AN IMPROVED ANALYSIS AND UNIFIED PERSPECTIVE ON DETERMINISTIC AND RANDOMIZED LOW-RANK MATRIX APPROXIMATION*

JAMES DEMMEL†, LAURA GRIGORI‡, AND ALEXANDER RUSCIANO‡§

Abstract. We introduce a Generalized LU Factorization (GLU) for low-rank matrix approximation. We relate this to past approaches and extensively analyze its approximation properties. The established deterministic guarantees are combined with sketching ensembles satisfying Johnson–Lindenstrauss properties to present complete bounds. Particularly good performance is shown for the subsampled randomized Hadamard transform (SRHT) ensemble. Moreover, the factorization is shown to unify and generalize many past algorithms, sometimes providing strictly better approximations. It also helps to explain the effect of sketching on the growth factor during Gaussian elimination.

Key words. low-rank approximation, spectrum preserving, kernel approximation, randomized algorithms, deterministic algorithms

MSC codes. 68Q25, 68R10, 68U05

DOI. 10.1137/21M1391316

1. Introduction. Many different problem domains produce matrices that can be approximated by a low-rank matrix. In some cases such as a divide-and-conquer approach to eigenproblems [2], there may be many large and small singular values separated by a gap. In other cases such as identifying a low-rank subspace from noisy data, we might expect there to be relatively few large singular values. Perhaps most generically in applied problems, there is no pronounced gap, but the spectrum still decays fairly quickly, and one might prefer to work with a more compact representation when computing quantities such as matrix-vector products.

We next define some related properties which can be of interest to these problems. The following definitions have appeared in the rank-revealing literature, such as in [24, 15, 11, 16] in similar forms. Here and later $A \in \mathbb{R}^{m \times n}$ is the matrix to be approximated, its singular values $\sigma_1 \geq \cdots \geq \sigma_{\min(m,n)}$ are sorted in descending order, and $A_k \in \mathbb{R}^{m \times n}$ is an approximation of A.

DEFINITION 1.1 (low-rank approximation). If A_k satisfies $||A - A_k||_2 \le \gamma \sigma_{k+1}(A)$ for some $\gamma \ge 1$, it is a (k, γ) low-rank approximation of A.

DEFINITION 1.2 (spectrum preserving). If A_k satisfies $\sigma_j(A) \geq \sigma_j(A_k) \geq \gamma^{-1}\sigma_j(A)$ for $1 \leq j \leq k$ and some $\gamma \geq 1$, it is (k, γ) spectrum preserving.

Many results in the rank-revealing literature [16, 14] use a strengthening of Definition 1.1, in which all singular values $\sigma_j(A-A_k)$ are bounded with respect to $\sigma_{k+j}(A)$.

^{*}Received by the editors January 12, 2021; accepted for publication (in revised form) October 18, 2022; published electronically May 11, 2023.

https://doi.org/10.1137/21M1391316

Funding: The work of the second author was supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme grant 810367.

[†]Department of Mathematics, Computer Science, UC Berkeley, Berkeley, CA 94720 USA (demmel@cs.berkeley.edu).

[‡]Sorbonne Université, Inria, CNRS Université de Paris, Laboratoire Jacques-Louis Lions, Paris, 75005 France (Laura.Grigori@inria.fr).

 $^{^{\}S}$ Department of Mathematics, UC Berkeley, Berkeley, CA 94709 USA (rusciano@berkeley.edu).

DEFINITION 1.3 (kernel approximation). If A_k satisfies $\sigma_{k+j}(A) \leq \sigma_j(A - A_k) \leq \gamma \sigma_{k+j}(A)$ for $1 \leq j \leq \min(m,n) - k$ and some $\gamma \geq 1$, it is a (k,γ) kernel approximation of A.

In all of these definitions, if we assume A_k is rank k, then $\gamma=1$ is optimal from the truncated SVD, so all methods can be compared with this standard. Different algorithms may end up representing A_k in different ways, but generally A_k is represented as a product of matrices which have at least one dimension much smaller than those of the original A. Note that in this work we do not require A_k to be rank k. Nevertheless, the rank of A_k will be chosen as a function of k in order to compete with the truncated SVD of rank k, and this motivates the choice of notation. For the choices made in this paper, our bounds always limit $\operatorname{rank}(A_k)$ to at most $\operatorname{rank}(A_k) = O(k \cdot \operatorname{polylog}(n))$. We also note that the approximation A_k can be truncated to be rank k and maintain Definition 1.1. This is a well-known strategy in the literature; for example, see section 6 of [32].

As we made the above definitions quite strong, and in particular we did not make any assumption on γ , we will not prove that all our results satisfy them exactly. Different algorithms can approach them to different degrees, and these definitions can be used as a measure of their quality. The bounds on singular values in Definitions 1.2 and 1.3 were first discussed in the context of deterministic rank-revealing factorizations, where in, e.g., [16] γ is a low degree polynomial in k and n. The algorithms in [3] can also be made to satisfy them. In Definition 1.2, inequality $\sigma_j(A) \geq \sigma_j(A_k)$ holds for any $A_k = PAQ$, where P and Q are such that $\sigma_{max}(P) \leq 1$ and $\sigma_{max}(Q) \leq 1$, e.g., if P and Q are orthogonal projections, by the multiplicative Weyl inequality (3.4). Some of our approximations A_k satisfy this, and some do not; we will not consider this inequality further. In Definition 1.3, inequality $\sigma_{k+j}(A) \leq \sigma_j(A-A_k)$ holds for any A_k of rank k, by the additive Weyl inequality (3.6). Some of our A_k satisfy this, and some do not; we will not consider this inequality further.

This paper has two main goals, both motivated by the history of low-rank factorizations. First, we show that many important low-rank factorizations can be viewed as a generalized LU factorization followed by setting the Schur complement equal to zero. We call this prototype algorithm **GLU**. Second, older research into low-rank factorizations bounded more quantities than recent results on randomized algorithms. In particular, Definitions 1.2 and 1.3 do not receive much discussion in randomized algorithms. However, approximating the singular values of A is useful for detecting a gap in the singular values and choosing accordingly the rank of the approximation. We will provide bounds on γ in each of Definition 1.1, Definition 1.2, and Definition 1.3 for **GLU**. In doing this, we first derive sharp deterministic bounds for truncated LU and QR factorizations in sections 3 and 5, and then in section 6 we complete the bounds by using properties of random matrix ensembles.

Our **GLU** approximation is essentially based on a truncated LU factorization that allows the leading block to be rectangular instead of square. Allowing the leading block to be rectangular enables much better low-rank approximation properties. Given the matrix $A \in \mathbb{R}^{m \times n}$, let A_{11} be the leading $l' \times l$ block which is assumed to have full column rank so that $l' \geq l$, and let $U \in \mathbb{R}^{m \times m}$ and $V \in \mathbb{R}^{n \times n}$ be invertible matrices. First we have an exact factorization of matrix A that is the natural generalization of an LU factorization,

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} I \\ A_{21}A_{11}^{+} & I \end{pmatrix} \begin{pmatrix} A_{11} & A_{12} \\ & \$(A_{11}) \end{pmatrix},$$

where A_{11}^+ is the Moore–Penrose pseudo-inverse of A_{11} and $\mathcal{S}(A_{11}) = A_{22} - A_{21} A_{11}^+ A_{12}$ denotes the generalized Schur complement (see, e.g., [5]). By applying the sketching matrices U and V and deleting the Schur complement, we get a low-rank factorization that can have remarkably good properties. Defining $\bar{A} = UAV = \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} \end{pmatrix}$, and assuming \bar{A}_{11} has full column rank,

$$(1.1) A_k := U^{-1} \begin{pmatrix} I \\ \bar{A}_{21} \bar{A}_{11}^+ \end{pmatrix} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \end{pmatrix} V^{-1}$$

is a complete mathematical description of our proposed **GLU** approximation. The inverses may look daunting at first because they are large matrices, but we will see that they are only tools to facilitate the analysis; actually the leading l' rows of U and leading l columns of V are the only parts required.

We have emphasized that \mathbf{GLU} factorization unifies many factorizations through appropriate choices of the settings of U, V. We believe other choices of U and V discussed in this paper are novel and practical, as we illustrate in main results Theorem 6.6 and Theorem 6.9. That said, this paper will not argue that these novel instantiations of \mathbf{GLU} should necessarily be adopted over similar methods like the low-rank factorization described in [7]. A comparison of the pros and cons is outside the scope of this work, and we mainly want to emphasize that our algorithm is practical and has a very general, transparent analysis.

The remainder of the introduction is divided into three sections for clarity. The first and second aim to highlight our contributions. The third gives notation we adopt.

1.1. Unifying approach. GLU generalizes past low-rank LU factorizations in two ways. First, it allows pre- and post-multiplication by matrices other than permutations. Second, it allows for rectangular A_{11} when computing Schur complements. Even without generalizing to rectangular A_{11} , GLU encompasses several well-known procedures. We provide examples to illustrate this and detailed derivations in section 4. Table 1 summarizes several deterministic and randomized approximation algorithms. It displays separately the case when $k \leq l = l'$ and the more general case when $k \leq l \leq l'$, and it cites existing as well as new bounds on the spectral and kernel approximations provided by these algorithms. Here we focus on identifying that existing deterministic and randomized algorithms depend on the same matrix factorizations.

In the case when $k \leq l = l'$, the rank-k approximation A_k becomes

$$A_k = U^{-1} \begin{pmatrix} I_l \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{pmatrix} (\bar{A}_{11} \quad \bar{A}_{12}) V^{-1}$$

$$= AV_1 (U_1 AV_1)^{-1} U_1 A,$$
(1.2)

where V_1 contains the leading l columns of V, U_1 contains the leading l rows of U, and $\bar{A} = UAV$. See (3.1) for more details. While $l \geq k$ is always the case, in applications l varies from being exactly k, as for deterministic algorithms, to being a polylog-factor larger than k for randomized algorithms. Now we define some notation we will use later. Let Q_1 be the orthogonal factor obtained from the thin QR decomposition of AV_1 , so Q_1 is of dimensions $m \times l$. Let $UQ_1 = \begin{pmatrix} \bar{Q}_{11} \\ \bar{Q}_{21} \end{pmatrix}$, where $\bar{Q}_{11} = U_1Q_1$ is $l \times l$. The approximation from (1.2) can be written as

(1.3)
$$A_k = AV_1(U_1AV_1)^{-1}U_1A$$
$$= Q_1(U_1Q_1)^{-1}U_1A.$$

Table 1

Summary of several deterministic and randomized algorithms for computing A_k , the low rank approximation of a matrix A of dimensions $m \times n$. U_1 is $l' \times m$, V_1 is $n \times l$, and Q_1 is the $m \times l$ orthogonal factor obtained from the thin QR decomposition of AV_1 . In the table, V_1 permutation and U_1 permutation refer to V_1 containing the leading l columns of a permutation matrix V, U_1 containing the leading l' rows of a permutation matrix U, respectively.

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	Existing algorithms: Instances of V_1, U_1 and the approximation A_k for $k \leq l = l'$,	
QR with column selection, $k = l = l'$ (a.k.a. strong rank revealing QR [16]) V_1 permutation, $U_1 = Q_1^T$, $A_k = Q_1Q_1^TA$; see (4.10) $new\ bounds\ on\ kernel\ approx.\ (Proposition\ 5.2)$ LU with column/row selection, $k = l = l'$ (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 , U_1 permutations V_1 , U_1 permutations v_1 , v_2 permutations v_2 permutations v_3 permutations v_4 permutation, v_4 permutations v_4 permutation, v_4	$A_k = AV_1(U_1AV_1)^{-1}U_1A$	
$(a.k.a. \ strong \ rank \ revealing \ QR \ [16]) \\ V_1 \ permutation, \ U_1 = Q_1^T, \ A_k = Q_1Q_1^TA; \ see \ (4.10) \\ new \ bounds \ on \ kernel \ approx. \ (Proposition \ 5.2) \\ LU \ with \ column/row \ selection, \ k = l = l' \\ (a.k.a. \ rank \ revealing \ LU, \ e.g., \ [14] \ or \ CUR) \\ V_1 \ with \ column/row \ selection, \ k = l = l' \\ (a.k.a. \ rank \ revealing \ LU, \ e.g., \ [14] \ or \ CUR) \\ N_1, \ U_1 \ permutations \\ new \ spectral, \ kernel \ bounds \ (Proposition \ 3.3) \\ Instances \ of \ V_1, \ U_1 \ and \ the \ approximation \ A_k^{ob} \ for \ k \leq l \leq l', \\ A_k^{ob} = AV_1(U_1AV_1)^+U_1A \\ Deterministic \ algorithms \ and \ bounds \\ ObliqueProj \ with \ column/row \ selection, \ k \leq l < l' \\ V_1, \ U_1 \ permutations \\ Nyström \ for \ A \ SPSD, \ k = l = l' \\ U_1 \ permutation, \ U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \\ U_1 \ permutation, \ U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \\ New \ algorithms: \ V_1, U_1 \ and \ the \ approximation \ A_k \ for \ k \leq l \leq l', \\ A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \ see \ (3.18) \\ Deterministic \ algorithm \ and \ bounds \\ GLU \ with \ column/row \ selection, \ k \leq l \leq l' \\ Randomized \ GLU, \ k \leq l$	Deterministic algorithms and bounds	Randomized algorithm and bounds
$V_1 \text{ permutation, } U_1 = Q_1^T, \ A_k = Q_1Q_1^TA; \text{ see } (4.10) \\ \text{new bounds on kernel approx. } (Proposition 5.2) \\ \hline \textbf{LU with column/row selection, } k = l = l' \\ \text{(a.k.a. rank revealing LU, e.g., } [14] \text{ or CUR}) \\ \hline \textbf{Eul with column/row selection, } k = l = l' \\ \text{(a.k.a. rank revealing LU, e.g., } [14] \text{ or CUR}) \\ \hline \textbf{Number the permutations} \\ \hline \textbf{Number the permutation} \\ \hline Number the p$	QR with column selection, $k = l = l'$	Randomized SVD, $k < l = l'$
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	(a.k.a. strong rank revealing QR [16])	(a.k.a. randomized QB, e.g., [17])
LU with column/row selection, $k = l = l'$ (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 (a.k.a. rank revealing LU, e.g., [14] or CUR) V_1 randomized SVD via row extraction [17], Randomized Randomized Randomized (Proposition 3.3) V_1 random, U_1 permutation; see (4.10) V_1 random, V_1 permutation and bounds V_1 Parmutation Algorithm and bounds V_1 Parmutation, V_1 permutation, V_2 Parmutation, V_3 Parmutation, V_4 Par	V_1 permutation, $U_1 = Q_1^T$, $A_k = Q_1 Q_1^T A$; see (4.10)	V_1 random, $U_1 = Q_1^T$, $A_k = Q_1 Q_1^T A$; see (4.10)
(a.k.a. rank revealing LU, e.g., [14] or CUR) (e.g. Randomized SVD via row extraction [17], RandomizedRowID [23]) $V_1, U_1 \text{ permutations} \qquad V_1 \text{ random}, U_1 \text{ permutation; see } (4.10) \\ new spectral, kernel bounds (Proposition 3.3) \\ \hline \textbf{Instances of } V_1, U_1 \text{ and the approximation } A_k^{ob} \text{ for } k \leq l \leq l', \\ A_k^{ob} = AV_1(U_1AV_1)^+U_1A \\ \hline \textbf{Deterministic algorithms and bounds} \\ \textbf{ObliqueProj with column/row selection}, k \leq l < l' \\ V_1, U_1 \text{ permutations} \\ \hline \textbf{Nyström for } A \text{ SPSD}, k = l = l' \\ U_1 \text{ permutation}, U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \\ \hline \textbf{New algorithms: } V_1, U_1 \text{ and the approximation } A_k \text{ for } k \leq l \leq l', \\ A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18) \\ \hline \textbf{Deterministic algorithm and bounds} \\ \hline \textbf{GLU with column/row selection}, k \leq l \leq l' \\ \hline \textbf{Randomized GLU}, k \leq l \leq $	new bounds on kernel approx. (Proposition 5.2)	new bounds on kernel approx (Cor. 6.15)
$ \begin{array}{c} \text{RandomizedRowID [23])} \\ V_1,U_1 \text{ permutations} & V_1 \text{ random},U_1 \text{ permutation; see } (4.10) \\ new spectral, kernel bounds (Proposition 3.3) & new spectral, kernel bounds (Proposition 3.3) \\ \hline \\ \textbf{Instances of } V_1,U_1 \text{ and the approximation } A_k^{ob} \text{ for } k \leq l \leq l', \\ A_k^{ob} = AV_1(U_1AV_1)^+U_1A \\ \textbf{Deterministic algorithms and bounds} & \textbf{Randomized algorithm and bounds} \\ \textbf{ObliqueProj with column/row selection}, k \leq l < l' \\ V_1,U_1 \text{ permutations} & \textbf{Randomized ObliqueProj}, k \leq l \leq l' \\ V_1,U_1 \text{ random} & \textbf{Clarkson and Woodruff } (\textbf{CW}), k \leq l \leq l' \\ V_1,U_1 \text{ based on CountSketch} \\ \textbf{Nyström for } A \text{ SPSD}, k = l = l' & \textbf{Randomized Nyström for } A \text{ SPSD}, k \leq l = l' \\ U_1 \text{ permutation}, U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T & U_1 \text{ random}, U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \\ \textbf{New algorithms: } V_1,U_1 \text{ and the approximation } A_k \text{ for } k \leq l \leq l', \\ A_k = [U_1^+(I-(U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18) \\ \textbf{Deterministic algorithm and bounds} & \textbf{Randomized algorithm and bounds} \\ \textbf{GLU with column/row selection}, k \leq l \leq l' & \textbf{Randomized GLU}, k \leq l \leq l' \\ \end{array}$	LU with column/row selection, $k = l = l'$	Randomized LU with row selection, $k < l = l'$
$V_1, U_1 \text{ permutations} \qquad V_1 \text{ random, } U_1 \text{ permutation; see } (4.10)$ $new \ spectral, \ kernel \ bounds \ (Proposition \ 3.3)$ $Instances \ of \ V_1, U_1 \ and \ the \ approximation \ A_k^{ob} \ for \ k \leq l \leq l',$ $A_k^{ob} = AV_1(U_1AV_1)^+U_1A$ $Deterministic \ algorithms \ and \ bounds$ $ObliqueProj \ with \ column/row \ selection, \ k \leq l < l'$ $V_1, U_1 \ permutations$ $Clarkson \ and \ Woodruff \ (CW), \ k \leq l \leq l'$ $V_1, U_1 \ based \ on \ CountSketch$ $Nyström \ for \ A \ SPSD, \ k = l = l'$ $Randomized \ Nyström \ for \ A \ SPSD, \ k \leq l = l'$ $U_1 \ permutation, \ U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \ U_1 \ random, \ U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T$ $New \ algorithms: \ V_1, U_1 \ and \ the \ approximation \ A_k \ for \ k \leq l \leq l',$ $A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \ see \ (3.18)$ $Deterministic \ algorithm \ and \ bounds$ $Randomized \ GLU, \ k \leq l \leq l'$ $Randomized \ GLU, \ k \leq l \leq l'$	(a.k.a. rank revealing LU, e.g., [14] or CUR)	(e.g. Randomized SVD via row extraction [17],
Instances of V_1, U_1 and the approximation A_k^{ob} for $k \leq l \leq l'$, $A_k^{ob} = AV_1(U_1AV_1)^+U_1A$ Deterministic algorithms and bounds ObliqueProj with column/row selection, $k \leq l < l'$ $V_1, U_1 \text{ permutations}$ Randomized ObliqueProj, $k \leq l \leq l'$ $V_1, U_1 \text{ random}$ Clarkson and Woodruff (CW), $k \leq l \leq l'$ $V_1, U_1 \text{ based on CountSketch}$ Nyström for A SPSD, $k = l = l'$ Randomized Nyström for A SPSD, $k \leq l = l'$ $U_1 \text{ permutation}, U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T$ $U_1 \text{ random}, U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T$ New algorithms: V_1, U_1 and the approximation A_k for $k \leq l \leq l'$, $A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18)$ Deterministic algorithm and bounds Randomized GLU, $k \leq l \leq l'$		RandomizedRowID [23])
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	V_1, U_1 permutations	V_1 random, U_1 permutation; see (4.10)
$A_k^{ob} = AV_1(U_1AV_1)^+U_1A$ Deterministic algorithms and bounds Randomized algorithm and bounds ObliqueProj with column/row selection, $k \le l < l'$ Randomized ObliqueProj, $k \le l \le l'$ V_1, U_1 permutations $V_1, U_1 \text{ random}$ Clarkson and Woodruff (CW), $k \le l \le l'$ V_1, U_1 based on CountSketch Nyström for A SPSD, $k = l = l'$ Randomized Nyström for A SPSD, $k \le l = l'$ U_1 permutation, $U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T$ U_1 random, $U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T$ New algorithms: V_1, U_1 and the approximation A_k for $k \le l \le l'$, $A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18)$ Deterministic algorithm and bounds Randomized algorithm and bounds GLU with column/row selection, $k \le l \le l'$ Randomized GLU, $k \le l \le l'$	new spectral, kernel bounds (Proposition 3.3)	new spectral, kernel bounds (Proposition 3.3)
$ \begin{array}{c} \textbf{ObliqueProj with column/row selection}, k \leq l < l' \\ V_1, U_1 \text{ permutations} \\ & V_1, U_1 \text{ random} \\ & & V_1, U_1 \text{ random} \\ & & & & Clarkson \text{ and Woodruff (CW)}, \ k \leq l \leq l' \\ & & & & & & & & & & & & & & & & & & $		
$V_1, U_1 \text{ permutations} \qquad V_1, U_1 \text{ random} \\ & \textbf{Clarkson and Woodruff (CW)}, \ k \leq l \leq l' \\ & V_1, U_1 \text{ based on CountSketch} \\ & \textbf{Nystr\"{o}m for A SPSD}, \ k = l = l' \\ & \textbf{Randomized Nystr\"{o}m for A SPSD}, \ k \leq l = l' \\ & U_1 \text{ permutation}, \ U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \ U_1 \text{ random}, \ U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \\ & \textbf{New algorithms: } \ V_1, U_1 \text{ and the approximation } A_k \text{ for } k \leq l \leq l', \\ & A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18) \\ & \textbf{Deterministic algorithm and bounds} & \textbf{Randomized algorithm and bounds} \\ & \textbf{GLU with column/row selection}, \ k \leq l \leq l' & \textbf{Randomized GLU}, \ k \leq l \leq l' \\ $	Deterministic algorithms and bounds	Randomized algorithm and bounds
	ObliqueProj with column/row selection, $k \le l < l'$	Randomized ObliqueProj, $k \leq l \leq l'$
	V_1, U_1 permutations	V_1, U_1 random
		Clarkson and Woodruff (CW), $k \le l \le l'$
$U_1 \text{ permutation, } U_1 = V_1^T, A_k = AV_1(V_1^TAV_1)^+(AV_1)^T \ U_1 \text{ random, } U_1 = V_1^T, \ A_k = AV_1(V_1^TAV_1)^+(AV_1)^T$ New algorithms: V_1, U_1 and the approximation A_k for $k \leq l \leq l'$, $A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18)$ Deterministic algorithm and bounds Randomized algorithm and bounds GLU with column/row selection, $k \leq l \leq l'$ Randomized GLU, $k \leq l \leq l'$		V_1, U_1 based on CountSketch
New algorithms: V_1, U_1 and the approximation A_k for $k \le l \le l'$, $A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18)$ Deterministic algorithm and bounds Randomized algorithm and bounds GLU with column/row selection, $k \le l \le l'$ Randomized GLU, $k \le l \le l'$		
$A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A]; \text{ see } (3.18)$ Deterministic algorithm and bounds Randomized algorithm and bounds GLU with column/row selection, $k \leq l \leq l'$ Randomized GLU, $k \leq l \leq l'$	U_1 permutation, $U_1 = V_1^T$, $A_k = AV_1(V_1^T A V_1)^+ (AV_1)^T$	U_1 random, $U_1 = V_1^T$, $A_k = AV_1(V_1^T A V_1)^+ (AV_1)^T$
GLU with column/row selection, $k \le l \le l'$ Randomized GLU, $k \le l \le l'$		
	Deterministic algorithm and bounds	Randomized algorithm and bounds
V_1, U_1 permutations V_1, U_1 random	GLU with column/row selection, $k \le l \le l'$	Randomized GLU, $k \le l \le l'$
	V_1, U_1 permutations	V_1, U_1 random
new spectral, kernel bounds (Proposition 3.3) new spectral, kernel bounds (Theorems 6.6, 6.9)	new spectral, kernel bounds (Proposition 3.3)	new spectral, kernel bounds (Theorems 6.6, 6.9)

Deterministic algorithms are typically based on truncated rank-revealing QR (RRQR) and rank-revealing LU (RRLU) factorizations. Both factorizations select k columns from the matrix A; that is, V is a column permutation matrix and AV_1 consists of the selected columns. In the case of a RRQR factorization, $U_1 = Q_1^T$, and the approximation becomes $A_k^o = Q_1Q_1^TA$, relying on an orthogonal projection. That is, letting $\mathcal{P}^o = AV_1(AV_1)^+$ and $k \leq l \leq l'$, the approximation is computed as

(1.4)
$$A_k^o = AV_1(AV_1)^+ A = Q_1 Q_1^T A = \mathcal{P}^o A,$$

(1.5) or equivalently
$$\mathcal{P}^o A(:,j) := \arg \min_{x \in range(AV_1)} \|x - A(:,j)\|_2$$
,

where A(:,j) denotes the jth column of A. Let $Q_1^T A = (R_{11} \ R_{12})$. The strong rank-revealing QR factorization [16] chooses V_1 such that $||R_{11}^{-1}R_{12}||_{max}$ is bounded by a small constant and the approximation A_k is spectrum preserving and a kernel approximation of A: γ in Definition 1.2 and Definition 1.3 is a low degree polynomial in n and k. The rank-revealing LU factorization selects k columns and k rows from the matrix A; that is, both U_1 and V_1 are formed by the leading rows, columns, respectively, of permutation matrices. For example, in LU CRTP [14] the columns

are selected by using a pivoting strategy on A referred to as tournament pivoting and based on strong RRQR, while the rows are selected by the same pivoting strategy applied to Q_1^T such that $||\bar{Q}_{21}\bar{Q}_{11}^{-1}||_{max}$ is bounded. The obtained approximation $A_k = AV_1(U_1AV_1)^{-1}U_1A$ is again spectrum preserving and a kernel approximation of A, with γ in Definition 1.2 and Definition 1.3 being a low degree polynomial in n and k (we note here that Definition 1.2 holds for the singular values of U_1AV_1 instead of the singular values of A_k). This factorization is also referred to as CUR, since AV_1 and U_1A correspond to columns and rows of the matrix A, respectively. Interpolative decomposition (ID), another popular approach, can also be described in this framework. We discuss it in more detail in section 4.

Several randomized algorithms rely on V_1 being a random matrix, typically based on Johnson-Lindenstrauss transforms or fast Johnson-Lindenstrauss transforms, such as the subsampled randomized Hadamard transform (SRHT) of Definition 6.7 introduced originally in [29]. The randomized QB, also referred to as randomized SVD (see, e.g., [17]), is obtained by choosing $U_1 = Q_1^T$ and corresponds to computing l steps of the QR factorization of UAV. The randomized SVD via row extraction is obtained by choosing U to be a row permutation such that $||\bar{Q}_{21}\bar{Q}_{11}^{-1}||_{max}$ is bounded. The obtained factorization corresponds to computing l steps of the LU factorization of UAV, and we refer to this as randomized LU with row selection. The permutation U can be obtained by computing the strong RRQR factorization of Q_1^T , or directly of $(AV_1)^T$ as in randomized ID; see, e.g., [23].

In the more general case when $k \leq l \leq l'$, which is the focus of this paper, the clean formulation of A_k described in (1.2) becomes a bit more complicated:

$$(1.6) A_k = [U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+][U_1A],$$

where U_1 and (U_1AV_1) are of dimensions $l' \times m$ and $l' \times l$, respectively. However, the algorithmic implementation is still straightforward and inexpensive. See (3.18) and Algorithm 3.1 for a detailed derivation. Subsection 1.2 summarizes properties of this novel approximation that are discussed in detail later in the paper. As in the previous case, U_1, V_1 can select deterministically or randomly columns, rows of A, respectively, or can be random matrices.

Several algorithms rely on an oblique projection to compute a low-rank approximation. One example is the popular approach introduced by Clarkson and Woodruff in [7], which we refer to as **CW**. Letting $\mathcal{P}^{ob} = AV_1(U_1AV_1)^+U_1$ and $k \leq l \leq l'$, we compute the approximation as

(1.7)
$$A_k^{ob} = AV_1(U_1AV_1)^+U_1A = \mathcal{P}^{ob}A,$$

(1.8) or equivalently
$$\mathcal{P}^{ob}A(:,j) := \arg\min_{x \in range(AV_1)} \|U_1(x - A(:,j))\|_2$$

where A(:,j) denotes the jth column of A. We show in Proposition 4.2 that this approximation is never more accurate, and can be less accurate than **GLU** approximation from (1.6) when l < l'. The Nyström method for symmetric positive semidefinite (SPSD) matrices is obtained by taking $U_1 = V_1^T$, with the approximation becoming $AV_1(V_1^TAV_1)^+(AV_1)^T$. In randomized Nyström, the matrix V_1 either can be chosen to randomly sample columns of A by using a uniform or a nonuniform importance sampling distribution or can be a random matrix. For more details see, e.g., [13].

We also note that the approximation A_k^{ob} from (1.7) is a special case of the approximation obtained by **GLU** when U_1, V_1 are such that

$$(1.9) (U_1AV_1)(U_1AV_1)^+U_1A = U_1A,$$

in which case A_k from (1.6) becomes equal to A_k^{ob} from (1.7). This condition is satisfied if and only if $range(U_1A) = range(U_1AV_1)$, since $(U_1AV_1)(U_1AV_1)^+$ is the projection onto $range(U_1AV_1)$. That is, V_1 spans the range of $(U_1A)^T$.

Given $C = U_1 A$ (or $C = U_1 (AA^T)A$, or any other matrix of the same dimensions) and $B = AV_1$ (again with other possibilities), another approximation for A is $A_k = B(B^+AC^+)C$, which is known to minimize $||A - BMC||_F$ over all possible choices of M; see homework 3.12 from [10]. When A and B are dense, it costs at least O(lmn) in general just to compute B^+A (ignoring Strassen-like algorithms), so asymptotically more than our approximations in (1.6) or (1.7). We will not consider it further.

1.2. Detailed bounds. GLU satisfies bounds at least as sharp as in the literature, and many are new. Proposition 3.3 gives bounds for the spectral and the kernel approximation provided by A_k from (1.6) for general U_1 and V_1 , and in section 6 properties of U_1 and V_1 specific to the algorithm are used to complete the bound. Both Proposition 3.3 and Proposition 5.2 provide new deterministic bounds not found in the literature. For example, Proposition 5.2 generalizes Theorem 9.1 of [17] to include singular values $\sigma_j(A)$ with j > 1. This generalization proves useful when analyzing Definition 1.3 for randomized algorithms, which we observe to be an advantage of GLU over CW, the approach introduced by Clarkson and Woodruff in [7].

Section 6 contains our new results after suitable random ensembles are chosen, that is, when V_1 and U_1 are random matrices. Extra attention is given to the SRHT ensemble of Definition 6.7, because the especially good bounds it can provide were not fully exploited in past literature. Using this ensemble, from Algorithm 3.1 for computing **GLU** we can see the number of arithmetic operations is $O(nm \log(l') + mll')$. Plugging in l' and l from Theorem 6.9, we can produce a low-rank approximation in $\tilde{O}(nm + k^2m\epsilon^{-3})$ time that relative to the squared error of the truncated SVD of rank k,

• approximates A with only $1 + O(\epsilon)$ times the squared Frobenius norm error,

(1.10)
$$||A - A_k||_F^2 = (1 + O(\epsilon))(\sigma_{k+1}^2 + \dots + \sigma_{\min(m,n)}^2),$$

• approximates A with only $O\left(1 + \frac{\epsilon \log(\min(m,n)/\delta)}{k \log(k/\delta)} \frac{\|\Sigma_2\|_F^2}{\|\Sigma_2\|_2^2}\right)$ times the squared spectral norm error,

$$||A - A_k||_2^2 = O\left(1 + \frac{\epsilon \log(\min(m, n)/\delta)}{k \log(k/\delta)} \frac{(\sigma_{k+1}^2 + \dots + \sigma_{\min(m, n)}^2)}{\sigma_{k+1}^2}\right) \sigma_{k+1}^2.$$

Here Σ_2 represents the trailing singular values, starting at k+1. This holds with probability $1-5\delta$, and l,l' grow polylogarithmically with δ , as in Remark 6.13. In other words, the algorithm we propose attains $\gamma=O(1)$ in Definition 1.1 for many families of A matrices encountered in practice with modest spectral decay (which makes the Frobenius norm not too much larger than the spectral norm). The same Theorem 6.9 shows this $\gamma=O(1)$ bound carries over to Definition 1.3. Further, Theorem 6.6 shows that Definition 1.2 is satisfied with $\gamma=O(\sqrt{\frac{k}{n}})$. To the best of our knowledge, no other work has found such a representation of a general matrix A in time o(nmk) satisfying these properties.

Our bounds have interesting implications for the growth factor of pre- and post-conditioned Gaussian elimination. Corollary 6.17 is a step towards a theoretical understanding of conditioning Gaussian elimination to avoid pivoting. Besides this, it

expands the classes of distributions for which pivoting is provably unnecessary to a class including Gaussian-distributed matrices. We pose an open question at the end, motivated by this analysis. For a more detailed discussion of the growth factor of Gaussian elimination and the implications of pivoting, see, e.g., [31] for an experimental study or [28] for its smoothed analysis.

- 1.3. Notation. As this paper is notation heavy, we first take a moment to collect some conventions we will use.
 - $A ext{ is } m \times n$.
 - Res(A,k) denotes A after setting its leading k-1 singular values to 0, so we restrict to the singular values starting at position k.
 - $[Q,R] = \mathbf{QR}(A)$ is the square QR decomposition of A, so Q is $m \times m$.
 - $[Q, R] = \mathbf{thin} \mathbf{QR}(A)$ is the thin QR decomposition of A, so Q is $m \times n$.
 - A^+ is the $n \times m$ Moore–Penrose pseudo-inverse. A_{11}^+ will refer to $(A_{11})^+$.
 - $[U, \Sigma, V] = \mathbf{SVD}(A)$ will be the full variant $(m \times n \Sigma)$ and with decreasing singular values. So U is $m \times m$, and V is $n \times n$.
 - $\Sigma_2 = \Sigma[k+1:,k+1:]$, the trailing singular values, starting at k+1.
 - $S(A_{11}) = A_{22} A_{21}A_{11}^{+}A_{12}$ is the Schur complement of A_{11} ; if the dimension of A_{11} is $l' \times l$, then $S(A_{11})$ is $(m-l') \times (n-l)$. Here $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$.

 • MATLAB-like notation to select submatrices; e.g., A[:k,:k] is the leading
 - $k \times k$ submatrix of A.
- 2. Related work. Low-rank matrix approximations have been extensively studied; hence this work is related to a large body of literature. Because of our emphasis on the LU factorization viewpoint, we should mention some work related to LU factorizations. Such papers providing information regarding Definitions 1.1, 1.2, and 1.3 are few, notably including perhaps the first [24], as well as later more efficient versions like [14]. These papers do not exploit randomness, however.

Exploiting randomness for low-rank factorizations has led to major speedups. Some literature in recent years has exploited this for LU factorizations, including perhaps most relevantly [30]. Their work has somewhat different goals, in that it seeks to find left and right permutation matrices, which makes it in some ways more Also, their paper only discusses spectral norm bounds on the residual. Interestingly, the fast version of their procedure (their Algorithm 4.4) uses an ensemble equivalent to the SRHT ensemble. The bounds we have in Theorem 6.9 are better for the spectral norm of the residual. Comparing our Theorem 6.9 with their Theorem 4.12, our approximation is always a factor on the order of \sqrt{n} more accurate, and a factor n more accurate when the spectrum decays sufficiently quickly. Our results utilizing the SRHT ensemble build on [4], which proved the SRHT ensemble has geometry preserving properties beyond those of the Johnson-Lindenstrauss transform properties. They used this fact to provide sharper spectral norm bounds on the residual for the randomized SVD approach to low-rank matrix approximation.

Outside of research into LU factorizations, many papers have focused on studying Johnson-Lindenstrauss embeddings. This has culminated in algorithms considered to run in nnz(A) time for many problems related to and including low-rank approximations. Notable such papers include [29, 7, 27]. For example, [27] uses the same factorization as [7], whose technical report we believe to be the first paper to use sketching from the left and right to speed up the algorithm. This body of literature has focused more on the properties of the random ensemble and less on the properties of the factorization itself. The error bound of these algorithms, given in terms of Frobenius norm of the residual, is similar to ours, $||A - A_k||_F^2 \le (1 + \epsilon) ||\Sigma_2||_F^2$ see,

e.g., [29, 7, 27] or [33] for a more detailed discussion. From this Frobenius norm error and Theorem 3.4 in [15], it is possible to derive a spectral norm error of the form $\|A - B\|_2^2 \le (1 + \epsilon \frac{\|\Sigma_2\|_F^2}{\|\Sigma_2\|_2^2}) \sigma_{k+1}^2$, as discussed in section 6 of [26]. But this error is weaker than the error in (1.11) attained by **GLU**. There are a few algorithms in the literature, such as the ones in [9, 8], that do achieve this stronger spectral error guarantee of the form $\|A - B\|_2^2 \le (1 + \frac{\epsilon}{k} \frac{\|\Sigma_2\|_F^2}{\|\Sigma_2\|_2^2}) \sigma_{k+1}^2$. In particular the work in [8] relies on sketching through a fast transform, while using the additional assumption that the matrix has a stable rank, that is, $\|A\|_F^2/\|A\|_2^2 \le k$. However, the low-rank approximation is obtained through an orthogonal projection as in (1.4). Hence, even if AV_1 can be computed at a lower cost, computing $Q_1^T A$ still requires 2mnk operations when A is dense (with non–Strassen-like algorithms), and thus these algorithms do not achieve o(nmk) complexity. They become expensive when k is modestly large, say a small power of n. An algorithm that has sublinear complexity and spectral error norm guarantee similar to ours is described in [26]. But it requires the matrix A to be positive semidefinite, and thus our algorithm is more general. We also note that the randomized low-rank community has typically not considered properties like spectrum preserving and kernel approximation (Definition 1.2 and Definition 1.3).

To date, procedures for the residual being within an ϵ factor as accurate as the truncated SVD with respect to the spectral norm do not gain any speed advantage by using fast Johnson–Lindenstrauss ensembles. This is because a repeated squaring must be used, and therefore structured sketching matrices have no advantage. Important work in this area includes [15] and [25].

This list is far from complete, and many different takes on the problem have been proposed which tangentially touch this paper; [17] and [33] are useful for finding more pointers into the literature.

3. Generalized LU factorization. Classically, as in [24] and [14], the rank-revealing LU factorization finds permutations U, V (usually stepwise over the procedure), forming $\bar{A} = UAV$, and LU factors \bar{A} but deletes the Schur complement after k steps. Thus,

$$\bar{A} = \left(\begin{array}{cc} I & 0 \\ \bar{A}_{21}\bar{A}_{11}^{-1} & I \end{array} \right) \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{12} \\ 0 & \mathcal{S}(\bar{A}_{11}) \end{array} \right) \approx \left(\begin{array}{cc} I \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{array} \right) \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{12} \end{array} \right) =: \bar{A}_k.$$

This naturally suggests the approximation $A \approx A_k := U^T \bar{A}_k V^T$. Letting V_1 be the first k columns of V and U_1 be the first k rows of U, some algebra (see Remark 3.4 for the more general case) shows the approximation to be $A \approx AV_1(U_1AV_1)^{-1}U_1A$.

This paper generalizes the rank-revealing LU factorization in two directions. First, we include other matrices on the left and right besides permutations. This allows for speedups through matrix sketching. Second, we generalize one step further by using rectangular Schur complements. This can greatly improve the quality of the low-rank approximation, as we will see in Proposition 3.3 and Theorem 6.6.

We describe this second modification in greater detail now. For the sake of analysis it will be convenient to let U, V be square matrices in the following discussion and subsequent Proposition 3.3. The relevant matrices are the $m \times n$ matrix A which we wish to approximate, the invertible $m \times m$ matrix U, and the invertible $n \times n$ matrix V. Now define

$$\bar{A} := UAV = \left(\begin{array}{cc} I_{l'} & 0 \\ \bar{A}_{21}\bar{A}_{11}^+ & I_{m-l'} \end{array} \right) \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{12} \\ 0 & \$(\bar{A}_{11}) \end{array} \right),$$

where this is valid when the $l' \times l$ block \bar{A}_{11} has full column rank so that $\bar{A}_{11}^+ \bar{A}_{11} = I$. In particular we are assuming $l' \ge l$. To help visualize the construction, the following depicts the block sizes:

$$\begin{split} \bar{A} &= \left(\begin{array}{cc} l', l & l', n-l \\ m-l', l & m-l', n-l \end{array} \right) \\ &= \left(\begin{array}{cc} l', l' & l', m-l' \\ m-l', l' & m-l', m-l' \end{array} \right) \left(\begin{array}{cc} l', l & l', n-l \\ m-l', l & m-l', n-l \end{array} \right). \end{split}$$

Deleting the $(m-l') \times (n-l)$ Schur complement and undoing the U,V factors gives the approximation we use as a definition,

$$(3.1) \hspace{1cm} A \approx A_k := U^{-1} \left(\begin{array}{c} I_{l'} \\ \bar{A}_{21} \bar{A}_{11}^+ \end{array} \right) \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{12} \end{array} \right) V^{-1}.$$

In (3.1), U and V are square, but for low-rank approximations this would be expensive. Only the leading l' rows of U and the leading l columns of V are actually required, but we find the square form helpful for the analysis. Accordingly for U, we assume that we may express

(3.2)
$$U = \begin{pmatrix} U_1 \\ U_2 \end{pmatrix}, U^{-1} = \begin{pmatrix} U_1^+ & U_2^T \end{pmatrix},$$

where U_2 has orthogonal rows spanning the orthogonal complement of U_1 . We also assume U_1 is full-rank so that U is invertible; any reasonable sketching matrix U_1 satisfies this property with probability 1. Similarly, we assume V may be expressed as

$$(3.3) V = \begin{pmatrix} V_1 & V_2 \end{pmatrix}, V^{-1} = \begin{pmatrix} V_1^+ \\ V_2^T \end{pmatrix},$$

where V_1 is full-rank and V_2 has orthogonal columns spanning the complement of the columns of V_1 . The assumptions on U_2 and V_2 are used later on in our theoretical results.

Singular values of a matrix product obey a well-known bound called the multiplicative Weyl inequality. We make use of this and its less known reverse version, so we give a proof here. Therefore, we state the inequality with a reference and prove its reverse version.

LEMMA 3.1. Say A is $m \times n$, and B is $n \times p$. For $1 \le k \le j$,

(3.4)
$$\sigma_i(AB) \le \sigma_{i-k+1}(A)\sigma_k(B).$$

Now assume for simplicity that $n \geq m \geq p$, both A, B are full rank, and $\operatorname{im}(B) \subset \ker(A)_{\perp}$. In other words, A is short-wide and B is tall-skinny, and the image of B is orthogonal to the kernel of A. Then for $1 \leq k \leq m-j$ and $j \leq p$, an inequality in the other direction is

(3.5)
$$\sigma_{m-k+1}(A)\sigma_{j+k-1}(B) \le \sigma_j(AB).$$

Besides these multiplicative inequalities, the additive Weyl inequality holds for any matrices A, B and $1 \le k, j \le n$, where n is the smaller of the row and column numbers, and says

(3.6)
$$\sigma_j(A+B) \le \sigma_{j-k+1}(A) + \sigma_k(B).$$

Proof. Inequalities (3.4) and (3.6) are well known. For example, see section 7.3, exercise 16 from [18].

We next prove (3.5). Let Σ_A, Σ_B be the square singular value matrices of A, B, respectively. Then AB is spectrally equivalent to $\Sigma_A U \Sigma_B$ for some $m \times p$ orthogonal matrix $U = V_1^T U_2$, with U_2 being the left singular matrix of B and V_1 being the right singular matrix of A. This U has orthonormal columns because it is norm preserving; $\operatorname{im}(U_2) \subset \operatorname{im}(V_1) = \ker(V_1^T)_{\perp}$, so if we let V extend V_1 to a square orthogonal matrix, then $\|V_1^T U_2 x\|_2 = \|V^T U_2 x\|_2 = \|x\|_2$. Σ_A is invertible based on the full rank assumption, and $U\Sigma_B$ is $m \times p$ with full column rank. Note that $(U\Sigma_B)^+ \Sigma_A^{-1}$ is a left inverse for $\Sigma_A U \Sigma_B$. Therefore, $(\Sigma_A U \Sigma_B)^+ = (U\Sigma_B)^+ \Sigma_A^{-1} P$, where P orthogonally projects onto $\operatorname{im}(\Sigma_A U \Sigma_B)$. Apply (3.4) to conclude $\sigma_j((\Sigma_A U \Sigma_B)^+) \leq \sigma_j((U\Sigma_B)^+ \Sigma_A^{-1})$. Combine this with another application of (3.4),

$$\sigma_{p-j+1}^{-1}(AB) = \sigma_{j}((AB)^{+}) \leq \sigma_{j}((U\Sigma_{B})^{+}\Sigma_{A}^{-1})$$

$$\leq \sigma_{j-k+1}((U\Sigma_{B})^{+})\sigma_{k}(\Sigma_{A}^{-1}) = \sigma_{p-(j-k+1)+1}^{-1}(V_{1}^{T}U_{2}\Sigma_{B})\sigma_{m-k+1}^{-1}(A)$$

$$= \sigma_{p-(j-k+1)+1}^{-1}(\Sigma_{B})\sigma_{m-k+1}^{-1}(A) = \sigma_{p-(j-k+1)+1}^{-1}(B)\sigma_{m-k+1}^{-1}(A).$$
(3.7)

We used that $V_1^T U_2$ is an orthogonal matrix to advance to line (3.7). Finally, reassign j = p - j + 1 to get the claimed (3.5).

Schur complements of rectangular blocks do not appear to be commonly used. The following derives an identity we require.

LEMMA 3.2. We continue to assume $l' \ge l$ and that \bar{A}_{11} has full column rank so that $\bar{A}_{11}^+\bar{A}_{11}=I$. Further introduce matrices U and V structured as explained in (3.2) and (3.3). Set $[Q,R] = \mathbf{Q}\mathbf{R}(AV)$ so that R is $m \times n$. Block R so that R_{11} is $l \times l$, and UQ so that $(UQ)_{11}$ is $l' \times l$. Then

$$S(\bar{A}_{11}) = S((UQ)_{11})R_{22}.$$

Proof. There is a factorization through a generalized LU factorization of A, in which the lower-triangular factor is the identity on the diagonal and the lower left factor is $\bar{A}_{21}\bar{A}_{11}^+$,

$$\bar{A} = \begin{pmatrix} I_{l'} \\ \bar{A}_{21}\bar{A}_{11}^{+} & I_{m-l'} \end{pmatrix} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \\ & \$(\bar{A}_{11}) \end{pmatrix}.$$

However, we could alternatively first use a QR factorization of AV followed by a generalized LU factorization of (UQ) (so that $(UQ)_{11}$ is $l' \times l$),

$$\begin{split} \bar{A} &= UAV = (UQ)R \\ &= \left(\begin{array}{cc} I_{l'} \\ (UQ)_{21}(UQ)_{11}^+ & I_{m-l'} \end{array} \right) \left(\begin{array}{cc} (UQ)_{11} & (UQ)_{12} \\ & \$((UQ)_{11}) \end{array} \right) \left(\begin{array}{cc} R_{11} & R_{12} \\ & R_{22} \end{array} \right) \\ (3.9) &= \left(\begin{array}{cc} I_{l'} \\ (UQ)_{21}(UQ)_{11}^+ & I_{m-l'} \end{array} \right) \left(\begin{array}{cc} (UQ)_{11}R_{11} & \dots \\ & \$((UQ)_{11})R_{22} \end{array} \right). \end{split}$$

The proof amounts to equating the blocks now between (3.8) and (3.9), but we provide a justification which essentially argues that the lower left block of the generalized LU factorization makes it unique.

The next proposition is critical for understanding the rank-revealing properties for **GLU**. It will combine with Proposition 5.2 to culminate in Theorem 6.6 and Theorem 6.9.

PROPOSITION 3.3. Let A be an $m \times n$ matrix, U and V be as in (3.2) and (3.3), $[Q,R] = \mathbf{Q}\mathbf{R}(AV)$, and finally $\bar{A} = UAV$. Block Q,R,A,\bar{A} as in Lemma 3.2; in particular, Q_{11} is $l' \times l$ and R_{11} is $l \times l$. Then the low-rank approximation suggested in (3.1), namely

$$A_k := U^{-1} \left(\begin{array}{c} I \\ \bar{A}_{21} \bar{A}_{11}^+ \end{array} \right) \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{12} \end{array} \right) V^{-1},$$

satisfies

$$(3.11) ||Res(A - A_k, j)||_F^2 \le ||Res(R_{22}, j)||_F^2 + ||(UQ)_{11}^+(UQ)_{12}^-(UQ)_{12}^-(Res(R_{22}, j))||_F^2,$$

$$||A - A_k||_2^2 \le ||R_{22}||_2^2 + ||(UQ)_{11}^+(UQ)_{12}R_{22}||_2^2,$$

(3.13)
$$\sigma_i^2(A - A_k) \le \|Res(R_{22}, j)\|_2^2 + \|(UQ)_{11}^+(UQ)_{12}Res(R_{22}, j)\|_2^2,$$

(3.14)
$$\sigma_i(A_k) \ge \sigma_i(A_k[:,:l']) = \sigma_i(AV_1V_1^+).$$

In the above, the relations for σ_j hold for $1 \le j \le \min(m,n) - k$. The relation for σ_i holds for $1 \le i \le k$.

Proof. The approximation loss in A_k is exactly the Schur complement $S(\bar{A}_{11})$. To establish this, we first do some matrix algebra. To start, we have

(3.15)
$$A_{k} = U^{-1} \begin{pmatrix} I \\ \bar{A}_{21}\bar{A}_{11}^{+} \end{pmatrix} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \end{pmatrix} V^{-1}$$
$$= U^{-1} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{21}\bar{A}_{11}^{+}\bar{A}_{12} \end{pmatrix} V^{-1}.$$

Next apply (3.2) from Lemma 3.2 to get $S(\bar{A}_{11}) = S((UQ)_{11})R_{22}$. From this and the fact that $U^{-1}\bar{A}V^{-1} = A$,

$$A - A_k = U^{-1} \begin{pmatrix} 0 \\ \$(\bar{A}_{11}) \end{pmatrix} V^{-1} = U^{-1} \begin{pmatrix} 0 \\ \$((UQ)_{11})R_{22} \end{pmatrix} V^{-1}$$

$$= \begin{pmatrix} 0 \\ U_2^T \$((UQ)_{11})R_{22}V_2^T \end{pmatrix}.$$

Now to get (3.10), recalling U_2 has orthonormal rows,

$$||A - A_k||_F^2 = ||S((UQ)_{11})R_{22}||_F^2 = ||[(UQ)_{22} - (UQ)_{21}(UQ)_{11}^+(UQ)_{12}]R_{22}||_F^2$$

$$= ||((UQ)_{21} (UQ)_{22})((UQ)_{11}^+(UQ)_{12}R_{22})||_F^2$$

$$\leq ||R_{22}||_F^2 + ||(UQ)_{11}^+(UQ)_{12}R_{22}||_F^2.$$

And for (3.12), similar steps produce

$$||A - A_k||_2^2 \le \left\| \left(\begin{array}{c} (UQ)_{11}^+(UQ)_{12}R_{22} \\ R_{22} \end{array} \right) \right\|_2^2 \le ||(UQ)_{11}^+(UQ)_{12}R_{22}||_2^2 + ||R_{22}||_2^2.$$

Even more generally, $\sigma_j(A - A_k) \leq \sigma_j(({}^{-(UQ)^+_{11}(UQ)_{12}R_{22}}))$ from the multiplicative Weyl inequality. Using this, and the additive Weyl inequality [18] in the second

inequality,

$$\begin{split} \sigma_{j+s-1}(A-A_k) &\leq \sigma_{j+s-1} \left(\left(\begin{array}{c} -(UQ)_{11}^+(UQ)_{12}(R_{22}-R_{22\mathrm{opt},j-1}+R_{22\mathrm{opt},j-1}) \\ R_{22}-R_{22\mathrm{opt},j-1}+R_{22\mathrm{opt},j-1} \end{array} \right) \right) \\ &\leq \sigma_s \left(\left(\begin{array}{c} -(UQ)_{11}^+(UQ)_{12}(R_{22}-R_{22\mathrm{opt},j-1}) \\ R_{22}-R_{22\mathrm{opt},j-1} \end{array} \right) \right) \\ &+ \sigma_j \left(\left(\begin{array}{c} -(UQ)_{11}^+(UQ)_{12}R_{22\mathrm{opt},j-1} \\ R_{22\mathrm{opt},j-1} \end{array} \right) \right) \\ &= \sigma_s \left(\left(\begin{array}{c} (UQ)_{11}^+(UQ)_{12}\mathrm{Res}(R_{22},j) \\ \mathrm{Res}(R_{22},j) \end{array} \right) \right), \end{split}$$

where we are letting $R_{22\text{opt},j-1}$ be the truncated SVD of rank j-1. In particular, this establishes

$$\begin{split} \sigma_{j}^{2}(A-A_{k}) &\leq \sigma_{1}^{2} \left(\left(\begin{array}{c} (UQ)_{11}^{+}(UQ)_{12} \mathrm{Res}(R_{22},j) \\ \mathrm{Res}(R_{22},j) \end{array} \right) \right) \\ &\leq \|(UQ)_{11}^{+}(UQ)_{12} \mathrm{Res}(R_{22},j)\|_{2}^{2} + \|\mathrm{Res}(R_{22},j)\|_{2}^{2}, \end{split}$$

and also by noting that the trailing $\min(m, n) - j$ singular values of $A - A_k$ are bound in this manner,

$$\|\operatorname{Res}(A - A_k, j)\|_F^2 \le \left\| \left(\begin{array}{c} (UQ)_{11}^+(UQ)_{12} \operatorname{Res}(R_{22}, j) \\ \operatorname{Res}(R_{22}, j) \end{array} \right) \right\|_F^2$$
$$= \|(UQ)_{11}^+(UQ)_{12} \operatorname{Res}(R_{22}, j)\|_F^2 + \|\operatorname{Res}(R_{22}, j)\|_F^2.$$

This completes (3.10)–(3.13). We proceed to the lower bound on $\sigma_i(A_k)$ claimed in (3.14). Let \bar{A}_k for the moment denote the middle matrix in (3.15), and $(\bar{A}_k)_1$ be the leading l columns of the matrix. Then because the rows of V_2^+ are orthogonal to the rows of V_1^+ , the rows of $(\bar{A}_k)_1V_1^+$ are orthogonal to those of $(\bar{A}_k)_1V_2^+$. Using this in the inequality step,

$$\sigma_i(A_k) = \sigma_i(U^{-1}\bar{A}_kV^{-1}) \ge \sigma_i(\left(\begin{array}{cc} U_1^+ & U_2^T \\ \end{array} \right) \cdot (\bar{A}_k)_1V_1^+) = \sigma_i(AV_1V_1^+). \qquad \qquad \Box$$

Remark 3.4. Recall the sizes $V_1 = V[:,:l]$, $U_1 = U[:l',:]$. When l' = l, the factorization in (3.1) can readily be rewritten in the more elegant form

(3.17)
$$A_k = AV_1(U_1AV_1)^{-1}U_1A.$$

One important feature of this is that only U_1, V_1 are actually needed to compute A_k . We will later see that the residual bounds in Proposition 3.3 can be computed with only U_1, V_1 , so it makes sense that we can find an analogue of (3.17) for l' > l. However, we actually need to set the rows $U_2 = U[l' + 1:,:]$ to be a basis for the orthogonal complement of the rows of U_1 in order to achieve this. Then $U^{-1} = [U_1^+, U_2^+]$, and we get a different form of (3.1) that is often faster to compute than (3.1),

$$A_{k} = U^{-1} \begin{pmatrix} I \\ \bar{A}_{21}\bar{A}_{11}^{+} \end{pmatrix} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \end{pmatrix} V^{-1} = \begin{pmatrix} U_{1}^{+} & U_{2}^{+} \end{pmatrix} \begin{pmatrix} I \\ \bar{A}_{21}\bar{A}_{11}^{+} \end{pmatrix} U_{1}A$$

$$= (U_{1}^{+} + U_{2}^{+}U_{2}AV_{1}(U_{1}AV_{1})^{+})U_{1}A$$

$$= \left[U_{1}^{+} + (I - U_{1}^{+}U_{1})AV_{1}(U_{1}AV_{1})^{+} \right] [U_{1}A]$$

$$= \left[U_{1}^{+} (I - (U_{1}AV_{1})(U_{1}AV_{1})^{+}) + (AV_{1})(U_{1}AV_{1})^{+} \right] [U_{1}A].$$

$$(3.18) = \left[U_{1}^{+} (I - (U_{1}AV_{1})(U_{1}AV_{1})^{+}) + (AV_{1})(U_{1}AV_{1})^{+} \right] [U_{1}A].$$

This final form should be viewed as a generalized LU factorization. The left factor is $m \times l'$, and the right factor (U_1A) is $l' \times n$. Also recall that U_1 is $l' \times m$, so the pseudo-inverse can be cheaply computed.

We summarize the factorization discussed above in (only partially specified because of U, V and the oversampling parameters l, l') Algorithm **GLU** (see Algorithm 3.1) and Algorithm **RLU** (see Algorithm 3.2). Recall that using square U, V was only to help with the theoretical guarantees. We also emphasize that Algorithm **RLU** is the special case of Algorithm **GLU** when the latter sets l = l'.

Algorithm 3.1 $[T,S] = \mathbf{GLU}(A,k)$. Generalized LU approximation computes a low-rank approximation $A \approx A_k = TS$, where T is a tall-skinny matrix and S is a short-wide matrix.

```
1: Input: target rank k, matrix A \in \mathbb{R}^{m \times n}

2: Output: T \in \mathbb{R}^{m \times l'}, S \in \mathbb{R}^{l' \times n}

3: Ensure: T = U_1^+(I - (U_1AV_1)(U_1AV_1)^+) + (AV_1)(U_1AV_1)^+, S = U_1A

4: Select oversampling parameters l' \geq l \geq k

5: Generate full-rank n \times l matrix V_1 and full-rank l' \times m matrix U_1

6: \hat{A} = U_1AV_1

7: T_1 = U_1^+(I - \hat{A}\hat{A}^+)

8: T_2 = AV_1

9: T_2 = T_2\hat{A}^+

10: T = T_1 + T_2

11: S = U_1A
```

Algorithm 3.2 $[T, \hat{A}, S] = \mathbf{RLU}(A)$. Rank-revealing LU computes a low-rank approximation $A \approx A_k = T\hat{A}^{-1}S$, where T is a tall-skinny matrix, S is a short-wide matrix, and \hat{A} is a small dense matrix.

```
1: Input: target rank k, matrix A \in \mathbb{R}^{m \times n}

2: Output: T \in \mathbb{R}^{m \times l}, S \in \mathbb{R}^{l \times n}, \hat{A} \in \mathbb{R}^{l \times l}

3: Ensure: T = AV_1, S = U_1A, \hat{A} = U_1AV_1

4: Select oversampling parameter l \geq k

5: Generate a full-rank n \times l matrix V_1 and a full-rank l \times m matrix U_1

6: T = AV_1

7: S = U_1A

8: \hat{A} = U_1T
```

The bounds in Proposition 3.3 are not fully developed, as V affects R and U affects UQ. In section 5 the R_{22} factor will be studied; it will be bound in terms of S^TV , where S is the right singular matrix of A. See Proposition 5.1 and the resulting Theorem 6.6. Section 6 describes how choosing suitable random ensembles for U_1, V_1 allows for the Frobenius norm of the residual to be arbitrarily close to that of the truncated SVD, as well as many other bounds. We therefore present what we consider to be our main results in section 6.

4. Relationship to other approaches. In this section we illustrate how GLU provides a general framework by proving the equivalence with Algorithm RandomizedLU_RowSelection and Algorithm RandomizedQB below and discussing re-

lations with interpolative decompositions. We also see a close connection with Algorithm **CW**, the approach from Clarkson and Woodruff [7]. We show that **CW** is never more accurate and can be less accurate than our approach when l' > l, and the same when l' = l.

We consider first the case when $k \leq l = l'$. For ease of understanding, we recall a few relations in this case. Let $\bar{A} = UAV$, where $A, \bar{A} \in \mathbb{R}^{m \times n}$, $U \in \mathbb{R}^{m \times m}$, and $V \in \mathbb{R}^{n \times n}$. Let $[Q, R] = \mathbf{QR}(AV)$, where $Q \in \mathbb{R}^{m \times m}$ and $R \in \mathbb{R}^{m \times n}$. Set $V_1 = V[:,:l]$, $U_1 = U[:l,:]$, $Q_1 = Q[:,:l]$, and $\bar{Q} = UQ$, and partition $\bar{Q} = \begin{pmatrix} \bar{Q}_{11} & \bar{Q}_{12} \\ \bar{Q}_{21} & \bar{Q}_{22} \end{pmatrix}$, $R = \begin{pmatrix} R_{11} & R_{12} \\ R_{22} \end{pmatrix}$ such that \bar{Q}_{11} and R_{11} are $l \times l$. It can be seen that $Q_1 R_{11}$ is the thin QR decomposition of AV_1 .

A decomposition of UAV can be obtained either directly from the LU factorization of UAV, or from the QR factorization of AV and then the LU factorization of UQ, as discussed previously in Lemma 3.2 in the general case $k \leq l \leq l'$. In the case discussed here, $k \leq l = l'$ and \bar{A}_{11} invertible, we obtain

$$(4.1) \bar{A} = UAV = \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} \end{pmatrix} = \begin{pmatrix} I \\ \bar{A}_{21}\bar{A}_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \\ S(\bar{A}_{11}) \end{pmatrix}$$
$$= (UQ)R = \begin{pmatrix} I \\ \bar{Q}_{21}\bar{Q}_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \bar{Q}_{11} & \bar{Q}_{12} \\ S(\bar{Q}_{11}) \end{pmatrix} \begin{pmatrix} R_{11} & R_{12} \\ R_{22} \end{pmatrix}$$
$$= \begin{pmatrix} I \\ \bar{Q}_{21}\bar{Q}_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} \bar{A}_{11} & \bar{A}_{12} \\ S(\bar{Q}_{11})R_{22} \end{pmatrix},$$

where, as shown in, e.g., [14], $\bar{Q}_{21}\bar{Q}_{11}^{-1} = \bar{A}_{21}\bar{A}_{11}^{-1}$, $S(\bar{A}_{11}) = S(\bar{Q}_{11})R_{22}$. We also have the following relations. Since $Q_1R_{11} = AV_1$, it follows that

$$(4.3) AV_1(U_1AV_1)^{-1}U_1A = Q_1R_{11}(U_1Q_1R_{11})^{-1}U_1A = Q_1(U_1Q_1)^{-1}U_1A.$$

Additionally,

$$(4.4) \hspace{1cm} AV_1(U_1AV_1)^{-1}U_1A = U^{-1}\begin{pmatrix} \bar{A}_{11} \\ \bar{A}_{21} \end{pmatrix} \bar{A}_{11}^{-1}U_1A = U^{-1}\begin{pmatrix} I_l \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{pmatrix} U_1A.$$

By dropping the Schur complement $S(\bar{A}_{11})$ from (4.1) and $S(\bar{Q}_{11})R_{22}$ from (4.2), and given (4.3) and (4.4), we obtain the approximations

$$(4.5) A_k = U^{-1} \begin{pmatrix} I_l \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{pmatrix} (\bar{A}_{11} \quad \bar{A}_{12}) V^{-1} = AV_1(U_1AV_1)^{-1}U_1A$$

$$(4.6) = U^{-1} \begin{pmatrix} I_l \\ \bar{Q}_{21}\bar{Q}_{11}^{-1} \end{pmatrix} (\bar{A}_{11} \quad \bar{A}_{12}) V^{-1} = Q_1 (U_1 Q_1)^{-1} U_1 A.$$

While we consider here infinite precision, we note that in finite precision, even for same U_1, V_1 , the approximations from (4.5) and (4.6) can be different.

Equations (4.5) and (4.6) reveal that, given V_1 , the permutation U and its submatrix U_1 can be obtained either from AV_1 or from the thin orthogonal factor Q_1 of AV_1 . In the first case, one approach consists in computing a strong RRQR factorization of $(AV_1)^T$,

(4.7)
$$(AV_1)^T U^T = \begin{pmatrix} \bar{A}_{11} \\ \bar{A}_{21} \end{pmatrix}^T = \tilde{Q} \begin{pmatrix} \tilde{R}_{11} & \tilde{R}_{12} \end{pmatrix},$$

Algorithm 4.1 $[T,S] = \mathbf{RandomizedQB}(A,k)$. Randomized QB approximation computes a low-rank approximation $A \approx A_k = TS$ where T is a tall-skinny matrix with orthonormal columns, and S is a short-wide matrix.

```
1: Input: target rank k, matrix A \in \mathbb{R}^{m \times n}
```

- 2: Output: orthogonal matrix $T \in \mathbb{R}^{m \times l}$, matrix $S \in l \times n$
- 3: **Ensure:** T has orthonormal columns, $S = T^T A$
- 4: Select the oversampling parameter $l \geq k$
- 5: Generate a full rank $n \times l$ matrix V_1
- 6: $\hat{A} = AV_1$
- 7: $[T, _] = \mathbf{thin} \mathbf{QR}(\hat{A})$
- 8: $S = T^T A$

such that $\|\tilde{R}_{11}^{-1}\tilde{R}_{12}\|_{max}$ is bounded. It can be seen that $\bar{A}_{21}\bar{A}_{11}^{-1}=\tilde{R}_{12}^T(\tilde{R}_{11}^{-1})^T$, and hence $\|\bar{A}_{21}\bar{A}_{11}^{-1}\|_{max}$ is also bounded. A detailed derivation can be found in [20], where this pivoting strategy is used to compute a block LU factorization with a growth factor smaller than the one of LU with partial pivoting. In the second case, U_1 is selected by computing a strong RRQR factorization of Q_1^T such that $\|\bar{Q}_{21}\bar{Q}_{11}^{-1}\|_{max}$ is bounded and the singular values of $\bar{Q}_{11} = U_1Q_1$ are upper bounded by 1 and lower bounded by 1/q(m,k), where q(m,k) is a low degree polynomial in m and k. For a derivation see [14], where this strategy is used in a deterministic algorithm referred to as LU CRTP, which computes a spectrum preserving and a kernel approximation of A.

We discuss now several instances of V_1 , U_1 that correspond to existing randomized algorithms in the literature. We choose to have a similar output for all the algorithms as $A_k = TS$, where $T \in \mathbb{R}^{m \times l}$ and $S \in \mathbb{R}^{l \times n}$.

Randomized QB presented in Algorithm 4.1 is also referred to in the literature as randomized SVD. The output is different in randomized SVD, though. It is obtained by computing the SVD of S and is provided as the product of an orthogonal matrix, a diagonal matrix, and a second orthogonal matrix. In our framework, this algorithm can be obtained by taking $U_1 = Q_1^T$, in which case it can be seen from (4.6) that the approximation becomes $Q_1Q_1^TA$.

Randomized LU with row selection is presented in Algorithm 4.2, in which V_1 is random and the row selection is obtained from Q_1 , the thin orthogonal factor of AV_1 , as in (4.6). The output is $T = U^T \binom{I_l}{Q_{21}Q_{11}^{-1}} = Q_1(U_1Q_1)^{-1}$ and $S = U_1A$, such that $||T||_{max}$ is bounded and S is formed by the selected rows of A. The algorithm could also return $T = Q_1$ and an orthogonal matrix, $S = (U_1Q_1)^{-1}U_1A$. We expect the computation of S to be numerically stable, since (U_1Q_1) is expected to be well conditioned. Obviously, the output can also be the product of a permutation matrix, a lower triangular matrix, and an upper triangular matrix, which would require computing the LU factorization of U_1AV_1 . Algorithm 4.2 corresponds to a randomized SVD with row selection algorithm described in section 5.2 of [17]. See discussion in [17] around (5.3).

Interpolative decompositions have three different versions, column ID, row ID, and double-sided ID. We refer the reader for a detailed discussion to [23]. For example, the column ID approximates A as CZ, where $C \in \mathbb{R}^{m \times l}$ is formed by l columns of A, $Z \in \mathbb{R}^{l \times n}$ contains an identity matrix, and $\|Z\|_{max} \leq 1$. As discussed in [23], such a decomposition can be obtained by computing l steps of a column pivoted QR factorization of A, AV = QR, where V is a permutation matrix. By letting

Algorithm 4.2 $[T,S] = \text{RandomizedLU_RowSelection}(A,k)$. Randomized LU with row selection approximation computes $A_k = TS$, performing sketching on the columns and selecting rows based on Q_1 , the thin orthogonal factor of AV_1 .

- 1: Input: target rank k, matrix $A \in \mathbb{R}^{m \times n}$
- 2: Output: $T \in \mathbb{R}^{m \times l}$ and $S \in \mathbb{R}^{l \times n}$
- 3: **Ensure:** $T = AV_1(U_1AV_1)^{-1} = U^T\begin{pmatrix} I_l \\ \bar{Q}_{21}\bar{Q}_{11}^{-1} \end{pmatrix} = Q_1(U_1Q_1)^{-1},$ $S = U_1A \in \mathbb{R}^{l \times n}$, where U is a permutation matrix, $U_1 = U[:l,:]$, and $Q_1 \in \mathbb{R}^{m \times l}$ has orthonormal columns, $||T||_{max}$ is bounded, and S is formed by rows of A
- 4: Select oversampling parameter $l \ge k$
- 5: Generate a full-rank $n \times l$ random matrix V_1
- 6: $[Q_1, _] = \mathbf{thin} \mathbf{QR}(AV_1)$
- 7: Permutation U is selected so that $UQ_1 = \bar{Q}_1 = \begin{pmatrix} \bar{Q}_{11} \\ \bar{Q}_{21} \end{pmatrix}$ results in $||\bar{Q}_{21}\bar{Q}_{11}^{-1}||_{max}$ being bounded by a small constant (see [22]). Here $U_1 = U[:l,:]$ and $\bar{Q}_{11} = U_1Q_1$
- 8: $T = U^T \begin{pmatrix} I_l \\ \bar{Q}_{21}\bar{Q}_{11}^{-1} \end{pmatrix}$; also note that $T = AV_1(U_1AV_1)^{-1}$
- 9: $S = U_1 A$

 $Q_1 = Q[:,:l], \ V_1 = V[:,:l], \ R_{11} = R[:l,:l], \ R_{12} = R[:l,l+1:],$ column ID is obtained as

$$(4.8) Q_1 Q_1^T A = Q_1 (R_{11} R_{12}) V^T = Q_1 R_{11} (I R_{11}^{-1} R_{12}) V^T$$

$$(4.9) = AV_1 \begin{pmatrix} I & R_{11}^{-1}R_{12} \end{pmatrix} V^T = CZ,$$

where $C = AV_1$ and $Z = \begin{pmatrix} I & R_{11}^{-1}R_{12} \end{pmatrix} V^T$. The elements of Z can be bounded by using strong RRQR to compute the column pivoted QR factorization of A, in which case $\|R_{11}^{-1}R_{12}\|_{max}$ is bounded by a small constant that can be chosen to be 1. A row ID approximates A as XR, where $X \in \mathbb{R}^{m \times l}$ contains an identity matrix, $\|X\|_{max} \leq 1$, and $R \in \mathbb{R}^{l \times n}$ is formed by l rows of A. It can be obtained by computing the column ID of A^T . However, we note that when U is a permutation matrix, the approximations from (4.5) and (4.6) also provide row IDs, since U_1A is formed by rows of A, and $\|\bar{A}_{21}\bar{A}_{11}^{-1}\|_{max}$ and $\|\bar{Q}_{21}\bar{Q}_{11}^{-1}\|_{max}$ can be bounded by 1. The randomized version is slightly different, and we discuss here **RandomizedRowID** in Algorithm 4.3, referred to as RandomizedID in Algorithm 15 from [23]. With our terminology, it corresponds to randomized LU with row selection, in which, as in (4.5), the rows are selected from AV_1 instead of its thin Q_1 factor.

The fact that these algorithms fit into the LU framework is simple, but it appears to have been overlooked in the literature. Therefore, it has its own proposition.

PROPOSITION 4.1. RandomizedLU_RowSelection is equivalent to RLU when the latter chooses the same V_1 and U_1 . RandomizedRowID is equivalent to RLU when the latter chooses the same V_1 and U_1 . RandomizedQB is equivalent to RLU when the latter chooses the same V_1 and $U_1 := T^T$.

Proof. The proof is mainly to recall the various definitions. First, Algorithm RandomizedLU_RowSelection produces $A_k = Q_1(U_1Q_1)^{-1}U_1A$. As claimed within the algorithm and (4.3), because $Q_1R_{11} = AV_1$ it follows that $AV_1(U_1AV_1)^{-1}U_1A = Q_1(U_1Q_1)^{-1}U_1A$. As $AV_1(U_1AV_1)^{-1}U_1A$ is the output factorization of Algorithm

Algorithm 4.3 [T, S] =**RandomizedRowID**(A, k). Randomized interpolative decomposition with row selection computes a randomized LU factorization $A_k = TS$, performing sketching on the columns and selecting rows based on AV_1 , where V_1 is a random matrix.

- 1: **Input:** target rank k, matrix $A \in \mathbb{R}^{m \times n}$
- 2: Output: $T \in \mathbb{R}^{m \times l}$ and $S \in \mathbb{R}^{l \times n}$
- 3: **Ensure:** $T = AV_1(U_1AV_1)^{-1} = U^T(\frac{I_l}{\bar{A}_{21}\bar{A}_{11}^{-1}})$, $S = U_1A \in \mathbb{R}^{l \times n}$, where U is a permutation matrix, $U_1 = U[:l,:]$, $UAV_1 = (\frac{\bar{A}_{11}}{\bar{A}_{21}})$, $||T||_{max}$ is bounded, and S is formed by rows of A
- 4: Select oversampling parameter $l \ge k$
- 5: Generate a full-rank $n \times l$ random matrix V_1
- 6: Permutation U is selected so that $UAV_1 = \begin{pmatrix} \bar{A}_{11} \\ \bar{A}_{21} \end{pmatrix} = \begin{pmatrix} I \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{pmatrix} \bar{A}_{11}$ results in $||\bar{A}_{21}\bar{A}_{11}^{-1}||_{max}$ being bounded by a small constant, e.g., by computing strong RRQR factorization of $(AV_1)^T$. Here $U_1 = U[:l,:]$ and $\bar{A}_{11} = U_1AV_1$.
- 7: $T = U^T \begin{pmatrix} I_l \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{pmatrix}$; also note that $T = AV_1(U_1AV_1)^{-1}$.
- 8: $S = U_1 A$

RLU, the factorizations agree. Second, consider Algorithm **RandomizedRowID**. As in (4.4), $AV_1(U_1AV_1)^{-1}U_1A = U^T\begin{pmatrix} I_l \\ \bar{A}_{21}\bar{A}_{11}^{-1} \end{pmatrix}U_1A$, which corresponds to the output of Algorithm **RandomizedRowID**.

We move on to the last equivalence. Recall that $[T, R] = \mathbf{thin} - \mathbf{QR}(AV_1)$. Selecting the same V_1 and $U_1 = T^T$, the random LU approximation given by Algorithm **RLU** would be

(4.10)
$$AV_1(T^TAV_1)^{-1}T^TA = TR(T^TTR)^{-1}T^TA = TT^TA,$$

which agrees with Algorithm RandomizedQB.

We consider now the more general case $k \leq l \leq l'$. The most popular approach involving sketching from the left and right is perhaps the method based on an oblique projection, in which the approximation is

(4.11)
$$A_k^{ob} = AV_1(U_1AV_1)^+U_1A.$$

This approximation is used in [7] and also described in the overview [33]. We refer to it as Algorithm **CW** after Clarkson and Woodruff. We now show that this approximation is never more accurate and can be less accurate than **GLU** when l' > l, and the same when l' = l.

PROPOSITION 4.2. Let $\bar{A} = UAV$, A_k be the output of Algorithm GLU, and A_k^{ob} be the output (4.11) of Algorithm CW. Then

$$\|A - A_k^{ob}\|_F^2 = \|A - A_k\|_F^2 + \|A_k - A_k^{ob}\|_F^2.$$

Proof. Similar to Remark 3.4,

$$\begin{split} A_k^{ob} &= AV_1(U_1AV_1)^+ U_1A = U^{-1} \left(\begin{array}{c} \bar{A}_{11} \\ \bar{A}_{21} \end{array} \right) \bar{A}_{11}^+ \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{12} \end{array} \right) V^{-1} \\ &= U^{-1} \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{11}\bar{A}_{11}^+ \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} - 8(\bar{A}_{11}) \end{array} \right) V^{-1}. \end{split}$$

From this calculation and the calculation leading to (3.16), it follows that

$$\begin{split} \|A - A_k^{ob}\|_F^2 &= \|U^{-1} \left[\bar{A} - \left(\begin{array}{cc} \bar{A}_{11} & \bar{A}_{11} \bar{A}_{11}^+ \bar{A}_{12} \\ \bar{A}_{21} & \bar{A}_{22} - 8(\bar{A}_{11}) \end{array} \right) \right] V^{-1}\|_F^2 \\ &= \| \left(\begin{array}{cc} 0 & U_1^+ (I - \bar{A}_{11} \bar{A}_{11}^+) \bar{A}_{12} \\ 0 & 8(\bar{A}_{11}) \end{array} \right) \|_F^2. \end{split}$$

This gives the Frobenius norm claim.

The work in [7] only considered the properties of the factorization (4.11) in the context of Johnson–Lindenstrauss transforms, specifically when l' is a poly-log factor larger than l, and not focusing on deterministic bounds. We note that the factorization in **CW** may be slightly cheaper to perform, although in typical settings (k is relatively small) the same term dominates the cost of both algorithms.

5. QR deterministic bounds. The following lemma is important in randomized low-rank approximation results. Our proof is novel, and (5.3), (5.4) significantly generalize past versions.

PROPOSITION 5.1. Let A be an $m \times n$ matrix with the SVD of A being $A = P\Sigma S^T$. As with Proposition 3.3, it is again convenient to suppose V is $n \times n$ with V as described in (3.3). Also let $[Q, R] = \mathbf{QR}(AV)$. Then block Q, R, S^TV, Σ using $Q_1 := Q[:,:l]$, $R_{11} = R[:l,:l]$, $(S^TV)_{11} = (S^TV)[:k,:l]$, and $\Sigma_1 = \Sigma[:k,:k]$, $\Sigma_2 = \Sigma[k+1:,k+1:]$. Then the nonzero singular values of $Q_1Q_1^TA - A$ are identical to those of R_{22} , i.e., for any $1 \le j \le \min(m,n) - l$

$$\sigma_j(Q_1Q_1^TA - A) = \sigma_j(R_{22}).$$

Moreover, assuming $(S^TV)_{11}$ has full row-rank (and therefore $k \leq l$), we have that

We may generalize this last equation with the goal of covering Definition 1.3. For any $1 \le j \le \min(m,n) - l$, there exists an $n \times n$ orthogonal matrix \tilde{S} independent of V, satisfying

(5.3)
$$\sigma_j^2(Q_1Q_1^TA - A) \le \|Res(\Sigma_2, j)\|_2^2 + \|Res(\Sigma_2, j)(\tilde{S}^TV)_{21}(\tilde{S}^TV)_{11}^+\|_{22}^2,$$

$$(5.4) \qquad \|Res(Q_1Q_1^TA - A,j)\|_F^2 \leq \|Res(\Sigma_2,j)\|_F^2 + \|Res(\Sigma_2,j)(\tilde{S}^TV)_{21}(\tilde{S}^TV)_{11}^+\|_F^2$$

with $(\tilde{S}^T V)_{11}$ being $k \times l$ as before.

Proof. We first observe by direct computation

$$\sigma_{j}(Q_{1}Q_{1}^{T}A-A)=\sigma_{j}(Q_{2}Q_{2}^{T}A)=\sigma_{j}(Q_{2}\left(\begin{array}{cc}0&R_{22}\end{array}\right)V^{-1})=\sigma_{j}(R_{22})$$

to establish the first claim. Next we invoke the fact that for the spectral and Frobenius norms, $Q_1Q_1^TA$ is the best approximation to A whose columns are in $\operatorname{im}(Q_1)$. For example, one can check that $Q_1Q_1^T$ satisfies the orthogonal projection properties with respect to these norms. Set $\bar{A} = P^TAV$. Then we explicitly propose an approximation \tilde{A}_k whose columns are contained in $\operatorname{im}(Q_1) = \operatorname{im}(AV_1)$, namely

$$\tilde{A}_k := P \left(\begin{array}{c} \bar{A}_{11} \\ \bar{A}_{12} \end{array} \right) \left(\begin{array}{cc} I & \bar{A}_{11}^+ \bar{A}_{12} \end{array} \right) V^{-1} = A V_1 \left(\begin{array}{cc} I & \bar{A}_{11}^+ \bar{A}_{12} \end{array} \right) V^{-1}.$$

Note that this is just using the **GLU** approximation, except in contrast to before, \bar{A}_{11} is $k \times l$, making it short and wide. Thus we may repeat the algebra around (3.16) in the first step,

$$||A - \tilde{A}_{k}||_{F}^{2} = ||\mathcal{S}((P^{T}AV)_{11})||_{F}^{2} = ||\mathcal{S}(\Sigma_{1}(S^{T}V)_{11})||_{F}^{2}$$

$$= ||\Sigma_{2}(S^{T}V)_{22} - \Sigma_{2}(S^{T}V)_{21}(\Sigma_{1}(S^{T}V)_{11})^{+}\Sigma_{1}(S^{T}V)_{12}||_{F}^{2}$$

$$= ||\Sigma_{2}(S^{T}V)_{22} - \Sigma_{2}(S^{T}V)_{21}(S^{T}V)_{11}^{+}(S^{T}V)_{12}||_{F}^{2}$$

$$\leq ||(\Sigma_{2} - \Sigma_{2}(S^{T}V)_{21}(S^{T}V)_{11}^{+})||_{F}^{2}$$

$$= ||\Sigma_{2}||_{F}^{2} + ||\Sigma_{2}(S^{T}V)_{21}(S^{T}V)_{11}^{+}||_{F}^{2}.$$

To be clear, S^TV was blocked so that $(S^TV)_{11}$ is $k \times l$. Note that we were able to distribute the pseudo-inverse in $(\Sigma_1(S^TV)_{11})^+$. In the generic case this follows from $(S^TV)_{11}$ having full row rank (this will be with probability 1 for suitably random V) and Σ_1 being invertible. If Σ_1 has trailing 0 values, the assumption of full row rank of $(S^TV)_{11}$ ensures $\operatorname{im}(AV_1) = \operatorname{im}(A)$ and therefore we can instead use $\tilde{A}_k := A$ to get the bound of 0.

For the spectral norm bound, the steps are the same, except the final equality becomes an inequality.

We actually are interested in the lower singular values as well, though, so we extend the proof. In the following, P_Y and P_{AV_1} project onto the complements of the images of Y and AV_1 , respectively; Y is rank j-1, and AV_1 is rank l. The same projection notation applies to the other projections. Using the additive Weyl inequality in the inequality step, similar to the use within Proposition 3.3, for $s \leq \min(m,n)-j$,

$$\begin{split} \sigma_{j+s-1}(Q_1Q_1^TA - A) &= \sigma_{j+s-1}(P_{AV_1}A) \\ &= \sigma_{j+s-1}(P_{AV_1}A - P_YP_{AV_1}A + P_YP_{AV_1}A) \\ &\leq \sigma_s(P_YP_{AV_1}A) = \sigma_s(P_{Y+AV_1}A) = \sigma_s(P_{\tilde{Y}}P_{V'}A). \end{split}$$

Additionally, in the third equality, $Y + AV_1$ is used to refer to direct sum of the images of Y and AV_1 , and the equality holds under the assumption these spaces are orthogonal. The fourth (last) equality holds when $\operatorname{im}(Y') \oplus \operatorname{im}(\tilde{Y}) = \operatorname{im}(Y) \oplus \operatorname{im}(AV_1)$, and $\operatorname{im}(Y')$ is orthogonal to $\operatorname{im}(\tilde{Y})$.

Now we make our choice of $Y + AV_1$. First let $P_1 = P[:,:j-1]$, the leading j-1 left singular vectors of A. Noting that $\operatorname{im}(P_1) \oplus \operatorname{im}((P_1P_1^TA - A)V_1)$ is rank l+j-1 and contains $\operatorname{im}(Q_1) = \operatorname{im}(AV_1)$, and that P_1 is orthogonal to $(A - P_1P_1^TA)V_1$, these are valid choices of Y' and \tilde{Y} , respectively. In summary,

(5.6)
$$\operatorname{im}(Y) \oplus \operatorname{im}(AV_1) = \operatorname{im}(P_1) \oplus \operatorname{im}((P_1P_1^TA - A)V_1) = \operatorname{im}(P_1) \oplus \operatorname{im}(AV_1),$$
$$Y = P_1 \text{ orthogonalized against } AV_1,$$
$$Y' = P_1,$$
$$\tilde{Y} = P_{P_1}AV_1.$$

We emphasize in (5.6) that the first two are orthogonal direct sums. This puts us essentially back into the situation surrounding (5.5). Indeed,

$$\sigma_{j+s-1}(Q_1Q_1^TA - A) \le \sigma_s(P_{\tilde{Y}}P_{Y'}A) = \sigma_s(P_{BV_1}B) = \sigma_s(B - \tilde{Q}_1\tilde{Q}_1^TB),$$

where B = Res(A, j), and \tilde{Q}_1 is an orthogonal matrix such that $\text{im}(\tilde{Q}_1) = \text{im}(BV_1)$. In particular, with s = 1,

$$\sigma_j(Q_1Q_1^TA - A) \le \sigma_1(\tilde{Q}_1\tilde{Q}_1^TB - B),$$

as well as, by comparing the singular values individually by varying s,

$$\|\operatorname{Res}(Q_1Q_1^TA - A, j)\|_F^2 \le \|\tilde{Q}_1\tilde{Q}_1^TB - B\|_F^2.$$

As a result, the right-hand sides we need to bound are the same as those bound when we established (5.1), and we may carry out the same steps as those around (5.5). The only change is A is replaced with B = Res(A, j), and accordingly Q_1 changes to have the same image as BV_1 . The effect of this is that the order of right singular vector matrix S changes; the leading j-1 singular values and singular vectors are removed. To capture this change, we may notationally let \tilde{S} be the permutation of the columns of S, moving the leading j-1 columns to the end. Then, in the spectral case with s=1,

$$\sigma_j^2(Q_1Q_1^TA - A) \le \|\text{Res}(\Sigma_2, j)\|_2^2 + \|\text{Res}(\Sigma_2, j)(\tilde{S}^TV)_{21}(\tilde{S}^TV)_{11}\|_2^2.$$

The Frobenius norm version follows similarly.

In (5.2) and (5.3), one could factor out the σ part to make the equations immediately take the form of Definition 1.1 and Definition 1.3. However, as in Theorem 6.9, the unfactored form can have advantages.

We conclude this section with a single proposition encompassing the QR factorization bounds we use.

PROPOSITION 5.2. Let A be an $m \times n$ matrix with the SVD of A being $A = P\Sigma S^T$. Set $[Q, R] = \mathbf{Q}\mathbf{R}(AV)$, where V is an $n \times n$ matrix as in (3.3). Then the singular values of $Q_1Q_1^TA - A$ are identical to those of $(m-l) \times (n-l)$ matrix R_{22} . Moreover, for $j \leq k$,

(5.7)
$$\sigma_j(A) \ge \sigma_j(Q_1 Q_1^T A) \ge \sigma_j(AV_1)\sigma_{\min}(V_1^+)$$

as well as for any given $j \leq \min(m,n) - l$, there is an orthogonal $n \times n$ matrix \tilde{S} independent of V such that

(5.8)
$$\sigma_j^2(R_{22}) \le \|Res(\Sigma_2, j)\|_2^2 + \|Res(\Sigma_2, j)(\tilde{S}^T V)_{21}(\tilde{S}^T V)_{11}^+\|_{22}^2$$

$$(5.9) \|Res(R_{22},j)\|_F^2 \le \|Res(\Sigma_2,j)\|_F^2 + \|Res(\Sigma_2,j)(\tilde{S}^TV)_{21}(\tilde{S}^TV)_{11}^+\|_F^2,$$

with $(\tilde{S}^T V)_{11}$ being $k \times l$ as before.

Proof. Excluding (5.7), the bounds are restatements of facts in Proposition 5.1. The upper bound is a consequence of the Weyl inequality,

$$\sigma_j(Q_1Q_1^T A) \le \sigma_1(Q_1Q_1^T)\sigma_j(A) = \sigma_j(A).$$

The lower bound follows from

$$\sigma_j(Q_1Q_1^TA) = \sigma_j\left(\left(\begin{array}{cc} R_{11} & R_{12} \\ 0 & 0 \end{array}\right)V^{-1}\right) \ge \sigma_j(R_{11}V_1^+) = \sigma_j(AV_1V_1^+),$$

and then from Weyl's inequality, Lemma 3.1. This bound is identical to that of Proposition 3.3.

6. Application of randomness. In this section, we combine our deterministic bounds with the past literature on sketching matrices. There are three applications. We first note that ensembles U and V used in Algorithm \mathbf{GLU} 's guarantees in Proposition 3.3 can be viewed through the oblivious subspace embedding property commonly used in literature. Second, we specialize the random ensemble to the subsampled randomized Hadamard transform (SRHT) introduced in [29] but whose analysis was strengthened in [4]. Our approach fits nicely with their work to give particularly strong operator norm bounds, but in asymptotically less time. Third, we specialize to the Gaussian ensemble to see an application to analyzing the growth factor in Gaussian elimination.

We begin by recalling a property associated with Johnson-Lindenstrauss embeddings. Different authors establish it in different ways, as in [29, 4, 7], but they all have found it necessary in providing sharp Frobenius bounds.

DEFINITION 6.1. We say U from \mathbb{R}^n to \mathbb{R}^s is (ϵ, δ, n) multiplication approximating if for any A, B having n rows,

$$||A^T U^T U B - A^T B||_F^2 \le \epsilon ||A||_F^2 ||B||_F^2,$$

with probability $1 - \delta$.

We also include a definition used consistently in the literature.

DEFINITION 6.2. An (k, ϵ, δ) oblivious subspace embedding (OSE) from \mathbb{R}^n to \mathbb{R}^s is a distribution $U \sim \mathbb{D}$ over $s \times n$ matrices. It must with probability $1 - \delta$ succeed in making

$$1 - \epsilon \le \sigma_{\min}^2(UQ) \le \sigma_{\max}^2(UQ) \le 1 + \epsilon$$

hold for any given orthogonal $n \times k$ matrix Q. We will assume $l \ge k$ and $\epsilon < 1/6$.

For Definition 6.2, there is a consequence we require. The first part is essentially Lemma 4.1 of [4], but we need to state it more generally.

LEMMA 6.3. Let U be an $s \times n$ matrix that is a (k, ϵ, δ) OSE from \mathbb{R}^n to \mathbb{R}^s , and let Q be an $(n \times k)$ orthogonal matrix. Provided $\epsilon < 1/6$, then with probability $1 - \delta$ both of the following hold:

$$\|(UQ)^+ - (UQ)^T\|_2 \le 3\epsilon,$$

 $\|U\|_2^2 = O\left(\frac{n}{k}\right),$

where in the second of these we require the additional assumption $\delta > 2e^{-k/5}$.

Proof. Let A = UQ. Then from Definition 6.2 and power series expansion, the singular values of A lie within $[\sqrt{1-\epsilon}, \sqrt{1+\epsilon}]$, and hence for simplicity we may say they lie within $[1-\epsilon, 1+\epsilon]$ with probability $1-\delta$. Let $l \times k$ diagonal matrix Σ contain these singular values. Therefore,

$$||A^{+} - A^{T}||_{2} = ||\Sigma^{T} - \Sigma^{+}||_{2} = \max_{i \le k} |\lambda_{i} - \lambda_{i}^{-1}|$$

$$\leq |1 - \epsilon - \frac{1}{1 - \epsilon}| \leq 3\epsilon,$$

where we have chosen to write the small extreme of σ_i ; the large extreme is identical.

For the second fact, let $V \leq \mathbb{R}^n$ be a uniformly distributed k-dimensional subspace with $\dim(V) = k$ independent of U; i.e., V is spanned by the first k columns of a Haar distributed matrix on \mathbb{R}^n independent of U. A consequence of Definition 6.2 is that $||Uv||_2 \leq 2$ with probability $1 - \delta$ holding uniformly for unit vectors v contained in V. Otherwise some fixed subspace V_0 would also fail to have this property with probability δ , violating Definition 6.2.

Now let x be the maximal right singular vector of U. The subsequent Lemma 6.4 gives $\sup_{v \in V, \|v\|_2 = 1} |\langle x, v \rangle| = \Omega(\sqrt{\frac{k}{n}})$ with probability $1 - \delta$. Next choose $v \in \operatorname{argmax}_{v \in V, \|v\|_2 = 1} |\langle x, v \rangle|$ to be a unit-vector with smallest angle with respect to x, and observe $\|Uv\|_2 = \Omega(\sqrt{\frac{k}{n}})\|Ux\|_2$. We conclude $\|U\|_2^2 = O\left(\frac{n}{k}\right)$ with probability $1 - \delta$. Otherwise this would contradict $\|Uv\|_2 \le 2$ holding with probability $1 - \delta$. \square

LEMMA 6.4. Let V be a k-dimensional uniformly distributed subspace of \mathbb{R}^n and $x \in \mathbb{R}^n$ be a unit vector drawn from a distribution independent of V. Then $\sup_{v \in V, \|v\|_2 = 1} |\langle x, v \rangle| = \Omega(\sqrt{\frac{k}{n}})$ with probability $1 - 2e^{-k/5}$.

Proof. We may assume $V = \operatorname{span}(e_1, \dots, e_k)$ and represent x as $\frac{(X_1, \dots, X_n)^T}{\sqrt{X_1^2 + \dots + X_n^2}}$ where X_i are i.i.d. variance $\frac{1}{n}$ Gaussians. Indeed, V can be taken to be the first k columns of Haar distributed orthogonal matrix \tilde{V} , and the WLOG assumption is equivalent to changing to the coordinates of \tilde{V} . As a result, we are interested in

$$\sup_{v \in V, \|v\|_2 = 1} |\langle x, v \rangle| = \frac{(X_1, \dots, X_n)}{\sqrt{X_1^2 + \dots + X_n^2}} \cdot \frac{(X_1, \dots, X_k, 0, \dots)^T}{\sqrt{X_1^2 + \dots + X_k^2}} = \frac{\sqrt{X_1^2 + \dots + X_k^2}}{\sqrt{X_1^2 + \dots + X_n^2}}.$$

Standard large-deviation bounds for chi-squared distribution, which is a subexponential random variable, can be used to lower bound this. We take bounds from [21, (4.3), (4.4)]. The right tail bound is

(6.1)
$$\mathbb{P}[X_1^2 + \dots + X_n^2 > 1 + 2\frac{\sqrt{\delta}}{\sqrt{n}} + 2\frac{\delta}{n}] \le e^{-\delta},$$

and the left tail bound is

$$(6.2) \mathbb{P}[X_1^2 + \dots + X_k^2 < \frac{k}{n} - 2\frac{\sqrt{k\delta}}{n}] \le e^{-\delta}.$$

From these and setting $\delta = k/5$ in (6.2) and $\delta = n/5$ in (6.1), we conclude that

$$\sup_{v \in V, ||v||_2 = 1} |\langle x, v \rangle| \ge \left(\frac{\frac{k}{n} - 2\frac{k}{\sqrt{5n}}}{1 + 2\frac{\sqrt{k}}{\sqrt{5n}} + 2\frac{k}{5n}} \right)^{.5} \ge \frac{1}{5} \sqrt{\frac{k}{n}}$$

holds with probability $1 - 2e^{-k/5}$.

The following lemma largely follows the steps of [4], but we have tried to abstract out the key probabilistic properties responsible in order to be more general. Another difference is that we also treat the spectral norm. It is a natural consequence of the prior lemmas and will bridge the gap between deterministic Proposition 3.3 and randomized Theorem 6.6. We do not attempt to tightly bound the constant coefficients.

LEMMA 6.5. Assume $l \times m$ matrix U is drawn from a distribution that is a $(k, \sqrt{\epsilon}, \delta)$ OSE from \mathbb{R}^m to \mathbb{R}^l . Let B be a fixed $(m-k) \times n$ matrix and $Q = [Q_1, Q_2]$

be a fixed orthogonal $m \times m$ matrix blocked so that Q_1 is $m \times k$. Then provided $\delta > 2e^{-k/5}$ and $\epsilon < 1/6$, with probability $1 - \delta$

$$||(UQ_1)^+(UQ_2)B||_2^2 = O\left(\frac{m}{k}\right)||B||_2^2.$$

Further assume that U is $(k, \sqrt{\epsilon}, \delta \frac{k}{n})$ OSE and $(\frac{\epsilon}{k}, \delta, m)$ multiplication approximating; then with probability at least $1 - 2\delta$,

$$||(UQ_1)^+(UQ_2)B||_F^2 = O(\epsilon) ||B||_F^2.$$

Proof. For the Frobenius bound, apply Lemma 6.3 in (6.3) and Definition 6.1 in (6.5) by noting that $Q_2^T Q_1 = 0$,

$$\|(UQ_1)^+(UQ_2)B\|_F^2 \le 2\|(UQ_1)^T(UQ_2)B\|_F^2 + 2\|((UQ_1)^+ - (UQ_1)^T)(UQ_2)B\|_F^2$$

$$(6.3) \leq 2\|Q_1^T U^T U Q_2 B\|_F^2 + 18\epsilon \|U Q_2 B\|_F^2$$

$$(6.4) \leq 2\|Q_1^T U^T U Q_2 B\|_F^2 + 36\epsilon \|B\|_F^2$$

(6.5)
$$\leq 2\frac{\epsilon}{k} \|Q_2 B\|_F^2 \|Q_1^T\|_F^2 + 36\epsilon \|B\|_F^2$$
$$\leq 2\epsilon \|B\|_F^2 + 36\epsilon \|B\|_F^2 = 38\epsilon \|B\|_F^2.$$

In the above, the step to (6.4) for $\|UQ_2B\|_F^2$ was obtained as follows. Let $C = Q_2B$, $C \in \mathbb{R}^{m \times n}$, and partition C into blocks of k columns. We refer to each such block as C_i , with $1 \le i \le n/k$. Without loss of generality we assume for simplicity that n divides k. Using Definition 6.2 and the thin QR decomposition $C_i = \tilde{Q}_i \tilde{R}_i$, where $\tilde{Q}_i \in \mathbb{R}^{m \times k}$, we obtain

$$||UC_i||_F^2 = ||U\tilde{Q}_i\tilde{R}_i||_F^2 \le ||U\tilde{Q}_i||_2^2 ||\tilde{R}_i||_F^2 \le (1 + \sqrt{\epsilon})||C_i||_F^2 \le 2||C_i||_F^2$$

with probability $1 - \delta \frac{k}{n}$. We further obtain

$$||UQ_2B||_F^2 = ||UC||_F^2 = \sum_{i=1}^{n/k} ||UC_i||_F \le 2||B||_F^2$$

with probability $1 - \delta$.

For the spectral bound, we may argue

$$||(UQ_{1})^{+}(UQ_{2})B||_{2}^{2} \leq ||(UQ_{1})^{+}(UQ_{2})||_{2}^{2} \cdot ||B||_{2}^{2}$$

$$\leq 2||(UQ_{1})^{T}(UQ_{2})||_{2}^{2} \cdot ||B||_{2}^{2}$$

$$+2||((UQ_{1})^{+} - (UQ_{1})^{T})(UQ_{2})||_{2}^{2} \cdot ||B||_{2}^{2}$$

$$\leq 2(1 + \sqrt{\epsilon})||UQ_{2}||_{2}^{2} \cdot ||B||_{2}^{2} + 18\epsilon||(UQ_{2})||_{2}^{2} \cdot ||B||_{2}^{2}$$

$$\leq 2(1 + \sqrt{\epsilon} + 9\epsilon)||U||_{2}^{2} \cdot ||Q_{2}||_{2}^{2} \cdot ||B||_{2}^{2}$$

$$= O\left(\frac{m}{k}\right)||B||_{2}^{2}.$$

In the former steps, we note in particular that (6.6) follows from Definition 6.2, and (6.7) follows from Lemma 6.3.

In the following, one of our main results, we continue with the notation of Propositions 3.3 and 5.2. We provide a bound on Definitions 1.1 and 1.3. While these bounds

do appear quite weak (often weaker than a naive Frobenius norm adaptation), we note that they match the guarantees of past literature for algorithms running in o(nmk) time, e.g., [14, 17, 30]. On the other hand, in Theorem 6.9 we notably achieve very sharp bounds on Definitions 1.1 and 1.3 by exploiting a special property of the SRHT ensemble.

THEOREM 6.6. Assume U_1 is drawn from a distribution that is an $(l, \sqrt{\epsilon}, \delta)$ OSE from \mathbb{R}^m into $\mathbb{R}^{l'}$. Similarly assume V_1^T is drawn from a distribution that is a $(k, \sqrt{\epsilon}, \delta)$ OSE from \mathbb{R}^n into \mathbb{R}^l . Then provided $\delta > 2e^{-l/5}$ and $\epsilon < 1/6$, with probability $1 - 2\delta$ for $j \leq k$,

$$\sigma_j(A_k) = \Omega\left(\sqrt{\frac{k}{n}}\right)\sigma_j(A).$$

Fixing a given $1 \le j \le \min(m, n) - l$, with probability $1 - 4\delta$ we also have

$$\sigma_j(A-A_k) = O\left(\sqrt{\frac{mn}{kl}}\right)\sigma_{k+j}(A).$$

If we additionally assume U_1 is drawn from a $(l, \sqrt{\epsilon}, \delta \frac{l}{n})$ OSE and a $(\frac{\epsilon}{l}, \delta, m)$ multiplication approximating and similarly V_1^T is drawn from a $(k, \sqrt{\epsilon}, \delta \frac{k}{m})$ and a $(\frac{\epsilon}{k}, \delta, n)$ multiplication approximating, then for a given $1 \le j \le \min(m, n) - l$,

$$||Res(A - A_k, j)||_F^2 = (1 + O(\epsilon))||Res(\Sigma_2, j)||_F^2$$

holds with probability $1-4\delta$.

Proof. We start with the Frobenius norm bound. The starting point is Proposition 5.2, which includes

$$\|\operatorname{Res}(R_{22},j)\|_F^2 \leq \|\operatorname{Res}(\Sigma_2,j)\|_F^2 + \|\operatorname{Res}(\Sigma_2,j)(\tilde{S}^TV)_{21}(\tilde{S}^TV)_{11}^+\|_F^2,$$

where V is an $n \times n$ matrix as in (3.3), $V_1 = V[:,:l]$, and \tilde{S} is an orthogonal $n \times n$ matrix independent of V. Let $\tilde{S} = [\tilde{S}_1, \tilde{S}_2]$ such that \tilde{S}_1 contains the first k columns of \tilde{S} . Hence $(\tilde{S}^TV)_{11} = \tilde{S}_1^TV_1$ and $(\tilde{S}^TV)_{21} = \tilde{S}_2^TV_1$. Then as V_1^T satisfies the Johnson–Lindenstrauss properties, apply Lemma 6.5 with $B = \text{Res}(\Sigma_2, j)^T$, $Q_1 = \tilde{S}_1$, $Q_2 = \tilde{S}_2$, and $U = V_1^T$ to conclude that for a given $1 \le j \le \min(m, n) - l$, $\|\text{Res}(R_{22}, j)\|_F^2 = (1 + O(\epsilon))\|\text{Res}(\Sigma_2, j)\|_F^2$ with probability $1 - 2\delta$. To complete the Frobenius bound, recall from Proposition 3.3 that

$$\|\operatorname{Res}(A-A_k,j)\|_F^2 \le \|\operatorname{Res}(R_{22},j)\|_F^2 + \|(UQ)_{11}^+(UQ)_{12}\operatorname{Res}(R_{22},j)\|_F^2,$$

and again apply Lemma 6.5, this time with $B = \text{Res}(R_{22}, j)$, to get $\|\text{Res}(A - A_k, j)\|_F^2 = (1 + O(\epsilon)) \|\text{Res}(\Sigma_2, j)\|_F^2$ with probability $1 - 4\delta$.

The spectral bound proceeds similarly, but using the spectral bounds of Proposition 5.2, Proposition 3.3, and Lemma 6.5 instead. Thus

$$\|\operatorname{Res}(R_{22},j)\|_{2}^{2} \leq \|\operatorname{Res}(\Sigma_{2},j)\|_{2}^{2} + \|\operatorname{Res}(\Sigma_{2},j)(\tilde{S}^{T}V)_{21}(\tilde{S}^{T}V)_{11}^{+}\|_{2}^{2} = O\left(\frac{n}{k}\right)\sigma_{j+k}^{2}(A).$$

And then using (3.13) of Proposition 3.3.

$$\sigma_j^2(A - A_k) \le \|\operatorname{Res}(R_{22}, j)\|_2^2 + \|(UQ)_{11}^+(UQ)_{12}\operatorname{Res}(R_{22}, j)\|_2^2 = O\left(\frac{mn}{kl}\right)\sigma_{j+k}^2(A),$$

which proves the spectral claim.

For the multiplicative lower bound on the singular values of A_k , from (3.14) and (5.7) in Proposition 3.3 and Proposition 5.2, respectively, it follows that for $j \leq k$,

$$\sigma_j(A_k) \ge \sigma_j(AV_1) \cdot \sigma_{\min}(V_1^+) = \Omega\left(\sqrt{\frac{k}{n}}\right) \sigma_j(AV_1) = \Omega\left(\sqrt{\frac{k}{n}}\right) \sigma_j(A).$$

This last step requires additional explanation. With $[P, \Sigma, S^T] = \mathbf{SVD}(A)$, let Σ_1 be formed by the first k singular values of A, and let S_1 be formed by the first k columns of S. We obtain

$$\sigma_j(AV_1) = \sigma_j(\Sigma S^T V_1) \ge \sigma_j(\Sigma_1 S_1^T V_1) \ge \sqrt{1 - \sqrt{\frac{1}{6}}} \sigma_j(\Sigma_1)$$

by using Definition 6.2.

Next, we specialize to the SRHT ensemble in order to see a case where the bounds of Definitions 1.1 and 1.3 are stronger than in Theorem 6.6.

DEFINITION 6.7. The SRHT ensemble embedding \mathbb{R}^n into \mathbb{R}^s is defined by generating $\sqrt{\frac{n}{s}}PHD$, where P is $s \times n$ selecting s rows, H is the normalized Hadamard transform, and D is an $n \times n$ diagonal matrix of uniformly random signs.

The key special additional property of the SRHT ensemble is from Lemma 4.8 of [4].

LEMMA 6.8. Let V^T be drawn from an SRHT of dimension $l \times n$. Then for $m \times n$ matrix A with rank ρ , with probability $1 - 2\delta$,

$$||AV||_2^2 \le 5||A||_2^2 + \frac{\ln(\rho/\delta)}{l}(||A||_F + \sqrt{8\ln(n/\delta)}||A||_2)^2.$$

Let $U \in \mathbb{R}^{l \times m}$ be drawn from SRHT ensembles, where l is chosen such that U satisfies $(k, \sqrt{\epsilon}, \delta \frac{k}{n})$ OSE and $(\frac{\epsilon}{k}, \delta, m)$ multiplication approximating properties, as in Lemma 6.5. From Lemma 4.1 of [4] we know that the SRHT with $l \geq 6\epsilon^{-1}(\sqrt{k} + 4\sqrt{\ln(\frac{max(m,n)}{\delta k})})^2\ln(n/\delta)$ rows is a $(k, \sqrt{\epsilon}, \delta \frac{k}{n})$ OSE (by substituting δ with $\delta \frac{k}{n}$). For the multiplication approximating property, by setting in Lemma 4.11 of [4] $R = 2\sqrt{\ln(3/\delta)}$ and $\delta = \delta/3$, we obtain that the number of rows should be $l \geq 4\epsilon^{-1}k(1+2\sqrt{\ln(3/\delta)})^2(1+\sqrt{8\ln(3m/\delta)})^2$. Hence both properties are satisfied by choosing $l = 6\epsilon^{-1}(\sqrt{6k\ln\frac{3}{\delta}} + 4\sqrt{\ln\frac{max(m,n)}{\delta k}})^2(1+\sqrt{8\ln\frac{max(3m,n)}{\delta}})^2$. We may substitute these parameters into Theorem 6.6, but numerous other ensembles could also be used. We have singled out the SRHT because it enjoys a remarkably good bound for the spectral norm approximation quality due to the prior lemma, but past work has not exploited this property fully. In particular, when the spectral norm and Frobenius norm are comparable (i.e., quickly decaying singular values), the quality is constant in the dimension rather than polynomial. Loosely speaking, as long as $\frac{\|A-A_k\|_F}{\|A-A_k\|_F} = O(\sqrt{k})$, then $\|A-A_k\|_2$ is around a constant factor from that of the k-truncated SVD. The theorem further strengthens this by proving the generalization to the lower singular values of $A-A_k$.

Theorem 6.9. Let U_1, V_1^T be drawn from SRHT ensembles with dimensions $l' \times m$, $n \times l$. We set $l \geq 6\epsilon^{-1}(\sqrt{6k\ln\frac{3}{\delta}} + 4\sqrt{\ln\frac{\max(m,n)}{\delta k}})^2(1+\sqrt{8\ln\frac{\max(3m,n)}{\delta}})^2$ and $l' \geq 6\epsilon^{-1}(\sqrt{6l\ln\frac{3}{\delta}} + 4\sqrt{\ln\frac{\max(m,n)}{\delta l}})^2(1+\sqrt{8\ln\frac{\max(3m,n)}{\delta}})^2$, where $\epsilon < 1/6$. Letting ρ be the rank of A, for simplicity assume $l' \geq \ln(m/\delta)\ln(\rho/\delta)$ and $l \geq \ln(n/\delta)\ln(\rho/\delta)$.

Then for any fixed $1 \le j \le \min(m, n) - l$, with probability $1 - 8\delta$ the approximation of A using GLU, A_k satisfies

$$\sigma_{j}^{2}(A - A_{k}) = O(1)\sigma_{k+j}^{2}(A) + O\left(\frac{\ln(\rho/\delta)}{l}\right) \|Res(A, k+j)\|_{F}^{2}$$

$$= O\left(1 + \frac{\epsilon \ln(\min(m, n)/\delta)}{k \ln(k/\delta)} \frac{\|Res(A, k+j)\|_{F}^{2}}{\sigma_{k+j}^{2}(A)}\right) \sigma_{k+j}^{2}(A).$$

Proof. It suffices to prove the first claim. Begin by using Proposition 5.2 and Lemma 6.3:

$$\sigma_j^2(R_{22}) \le \|\operatorname{Res}(\Sigma_2, j)\|_2^2 + \|\operatorname{Res}(\Sigma_2, j)(\tilde{S}^T V)_{21}(\tilde{S}^T V)_{11}^+\|_2^2$$

$$\le \|\operatorname{Res}(\Sigma_2, j)\|_2^2 + 2\|\operatorname{Res}(\Sigma_2, j)(\tilde{S}^T V)_{21}\|_2^2,$$

with probability $1 - \delta$. Next apply Lemma 6.8 to the second term to get

(6.8)

$$\sigma_{j}^{2}(R_{22}) \leq O\left(1 + \frac{\ln(\rho/\delta)\ln(n/\delta)}{l}\right) \|\operatorname{Res}(\Sigma_{2}, j)\|_{2}^{2} + O\left(\frac{\ln(\rho/\delta)}{l}\right) \|\operatorname{Res}(\Sigma_{2}, j)\|_{F}^{2}$$

$$(6.9) \qquad = O(1) \|\operatorname{Res}(\Sigma_{2}, j)\|_{2}^{2} + O\left(\frac{\ln(\rho/\delta)}{l}\right) \|\operatorname{Res}(\Sigma_{2}, j)\|_{F}^{2},$$

where ρ is the rank of A, with probability $1 - 2\delta$. Continuing from the result of Proposition 3.3, with probability $1 - \delta$ we obtain

$$\sigma_j^2(A - A_k) \le \|\operatorname{Res}(R_{22}, j)\|_2^2 + \|(U_1 Q_1)^+ (U_1 Q_2) \operatorname{Res}(R_{22}, j)\|_2^2 \le \|\operatorname{Res}(R_{22}, j)\|_2^2 + 6\|(U_1 Q_2) \operatorname{Res}(R_{22}, j)\|_2^2.$$

From Theorem 6.6 we know $\|\operatorname{Res}(R_{22},j)\|_F^2 \leq (1+O(\epsilon))\|\operatorname{Res}(\Sigma_2,j)\|_F^2$ with probability $1-2\delta$, because the SRHT with the parameter settings specified for l and l' satisfies the multiplication approximation and OSE properties. Thus repeating the same steps using Lemma 6.8 to complete the proof for the first bound, with probability $1-2\delta$, we obtain

$$\begin{split} \sigma_j^2(A-A_k) &\leq \|\mathrm{Res}(R_{22},j)\|_2^2 + 6\|(U_1Q_2)\mathrm{Res}(R_{22},j)\|_2^2 \\ &\leq O\left(1 + \frac{\ln(\rho/\delta)\ln(m/\delta)}{l'}\right) \|\mathrm{Res}(R_{22},j)\|_2^2 \\ &+ O\left(\frac{\ln(\rho/\delta)}{l'}\right) \|\mathrm{Res}(R_{22},j)\|_F^2 \\ &= O(1) \|\mathrm{Res}(R_{22},j)\|_2^2 + O\left(\frac{\ln(\rho/\delta)}{l'}\right) \|\mathrm{Res}(R_{22},j)\|_F^2 \,. \end{split}$$

By using the bounds on $\sigma_j(R_{22})$ from (6.8) and the fact that $\|\text{Res}(R_{22},j)\|_2 = \sigma_j(R_{22})$, we further obtain

$$\sigma_j^2(A - A_k) \le O(1)\sigma_{k+j}^2(A) + O\left(\frac{\ln(\rho/\delta)}{l}\right) \|\operatorname{Res}(\Sigma_2, j)\|_F^2.$$

A few remarks are in order.

Remark 6.10. First, the SRHT ensemble is only defined for powers of 2. This is not a theoretical issue because matrices can be padded. However, as discussed in [4] there

are orthogonal ensembles related to the SRHT, namely the discrete cosine transform and Hartley transform, for which the key probabilistic requirement in Lemma 6.8 carries over, so this corollary also carries over.

Remark 6.11. Second, we consider much of the work in this section as adapting [4] to algorithm **GLU** which sketches A's columns and rows and proves a spectral norm bound comparable to the above. Their work does not specify how to proceed after finding $A \approx Q_1 Q_1^T A$ and therefore follows **RandomizedQB**. Therefore, if one follows their approach, creating a compressed representation of A would still require O(nmk) time because $Q_1^T A$ must be computed. We state the relevant part of their theorem here to provide context.

THEOREM 6.12 ([4, Thm. 2.1]). Let $A \in \mathbb{R}^{m \times n}$ have rank ρ and n a power of 2. Fix an integer k satisfying $2 \le k < \rho$. Let $0 < \epsilon < 1/3$ and $0 < \delta < 1$. Let $Y = AV^T$, where $V \in \mathbb{R}^{r \times n}$ is drawn from the SRHT ensemble with $r = 6\epsilon^{-1}(\sqrt{k} + \sqrt{8\ln(n/\delta)})^2)\ln(k/\delta)$). Then with probability $1 - 5\delta$

$$||A - YY^{+}A||_{2} \le \left(4 + \sqrt{\frac{3\ln(n/\delta)\ln(\rho/\delta)}{r}}\right) ||\Sigma_{2}||_{2} + \sqrt{\frac{3\ln(\rho/\delta)}{r}} ||\Sigma_{2}||_{F}$$

From this we see that our Theorem 6.9 has qualitatively the same accuracy guarantee on the residual error. For many types of matrices A, in particular for those with fast spectral decay, Theorem 6.9 will be within a constant factor of the rank k truncated SVD's spectral approximation error.

Remark 6.13. Let us consider the computational cost of computing the \mathbf{GLU} approximation of A through Theorem 6.9, storing the result in the form of (1.6), following Algorithm 3.1.

Simply by following the algorithmic description, we see the largest cost terms are $O(nm \log(l') + mll')$. We present a short table tabulating this.

$$\begin{array}{ll} \hat{A} = U_1(AV_1) & O(nm\log(l)) \\ T_1 = U_1^+(I - \hat{A}\hat{A}^+) & O(ml'\log(m) + ll'^2) \text{ because} \\ & U_1^+ = \sqrt{\frac{l'}{m}}(PHD)^T = \sqrt{\frac{l'}{m}}DHP^T \\ T_2 = AV_1 & \text{Stored from first step} \\ T_2 = T_2\hat{A}^+ & O(mll') \\ T = T_1 + T_2 & O(ml') \\ S = U_1A & O(mn\log(l')) \end{array}$$

Specializing as in the theorem, we additionally required $l \geq 6\epsilon^{-1}(\sqrt{6k\ln\frac{3}{\delta}} + 4\sqrt{\ln\frac{\max(m,n)}{\delta k}})^2(1+\sqrt{8\ln\frac{\max(3m,n)}{\delta}})^2$ and $l' \geq 6\epsilon^{-1}(\sqrt{6l\ln\frac{3}{\delta}} + 4\sqrt{\ln\frac{\max(m,n)}{\delta l}})^2(1+\sqrt{8\ln\frac{\max(3m,n)}{\delta}})^2$. Using these bounds on l and l', we say the runtime is $\tilde{O}(nm+k^2m\epsilon^{-3})$. Various poly-log factors are hidden here, involving n, m, k, δ . In more detail, plugging l and l' into the prior complexity bound and assuming m < n so that $l' = O(\epsilon^{-1}(l\ln(3/\delta) + \ln(n/(\delta k)))\ln(n/\delta))$ and $l = O(\epsilon^{-1}(k\ln(3/\delta) + \ln(n/(\delta k)))\ln(n/\delta))$, we get the Big-Oh of

$$nm\log\left(\epsilon^{-1}(l\ln(3/\delta) + \ln(n/(\delta k)))\ln(n/\delta)\right) + m\epsilon^{-3}(k\ln(3/\delta) + \ln(n/(\delta k)))^2\ln^3(n/\delta)\ln(3/\delta).$$

Note that in the runtime bound, because there is asymmetry between m and n, it turns out to be faster if m < n and thus A is short-wide. If this is not the case for A, then one could simply run the algorithm on A^T .

Remark 6.14. As stated, Theorem 6.9 provides bounds for the **GLU** with sketching from the left and right. We noted in the prior remark how this retains the performance of [4] while increasing the speed. We could stop the analysis at (6.8), and also borrow the bounds already found in Theorem 6.6 and Proposition 5.2. Then we obtain new bounds for the randomized QB factorization,

COROLLARY 6.15. Let $n \times l$ matrix V_1^T be drawn from an SRHT ensemble, $l \ge 6\epsilon^{-1}(\sqrt{6k\ln\frac{3}{\delta}} + 4\sqrt{\ln\frac{\max(m,n)}{\delta k}})^2(1+\sqrt{8\ln\frac{\max(3m,n)}{\delta}})^2$, and for simplicity assume $l \ge \ln(n/\delta)\ln(\rho/\delta)$. Then we have

$$||Res(R_{22}, j)||_F^2 \le (1 + O(\epsilon))||Res(A, k + j)||_F^2$$

with probability $1-2\delta$, as well as

$$\sigma_j^2(R_{22}) \le O(1)\sigma_{k+j}^2(A) + O\left(\frac{\ln(\rho/\delta)}{l}\right) ||Res(A, k+j)||_F^2$$

for $1 \le j \le \min(m, n) - l$ with probability $1 - 3\delta$ for a particular j. We also have upper and lower bounds on the largest singular values, as for $1 \le j \le k$,

$$\sigma_j(A) \geq \sigma_j(Q_1Q_1^TA) = \Omega\left(\sqrt{\frac{k}{n}}\right)\sigma_j(A)$$

holds with probability $1-2\delta$. Actually, using the deterministic bound of [15] found in (4.7),

$$\sigma_j(A) \ge \sigma_j(Q_1 Q_1^T A) \ge \frac{\sigma_j(A)}{1 + O(\sqrt{\frac{n}{k}}) \frac{\sigma_{k+1}(A)}{\sigma_j(A)}}$$

holds with probability $1 - \delta$.

We move on to the third application, controlling the growth factor during Gaussian elimination by right and left multiplication by square random matrices. The theoretical result we establish is that the growth factor is well behaved if we multiply by square Gaussian random matrices. Note that the bounds in Propositions 3.3 and 5.1 will in this case be the same for Gaussian random matrices as for Haar random matrices, because they differ by lower and upper triangular factors and U_1, V_1 are now square. We make use of bounds proven for the Haar ensemble. The work [11], which viewed the problem in terms of the Haar ensemble, required a randomized QB factorization as a subroutine to compute the generalized Schur decomposition of the matrix by a divide-and-conquer approach. This required a bound on the smallest singular value of the $k \times k$ minors. Eventually a tight bound on these was given in [12] by means of the exact probability distribution, which we will use.

As pointed out in [3], Theorem 3.2 and Lemma 3.5 of [12] give an exact density of the smallest singular value of a Haar minor. Analyzing this formula gives the following bound, which is sharp up to a constant in the primary range of interest, $\sigma_{\min} = O(\frac{1}{\sqrt{k(n-k)}})$.

LEMMA 6.16. Let
$$\delta > 0$$
, $k, (n - k) > 30$; then $\mathbb{P}[\sigma_{\min} \le \frac{\delta}{\sqrt{k(n - k)}}] \le 2.02\delta$.

We will define the ℓ_2 growth factors of \bar{A} as $\rho_U(\bar{A}) := \max_p \|\mathbb{S}_p(\bar{A})\|_2 / \|\bar{A}\|_2$ and $\rho_L(\bar{A}) := \max_p \|\bar{A}_{21}\bar{A}_{11}^{-1}\|_2$, where \mathbb{S}_p is the Schur complement of the top $p \times p$ block. From Proposition 3.3, (3.2), and Proposition 5.1 it is not difficult to see that both are bounded as

$$\begin{split} \rho_U(\bar{A}), \rho_L(\bar{A}) &\leq \max_j \left[\|X[:j,:j]^{-1}\|_2 \|R[j+1:,j+1:]\|_2 / \|\bar{A}\|_2 \right] \\ &\leq \max_j \left[\|(UQ)[:j,:j]^{-1}\|_2 \|(S^TV)[j+1:,j+1:]^{-1}\|_2 \right] \\ &= \max_j \left[\sigma_{\min}^{-1}((UQ)[:j,:j]) \sigma_{\min}^{-1}((S^TV)[:j,:j]) \right]. \end{split}$$

Note that ρ_U and ρ_L control what is typically called the growth factor of \bar{A} . The growth factor is the largest magnitude entry appearing in the matrices L, U returned by Gaussian elimination. This is because of norm equivalence, with the operator and max-element norm differing by at most a factor of \sqrt{n} . Therefore, our ℓ_2 growth factors are equivalent for the purpose of proving stability.

COROLLARY 6.17. Suppose we want to solve Ax = b by Gaussian elimination, and we pre-condition, post-condition A by Haar distributed matrices U, V. That is, we solve UAVx' = Ub and output V^Tx' . Then the U and L ℓ_2 growth factors introduced above satisfy

$$\mathbb{E}[\log(\max(\rho_U(\bar{A}), \rho_L(\bar{A})))] = O(\log(n)).$$

Proof. Because U and V are Haar, the matrices UQ and S^TV in Propositions 3.3 and 5.1 are Haar distributed. Apply Lemma 6.16 to the minors (call them generically M) of UQ and S^TV with size in the range [30, n-30],

$$\mathbb{P}[\sigma_{\min}^{-1}(M) > n^{2+a}] < 2.02n^{-1-a}.$$

To control all minors in this range, simply perform a union bound over all < 2n minors being considered. Let B_1 be the inverse of the smallest singular value of the minors in range [30, n-30] of UQ and S^TV . Then $\mathbb{P}[B_1 \ge n^{2+a}] \le 4.04n^{-a}$. Setting a = x-2, this is $\mathbb{P}[\log_n(B_1) \ge x] \le 4.04n^{2-x}$.

To deal with the minors in range [0,30], we cite a result in random matrix theory which says that these minors scaled by \sqrt{n} approach a matrix of i.i.d. N(0,1) random variables. The convergence is with respect to total variation distance; see [19]. Let B_2 be the inverse of the smallest singular value of these 60 minors. For the claimed result, what matters is $\mathbb{E}[\log_n(B_2)] = C_1'$ for some constant C_1' . This is apparent from work similar to [12] but for Gaussian matrices; see, for example, the bound on the condition number in [6].

Combining the bounds for B_1 and B_2 , and denoting the minors (UQ)[:j,:j] and $(S^TV)[:j,:j]$ as $(UQ)_j$ and $(S^TV)_j$ respectively, we obtain

$$\mathbb{E}[\log_n(\max(\rho_U(\bar{A}), \rho_L(\bar{A})))] \leq \mathbb{E}[\log_n(\max_j) \\ [\sigma_{\min}^{-1}((UQ)[:j,:j])\sigma_{\min}^{-1}((S^TV)[:j,:j]))] \\ \leq \mathbb{E}[\log_n(B_1)] + \mathbb{E}[\log_n(B_2)] \\ \leq C_1' + \int_0^2 1 dx + 4.04 \int_2^\infty n^{2-x} dx \\ \leq C,$$

where we have used the fact that $\mathbb{E}[\log_n(B_2)] \leq \int_0^\infty \mathbb{P}[\log_n(B_2) > x] dx$. Since we have shown $\mathbb{E}[\log_n(\max(\rho_U(\bar{A}), \rho_L(\bar{A})))] \leq C$, the statement of the corollary follows.

Of course, it is impractical to use a Gaussian or Haar matrix to condition a matrix in this context. We might as well then solve the system by means of QR factorization.

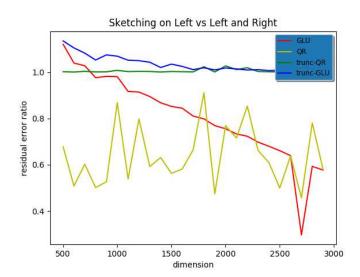


Fig. 1. The QR follows Algorithm 4.1 with the sketching ensemble being SRHT with the value of l' displayed on the x axis. The **GLU** follows Algorithm 3.1, again using SRHT. The matrix A is held fixed as being the diagonal matrix of dimension 3000, with entries $a_{ii} = (1-i/n)^{20 \cdot \ln(n)}$ to give some spectral decay. To reduce the impact of variance, we repeat the procedures 10 times for every value of l' and record the median on the plot. The varying of l' affects the **GLU** and its truncated version but does not affect the QR algorithms.

However, this sheds light on the strategy of using conditioners to avoid pivoting during Gaussian elimination. This has been popularized in work such as [1]. The theoretical support of such work has been lacking. Corollary 6.17 is the first theoretical result we are aware of that shows that a random conditioner can be used to provably avoid the need to pivot.

It also could be considered a generalization of the well-known fact that Gaussian random matrices have low pivot growth during Gaussian elimination. Indeed, we have shown that this is the case for any distribution of singular values—not just that of the Gaussian random matrix. The most interesting question still remains if faster conditioners can be used to make the approach both theoretically and practically sound for all matrices A. More concretely, we pose the following question.

Remark 6.18. Is there a random matrix ensemble S such that SA can be computed quickly, but also $\sigma_{\min}((SA)[:k,:k]) = O(\frac{1}{\operatorname{poly}(n)})$ when A is an orthogonal matrix?

7. Experiments. In this section, we present two numerical experiments. In general, probabilistic upper bounds are considered to be pessimistic, acting more to help guide practical settings. Therefore, we do not directly compare experiments with the bounds established in this paper. However, qualitative insights can still be gained from our bounds.

First, the QR approach in Algorithm 4.1 is the most efficient in its use of randomness, as it does not sketch from the left and right, but rather just from the right. The works [4, 17, 15] all adopt or include discussions of this approach. For small values of k it is likely to be the fastest and most accurate approach. Indeed, that is the case for the relatively small values of n and k we use in the experiments we present. However, as k becomes larger, sketching from both left and right becomes computationally faster, as is done in Algorithm 3.1. In order to illustrate the accuracy loss caused by sketching from both sides, we present Figure 1. Note that the variants that

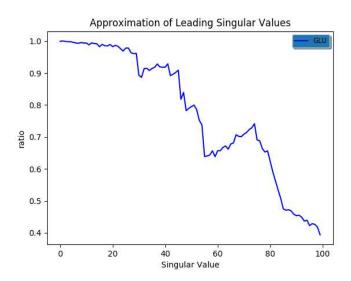


Fig. 2. The ensemble used for **GLU** is again the SRHT. We have A similar to the last experiment, this time diagonal with n=3000 and $a_{ii}=(1-i/n)^{20\cdot \ln(n)}$. On the y-axis is plotted $\sigma_j(A_k)/\sigma_j(A)$, with j being given on the x-axis, where A_k is the untruncated approximation.

only sketch from the left are more accurate than their counterparts. Also, naturally, truncating the approximation to rank k (following the procedure of [32]) reduces the accuracy. The target approximation rank is k = 20, and the y-axis is spectral residual error normalized by σ_{21} . The value of l is held fixed at l = 100, and the x-axis varies l' from 500 to 2500, in increments of 100.

In reading the plot, we note that the approximation error of **GLU** at l' = 500, although not competitive with QR, is within a factor of 2. In comparison, the null approximation with a trivial matrix of zeros would have yielded about $1/.3247 \approx 3$ on the y-axis. Naturally, as l' increases, **GLU** eventually becomes indistinguishable from QR in terms of approximation quality.

As an additional experiment to study the approximation of the leading singular values that we provided bounds for in Theorem 6.6, we present Figure 2. This time, the untruncated **GLU** algorithm is run once with l=100, l'=500. Thus we show the ratio of the singular values of the output A_k of Algorithm 3.1 to that of A for a single sample run. The first interesting thing to note is that the under-approximation loss is very close to monotonically increasing. The second interesting thing we notice is the relatively large drop around the 45th singular value. As the target approximation rank is 20 as in Figure 1, it makes sense that the approximation gets gradually worse. However, the leading 20 singular values are well approximated, with the ratio being close to 1.

8. Conclusion. We have provided a thorough analysis of a new low-rank approximation procedure **GLU**. Along the way, we have seen that it is closely related to many different past approaches. In the dense case, our procedure is as fast as past approaches to within a log factor, and comes with spectral and Frobenius norm bounds on the residual, as well as multiplicative bounds for the other singular values.

For future work, Remark 6.18 seems useful and interesting. Finding applications which particularly benefit from the speed and accuracy guarantees of our procedure is also of interest.

Acknowledgments. We would like to thank the two anonymous referees for their constructive comments that helped us improve our manuscript. We would also like to thank Oleg Balabanov for pointing out to us the connection to oblique projections and other helpful discussions.

REFERENCES

- M. BABOULIN, X. S. LI, AND F.-H. ROUET, Using random butterfly transformations to avoid pivoting in sparse direct methods, in Proceedings of VECPAR 2014, Springer, New York, 2014, pp. 135–144.
- Z. Bai, J. Demmel, and M. Gu, An inverse free parallel spectral divide and conquer algorithm for nonsymmetric eigenproblems, Numer. Math., 76 (1997), pp. 279–308, https://doi.org/ 10.1007/s002110050264.
- [3] G. Ballard, J. Demmel, I. Dumitriu, and A. Rusciano, A Generalized Rank-Revealing Factorization, preprint, https://arxiv.org/abs/arXiv:1909.06524v1, 2019.
- [4] C. BOUTSIDIS AND A. GITTENS, Improved matrix algorithms via the subsampled randomized Hadamard transform, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 1301–1340, https://doi.org/10.1137/120874540.
- [5] D. CARLSON, E. HAYNSWORTH, AND T. MARKHAM, A generalization of the Schur complement by means of the Moore-Penrose inverse, SIAM J. Appl. Math., 26 (1974), pp. 169–175, https://doi.org/10.1137/0126013.
- [6] Z. CHEN AND J. J. DONGARRA, Condition numbers of Gaussian random matrices, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 603–620, https://doi.org/10.1137/040616413.
- [7] K. L. CLARKSON AND D. P. WOODRUFF, Low-rank approximation and regression in input sparsity time, J. ACM, 63 (2017), 54, https://doi.org/10.1145/3019134.
- [8] M. B. COHEN, S. ELDER, C. MUSCO, C. MUSCO, AND M. PERSU, Dimensionality reduction for k-means clustering and low rank approximation, in Proceedings of the 47th Annual ACM Symposium on Theory of Computing, STOC '15, ACM, New York, 2015, pp. 163–172, https://doi.org/10.1145/2746539.2746569.
- [9] M. B. COHEN, J. NELSON, AND D. P. WOODRUFF, Optimal approximate matrix product in terms of stable rank, in Proceedings of the 43rd International Colloquium on Automata, Languages, and Programming (ICALP 2016), Vol. 55, 2016, 11, https://doi.org/10.4230/ LIPIcs.ICALP.2016.11.
- [10] J. DEMMEL, Applied Numerical Linear Algebra, SIAM, Philadelphia, 1997, https://doi.org/ 10.1137/1.9781611971446.
- [11] J. DEMMEL, I. DUMITRIU, AND O. HOLTZ, Fast linear algebra is stable, Numer. Math., 108 (2007), pp. 59–91, https://doi.org/10.1007/s00211-007-0114-x.
- [12] I. Dumitriu, Smallest eigenvalue distributions for two classes of β -Jacobi ensembles, J. Math. Phys., 53 (2012), 103301.
- [13] A. GITTENS AND M. W. MAHONEY, Revisiting the Nyström method for improved large-scale machine learning, J. Mach. Learn. Res., 17 (2016), pp. 3977–4041.
- [14] L. GRIGORI, S. CAYROLS, AND J. W. DEMMEL, Low rank approximation of a sparse matrix based on LU factorization with column and row tournament pivoting, SIAM J. Sci. Comput., 40 (2018), pp. C181–C209, https://doi.org/10.1137/16M1074527.
- [15] M. Gu, Subspace iteration randomization and singular value problems, SIAM J. Sci. Comput., 37 (2015), pp. A1139–A1173, https://doi.org/10.1137/130938700.
- [16] 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, https://doi.org/10.1137/ 0917055.
- [17] 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, https://doi.org/10.1137/090771806.
- [18] R. A. HORN AND C. R. JOHNSON, Matrix Analysis, 2nd ed., Cambridge University Press, Cambridge, UK, 2013.
- [19] T. JIANG, Maxima of entries of Haar distributed matrices, Probab. Theory Related Fields, 131 (2005), pp. 121–144, https://doi.org/10.1007/s00440-004-0376-5.
- [20] A. Khabou, J. Demmel, L. Grigori, and M. Gu, LU factorization with panel rank revealing pivoting and its communication avoiding version, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 1401–1429, https://doi.org/10.1137/120863691.
- [21] B. LAURENT AND P. MASSART, Adaptive estimation of a quadratic functional by model selection, Ann. Statist., 28 (2000), pp. 1302–1338, https://doi.org/10.1214/aos/1015957395.

- [22] M. W. Mahoney and P. Drineas, CUR matrix decompositions for improved data analysis, Proc. Natl. Acad. Sci. USA, 106 (2009), pp. 697–702.
- [23] P.-G. MARTINSSON AND J. A. TROPP, Randomized numerical linear algebra: Foundations and algorithms, Acta Numer., 29 (2020), pp. 403–572, https://doi.org/10.1017/ S0962492920000021.
- [24] L. MIRANIAN AND M. Gu, Strong rank revealing LU factorizations, Linear Algebra Appl., 367 (2003), pp. 1–16.
- [25] C. Musco and C. Musco, Randomized block Krylov methods for stronger and faster approximate singular value decomposition, in Proceedings of the 28th International Conference on Neural Information Processing Systems, Vol. 1, NIPS'15, MIT Press, Cambridge, MA, 2015, pp. 1396–1404.
- [26] C. Musco and D. P. Woodruff, Sublinear time low-rank approximation of positive semidefinite matrices, in Proceedings of the 2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS), IEEE, Washington, DC, 2017, pp. 672–683, https://doi.org/10.1109/FOCS.2017.68.
- [27] J. NELSON AND H. L. NGUYEN, OSNAP: Faster numerical linear algebra algorithms via sparser subspace embeddings, in Proceedings of the 2013 IEEE 54th Annual Symposium on Foundations of Computer Science, IEEE, Washington, DC, 2013, pp. 117–126, https://doi.org/10.1109/FOCS.2013.21.
- [28] A. SANKAR, D. A. SPIELMAN, AND S.-H. TENG, Smoothed analysis of the condition numbers and growth factors of matrices, SIAM J. Matrix Anal. Appl., 28 (2006), pp. 446–476, https://doi.org/10.1137/S0895479803436202.
- [29] T. SARLOS, Improved approximation algorithms for large matrices via random projections, in Proceedings of the 2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS'06), IEEE, Washington, DC, 2006, pp. 143–152.
- [30] G. SHABAT, Y. SHMUELI, Y. AIZENBUD, AND A. AVERBUCH, Randomized LU decomposition, Appl. Comput. Harmon. Anal., 44 (2018), pp. 246–272, https://doi.org/10.1016/ j.acha.2016.04.006.
- [31] L. N. Trefethen and R. S. Schreiber, Average-case stability of Gaussian elimination, SIAM J. Matrix Anal. Appl., 11 (1990), pp. 335–360, https://doi.org/10.1137/0611023.
- [32] 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, https://doi.org/10.1137/17M1111590.
- [33] D. P. WOODRUFF, Sketching as a tool for numerical linear algebra, Found. Trends Theor. Comput. Sci., 10 (2014), pp. 1–157, https://doi.org/10.1561/040000066.