 1 : ib)$ 
12:    $\mathbf{Y}_i = \mathbf{G}(:, (i-1)b + 1 : ib) - \mathbf{Q}(\mathbf{B}\mathbf{\Omega}_i)$ 
13:    $[\mathbf{Q}_i, \mathbf{R}_i] = \text{qr}(\mathbf{Y}_i)$ 
14:    $[\mathbf{Q}_i, \tilde{\mathbf{R}}_i] = \text{qr}(\mathbf{Q}_i - \mathbf{Q}(\mathbf{Q}^T\mathbf{Q}_i))$ 
15:    $\mathbf{R}_i = \tilde{\mathbf{R}}_i\mathbf{R}_i$ 
16:    $\mathbf{B}_i = \mathbf{R}_i^{-T}(\mathbf{H}(:, (i-1)b + 1 : ib)^T - \mathbf{Y}_i^T\mathbf{Q}\mathbf{B} - \mathbf{\Omega}_i^T\mathbf{B}^T\mathbf{B})$ 
17:    $\mathbf{Q} = [\mathbf{Q}, \mathbf{Q}_i]$ 
18:    $\mathbf{B} = \begin{bmatrix} \mathbf{B} \\ \mathbf{B}_i \end{bmatrix}$ 
19:    $E = E - \|\mathbf{B}_i\|_F^2$ 
20:   if  $\sqrt{E} < \epsilon$  then stop
21: end for

```

---

In Algorithm 4, a sufficient large value of  $\tilde{l}$  should be set according to problem-specific experience and the concern of computing time. If the set  $\tilde{l}$  is not large enough for attaining the specified accuracy criterion, we need to regenerate the  $\mathbf{\Omega}$  matrix and rerun the algorithm to collect additional columns/rows of  $\mathbf{Q}$  and  $\mathbf{B}$ . This situation and the power scheme both increase the number of passes over  $\mathbf{A}$ . But compared to other algorithms for the fixed-precision problem, this fixed-precision randQB\_FP algorithm involves much fewer passes over  $\mathbf{A}$ .

**5. Numerical results.** In this section we compare the proposed algorithms against several existing algorithms in terms of execution time, memory usage, and accuracy. All experiments are carried out on a Linux server with two 12-core Intel Xeon E5-2630 CPUs at 2.30 GHz and 32GB RAM. For comparison of speed, the proposed algorithms have been implemented in C based on the codes shared by the authors of [1, 16]. The program is coded with OpenMP derivatives and compiled with the Intel ICC compiler with MKL libraries [17], to take full advantage of the

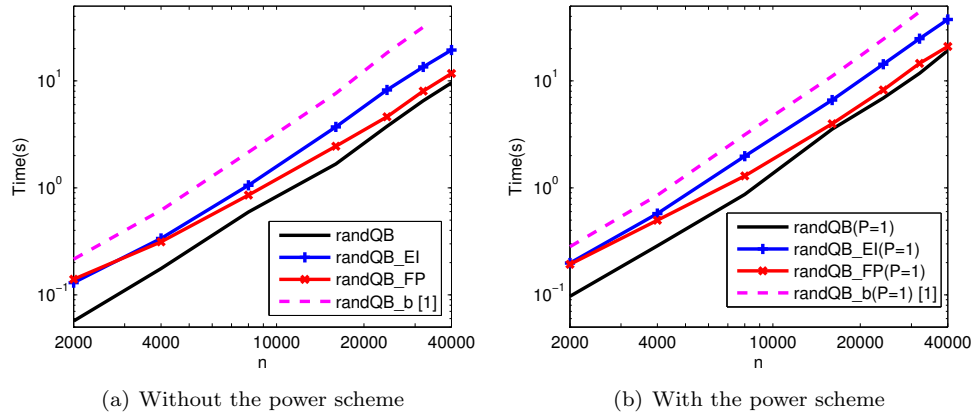


FIG. 3. Runtime of the randomized QB-factorization algorithms for dense matrices ( $l = 200$ ).

multicore CPUs. The QR factorization and other basic linear algebra operations are implemented through LAPACK routines which are automatically executed in parallel.

**5.1. Comparison of speed.** We compute the QB factorization of an  $n \times n$  matrix  $\mathbf{A}$ . The matrix is stored in RAM, though the pass-efficient algorithm `randQB_FP` is suitable for handling a large matrix stored on hard disk. For the latter situation, the runtime advantage of `randQB_FP` can be estimated by the number of passes over  $\mathbf{A}$ . Here we want to show that `randQB_FP` has little sacrifice while owning the pass-efficient property. Notice the singular value distribution of the matrix is immaterial for this runtime comparison. Four algorithms are compared:

- The `randQB` algorithm in Figure 1(a);
- The `randQB_b` algorithm in Figure 1(b), obtained from [16];
- The `randQB_EI` algorithm presented in section 3;
- The `randQB_FP` algorithm presented in section 4.

We compare their speed using both dense and sparse matrices, both as a function of the dimension of the matrix and the parameter  $l$  denoting the number of the output  $\mathbf{Q}$ 's columns. The block size is  $b = 20$  for the `randQB_b`, `randQB_EI`, and `randQB_FP` algorithms. For each runtime measurement, the average time over 20 runs is reported. Notice that the compared `randQB_b` algorithm is an efficient parallel implementation open-sourced in [16], also based on Intel MKL libraries.

In the first experiment we test the algorithms on dense matrices of varying size.  $n$  ranges from 2,000 to 40,000. The value of  $l$  is always 200. The results are shown in Figure 3 for the situations without and with the power scheme. The data of the blocked `randQB` algorithm for the matrix with  $n = 40,000$  are not available due to the unreasonably long runtime of the program from [16]. From the results in Figure 3(a), we see that the `randQB_EI` and `randQB_FP` algorithms are 2.4X (13.47s vs. 31.78s) and 4.0X (8.01s vs. 31.78s) faster than the implementation of `randQB_b` algorithm, respectively, when  $n = 32,000$ . If the power scheme is imposed, the acceleration ratios decrease to 1.8X and 3.0X, respectively, which are still remarkable. The `randQB` algorithm has the fastest computational speed, but its advantage over the `randQB_FP` algorithm becomes marginal when the matrix size is large.

Here we only show the runtime results with the power parameter  $P = 1$ , as for many applications this already achieves sufficient accuracy.

The memory costs for some large matrices are listed in Table 2. For the `randQB`, `randQB_EI`, and `randQB_FP` algorithms, the memory cost is mainly due to storing

TABLE 2

The memory usage of the randomized QB-factorization algorithms for dense matrices ( $l = 200$ ).

$n$	randQB	randQB_EI	randQB_FP	randQB_b [1]
16,000	2,308 MB	2,303 MB	2,357 MB	6,237 MB
24,000	4,792 MB	4,796 MB	4,873 MB	13,618 MB
32,000	8,253 MB	8,253 MB	8,356 MB	23,931 MB
40,000	12,694 MB	12,581 MB	12,714 MB	N.A.

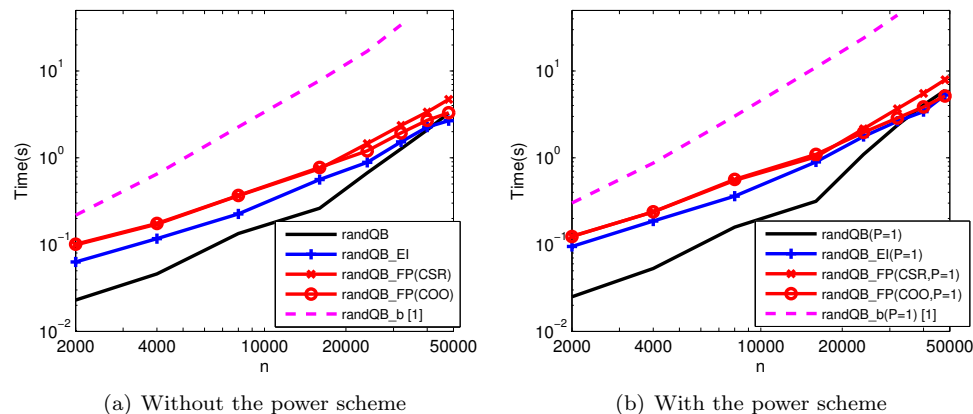


FIG. 4. Runtime of the randomized QB-factorization algorithms for sparse matrices ( $l = 200$ ).

matrix  $\mathbf{A}$ . For the randQB\_b algorithm, it needs additional memory to store the residual matrix and the product of  $\mathbf{QB}$ . So, the proposed algorithms consume about 1/3 of that used by the blocked randQB algorithm. If we allow that  $\mathbf{A}$  can be overwritten by the residual matrix, the memory cost of the randQB\_b algorithm [1] can be reduced, but will still be 2X larger than the proposed algorithms. For the largest case with  $n = 40,000$ , the randQB\_b algorithm actually requests more memory than the size of RAM ( $\sim 32$  GB), which explains the aforementioned long runtime of randQB\_b.

The second experiment is about the algorithms' efficiency for sparse matrices. We generate sparse matrices with roughly 0.3% nonzero elements. They are stored in compressed sparse row (CSR) format [18]. The runtimes of the algorithms are shown in Figure 4. The results of the randQB\_b algorithm for the matrices with  $n \geq 40,000$  are not available due to unreasonably long runtime. In contrast, it only takes a couple of seconds for the other algorithms to process the largest matrix with  $n = 48,000$ . We see that the proposed algorithms take usage of the sparsity, while the blocked randQB algorithm cannot. The speedup ratios of the former to the latter increase as the matrix size increases. For  $n = 32,000$ , the randQB\_EI and randQB\_FP algorithms are more than 22X and 14X faster than the implementation of randQB\_b algorithm, respectively. Different from the situation for dense matrices, the randQB\_EI algorithm becomes faster than randQB\_FP. This implies that lumping the multiplications of a sparse matrix all together brings less benefit. And, randQB\_EI could run faster than the randQB algorithm for a matrix size over 48,000. This can be explained by the comparison of (10) and (12), the inefficiency of orthogonalizing the whole matrix of  $l$  columns, and that the sparse matrix removes the benefit of BLAS-3 operation to randQB. Another interesting phenomenon is that if we instead store a large sparse matrix with the COO (coordinate) format, the runtime of the

TABLE 3

The memory usage of the randomized QB-factorization algorithms for sparse matrices ( $l = 200$ ).

$n$	randQB	randQB_EI	randQB_FP	randQB_b [1]
16,000	162 MB	174 MB	223 MB	6,153 MB
24,000	232 MB	239 MB	312 MB	13,572 MB
32,000	293 MB	303 MB	402 MB	23,917 MB
40,000	338 MB	359 MB	488 MB	N.A.
48,000	405 MB	426 MB	582 MB	N.A.

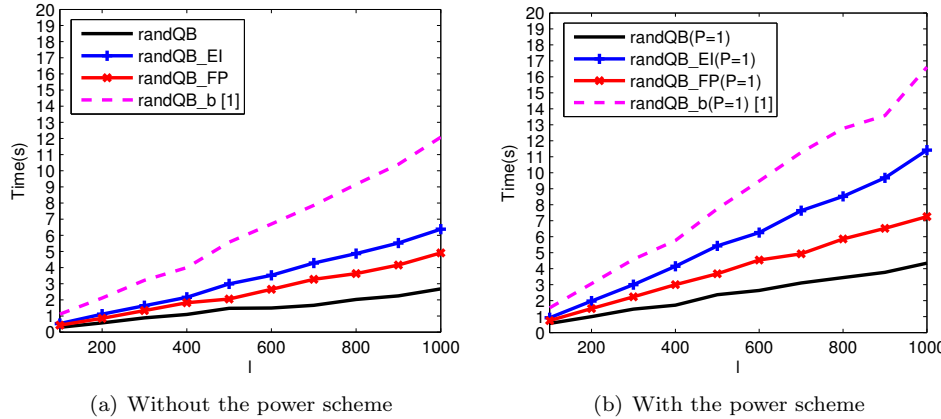


FIG. 5. Runtime of the randomized QB-factorization algorithms for a dense matrix ( $n = 8,000$ ) with varying value of  $l$ .

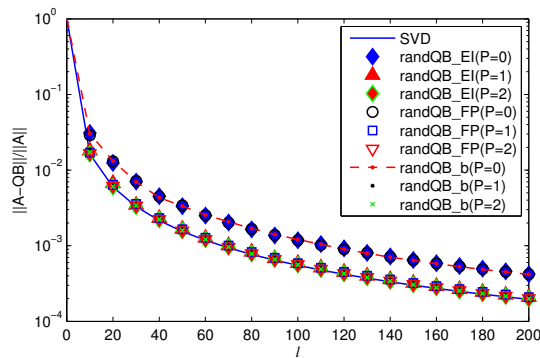
randQB\_FP algorithm can be reduced by 30%. This means that the COO format is more adaptive to parallel computing. If we set  $P = 1$ , similar observations regarding the experimental results can be drawn, as shown in Figure 4(b).

The memory cost of these algorithms is listed in Table 3, from which we see more prominent memory saving of the proposed algorithms over the randQB\_b algorithm. While compared with randQB, the proposed algorithms consume comparable memory.

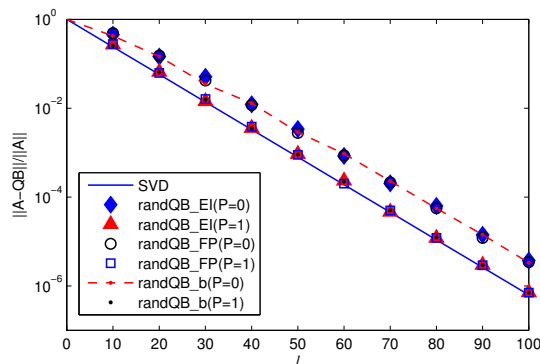
Last, we test a dense matrix with size  $n = 8,000$ , and vary the value of  $l$ . The trends of the runtime are plotted in Figure 5. It shows that the randQB\_EI and randQB\_FP algorithms without the power scheme are about 1.9X and 2.5X faster than randQB\_b, respectively. If the power scheme is imposed, the speedup ratios to randQB\_b decrease, but randQB\_FP is still more than 2X faster than randQB\_b.

**5.2. Comparison of accuracy.** Three kinds of matrices are tested standing for different distribution patterns of singular values:

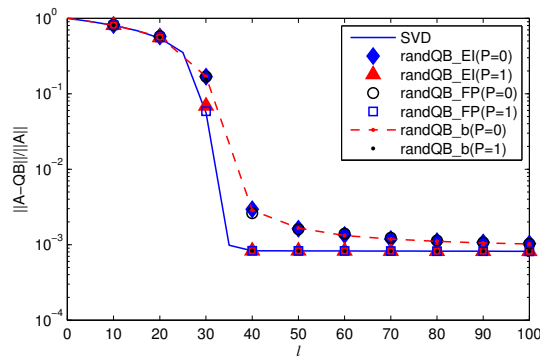
- Matrix 1 (slow decay):  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are randomly drawn matrices with orthonormal columns, and the diagonal matrix  $\mathbf{\Sigma}$  has diagonal elements  $\sigma_j = 1/j^2$ .
- Matrix 2 (fast decay):  $\mathbf{A}$  is formed just like Matrix 1, but the diagonal elements of  $\mathbf{\Sigma}$  is given by  $\sigma_j = e^{-j/7}$ . It reflects a fast decay of singular values.
- Matrix 3 (S-shape decay):  $\mathbf{A}$  is built in the same manner as Matrices 1 and 2, but the diagonal elements of  $\mathbf{\Sigma}$  are given by  $\sigma_j = 0.0001 + (1 + e^{j-30})^{-1}$ . It makes the singular values first hover around 1, then decay rapidly, and finally level out at about 0.0001.



(a) Errors of Matrix 1 whose singular values decay slowly



(b) Errors of Matrix 2 whose singular values decay rapidly



(c) Errors of Matrix 3 with S-shape decay of singular values

FIG. 6. Errors of different algorithms for approximating the test matrices ( $n = 2,000$ ,  $b = 10$ ).

For each kind, we generate a  $2,000 \times 2,000$  matrix, for which we compare the errors of the proposed techniques and the blocked randQB scheme [1] for varying  $l$  values. The results are shown in Figure 6, where we see that the proposed techniques have just the same accuracy as the blocked randQB algorithm. If we use the power scheme, even with a power parameter as small as  $P = 1$ , the errors of the randQB\_EI and randQB\_FP algorithms are remarkably reduced. And, the power schemes with  $P = 1$  and  $P = 2$  produce indistinguishable results for the tested matrices. Both are extremely close to the optimal results from SVD.



**A single-pass algorithm for low-rank approximation**

**Input:**  $A, l$ .

**Output:**  $Q, \tilde{Q}, \hat{B}$  such that  $A \approx Q\hat{B}\tilde{Q}^T$ .

(1)  $\Omega = \text{randn}(n, l); \tilde{\Omega} = \text{randn}(m, l)$ ;

(2)  $Y = A\Omega; \tilde{Y} = A^T\tilde{\Omega}$ ;

(3)  $Q = \text{orth}(Y); \tilde{Q} = \text{orth}(\tilde{Y})$ ;

(4)  $\hat{B} = (\tilde{\Omega}^T Q)^{-1} \tilde{Y}^T \tilde{Q}$ .

FIG. 7. An existing single-pass algorithm for low-rank matrix factorization [2].

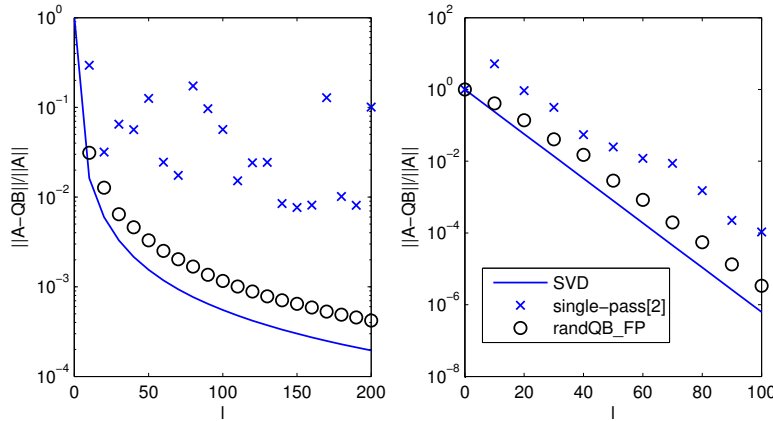


FIG. 8. The approximation errors of the two single-pass algorithms and truncated SVD for Matrix 1 (left) and Matrix 2 (right).

**5.3. Performance of the single-pass algorithm.** Without the power scheme, the `randQB_FP` algorithm is a single-pass algorithm (see Algorithm 3). This is because  $G = A\Omega$  and  $H = A^T G$  can be executed through one pass over matrix  $A$ , providing that  $A$  is in the row-major format. Another single-pass algorithm was proposed in [2], as a remedy to the `randQB` algorithm. It is shown in Figure 7, whose step (3) produces matrices  $Q$  and  $\tilde{Q}$  such that  $A \approx QQ^T A\tilde{Q}\tilde{Q}^T$ . Then, a small matrix  $\hat{B} = Q^T A\tilde{Q}$  is approximately solved in step (4), because  $\tilde{Q}^T \tilde{Y} \approx \tilde{Q}^T A^T Q Q^T \tilde{\Omega} = \hat{B}^T Q^T \tilde{\Omega}$ . This single-pass algorithm corresponds to the low-rank factorization in the form of  $A \approx Q\hat{B}\tilde{Q}^T$ . Obviously, it includes more approximations and is not equivalent to the `randQB` algorithm. In contrast, Algorithm 3 is mathematically equivalent to `randQB` (see Proposition 2) and is supposed to be more accurate. With Matrices 1 and 2 from section 5.2, we can compare the accuracy of the both algorithms, whose results are plotted in Figure 8.

From Figure 8 we see that the approximation error of the single-pass algorithm in [2] is often one order of magnitude larger than that of our `randQB_FP` based algorithm. Actually, it does not even decrease as the rank of the approximation matrix increases. We also calculate the top 50 singular values, and the oversampling with  $s = 10$  is applied to the both algorithms. The results are shown in Figure 9, along with those obtained from the `randQB` algorithm, where the results of `randQB_FP` and `randQB` are indistinguishable. For the matrix with slow decay of singular value the result from `randQB_FP` shows moderate accuracy on the top singular values, whose error is usually orders of magnitude smaller than that of the single-pass algorithm in [2].

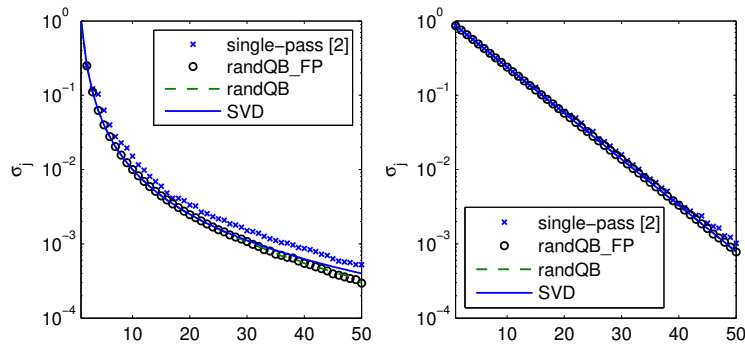


FIG. 9. The computational results of top singular values for Matrix 1 (left) and Matrix 2 (right).

TABLE 4

The results on solving the fixed-precision problems based on three test matrices.

Matrix type	$\epsilon$	randQB_EI			randQB_FP		Truncated SVD		RangeFinder [2]	
		Rank	Time(s)	Error	Rank	Time(s)	Rank	Time(s)	Rank	Time(s)
Matrix 1	1e-2	<b>15</b>	1.19	9.3e-3	<b>15</b>	3.13	<b>15</b>	123	115	1.1
	1e-4	<b>327</b>	8.29	9.98e-5	<b>328</b>	3.71	<b>313</b>	123	2,084	29.1
Matrix 2	1e-4	<b>66</b>	2.16	8.37e-5	<b>66</b>	2.73	<b>65</b>	115	101	1.1
	1e-5	<b>82</b>	2.68	8.86e-6	<b>82</b>	3.17	<b>81</b>	115	113	1.1
Matrix 3	1e-2	<b>33</b>	1.56	4.1e-3	<b>33</b>	3.62	<b>32</b>	126	3,618	87.8
	1.5e-3	<b>1,588</b>	18.7	1.499e-3	<b>1,587</b>	15.8	<b>1,587</b>	126	7,916	379

**5.4. Results of solving the fixed-precision problems.** In this subsection we test the proposed algorithms with some fixed-precision problems. The optimal solution is the factorization with the smallest sizes of  $\mathbf{Q}$  and  $\mathbf{B}$ , which corresponds to a lower amount of subsequent computation. It can be achieved by first calculating SVD of the input matrix  $\mathbf{A}$ , and then checking  $(\sum_{j=k+1}^{\min(m,n)} \sigma_j^2)^{1/2}$ , where  $\sigma_j$  is  $\mathbf{A}$ 's  $j$ th singular value, to determine the smallest  $k$  satisfying the accuracy criterion. Here, we always consider the accuracy criterion with a relative tolerance:  $\|\mathbf{A} - \mathbf{QB}\|_F < \epsilon \|\mathbf{A}\|_F$ .

The proposed algorithms are compared with the SVD based method and the *adaptive randomized range finder* (Algorithm 4.2 in [2]). The row-by-row calculation scheme mentioned in section 3 is implemented into our algorithms. The experiments are carried out with MATLAB on the aforementioned Linux server. The built-in commands like “svd,” “qr,” etc., are employed, which naturally take advantage of parallel computing.

For the randQB\_FP algorithm, we empirically set  $\tilde{l} = 50b$ . With a suitable value of  $b$ , this produces a large enough  $\tilde{l}$  for attaining the accuracy criteria in the experiments. A more sophisticated approach for setting  $\tilde{l}$  and  $b$  can be investigated in the future.

We first construct the three kinds of matrices in section 5.2, each of  $8,000 \times 8,000$  size, and test them with the four methods. Their results are shown in Table 4. For randQB\_EI and randQB\_FP, the power scheme with  $P = 1$  is used. The block size is set to  $b = 10$  in all tests, except the last one for which  $b = 40$ . In Table 4, “ $\epsilon$ ” stands for the threshold for relative error, and “error” means the relative Frobenius-norm error of the produced QB factorization. From the table we see that the results of randQB\_EI and randQB\_FP algorithms all satisfy the set accuracy demands. And, the corresponding ranks (i.e., the number of columns in  $\mathbf{Q}$ ) are

TABLE 5

The results on solving the fixed-precision problems based on two real data.

Matrix	$\epsilon$	Parameters	randQB_EI		randQB_FP		Truncated SVD		RangeFinder [2]		
			Rank	Time(s)	Error	Rank	Time(s)	Rank	Time(s)	Rank	Time(s)
Image	0.1	$P=1, b=10$	468	8.1	0.0999	471	3.25	<b>426</b>	44.2	2,913	79.0
		$P=1, b=20$	468	4.23	0.0999	472	4.44				
		$P=2, b=10$	<b>441</b>	9.98	0.0999	<b>443</b>	3.47				
		$P=2, b=20$	441	5.76	0.0999	443	7.26				
AMiner	0.5	$P=1, b=50$	2440	108	0.4999	2,449	143	<b>2,115</b>	1,049	8,018	399
		$P=2, b=50$	<b>2,229</b>	134	0.4999	<b>2,242</b>	205				

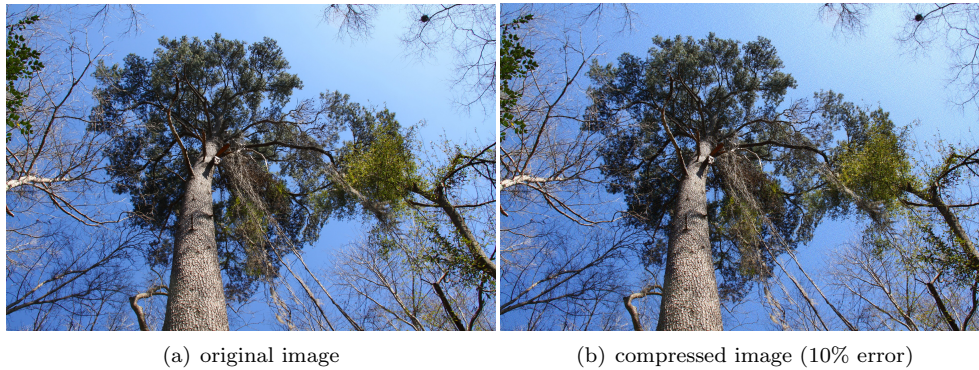


FIG. 10. The original image and the compressed image obtained with the proposed algorithm.

very close to the optimal values from the SVD based approach. As for the runtime, the proposed algorithms are usually several tens times faster than SVD. Notice that our MATLAB programs are less optimized than the built-in `svd` command. So, more significant speedup could be expected for the implementation in C. Although the adaptive range finder is built on a theory with spectral norm of matrix, in our experiments it always produces a QB factorization satisfying the accuracy demand in Frobenius norm. However, its results (factor matrices) are much larger than necessary.

We then test the algorithms with two real data. One is from a scenic image, and the other is from an information retrieval application “AMiner” [19]. The colored image is represented by a  $9,504 \times 4,752$  matrix. The other is an  $8,130 \times 100,000$  keyword-person matrix produced with the term frequency and inverse document frequency model [20]. This sparse matrix has about 0.2% nonzero elements. The computational results are listed in Table 5, with different power parameters and block sizes. They again validate that the proposed algorithms can automatically satisfy the accuracy criterion. And, with  $P = 2$  the result of rank is substantially reduced, approaching the optimal value. For the same power scheme, setting larger block size  $b$  we can reduce the runtime of `randQB_EI`. In contrast, the runtime of `randQB_FP` increases with the block size, as we have set  $\tilde{l} = 50b$ . Notice that with the relative error tolerance  $\epsilon = 0.1$ , the image is largely compressed ( $\sim 7X$  size reduction), with little loss of quality (see Figure 10). And, the singular value of the AMiner matrix decays very slowly, but even with large approximation error its low-rank approximation could bring improved performance of information retrieval (cf. [20], section 11.3).

For the second data, which is a large sparse matrix, the Krylov subspace iterative method `svds` [18, 20] is also tested. However, it costs 2,281 seconds for computing the first 1,000 singular values/vectors. It is much slower than executing `svd` to the matrix's dense version. Besides, `svd` requests more than 20 GB memory, while the proposed randomized algorithms only costs 3 GB memory or so for this case.

**6. Conclusions.** Efficient techniques are proposed for the fixed-precision low-rank approximation of large matrices. Our contributions are as follows.

- A simple and accurate error indicator in the Frobenius norm is proposed, which enables efficient rank determination and can be used in the blocked `randQB` algorithm [1] and other incremental QB-form factorization algorithms (like that in [14]). We have proved its accuracy and validity for the problems with relative accuracy tolerance larger than  $2.1 \times 10^{-7}$ . Numerical experiments on large dense and sparse matrices have shown that the proposed rank determination scheme brings several to several tens times speedup and memory saving to the blocked `randQB` algorithm, without loss of accuracy.
- Based on the blocked `randQB` algorithm, we propose a pass-efficient algorithm called `randQB_FP`. It is mathematically equivalent to the blocked `randQB` algorithm, but reduces the passes over matrix  $\mathbf{A}$  to the fewest. The `randQB_FP` algorithm also suits the fixed-precision problem and can derive a single-pass algorithm under a certain condition. Numerical results have validated the efficiency and accuracy of the `randQB_FP` algorithm and shown that the derived single-pass algorithm is much more accurate than an existing counterpart.
- Real data are tested to demonstrate the effectiveness of the proposed algorithms for the fixed-precision problem. Compared with the adaptive range finder approach [2], the proposed algorithms run faster and produce much smaller factor matrices while attaining the accuracy criterion.

Future work includes extending and applying the proposed algorithms to more practical data mining and machine learning scenarios.

**Acknowledgments.** The authors thank Prof. P.-G. Martinsson for sharing the source code of the blocked `randQB` algorithm, which makes this work possible. The authors are also grateful to Prof. Jie Tang with Tsinghua University for providing the data of Aminer.

#### REFERENCES

- [1] P.-G. MARTINSSON AND S. VORONIN, *A randomized blocked algorithm for efficiently computing rank-revealing factorizations of matrices*, SIAM J. Sci. Comput., 38 (2016), pp. S485–S507.
- [2] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Rev., 53 (2011), pp. 217–288.
- [3] P. DRINEAS AND M. W. MAHONEY, *RandNLA: Randomized numerical linear algebra*, Commun. ACM, 59 (2016), pp. 80–90.
- [4] S. VORONIN AND P.-G. MARTINSSON, *RSVDPACK: Subroutines for Computing Partial Singular Value Decompositions Via Randomized Sampling on Single Core, Multi Core, and GPU Architectures*, preprint, arXiv:1502.05366v3 [math.NA], 2016.
- [5] S. ERIKSSON-BIQUE, M. SOLBRIG, M. STEFANELLI, S. WARKENTIN, R. ABBEY, AND I. C. F. IPSEN, *Importance sampling for a Monte Carlo matrix multiplication algorithm, with application to information retrieval*, SIAM J. Sci. Comput., 33 (2011), pp. 1689–1706.
- [6] P. DRINEAS, R. KANNAN, AND M. W. MAHONEY, *Fast Monte Carlo algorithms for matrices II: Computing a low-rank approximation to a matrix*, SIAM J. Sci. Comput., 36 (2006), pp. 158–183.

- [7] V. ROKHLIN, A. SZLAM, AND M. TYGERT, *A randomized algorithm for principal component analysis*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 1100–1124.
- [8] H. JI AND Y. LI, *GPU accelerated randomized singular value decomposition and its application in image compression*, in Proceedings of the Modeling, Simulation, and Visualization Capstone Conference, Suffolk, VA, 2014.
- [9] T. MARY, I. YAMAZAKI, J. KURZAK, P. LUSZCZEK, S. TOMOV, AND J. DONGARRA, *Performance of random sampling for computing low-rank approximations of a dense matrix on GPUs*, in Proceedings of SC'2015, Austin, TX, 2015.
- [10] M. W. MAHONEY, *Randomized algorithms for matrices and data*, Found. Trends Mach. Learn., 3 (2011), pp. 123–224.
- [11] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 3rd ed., Johns Hopkins University Press, Baltimore, MD, 1996.
- [12] M. GU AND S. C. EISENSTAT, *Efficient algorithms for computing a strong rank-revealing QR factorization*, SIAM J. Sci. Comput., 17 (1996), pp. 848–869.
- [13] P. G. MARTINSSON, G. QUINTANA-ORTI, N. HEAVNER, AND R. VAN DE GEIJN, *Householder QR factorization with randomization for column pivoting (HQRRP)*, SIAM J. Sci. Comput., 39 (2017), pp. C96–C115.
- [14] J. A. DUERSCH AND M. GU, *Randomized QR with column pivoting*, SIAM J. Sci. Comput., 39 (2017), pp. C263–C291.
- [15] C. ECKART AND G. YOUNG, *The approximation of one matrix by another of lower rank*, Psychometrika, 1 (1936), pp. 211–218.
- [16] S. VORONIN AND P. G. MARTINSSON, RandQR, [http://amath.colorado.edu/faculty/martinss/main\\_codes.html](http://amath.colorado.edu/faculty/martinss/main_codes.html), 2015.
- [17] Intel Parallel Studio XE Cluster Edition for Linux, <https://software.intel.com/en-us/intel-parallel-studio-xe>, 2017.
- [18] Y. SAAD, *Iterative Methods for Sparse Linear Systems*, SIAM, Philadelphia, 2003.
- [19] *AMiner*, <http://www.aminer.org>.
- [20] L. ELDEN, *Matrix Methods in Data Mining and Pattern Recognition*, SIAM, Philadelphia, 2007.
- [21] R. LEHOUCQ, D. SORENSEN, AND C. YANG, *ARPACK User's Guide: Solution of Large-Scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*, SIAM, Philadelphia, 1998.
- [22] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for low-rank matrix approximation*, SIAM J. Matrix Anal. Appl., 38 (2017), pp. 1454–1485.
- [23] H. LI, G. C. LINDERMAN, A. SZLAM, K. P. STANTON, Y. KLUGER, AND M. TYGERT, *Algorithm 971: An implementation of a randomized algorithm for principal component analysis*, ACM Trans Math Software, 43 (2017), 28.
- [24] W. YU AND X. WANG, *Advanced Field-Solver Techniques for RC Extraction of Integrated Circuits*, Springer, New York, 2014.