

Low-Precision QR Factorization: Analysis, Algorithms, and Applications

L. Minal Yung, Alyson Fox, and Geoffrey Sanders

August 26, 2019

1 Introduction

Development of new graphics processing units (GPUs) that perform half precision arithmetic up to 16 times faster than double precision arithmetic motivates use of low precision computations whenever possible. GPUs were initially designed to support computer graphics and evolved to ~~continue supporting that sole utility prior to the emergence of~~ general purpose GPUs (GPGPUs) ~~in recent years. GPGPUs were largely fueled by growing interest in deep learning, which requires~~ dealing with large amounts of data that can be processed in parallel and can benefit from efficiently harnessing the highly parallelizable and wide memory bandwidth structure of GPUs. Some drawbacks of GPUs are that they are limited to simple tasks, and have a narrower range of capabilities in comparison to central processing units (CPUs). However, there are plenty of applications and algorithms in scientific computing that can potentially take advantage of GPUs.

Training deep artificial neural networks with various tasks such as multiplication [1] and storage [2], assigned to lower precision arithmetic can easily be implemented on GPUs, and these works feature methods that are already in common practice in data science applications. The accuracy of a numerical algorithm depends on several factors including numerical stability and well-conditionedness of the problem, both of which may be sensitive to rounding errors. Low precision floats use fewer bits than high precision floats to represent the real numbers, and naturally incur larger rounding errors. Therefore, error attributed to round-off may have a larger influence over the total error when using low precision, and some standard algorithms that are in wide use may no longer be numerically stable when using half precision floating arithmetic and storage.

The low precision computing environments that we consider are designed to imitate those of new GPUs, which employ multiple precision types. For example, Tesla V100's Tensor Cores perform matrix-multiply-and-accumulate of half precision input data with full-precision products and single precision summation accumulate [cite]. The existing rounding error analyses only allow a uniform floating point precision for storage and operations, and cannot accurately represent mixed-precision settings such as TensorCores. Therefore, using existing error analyses by rounding up to the higher precision within the mixed-precision setting would be too optimistic, and rounding down to the lowest precision within the mixed-precision setting would be too pessimistic. One form of battling the pessimistic nature of deterministic error bounds is probabilistic error bounds, but these also suffer from being restricted to uniform precision procedures. In this work, we develop a framework

for mixed-precision rounding error analysis, and explore half-precision QR factorization for data and graph analysis applications.

Our findings are that the new mixed-precision error analysis produces tighter error bounds, and block QR algorithms are able to operate in low precision more robustly than non-block techniques.

Given a matrix $A \in \mathbb{R}^{m \times n}$ for $m \geq n$, we consider performing the QR factorization, where

$$A = QR, \quad Q \in \mathbb{R}^{m \times m}, \quad R \in \mathbb{R}^{m \times n}.$$

Q is orthogonal, $Q^T Q = I_{m \times m}$, and R is upper-trapezoidal, $R_{ij} = 0$ for $i > j$.

The above formulation is a full QR factorization whereas a more efficient *thin* QR factorization results in $Q \in \mathbb{R}^{m \times n}$ and $R \in \mathbb{R}^{n \times n}$.

$$A = QR = [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0_{m-n \times n} \end{bmatrix} = Q_1 R_1$$

Here, $Q_1 R_1$ is the *thin* QR factorization, where the columns of Q_1 are orthonormal, and R_1 is upper-triangular. In many applications, computing the *thin* decomposition requires less computation and is sufficient in performance.

In section 2, we will give an overview of the modern developments in hardware that motivates rounding error analysis that supports multiple precision types and present a set of error analysis tools. The HQR factorization algorithm and a mixed-precision rounding error analysis of its implementation is in section 3. In section 4, we present the TSQR algorithm as well as numerical experiments that show that TSQR can be useful in low precision environments. Section 5 explores the use of low and mixed precision QR algorithms as subroutines for two applications, spectral clustering and sparse regression in the context of discovery of equations.

1.1 Notation

While important definitions are stated explicitly in the text, Table 1 serves to establish basic notation.

2 Floating Point Numbers and Error Analysis Tools

2.1 Modern GPU Hardware

2.2 Representation of Real Numbers

Consider floating point number systems that are defined by

$$\text{significand} \times \text{base}^{\text{exponent}} \quad (1)$$

This is the generic form of floating point representations, including the IEEE 754 Standard which was established in 1985 and has been been accepted and followed by most modern machines since. Let $F \subset \mathbb{R}$ denote the space of some floating point number system with base $b \in \mathbb{N}$, precision $t \in \mathbb{N}$, significand $\mu \in \mathbb{N}$, and exponent range $[q_{\min}, q_{\max}] \subset \mathbb{Z}$. Then every element y in F can be written as

$$y = \pm \mu \times b^{q-t} \quad (2)$$

First time introduced, I would move up to second sentence

(44 bits)

(the 16 bits)

However, in recent years where developed, which was

List a few more applications that would benefit from mixed precision. What is this? Inform the reader

How many bits?

double? what do you mean? unclear to reader

cite

cite

✓

Symbol(s)	Definition(s)	Section(s)
Q	Orthogonal factor of matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$; m -by- n (full) or m -by- n (thin)	1
R	Upper triangular or trapezoidal factor of matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$; m -by- n (full) or n -by- n (thin)	1
$\mathbf{A}^{(k)}$	Matrix \mathbf{A} after k Householder transformations.	3.1
$\mathbf{fl}(\mathbf{x})$, \mathbf{x}	Quantity \mathbf{x} calculated from floating point operations	2
b, t, μ, η	Base/precision/mantissa/exponent bits	2
Inf	Values outside the range of representable numbers	3.1.2
k	Number of FLOPs	2
u_q	Unit round-off for precision t and base b : $\frac{1}{2}b^{t-1}$	2
δ_q	Quantity bounded by: $ \delta_q < u_q$	2
$\gamma_q^{(k)}, \theta_q^{(k)}$	Quantity bounded by: $ \theta_q^{(k)} \leq \gamma_q^{(k)}$	2
\mathbf{x}, \mathbf{A}	Vector, matrix	2
m, n	Number of rows, columns of matrix, or length of vector	1
i, j	Row, column index of matrix or vector	3.1
$\ \mathbf{x}\ _2, \ \mathbf{A}\ _2$	Vector, operator 2-norm	3
$ c , \mathbf{x} , \mathbf{A} $	Absolute value of constant, all elements of vector, matrix	3
x_i	i^{th} element of vector \mathbf{x}	3.1
$\mathbf{A}(a:b:, c:d)$, $\mathbf{A}[:, c:d]$	Rows a to b , columns c to d of matrix \mathbf{A}	3.1
\hat{e}_i	Cardinal vector	3
$\mathbf{0}_{m \times n}$, \mathbf{I}_n	m -by- n zero matrix, n -by- n identity matrix	1
$\mathbf{I}_{m \times n}$	$(\mathbf{I}_n - \mathbf{0}_{n \times (m-n)})^\top$	4
$\mathbf{P}_v, \mathbf{P}_t$	Householder transformation defined by \mathbf{v} , i^{th} Householder transformation in HQR	3.1
u_s, u_p, u_w	Unit round-off for sum, product, and storage (write)	2.3

Table 1: Basic definitions

where μ is any integer in $[0, b^t - 1]$, and η is an integer in $[\eta_{\min}, \eta_{\max}]$. While base, precision, and exponent range are fixed and define a floating point number system, the sign, significand, and exponent uniquely identifies a single number within that system. As the significand is always an integer in $[0, b^{t-1}]$, it can be combined with the precision and represent a fraction between 0 and 1. $\frac{\mu}{b^t}$. Table 2 shows IEEE 754 floating point number types described by the same parameters as in Equation 2, and the remainder of this paper will discuss floating point arithmetic as is defined by the IEEE 754 Standard.

Name	b	t	# of exponent bits	η_{\min}	η_{\max}	u
IEEE754 half	2	11	5	-15	16	4.883e-04
IEEE754 single	2	24	8	-127	128	5.960e-08
IEEE754 double	2	53	11	-1023	1024	1.110e-16

Table 2: IEEE754 formats and their primary attributes.

Although operations we use on \mathbb{R} cannot be replicated exactly due to the finite cardinality of \mathbb{F} , we can still approximate the accuracy of analogous floating point operations. We adopt

the rounding error analysis tools developed in [3], which allow a relatively simple framework for formulating error bounds for complex linear algebra operations. A short analysis of floating point operations (cf. Theorem 2.2 [3]) shows that the relative error is controlled by the unit round-off, $u = \frac{1}{2}b^{t-1}$.

Let op be any basic operation between 2 floating point numbers from the set $\text{OP} = \{+, -, \times, \div\}$. The true value $(x \text{ op } y)$ lies in \mathbb{R} and it is rounded using some conversion to a floating point number, $\mathbf{fl}(x \text{ op } y)$, admitting a rounding error. The IEEE 754 Standard requires *correct rounding*, which rounds the exact solution $(x \text{ op } y)$ to the closest floating point number, and in case of a tie, to the floating point number that has a mantissa ending in an even number. *Correct rounding* gives us an assumption for the error model where a single basic floating point operation yields a relative error, δ , bounded in the following sense:

$$\mathbf{fl}(x \text{ op } y) = (1 + \delta)(x \text{ op } y), \quad |\delta| \leq u, \quad \text{op} \in \{+, -, \times, \div\}. \quad (3)$$

We use Equation 3 as a building block in accumulating errors from successive floating point operations (FLOPs). For example, consider computing $x + y + z$ where $x, y, z \in \mathbb{R}$. Let's assume that the machine can only compute one operation at a time, and we take the convention of computing the left-most operation first. Then there is a rounding error in computing $s_1 := \mathbf{fl}(x + y) = (1 + \delta)(x + y)$, and another rounding error in computing $s_2 := \mathbf{fl}(s_1 + z) = (1 + \delta)(s_1 + z)$, where $|\delta|, |\hat{\delta}| < u$. Let's expand the final result:

$$\mathbf{fl}(x + y + z) = \mathbf{fl}(\mathbf{fl}(x + y) + z) = \mathbf{fl}((1 + \delta)(x + y) + z) \quad (4)$$

$$= (1 + \delta)((1 + \delta)(x + y) + z) \quad (5)$$

$$= (1 + \delta)(1 + \delta)(x + y) + (1 + \delta)z \quad (6)$$

We can see that keeping track of rounding errors from each operation can quickly grow to be challenging, even with just two successive operations. A way of simplifying complicated expressions like Equation 6 is crucial in developing error analyses for complex linear algebra operations. Lemma 2.1 introduces one such convenient and elegant bound that simplifies accumulated rounding errors.

Lemma 2.1 (Lemma 3.1 [3]). *Let $|\delta_i| < u$ and $\mu_i \in \{-1, +1\}$ for $i = 1, \dots, k$, and $ku < 1$. Then,*

$$\prod_{i=1}^k (1 + \delta_i)^{\mu_i} = 1 + \theta^{(k)} \quad (7)$$

where

$$|\theta^{(k)}| \leq \frac{ku}{1 - ku} =: \gamma^{(k)}, \quad (8)$$

In other words, $\theta^{(k)}$ represents the accumulation of rounding errors from k successive operations, and it is bounded by $\gamma^{(k)}$. Allowing $\theta^{(k)}$'s to be any arbitrary value within the corresponding $\gamma^{(k)}$ bounds further aids in keeping a clear, simple error analysis. Applying this lemma to our example of adding three numbers results in:

$$\mathbf{fl}(x + y + z) = (1 + \delta)(1 + \delta)(x + y) + (1 + \delta)z = (1 + \theta^{(2)})(x + y) + (1 + \theta^{(1)})z. \quad (9)$$

Since $|\theta^{(1)}| \leq \gamma^{(1)} < \gamma^{(2)}$, we can further simplify Equation 9 to

$$\mathbf{fl}(x + y + z) = (1 + \hat{\theta}^{(2)})(x + y + z), \quad \text{where } |\hat{\theta}^{(2)}| \leq \gamma^{(2)}. \quad (10)$$

Typically, error bounds formed in the fashion of Equation 10 converted to relative errors in order to put the error magnitudes in perspective. In our example, for nonzero $(x + y + z)$, we have:

$$\frac{|(x + y + z) - \mathbf{fl}(x + y + z)|}{|x + y + z|} \leq \gamma^{(2)} \quad (11)$$

Although Lemma 2.1 only requires $ku < 1$, we actually need $ku < \frac{1}{2}$, which implies $\gamma^{(k)} < 1$, in order to maintain a meaningful relative error bound. While this assumption, $\gamma^{(k)} < 1$, is easily satisfied by fairly large k in higher precision floating point numbers, it is a problem even for small k in lower precision floating point numbers. Table 3 shows the maximum value of k that still guarantees a relative error below 100% ($\gamma^{(k)} < 1$). Thus, accumulated rounding errors in lower precision types

precision	u	$k = \text{argmax}^{(k)}(\gamma^{(k)} < 1)$
half	4.883e-04	512
single	5.960e-08	$\approx 4.194e06$
double	1.110e-16	$\approx 2.252e15$

Table 3: Upper limits of meaningful relative error bounds in the $\gamma^{(k)}$ notation.

grow unstable very quickly and with fewer operations in comparison to higher precision types. As k represents the number of FLOPs, this restricts low-precision floating point operations to smaller problem sizes and lower complexity algorithms.

To clearly illustrate how this restricts rounding error analysis in half precision, we now consider performing the dot product of two vectors. A forward error bound for dot products is:

$$\frac{|\mathbf{x}^\top \mathbf{y} - \mathbf{fl}(\mathbf{x}^\top \mathbf{y})|}{|\mathbf{x}^\top \mathbf{y}|} \leq \gamma^{(m)}, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^m. \quad (12)$$

where details and proof for this statement can be found in Section 3.4 of [3]. While this result does not guarantee a high relative accuracy when $|\mathbf{x}^\top \mathbf{y}| \ll |\mathbf{x}|^\top |\mathbf{y}|$, high relative accuracy is expected in some special cases. For example, let $\mathbf{x} = \mathbf{y}$. Then we have exactly $|\mathbf{x}^\top \mathbf{x}| = |\mathbf{x}|^\top |\mathbf{x}| = \|\mathbf{x}\|_2^2$. This leads to

$$\left| \frac{\|\mathbf{x}\|_2^2 - \mathbf{fl}(\|\mathbf{x}\|_2^2)}{\|\mathbf{x}\|_2^2} \right| \leq \gamma_p^{(d+2)} \quad (13)$$

Since vectors of length m accumulate rounding errors that are bounded by $\gamma^{(m)}$, the worst-case relative error bound for a dot product of vectors of length 512 is already at 100% ($\gamma_{\text{half}}^{(512)} = 1$).

We present a simple numerical experiment that shows that the standard deterministic error bound is too pessimistic, and cannot be practically used to approximate rounding error for half-precision arithmetic. In this experiment, we generated 2 million random half-precision vectors of length 512 from two random distributions: the standard normal distribution, $N(0, 1)$, and the uniform distribution over $(0, 1)$. Half precision arithmetic was simulated by calling Algorithm 1 for every multiplication and summation step required in calculating the dot product, $\mathbf{fl}(\mathbf{x}^\top \mathbf{y})$.

Algorithm 1: $\mathbf{z}_{\text{half}} = \text{simHalf}(f, \mathbf{x}_{\text{half}}, \mathbf{y}_{\text{half}})$ Simulate function $f \in \{\text{OP}\} \cup \{\text{dot_product}\}$ in half precision arithmetic given input variables \mathbf{x}, \mathbf{y} . Function **castup** converts half precision floats to single precision floats, and **castdown** converts single precision floats to half precision floats by rounding to the nearest half precision float.

Input: $\mathbf{x}_{\text{half}}, \mathbf{y}_{\text{half}} \in \mathbb{F}_{\text{half}}^m, f: \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}^n$
Output: $\mathbf{fl}(f(\mathbf{x}_{\text{half}}, \mathbf{y}_{\text{half}})) \in \mathbb{F}_{\text{half}}^n$

```

1  $\mathbf{x}_{\text{single}}, \mathbf{y}_{\text{single}} \leftarrow \text{castup}([\mathbf{x}_{\text{half}}, \mathbf{y}_{\text{half}}])$ 
2  $\mathbf{z}_{\text{single}} \leftarrow \mathbf{fl}(f(\mathbf{x}_{\text{single}}, \mathbf{y}_{\text{single}}))$ 
3  $\mathbf{z}_{\text{half}} \leftarrow \text{castdown}(\mathbf{z}_{\text{single}})$ 
4 return  $\mathbf{z}_{\text{half}}$ 
```

The casting up step is exact since all half precision numbers can be exactly represented in single precision, $\mathbb{F}_{\text{half}} \subset \mathbb{F}_{\text{single}}$; the second step incurs a rounding error from a single precision arithmetic operation, and the casting down step incurs a rounding error from casting down to half precision. Note that using Algorithm 1 for any operation in OP results in simulating half precision arithmetic, whereas using it with the dot product results in simulating mixed precision arithmetic instead. The relative error in this experiment is formulated as the left hand side of the inequality in Equation 12, where all operations outside of calculating $\mathbf{fl}(\mathbf{x}^\top \mathbf{y})$ is executed by casting up to double precision format and using double precision arithmetic. Table 4 shows statistics from computing the relative error for simulated half precision dot products of 512-length random vectors.

Random Distribution	Average	Standard deviation	Maximum
Standard normal	1.627e-04	1.640e-04	2.838e-03
Uniform (0, 1)	2.599e-03	1.854e-03	1.399e-02

Table 4: Statistics from dot product backward relative error in for 512-length vectors stored in half-precision and computed in simulated half-precision.

We see that the inner products of vectors sampled from the standard normal distribution have backward relative errors that don't deviate much from the unit round-off (4.883e-4), whereas the vectors sampled from the uniform distribution tend to accumulate larger errors. Even so, the theoretical upper error bound of 100% is too pessimistic, and it is difficult to predict the kind of results this experiment shows. Recent works in developing probabilistic bounds on rounding errors of floating point operations have shown that the inner product relative backward error for the conditions used for this experiment is bounded by $5.466e-2$ with probability 0.99.

Most importantly, no rounding error bounds (deterministic and probabilistic) allow flexibility in the precision types used for different operations. This restriction is the biggest obstacle in gaining an understanding of rounding errors to expect from computations done on emerging hardware that support mixed-precision such as GPUs that employ mixed-precision arithmetic. In this paper, we develop a mixed-precision error analysis that allows multiple precision types and is based on the rounding error analysis framework established in [3].

Lemma 2.2 shows rules from Lemma 3.3 in [3] that summarize how to accumulate errors represented by θ 's and γ 's.

Lemma 2.2. For any positive integer k , let $\theta^{(k)}$ denote a quantity bounded according to $|\theta^{(k)}| \leq \frac{1-u}{1-ku} =: \gamma^{(k)}$. The following relations hold for positive integers i, j , and nonnegative integer k .

Arithmetic operations between $\theta^{(k)}$'s.

$$(1 + \theta^{(k)})(1 + \theta^{(j)}) = (1 + \theta^{(k+j)})$$

$$\frac{1 + \theta^{(k)}}{1 + \theta^{(j)}} = \begin{cases} 1 + \theta^{(k-j)}, & j \leq k \\ 1 + \theta^{(k+2j)}, & j > k \end{cases}$$

Operations on γ 's:

$$\gamma_i^{(k)} \gamma_j^{(j)} \leq \gamma_{\min(k,j)}^{(k+j)} \quad \text{for } \max(k,j)u \leq \frac{1}{2}$$

$$n\gamma_i^{(k)} \leq \gamma_i^{(nk)}, \quad \text{for } n \leq \frac{1}{uk}$$

$$\gamma_i^{(k)} + u \leq \gamma_i^{(k+1)}$$

$$\gamma_i^{(k)} + \gamma_j^{(j)} + \gamma_i^{(k)} \gamma_j^{(j)} \leq \gamma_i^{(k+j)}$$

In Lemma 2.3, we present modified versions of the rules in Lemma 2.2. This mixed-precision error analysis relies on the framework given by Lemma 2.1, which best allows us to keep a simple analysis. These relations allow us to easily accumulate errors in terms of θ 's and γ 's, and aid in writing clear and simpler error analyses. The modifications support multiple precision types, whereas Lemma 2.2 assumes that the same precision is used in all operations.

We distinguish between the different precision types using subscripts — these types include products (p), sums (s), and storage formats (w).

Lemma 2.3. For any nonnegative integer k and some precision q , let $\theta_q^{(k)}$ denote a quantity bounded according to $|\theta_q^{(k)}| \leq \frac{k u_q}{1 - u_q} =: \gamma_q^{(k)}$. The following relations hold for two precisions s and p , positive integers, j, j_p , non-negative integers k , and k_p , and $c > 0$.

$$(1 + \theta_p^{(k_p)})(1 + \theta_p^{(j_p)})(1 + \theta_s^{(k)}) = (1 + \theta_p^{(k_p+j_p)})(1 + \theta_s^{(k+j)}) \quad (14)$$

$$\frac{(1 + \theta_p^{(k_p)})(1 + \theta_s^{(k)})}{(1 + \theta_p^{(j_p)})(1 + \theta_s^{(j)})} = \begin{cases} (1 + \theta_s^{(k+j)}) / (1 + \theta_p^{(k_p+j_p)}), & j_s \leq k_s, j_p \leq k_p \\ (1 + \theta_s^{(k+(2j_s))}) / (1 + \theta_p^{(k_p+j_p)}), & j_s \leq k_s, j_p > k_p \\ (1 + \theta_s^{(k+j)}) / (1 + \theta_p^{(k_p+2j_p)}), & j_s > k_s, j_p \leq k_p \\ (1 + \theta_s^{(k+(2j_s))}) / (1 + \theta_p^{(k_p+2j_p)}), & j_s > k_s, j_p > k_p \end{cases} \quad (15)$$

Without loss of generality, let $1 \gg u_p \gg u_s > 0$. Let d , a nonnegative integer, and $c \in [0, \frac{u_s}{u_p}]$ be numbers that satisfy $k u_s = d u_p + r u_s$. Alternatively, d can be defined by $d := \lfloor \frac{k u_s}{u_p} \rfloor$.

$$\gamma_s^{(k_s)} \gamma_p^{(k_p)} \leq \gamma_p^{(k_p)} \quad \text{for } k_p u_p \leq \frac{1}{2} \quad (16)$$

$$\gamma_s^{(k_s)} + u_p \leq \gamma_p^{(d+2)} \quad (17)$$

$$\gamma_p^{(k_p)} + u_s \leq \gamma_p^{(k_p+1)} \quad (18)$$

$$\gamma_p^{(k_p)} + \gamma_s^{(k_s)} + \frac{\gamma_p^{(k_p)} \gamma_s^{(k_s)}}{\gamma_p^{(k_p)}} \leq \gamma_p^{(k_p+d+1)} \quad (19)$$

A proof for Equation 19 is shown in Appendix A.3.1. We use these principles to establish a mixed-precision rounding error analysis for computing the dot product, which is crucial in many linear algebra routines such as the QR factorization.

2.3 Inner product Mixed-Precision error

We will see in Section 3 that the inner product is a building block of the HQR factorization (HQR) algorithm, which was introduced in [4]. More generally, it is used widely in most linear algebra tools such as matrix-vector multiply and projections. Thus, we will generalize classic round-off error analysis of inner products to algorithms that may employ different precision types to different operations. Specifically, we consider performing an inner product with the storage precision, u_w , being lower than the summation precision, u_s . This is designed to provide a more accurate rounding error analysis of mixed-precision floating point operations present in recent GPU technologies such as NVIDIA's TensorCore. Currently, TensorCore computes the inner product of vectors stored in half-precision by employing full precision multiplications and a single-precision accumulator. As the majority of rounding errors from computing inner products occur during summation (c.f. Section 3.1, [9]), the single precision accumulator immensely reduces the error in comparison to using only half-precision operations. This increase in accuracy combined with its speedy performance motivates us to: 1) study how to best utilize mixed-precision arithmetic in algorithms, and 2) to develop more accurate error analyses appropriate for mixed-precision algorithms.

Lemma 2.4 and Corollary 2.5 present two mixed-precision forward error bounds for inner products, which show a tighter bound than the existing error bounds. In both cases, we assume storage in the lowest precision with round-off value, u_w , and summation performed with a higher precision with round-off value, u_s , and let $d \approx \min(u_s/u_w)$, where m is the length of the vectors. Although there are additional differing assumptions in these two lemmas, results from both show a strong dependence on d .

Lemma 2.4. Let w , p , and s each represent floating point precisions for storage, product, and summation, where the varying precisions are defined by their unit round-off values denoted by u_w , u_p , and u_s . Let $\mathbf{x}, \mathbf{y} \in \mathbb{F}_w^m$ be two arbitrary vectors stored in w precision. If an inner product performs multiplications in precision p , and addition of the products using precision s , then

$$\|(\mathbf{x}^T \mathbf{y}) - (\mathbf{x} + \Delta \mathbf{x}) \mathbf{y} = \mathbf{x} (\mathbf{y} + \Delta \mathbf{y})\| \quad (20)$$

where $|\Delta \mathbf{x}| \leq \gamma_p^{(i, m-1)} |\mathbf{x}|$, $|\Delta \mathbf{y}| \leq \gamma_p^{(i, m-1)} |\mathbf{y}|$ componentwise, and

$$\gamma_p^{(i, m-1)} := (1 + u_p)(1 + \gamma_s^{(m-1)}) - 1$$

This result is then stored in precision w and if we further assume that $u_w = u_p > u_s$, then $|\Delta \mathbf{x}| \leq \gamma_w^{(d+2)} |\mathbf{x}|$ and $|\Delta \mathbf{y}| \leq \gamma_w^{(d+2)} |\mathbf{y}|$ where $d := \lfloor \frac{(m-1)u_s}{u_w} \rfloor$.

Corollary 2.5 presents another mixed-precision forward error bound for mixed-precision inner products with additional constraints. Here, we assume that the vectors are being stored in a lower precision than the precision types being used for multiplications and additions. This scenario is similar to how TensorCore technology works in GPUs.

Corollary 2.5. In addition to the assumptions in Lemma 2.4, assume $1 \gg u_w \gg u_s > 0$, and for any two numbers x, y in \mathbb{F}_w , their product xy is in \mathbb{F}_s . Let $\mathbf{x}, \mathbf{y} \in \mathbb{F}_w^m$ be two arbitrary vectors stored in w precision. If an inner product performs multiplications in full precision, and addition of the products using precision s , then

$$\|(\mathbf{x}^T \mathbf{y}) - (\mathbf{x} + \Delta \mathbf{x}) \mathbf{y} = \mathbf{x} (\mathbf{y} + \Delta \mathbf{y})\| \quad (21)$$

where $|\Delta \mathbf{x}| \leq \gamma_w^{(d+1)} |\mathbf{x}|$, $|\Delta \mathbf{y}| \leq \gamma_w^{(d+1)} |\mathbf{y}|$ componentwise, and $d := \lfloor \frac{(m-1)u_s}{u_w} \rfloor$.

Proofs for Lemma 2.4 and Corollary 2.5 are shown in Appendix A.3.1. The analyses for these two differ only in the type of mixed-precision arithmetic performed within the inner product subroutine, and the difference is revealed to result in either $\gamma_w^{(d(1))}$ or $\gamma_w^{(d(2))}$. For the rest of this paper, we will refer to the forward error bound for the inner product as $\gamma_w^{d/2}$ for $d = 1, 2$ to generalize the analysis for varying assumptions. This simplification allows us to use the same analysis for the remaining steps of the HQR algorithm presented in the following sections.

3 Mixed-Precision HQR Factorization

The HQR algorithm uses Householder transformations to zero out elements below the diagonal of a matrix. We present this as zeroing out all but the first element of some vector, $\mathbf{x} \in \mathbb{R}^m$.

Lemma 3.1. *Given vector $\mathbf{x} \in \mathbb{R}^m$, there exist Householder vector \mathbf{v} and Householder transformation matrix \mathbf{P}_v such that $\mathbf{P}_v \mathbf{x}$ zeros out \mathbf{x} below the first element.*

$$\begin{aligned} \sigma &= -\text{sign}(x_1)\|\mathbf{x}\|_2, \quad \mathbf{v} = \mathbf{x} - \sigma\mathbf{e}_1, \\ \beta &= \frac{2}{\mathbf{v}^T \mathbf{v}} = -\frac{1}{\sigma v_1}, \quad \mathbf{P}_v = \mathbf{I}_m - \beta \mathbf{v} \mathbf{v}^T \end{aligned} \quad (22)$$

The transformed vector, $\mathbf{P}_v \mathbf{x}$, has the same 2-norm as \mathbf{x} since Householder transformations are orthogonal.

$$\mathbf{P}_v \mathbf{x} = \sigma \mathbf{e}_1 \quad (23)$$

In addition, \mathbf{P}_v is symmetric and orthogonal, $\mathbf{P}_v = \mathbf{P}_v^T = \mathbf{P}_v^{-1}$, and therefore involutory, $\mathbf{P}_v^2 = \mathbf{I}$.

3.1 HQR Factorization Algorithm

Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and Lemma 3.1, HQR is done by repeating the following processes. For $i = 1, 2, \dots, n$,

Step 1) Find and store the Householder vector (\mathbf{v}) and constant (β) that zeros out the i^{th} column of \mathbf{A} beneath the i^{th} element.

Step 2) Apply the Householder transformation to the bottom right partition ($\mathbf{A}[i:m, i:n]$) of the matrix.

until only an upper triangular matrix remains.

Consider the following 4-by-3 matrix example adapted from [3]. Let \mathbf{P}_i represent the i^{th} Householder transformation of this algorithm.

$$\begin{aligned} \mathbf{A} = \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} &\xrightarrow{\text{apply } \mathbf{P}_1 \text{ to } \mathbf{A}^{(1)} = \mathbf{A}} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix} \xrightarrow{\text{apply } \mathbf{P}_2 \text{ to } (\mathbf{A}^{(1)} - \mathbf{P}_1 \mathbf{A})} \\ \\ \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \end{bmatrix} &\xrightarrow{\text{apply } \mathbf{P}_3 \text{ to } (\mathbf{A}^{(2)} - \mathbf{P}_2 \mathbf{P}_1 \mathbf{A})} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \end{bmatrix} \end{aligned}$$

Since the final matrix $\mathbf{P}_1 \mathbf{P}_2 \mathbf{P}_3 \mathbf{A}$ is upper-triangular, this is the R factor of the QR decomposition. Set $\mathbf{Q}^T := \mathbf{P}_1 \mathbf{P}_2 \mathbf{P}_3$. Then we can formulate Q via

$$\mathbf{Q} = (\mathbf{P}_1 \mathbf{P}_2 \mathbf{P}_3)^T = \mathbf{P}_1^T \mathbf{P}_2^T \mathbf{P}_3^T = \mathbf{P}_1 \mathbf{P}_2 \mathbf{P}_3,$$

where the last equality results from the symmetric property of \mathbf{P}_i 's.

Returning to the general case, we have

$$\mathbf{Q}_{\text{full}} = \mathbf{P}_1 \cdots \mathbf{P}_n \quad \text{and} \quad \mathbf{R}_{\text{full}} = \mathbf{Q}^T \mathbf{A} = \mathbf{P}_n \cdots \mathbf{P}_1 \mathbf{A}, \quad (24)$$

for the orthogonal factor in a full QR factorization, and

$$\mathbf{Q}_{\text{thin}} = \mathbf{P}_1 \cdots \mathbf{P}_n \mathbf{I}_{m \times n} \quad \text{and} \quad \mathbf{R}_{\text{thin}} = \mathbf{I}_{m \times n}^T \mathbf{Q}^T \mathbf{A} = \mathbf{I}_{m \times n}^T \mathbf{P}_n \cdots \mathbf{P}_1 \mathbf{A}, \quad (25)$$

3.1.1 HQR Factorization Implementation

The Householder transformation is implemented by a series of inner and outer products, since Householder matrices are rank-1 updates of the identity. This is much less costly than forming \mathbf{P}_v , then performing matrix-vector or matrix-matrix multiplications. For some $\mathbf{P}_v = \mathbf{I} - \beta \mathbf{v} \mathbf{v}^T$, we result in the following computation.

$$\mathbf{P}_v \mathbf{x} = (\mathbf{I} - \beta \mathbf{v} \mathbf{v}^T) \mathbf{x} = \mathbf{x} - (\beta \mathbf{v}^T \mathbf{x}) \mathbf{v} \quad (26)$$

The routine shown in Equation 26 is used in forming R and Q. Given a vector $\mathbf{x} \in \mathbb{R}^m$, Algorithm 2 calculates the Householder constant β and Householder vector \mathbf{v} , that zero out \mathbf{x} below the first element, and also returns σ . Algorithm 3 is the HQR algorithm where information necessary to build Q is returned instead of explicitly forming Q—the Householder vector and constant at the k^{th} step is stored as the k^{th} column of matrix $\mathbf{V} \in \mathbb{R}^{m \times n}$ and the k^{th} element of vector $\beta \in \mathbb{R}^n$.

Finally, the Q factor can be built using Algorithm 4. While this algorithm shows how to left multiply Q to any input matrix B given V and β , putting in $\mathbf{B} \equiv \mathbf{I}_{m \times n}$ will yield \mathbf{Q}_{thin} .

Algorithm 2: ($\beta, \mathbf{v}, \sigma = \text{hh_vec}(\mathbf{x})$). Given a vector $\mathbf{x} \in \mathbb{R}^m$, return the Householder vector \mathbf{v} , a Householder constant β , and σ such that $(\mathbf{I} - \beta \mathbf{v} \mathbf{v}^T) \mathbf{x} = \sigma \mathbf{e}_1$, and $\mathbf{v}_1 = 1$.

Input: $\mathbf{x} \in \mathbb{R}^m$
Output: $\mathbf{v} \in \mathbb{R}^m$, and $\sigma, \beta \in \mathbb{R}$ such that $(\mathbf{I} - \beta \mathbf{v} \mathbf{v}^T) \mathbf{x} = \pm \|\mathbf{x}\|_2 \mathbf{e}_1 = \sigma \mathbf{e}_1$
(σ is chosen the sign of x_1 to avoid cancellation of x_1 , the so-called standard in LAPACK, LAPACK procedure 13). Then return $\beta, \mathbf{v}, \sigma$.

```

1  $\mathbf{v} \leftarrow \mathbf{x}$ 
2  $\sigma \leftarrow -\text{sign}(x_1)\|\mathbf{x}\|_2$ 
3  $\mathbf{v}_1 \leftarrow x_1 = \sigma$ 
4  $\beta \leftarrow -\frac{2}{\sigma}$ 
5  $\mathbf{v} \leftarrow \frac{1}{v_1} \mathbf{v}$ 
6 return  $\beta, \mathbf{v}, \sigma$ 

```

3.1.2 Normalization of Householder Vectors

Equation 22 gives a single Householder transformation matrix \mathbf{P}_v for all \mathbf{v}^i in $\text{Span}(\mathbf{v})$, which allows for many different ways of normalizing the Householder vectors as well as the choice of not

Algorithm 3: $V, \beta, \text{qr}(A)$ Given a matrix $A \in \mathbb{R}^{m \times n}$ where $m \geq n$, return matrix $V \in \mathbb{R}^{m \times n}$, vector $\beta \in \mathbb{R}^n$, and upper triangular matrix R . An orthogonal matrix Q can be generated from V and β , and $QR = A$.

Input: $A \in \mathbb{R}^{m \times n}$ where $m \geq n$
Output: V, β, R

```

1  $V, \beta \leftarrow 0_{m \times n}, 0_m$ 
2 for  $i = 1 : n$  do
3    $v, \beta, \sigma \leftarrow \text{hh\_vv}(A[i : \text{end}, i])$ 
4    $V[i : \text{end}, i], \beta, A[i, i] \leftarrow v, \beta, \sigma$  // Store the Householder vectors and constants
5   // The next two steps update A
6    $A[i + 1 : \text{end}, i] \leftarrow \text{zeros}(m - i)$ 
7    $A[i : \text{end}, i + 1 : \text{end}] \leftarrow A[i : \text{end}, i + 1 : \text{end}] - \beta v v^T A[i : \text{end}, i + 1 : \text{end}]$ 
8 return  $V, \beta, A[1 : n, 1 : n]$ 

```

Algorithm 4: $QB \leftarrow \text{hh_mult}(V, B)$ Given a set of householder vectors $\{v_i\}_{i=1}^n$ and their corresponding constants $\{\beta_i\}_{i=1}^n$, compute $P_1 = P_n B$, where $P_i = I - \beta_i v_i v_i^T$.

Input: $V \in \mathbb{R}^{m \times n}$, $\beta \in \mathbb{R}^n$ where $m \geq n$, $B \in \mathbb{R}^{m \times d}$
Output: QB

```

1 for  $i = 1 : n$  do
2    $B \leftarrow B - \beta v_i (v_i^T B)$ 
3 return  $B$ 

```

normalizing them. However, this equivalence ($P_v = P_{\tilde{v}}$ for all $v^T \in \text{Span}(v)$) is not guaranteed due to rounding errors when using floating point numbers and operations. When using high precision floating point numbers such as double-precision floats, rounding errors that accumulate from the normalization of Householder vectors rarely and barely contribute to the overall stability of the HQR algorithm performed. In contrast, lower precision floating point numbers with limited dynamic range may be more sensitive to the un/normalization choice. For example, if we leave the Householder vectors unnormalized while using half-precision, it is possible to accumulate 1mf's in inner products of "large" vectors. As a result, picking a normalization scheme for v is important in low-precision calculations. Some methods and reasons for the normalization of v are as follows:

- Set the first element of v , v_1 , as 1 for efficient storage of many Householder vectors.
- Set the 2-norm of v to $\sqrt{2}$ to always have $\beta = 1$.
- Set the 2-norm of v to 1 to prevent extremely large values, and to always have $\beta = 2$.

The LINPACK implementation of the HQR factorization uses +14(-14) the first method of normalizing via setting v_1 to 1, and is shown in Algorithm 2. The first normalizing method adds an extra rounding error to β and v each, whereas the remaining methods incur no rounding error in forming $\beta = 1$ and 2 can be represented exactly.

3.2 Rounding Error Analysis

We present an error analysis for the HQR factorization where all inner products are performed with mixed-precision, and all other calculations are done in the storage precision, w .

Assumption 3.2 lays out the generalized mixed-precision inner product we will be using over and over again in the remainder of this paper.

Assumption 3.2. Let w , p , and s each represent floating point precisions for storage, product, and summation, where the varying precisions are defined by their unit round-off values denoted by u_w , u_p , and u_s , and we can assume $1 \gg u_w \gg u_s$, and $u_p \in (0, u_w]$. Within the inner product subroutine, products are done in precision p and summation is done in precision s , and the result stored in precision w . All operations other than inner products are done in the storage precision, w .

3.2.1 Error analysis for forming Householder Vector and Constant

Calculating the Householder vector and constant is a major routine for the HQR factorization.

Error analysis for v : In this section, we show how to bound the error when employing the mixed precision dot product procedure for Algorithm 2. We begin by extending the inner-product error shown in Lemmas 2.4 and 2.5 to the 2-norm error.

Lemma 3.3 (2-norm round-off error). Consider a mixed-precision scheme as is outlined in Assumption 3.2. Let $x \in \mathbb{F}_w^m$ be an arbitrary n -length vector stored in w precision. The forward error bound for computing the 2-norm of x is

$$\|(\|x\|_2)\| = (1 + \theta_w^{d+1}) \|x\|_2, \quad (27)$$

where $|\theta_w^{d+1}| \leq \gamma_w^{(d+1)} |x|$ for $d \in \{1, 2\}$ and $d := \lfloor \frac{(m-1)m}{u_w} \rfloor$.

There is no error incurred in evaluating the sign of a number or flipping the sign. Therefore, the error bound for computing $\sigma = -\text{sign}(x_1)\|x\|_2$ is exactly the same as that for the 2-norm.

$$\|(\sigma - \hat{\sigma})\|_2 = \|\text{sign}(x_1)\|_2 \|x\|_2 = \sigma + \Delta\sigma, \quad |\Delta\sigma| \leq 2^{-(t+1)}\|\sigma\|_2 \quad (2s)$$

Let \hat{v}_i be the approximate value v_i and $(\hat{v}_i = x_i - \sigma)$. We can now show the round-off error for \hat{v}_i and v_i where $i = 2, \dots, n$. Then the round-off errors for \hat{v}_i and v_i are

$$\|v_1\|_2 = v_1 = \hat{v}_1 + \Delta v_1$$

$$\|v_1 - \sigma\|_2 = (1 + \delta_n)(\sigma + \Delta\sigma) = (1 + \theta^{(1,2,t+1,2)})\hat{v}_1$$

and

$$\|v_i\|_2 = v_i = \hat{v}_i + \Delta v_i = (1 + \delta_n) \frac{\hat{v}_i + \Delta v_i}{x_i} = (1 + \theta^{(1,2,t+1,2)})\hat{v}_i.$$

The above equalities (as opposed to inequalities) are permitted since θ values are allowed to be flexible within the corresponding γ bounds.

Error analysis for β : Now we show the derivation of round-off error for the Householder computation.

$$\begin{aligned} \beta &= \beta - \left(\frac{v}{v_1} - (1 + \delta_n) \frac{\hat{v}_1 + \Delta v_1}{v_1} \right) \\ &= -(1 + \theta^{(1)}) \frac{(1 + \theta^{(2)})}{(1 + \theta^{(1,2,t+1,2)})v_1} \\ &= (1 + \theta^{(2,t+1,2,t+1,1)})\beta \\ &= (1 + \theta^{(2,t+1,2,t+1,2)})\beta. \end{aligned}$$

where $t = 1$ or $t = 2$ depending on which mixed-precision inner product procedure was used. These two results are formalized in Lemma 3.4 below.

Lemma 3.4. Given $x \in \mathbb{R}^m$, consider the constructions of $\beta \in \mathbb{R}$ and $v \in \mathbb{R}^m$ such that $P_v x = \sigma v$ (cf. Lemma 3.1) by using Algorithm 2. Then the forward error of forming v and β with the floating point arithmetic with the mixed-precision scheme outlined in Assumption 3.2 are

$$\|v\|_2 = (1 + \theta^{(1,2,t+1,2,t+1,2)})\|v\|_2, \text{ and } \beta = (1 + \theta^{(2,t+1,2,t+1,2)})\beta,$$

$$\text{where } t \in \{1, 2\} \text{ and } \theta = \frac{n}{m-1}.$$

3.2.2 Applying a Single Householder Transformation

Applying a Householder transformation is implemented by a series of inner and outer products, since Householder matrices are rank-1 updates of the identity. This is much less costly than forming P_v , then performing matrix-vector or matrix-matrix multiplications. For some $P_v = I - \beta vv^T$, we result in the following computation.

$$P_v x = (I - \beta vv^T)x = x - \beta v v^T x \quad (29)$$

13

14

Applying P_v to zero out the target column of a matrix. Let $x \in \mathbb{R}^m$ be the target column we wish to zero out (together with the first element). Recall that we chose a sparse v such that $P_v x = \sigma v$, As a result, the only error lies in the first element, σ , and that is shown in Equation 2s. Note that the normalization choice of v does not impact the Householder transformation matrix (P_v) nor its action on x , $P_v x$.

Applying P_v to the remaining columns of the matrix. Now, let x and v have no special relationship, as v was constructed given some pivoting column. Set $w = \beta v^T x$. Note that x is exact, whereas v and β were still computed with floating point operations. The errors incurred from computing v and β need to be included as well as the new rounding errors accumulating from the action of applying P_v to a column.

$$\|(\hat{v}^T x)\|_2 = (1 + \theta^{(1,t)})\|v + \Delta v\|_2^T x,$$

where $\theta^{(1,t)}$ is incurred from the action of a dot product.

$$= (1 + \theta^{(1,t+1,2,t+1,2)})\|v\|_2^T x$$

$$= (1 + \theta^{(1,t+1,2,t+1,2)})v^T x$$

Now we can form $\|w\|_2$.

$$w = (1 + \theta^{(2)})\beta(1 + \Delta\beta)(1 + \theta^{(1,t+1,2,t+1,2)})v^T xw$$

where $\theta^{(2)}$ results from multiplying β and $v^T x$ to w .

$$= (1 + \theta^{(2)})\beta(1 + \theta^{(1,t+1,2,t+1,2)})\beta(1 + \theta^{(1,t+1,2,t+1,2)})v^T xw$$

$$= (1 + \theta^{(2,t+1,2,t+1,2)})w$$

Finally, we can add in the vector subtraction operation and complete the rounding error analysis of applying a Householder transformation to any vector.

$$\|(\hat{P}_v x)\|_2 = \|x - w\|_2 = (1 + \delta_n)(1 + \theta^{(1,t+1,2,t+1,2)})w$$

$$= (1 + \theta^{(1,t+1,2,t+1,2)})\|P_v x\|_2 \quad (31)$$

$$= (1 + \delta_n)\|P_v x\|_2 \leq \gamma_{\text{float}(t,1)}^2 \quad (32)$$

Details behind the matrix norm error bound in Equation 32 is shown in the appendix. Constructing both Q and R relies on applying Householder transformations in the above two ways: 1) to zero out below the diagonal of a target column, and 2) to update the bottom right submatrix. We now have the tools to formulate the forward error bound on Q and R calculated from the HQR factorization.

3.2.3 HQR Factorization Forward Error Analysis

Consider a thin QR factorization where $A \in \mathbb{R}^{m \times n}$ for $m \geq n$, we have $Q \in \mathbb{R}^{m \times m}$ and $R \in \mathbb{R}^{n \times n}$. The pseudocode in Section 3 shows that each succeeding Householder transformation is applied to a smaller lower right submatrix each time.

Rounding error bound plots for the HQR Factorization

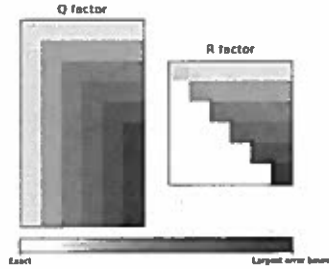


Figure 1: Grayscale representation of distribution of rounding errors bounds for the HQR algorithm.

$d = \lfloor \frac{(m-1)u_c}{u_m} \rfloor$, and $c = 1$ or $c = 2$. Then we have rowwise forward error bounds

$$\tilde{R} = R + \Delta R = P_n - P_1 A, \quad (33)$$

$$Q = Q + \Delta Q = P_1 - P_n I, \quad (34)$$

$$(35)$$

where

$$\|\Delta Q\|_F \leq n^{3/2} \gamma_w^{(6d+6c+13)}, \quad (36)$$

and for column j in $\{1, \dots, n\}$,

$$\|\Delta R[:, j]\|_2 \leq j \gamma_w^{(6d+6c+13)} \|A[:, j]\|_2. \quad (37)$$

We can also form a backward error. Let $A + \Delta A = \tilde{Q}\tilde{R}$, where Q and R are obtained via Algorithm 3. Then,

$$\|\Delta A\|_F \leq n^{3/2} \gamma_w^{(6d+6c+13)} \|A\|_F. \quad (38)$$

3.2.4 HQR Comparison to Uniform Precision Analysis

The mixed-precision segments of the analysis behind Theorem 3.5 derive from the mixed-precision inner product scheme outlined in Assumption 3.2, and is propagated to form the error bounds for a single Householder transformation as is shown in Equation 32. All steps to form the error bounds in Theorem 3.5 from the error bound for a single Householder transformation (Equation 32) directly

Consequently, rounding error bounds for each element of R and Q can be specifically computed by its location within the matrices, as is displayed in Figure 1 for a 10-by-6 example. Instead of continuing with a component-wise analysis of how accumulated rounding errors are distributed by HQR, we transition into normwise error analyses. To do this, we use the analysis in shown in the preceding section (summarized in Equation 32) to implicitly form the matrix norm error of the Householder transformation matrix, P_v . Then, we use the result of Lemma 3.7 in [3] to get a normwise bound on the perturbation effect of multiple matrix multiplications. This result is shown in Theorem 3.5, and the proof is shown extensively in A.3.

Theorem 3.5. Let $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ have full rank, n . Let $Q \in \mathbb{R}^{m \times m}$ and $R \in \mathbb{R}^{n \times n}$ be the thin QR factors of A obtained via the HQR algorithm with a mixed-precision scheme as is outlined in Assumption 3.2. Let

follow the analyses in Section 19.3 of [3]. In these steps, we generalize the single Householder transformation error bound,

$$\|(\tilde{P}_v x) - (P_v + \Delta P_v)x\|_F \leq \epsilon, \quad (39)$$

for some small quantity $0 < \epsilon \ll 1$, and propagate it through the for-loop in Algorithm 3. This then results in forward error bound coefficients m or $n^{3/2}\epsilon$. Since this ϵ value remains constant, the rounding error analysis for both mixed-precision and uniform-precision schemes are essentially the same with different values for ϵ . The uniform precision equivalent of Equation 32 is shown in Equation 40.

$$\|(\tilde{P}_v x) - (P_v + \Delta P_v)x\|_F \leq \zeta^{(m)}, \quad (40)$$

which is derived in detail in [3]. Therefore, we only need to compare $\gamma^{(6d+6c+13)}$ against $\zeta^{(m)}$, where ϵ is a small integer. Although d relies on both m and the precisions u and u_c , we can generally assume that $cm \gg (6d+6c+13)$ in most mixed-precision settings. Therefore, the new bounds in Theorem 3.5 are much tighter than the existing ones, and more accurately describe the kind of rounding error accumulated in mixed-precision computational settings.

4 Tall-and-Skinny QR

Some important problems that require QR factorizations of overdetermined systems include least squares problems, eigenvalue problems, low rank approximations, as well as other matrix decompositions. Although Tall-and-Skinny QR (TSQR) broadly refers to row-block QR factorization methods, we will discuss a specific variant of TSQR which is also known as the AllReduce algorithm [5]. In this paper, the TSQR/AllReduce algorithm refers to the most parallel variant of all row-block QR factorization algorithms discussed in [6]. A detailed description and rounding error analysis of this algorithm can be found in [5], and we present a pseudocode for the algorithm in Algorithm 5. Our initial interest in this algorithm came from its parallelizable nature, which is particularly suitable to implementation on GPUs and could lead to speed-ups. Additionally, our numerical simulations (discussed in Section 4.3) show that TSQR can not only increase the speed but also outperform the traditional HQR factorization in low precisions.

4.1 TSQR/AllReduce Algorithm

Algorithm 5 takes a tall-and-skinny matrix, A and organizes it into row-blocks. HQR factorization is performed on each of those blocks, and pairs of R factors are combined to form the next set of A matrices to be QR factorized. This is repeated until only a single R factor remains, and the Q factor is built from all of the Householder constants and vectors stored at each level. The most gains from parallelizing can be made in the initial level where the maximum number of independent HQR factorizations occur. For any given tall-and-skinny matrix, more than one configuration of this algorithm may be given. A 1000-by-100 matrix can be partitioned into 2, 4, 8, or 16 initial row-blocks, but a 1000-by-700 matrix can only be partitioned into 2 initial blocks. The choice in the initial partition may be given by the number of nodes available in the shape of the matrix, and determine the recursion level which we call L . Our numerical experiments show that this choice has an impact in the accuracy of the QR factorization.

A L -level refers to number of recursions in a particular TSQR implementation. An L -level TSQR algorithm partitions the original matrix into 2^L submatrices in the initial or 0^{th} level of the algorithm, and 2^{L-i} QR factorizations are performed in level i for $i = 1, \dots, L$. The set of matrices

E.g.,

we refer to each as the ...

but typically is determined by the # of nodes available or the shape of the matrix, e.g. 1000x100

that are QR factorized at each level i are called $A_i^{(j)}$ for $j = 1, \dots, 2^{L-i}$, where superscript i corresponds to the level, and the subscript j indexes the row blocks within level i . In the following sections, Algorithm 5 (tsqr) will find a TSQR factorization of a matrix $A \in \mathbb{R}^{m \times n}$ where $m \gg n$. The inline function qr refers to Algorithm 3, hh_mult is Algorithm 4, and we use Algorithm 2 as a subroutine of qr.

4.1.1 TSQR Notation

Due to the multi-level nature of the TSQR algorithm, we will introduce new notation to avoid confusion. In the final task of constructing Q , we need to aggregate the Q factors from each block at each level. We partition each $Q_i^{(j)}$ factor from level i and apply $Q_i^{(j-1)}$ factors from level $i-1$ to them. We write $Q_i^{(j)} = [\tilde{Q}_i^{(j,1)} \tilde{Q}_i^{(j,2)}]^T$ where the partition (approximately) halves $Q_i^{(j)}$. The functions $\alpha(j)$ and $\phi(j)$ are defined such that $Q_i^{(j)}$ is applied to $\tilde{Q}_{\alpha(j)\phi(j)}^{(i+1)}$. For $j = 1, \dots, 2^{L-i}$ at level i , we need $j = 2(\alpha(j)-1) + \phi(j)$ where $\alpha(j) = \lfloor \frac{j}{2} \rfloor$ and $\phi(j) = 2+j-2\alpha(j)$. Section 4.1.2 shows full linear algebra details for a single-level ($L = 1$, 2 initial blocks) example. The reconstruction of Q can be implemented more efficiently (c.f. [7]), but the reconstruction method in Algorithm 5 is presented for a clear, straightforward explanation.

4.1.2 Single-level Example

In the single-level version of this algorithm, we first bisect A into $A_1^{(0)}$ and $A_2^{(0)}$ and compute the QR factorization of each of these submatrices. We combine the resulting upper-triangular matrices i.e., $A_1^{(1)} = \begin{bmatrix} R_1^{(0)} \\ R_2^{(0)} \end{bmatrix}$ which we QR factorize and repeat.

$$A = \begin{bmatrix} A_1^{(0)} \\ A_2^{(0)} \end{bmatrix} = \begin{bmatrix} Q_1^{(0)} R_1^{(0)} \\ Q_2^{(0)} R_2^{(0)} \end{bmatrix} = \begin{bmatrix} Q_1^{(0)} & 0 \\ 0 & Q_2^{(0)} \end{bmatrix} \begin{bmatrix} R_1^{(0)} \\ R_2^{(0)} \end{bmatrix} = \begin{bmatrix} Q_1^{(0)} & 0 \\ 0 & Q_2^{(0)} \end{bmatrix} A_1^{(1)} = \begin{bmatrix} Q_1^{(0)} & 0 \\ 0 & Q_2^{(0)} \end{bmatrix} Q_1^{(1)} R.$$

Whereas the R factor of $A_1^{(1)}$ is the final R factor of the QR factorization of the original matrix, A , we still need to construct Q . Dissecting $Q_1^{(1)}$ into two submatrices ($\tilde{Q}_1^{(1,1)}$ and $\tilde{Q}_1^{(1,2)}$) allows us to write and compute the product more compactly.

$$Q = \begin{bmatrix} Q_1^{(0)} & 0 \\ 0 & Q_2^{(0)} \end{bmatrix} Q_1^{(1)} = \begin{bmatrix} Q_1^{(0)} & 0 \\ 0 & Q_2^{(0)} \end{bmatrix} \begin{bmatrix} \tilde{Q}_1^{(1,1)} \\ \tilde{Q}_1^{(1,2)} \end{bmatrix} = \begin{bmatrix} Q_1^{(0)} \tilde{Q}_1^{(1,1)} \\ Q_2^{(0)} \tilde{Q}_1^{(1,2)} \end{bmatrix}.$$

Algorithm 5 takes a tall-and-skinny matrix A and level L and finds a QR factorization by initially partitioning A into 2^L row blocks. This initial includes the building of the final Q .

4.2 TSQR Rounding Error Analysis

The TSQR algorithm presented in Algorithm 5 is a divide-and-conquer strategy for the QR factorization that uses the HQR for the subproblems. Divide-and-conquer methods can naturally be implemented in parallel and accumulate less rounding errors. For example, the single-level TSQR decomposition of a tall-and-skinny matrix, A , requires $O(n^2 \log_2(\frac{m}{n}))$ FLOPs of matrices of sizes $\lfloor \log_2(\frac{m}{n}) \rfloor$ -by- n , $\lfloor \log_2(\frac{m}{n}) \rfloor$ -by- n , and $2n$ -by- n . The single-level TSQR strictly uses more FLOPs, but the dot

Algorithm 5: $Q, R = \text{tsqr}(A, L)$. Finds a QR factorization of a tall, skinny matrix, A .

Input: $A \in \mathbb{R}^{m \times n}$ where $m \gg n$, $L \leq \lfloor \log_2(\frac{m}{n}) \rfloor$, and 2^L is the initial number of blocks.
Output: $Q \in \mathbb{R}^{m \times n}$, $R \in \mathbb{R}^{n \times n}$ such that $QR = A$.

```

1  $h \leftarrow \lfloor \frac{m}{2^L} \rfloor$  // Number of rows for all but the last block
2  $r \leftarrow m - (2^L - 1)h$  // Number of rows for the last block
3 // Split  $A$  into  $2^L$  blocks, each that have  $h$  rows
4 for  $j = 1 : 2^L - 1$  do
5    $A_j^{(0)} \leftarrow A[(j-1)h + 1 : jh, :]$ 
6  $A_{2^L}^{(0)} \leftarrow A[(2^L-1)h + 1 : m, :]$  // Last block may have more rows
7 // Store Householder vectors in columns of matrix  $V$ , Householder exponents
  in components of vector  $\beta$ , and set up the partitioning
8 for  $i = 0 : L - 1$  do
9   // The inner loop can be parallelized
10  for  $j = 1 : 2^{L-i}$  do
11     $V_{2j-1}^{(i)}, \beta_{2j-1}^{(i)}, R_{2j-1}^{(i)} \leftarrow \text{qr}(A_{2j-1}^{(i)})$ 
12     $V_{2j}^{(i)}, \beta_{2j}^{(i)}, R_{2j}^{(i)} \leftarrow \text{qr}(A_{2j}^{(i)})$ 
13     $A_j^{(i+1)} \leftarrow \begin{bmatrix} R_{2j-1}^{(i)} \\ R_{2j}^{(i)} \end{bmatrix}$ 
14  // At the bottom-most level, get the final  $R$  factor
15   $V_1^{(L)}, \beta_1^{(L)}, R \leftarrow \text{qr}(A_1^{(L)})$ 
16   $Q_1^{(L)} \leftarrow \text{hh\_mult}(V_1^{(L)}, I_{2n \times n})$ 
17  // Compute  $Q$  by applying stored Householder vectors,  $A$ , and  $Q_1^{(L)}$ 
18  for  $i = L - 1 : 1$  do
19    for  $j = 1 : 2^{L-i}$  do
20       $Q_j^{(i)} \leftarrow \text{hh\_mult}(V_j^{(i)}, \begin{bmatrix} \tilde{Q}_{\alpha(j)\phi(j)}^{(i+1)} \\ 0_{n \times n} \end{bmatrix})$ 
21  // At the top-most level, construct the final  $Q$  by the final  $Q$  factors
22   $Q \leftarrow []$ 
23  for  $j = 1 : 2^L$  do
24     $Q \leftarrow [\text{hh\_mult}(V_j^{(0)}, \begin{bmatrix} \tilde{Q}_{\alpha(j)\phi(j)}^{(1)} \\ 0_{h \times n} \end{bmatrix}) \quad Q]$ 
25 return  $Q, R$ 

```

Way too much remove
we can remove them where you can.
We is typically best used only in the numerical section

new will introduce new notation

I haven't deck notation...

however the final Q still need to be constructed

More generally

why only 3? why give invitation.

product subroutines may accumulate smaller rounding errors and certainly have smaller upper bounds since they are performed on shorter vectors and lead to a more accurate solution overall. These concepts are elaborated in [9], where the rounding error analysis of TSQR is shown in detail in [9]. We summarize the main results in Theorem 4.1.

Theorem 4.1. Let $A \in \mathbb{R}^{m \times n}$ with $m \geq n$ and $Q \in \mathbb{R}^{m \times m}$ and $R \in \mathbb{R}^{n \times n}$ be the thin QR factors of A obtained via Algorithm 3. Then we have rigorous forward error bounds

$$A = A + \Delta A = QR + \Delta R, \\ Q = Q + \Delta Q,$$

where

$$\|\Delta R\|_F \leq \left[n^2 \frac{\kappa}{2} + (1 + n^2 \frac{\kappa}{2}) \left\{ (1 + n^2 \kappa)^L - 1 \right\} \right] \|A\|_F, \\ \|\Delta Q\|_F \leq \sqrt{n} \left[(1 + n^2 \frac{\kappa}{2}) (1 + n^2 \kappa)^L - 1 \right]. \quad (42)$$

Furthermore, if we assume $n^2 \frac{\kappa}{2}, n^2 \kappa \ll 1$, the coefficient for $\|A\|_F$ in Equations 41 can be approximated as

$$\left[n^2 \frac{\kappa}{2} + (1 + n^2 \frac{\kappa}{2}) \left\{ (1 + n^2 \kappa)^L - 1 \right\} \right] \approx n^2 \frac{\kappa}{2} + Ln^2 \kappa, \quad (43)$$

and the right hand side of Equation 42 can be approximated as

$$\sqrt{n} \left[(1 + n^2 \frac{\kappa}{2}) (1 + n^2 \kappa)^L - 1 \right] \approx \sqrt{n} \left(n^2 \frac{\kappa}{2} + Ln^2 \kappa \right) \quad (44)$$

We can also form a backward error, where $A + \Delta A \text{TSQR} = QR$ and both Q and R are obtained via Algorithm 3. Then,

$$\|\Delta A \text{TSQR}\|_F = \|Q\Delta R + \Delta Q R\|_F \approx \sqrt{n} \left(n^2 \frac{\kappa}{2} + Ln^2 \kappa \right) \|A\|_F \quad (45)$$

In Section 3.2.4, we discussed how the steps of the HQR algorithm result as some multiple of ϵ , where ϵ is the error bound for a single Householder transformation and is described in Equation 39.

Similarly, the analysis behind Theorem 4.1 can be generalized via defining ϵ_1 to be the error bound for a single Householder transformation corresponding to the vector length at the initial level 0, and ϵ_2 to be the error bound for a Householder transformation corresponding to vector length in all deeper levels, $2n$. This generalization leads to the error bound coefficients

$$m_1 + Ln_2 \quad \text{for } \|\Delta Q\|_F, \|\Delta A \text{TSQR}\|_F, \\ \sqrt{n}(m_1 + Ln_2) \quad \text{for } \|\Delta R\|_F, \|\Delta A\|_F. \quad (46)$$

and in a uniform precision setting, these correspond to

$$\epsilon_1 = \gamma_{\text{HQR}}^{(m)} \quad \text{and} \quad \epsilon_2 = \gamma_{\text{HQR}}^{(2n)}, \quad (48)$$

and in the mixed precision setting outlined in Assumption 3.2, they correspond to

$$\epsilon_1 = \gamma_{\text{HQR}}^{(m)} \epsilon_{\text{HQR}}^{(1)} \quad d_1 = \left[\frac{m}{2} - 1 \right] \epsilon_{\text{HQR}}^{(1)} \quad \text{and} \quad (49)$$

$$\epsilon_2 = \gamma_{\text{HQR}}^{(2n)} \epsilon_{\text{HQR}}^{(1)} \quad d_2 = \left[\frac{2n}{2} - 1 \right] \epsilon_{\text{HQR}}^{(1)}. \quad (50)$$

where $d_1 \approx \delta d_2 = 1$ respectively

In both settings, we see that increasing L may decrease ϵ_{HQR} still increase the overall bound, the larger L , still could have an adverse effect on the coefficients in Theorem 4.1. This trade-off is precisely the balance between the size of initial blocks and the number of levels in the TSQR algorithm, and an optimal TSQR scheme would ideally minimize ϵ_1 and ϵ_2 with the choice of L . These error bounds are stated in detail in Section 4.2.1 following Section.

following Section.

4.2.1 HQR and TSQR error bound comparison

We compare the error bounds for HQR and TSQR algorithms. Let's consider the larger error bounds in the uniform precision equivalents of Theorems 4.5 and 4.1, which are the bounds of ΔQ and ΔA . In order for the a meaningful TSQR error bound to outperform the bound for the HQR algorithm, we need integers $m, n > 0$, and $L \geq 0$ such that,

$$1 \gg n^{1/2} \gamma_{\text{HQR}}^{(m)} \gg n^{1/2} \gamma_{\text{TSQR}}^{(2n)} + L \gamma_{\text{HQR}}^{(2n)}$$

If we assume $\frac{m}{2} = 2n$, the HQR bound is $\frac{L+1}{2}$ larger than the bound for TSQR with L levels. For example, in single precision, a HQR of a $2^{15} \times 2^{15}$ matrix results in an upper bound relative back-

ward error $\|A - QR\|_F / \|A\|_F$ of ≈ 1.002 , but a TSQR with $L = 8$ is bounded by $\approx 3.516e-02$. This exemplifies a case in which stability is not guaranteed in HQR, but is guaranteed in TSQR. Even in the worst-case scenario, we can consider $2^{20} \times 2^{12}$ matrix, HQR factorization is performed with double precision. The error bound for HQR is $1.686e-7$ whereas the error bound for TSQR with 12 levels, which is $5.351e-10$. In general, we can conjecture that values of L that can make $m/2$ and $2Ln$ much smaller than m , should produce a TSQR that outperforms HQR in worst-case scenarios, at least in uniform precision settings. However, the range of matrix sizes that TSQR can accommodate decreases as L grows larger. Figure 2 shows the matrix sizes HQR, 2-level TSQR and 4-level TSQR can accommodate as well as their respective error bounds.

remains stable.

Let's consider a mixed-precision setting such as in Assumption 3.2, and we assume $m_p = m$, so that $\epsilon_2 = 2$. In order for the a meaningful TSQR error bound to outperform the bound for the HQR algorithm, we now need integers $m, n > 0$, and $L \geq 0$ such that,

$$1 \gg n^{1/2} \gamma_{\text{HQR}}^{(m)} \epsilon_{\text{HQR}}^{(1)} \gg n^{1/2} \gamma_{\text{TSQR}}^{(2n)} + L \gamma_{\text{HQR}}^{(2n)} \epsilon_{\text{HQR}}^{(1)}$$

where $d = \left\lfloor \frac{m-1}{2} \right\rfloor$, $d_1 = \left\lfloor \frac{m}{2} - 1 \right\rfloor \epsilon_{\text{HQR}}^{(1)}$ and $d_2 = \left\lfloor \frac{2n-1}{2} \right\rfloor \epsilon_{\text{HQR}}^{(1)}$.

In contrast to the analysis for uniform precision settings, larger L values do not necessarily reduce the error bounds of TSQR. While large L can imply $m \gg m/2$ and $2Ln$, it doesn't always lead to $d \gg d_1 + Ld_2$. Although the theoretical error bounds do not give a clear indication of the worst-case performances of HQR and TSQR in mixed-precision settings, our numerical simulations show that TSQR can still outperform HQR in ill-conditioned matrices. These experiments are discussed in detail in the next section.

4.3 Numerical Experiment

In Section 3.2.4, we showed that conditions exist where TSQR could outperform HQR, and that these conditions were hard to identify in mixed-precision settings. An empirical comparison of these two QR factorization algorithms in double precision can be found in [9] where they conclude that deeper TSQR tends to produce more accurate QR factorizations than HQR. However, using TSQR with deep levels (large L) can actually start to perform worse than TSQR with shallow

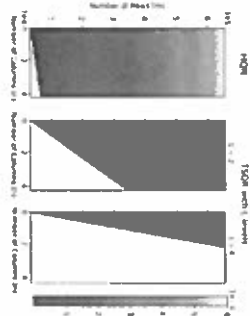


Figure 2: Non-white space indicates allowable matrix sizes for each scheme and value represents error bounds for $\|A\|_F$ in double precision. (for uniform precision only?) include in figure description

4.3.1 Experiment Details

We used Julia v1.0.4 for all of the numerical simulations. The programming language allows half precision storage as well as castup and castdown operations to and from single and double precision, but has no half precision arithmetic. Therefore, we relied on using Algorithm 1 for $f \in \text{OPU}$ (dot-product) to simulate half and mixed-precision arithmetic operations. For HQR, we created a mixed precision version of the LAPACK routine xGECQRF where the dot product subroutine was approximated by $\text{dot}(x, y) \approx \text{dot}(x, y) \cdot 2^{\text{exp}(\text{exponent}(x))}$ to simulate the mixed-precision setting described in Assumption 3.2 with $\alpha_p = 0$ (which implies $\alpha = 1$), and used Algorithm 1 on all other basic operations in OP to simulate half-storage precision arithmetic. This HQR was then used as a subroutine of TSQR as well. All in all, our experiments nearly replicated the mixed-precision setting we assigned for the error analysis in Sections 3.2 and 4.2.

We generated random matrices of size 4000×100 and computed their HQR and TSQR for $L = 1, \dots, 6$ in a mixed-precision setting that simulates Assumption 3.2 with $\alpha = 1$. The relative backward error $\|QR - A\|_F / \|A\|_F$ was computed by casting up Q , R , and A to double precision to compute $\|QR - A\|_F / \|A\|_F$. The mixed-precision HQR error bounds $n^2 \epsilon_{\text{HQR}}^{\text{mixed}}$ and $n^2 \epsilon_{\text{TSQR}}^{\text{mixed}}$ for $m = 4000$ and $n = 100$ are 0.586 and 9.364 respectively, and the mixed-precision TSQR bounds for $L = 1, \dots, 5$ are even larger. This indicates that our error bounds do not guarantee stability.

4.3.2 Results

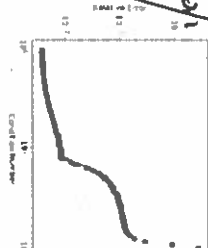


Figure 3: HQR errors for matrices with varying condition numbers.

and all levels of TSQR are stable, but degrading the errors. Another trend occurs for matrices with higher condition numbers, where single-level and 2-level TSQR yield smaller errors than HQR. In these cases, TSQR with 4 or more levels have errors similar to or worse than 2-level TSQR, but their errors tend to not rise above the HQR errors.

These results suggest that TSQR can significantly outperform HQR even in mixed-precision settings, and particularly when HQR is unstable due to high condition numbers. Although this experiment focused on condition numbers, identifying other properties that point to better performance of TSQR than HQR can further broaden the potential use of mixed-precision TSQR in applications.

5 Applications

Many applications in scientific computing typically employ double precision when lower precision may actually be sufficient. However, identifying which applications can tolerate lower precision and still achieve sufficient results is a challenging task that wasn't previously relevant prior to the emergence of new technology that supports low precision. Since low and mixed precision settings benefit from speedups of higher storage applications that possess large amounts of data

Can it make sense to use this as expected on new seeds?

$$A_n = \frac{Q(nE + I)}{\|Q(nE + I)\|_F} \quad (11)$$

where $Q \in \mathbb{R}^{m \times n}$ is a random orthogonal matrix and $E \in \mathbb{R}^{n \times n}$ is the matrix of 1's. The random orthogonal matrix Q is generated by taking a QR factorization of an iid 4000×100 matrix sampled from $\text{Exp}(1/n)$, and we used the built-in QR factorization function in Julia. By construction A_n has 2-norm condition number not $n + 1$. By varying n from $1e-4$ to 1, we varied the condition number from 1.1 to 101, and we generated 10 samples for each value of n .

We computed the relative backward error, $\frac{\|QR - A\|_F}{\|A\|_F}$, of the traditional HQR algorithm (Figure 3) as well as the errors of TSQR with $L = 1, \dots, 5$ (Figure 4). Figure 3 shows the errors for HQR plotted against the theoretical condition number, and Figure 4 shows the errors for HQR and TSQR for our test matrices where $L = 0$ is the HQR error bound. The 5 correspond to TSQR with L levels. The black markers on the vertical axis of Figure 4 correspond to the HQR errors shown in Figure 3.

Figure 3 shows that the error of the QR factorization worsens as the condition number increases, and there is a non-monotonically increasing trend between errors and condition numbers. Figure 4 reveals two different trends for the errors as we compare HQR to TSQR with only one level, then to deeper levels of TSQR. One trend occurs for matrices with smaller condition numbers, where HQR

not really true, be careful. Flops were more expensive that data storage so it wasn't considered. (continued)

Figure 4 shows that the error of the QR factorization worsens as the condition number increases, and there is a non-monotonically increasing trend between errors and condition numbers. Figure 4 reveals two different trends for the errors as we compare HQR to TSQR with only one level, then to deeper levels of TSQR. One trend occurs for matrices with smaller condition numbers, where HQR

which are derived from
Block stochastic matrices.

eigenspaces, DISCAN to partition the embeddings of the nodes onto these eigenspaces, and precision and recall to evaluate clustering performance. We used a static graph of 2000 nodes with 10 known true partitions for the Graph Challenge [10]. The graphs we used were undirected and unweighted; the only elements in the adjacency matrices were 0's and 1's which can easily be represented in half, single, and double precision floats. For $r = 1, \dots, 10$, let $Y_{\text{node},r} \in \mathbb{R}^{n \times r}$ be the half precision storage of the r^{th} random matrix. Since any half precision float can be exactly represented in single and double precisions, $Y_{\text{node},r}$'s can be easily cast up to single and double precisions, $Y_{\text{node},r}$ and $Y_{\text{node},r}$. We performed mixed-precision HQB within subspace iteration initialized by $Y_{\text{node},r}$ and uniform precision HQB for subspace iteration initialized by $Y_{\text{node},r}$'s and $Y_{\text{node},r}$'s. For trial $i = 1, \dots, 10$, we repeated the following steps:

Step 1. Identifying an orthogonal basis of dimension 10 with subspace iteration using the appropriate HQB routine for $Y_{\text{node},r}$, $Y_{\text{node},r}$, and $Y_{\text{node},r}$.

Step 2. Apply DISCAN to the output matrices of previous step to cluster most nodes into communities, and the remaining nodes as outliers.

Step 3. ~~error~~ ^{using} we measure clustering performances of DISCAN on the three different precision subspace iteration embeddings with precision and recall.

Subspace Iteration Results. Figure 5 shows the eigenspace error, $\|Y - QQ^T Y\|_F / \|Y\|_F$, from the subspace iteration step of our trial. The r values were set to 10, single and double for the uniform precision HQB, and 10, single and double for the mixed-precision HQB, and the solid lines are plotted to show the unit round-off values. The uniform precision implementations of subspace iterations reached their stopping criterion set by r , and the mixed-precision implementation fluctuated close to r but ~~could not reach it~~. The convergence rate was approximately the same across the three different implementations, which suggests that the lower precision routines (mixed-precision HQB or uniform single precision HQB) can be used as proxy preconditioners to the high precision routine. In addition, if double precision eigenspace error is not necessary to achieve sufficient clustering results, we can simply use the lower precision HQB subspace iterations ~~as preconditioners~~.

Clustering Results. Table 5 shows the precision and recall results from the 10 trials for each-precision implementation. The DISCAN algorithm and the calculation of precision and recall were computed in double precision, and the variance in precision and recall values for these 10 trials were in the range of $1e-6$. The results show a discernible difference in the precision and recall, which suggests that a lower precision HQB within subspace iteration can still lead to a sufficiently accurate clustering. Possible future works include trying different clustering methods for Step 3. *I don't think this is necessary*

HQB scheme	Precision	Recall
mixed-precision	0.9822	0.9393
single-precision	0.9817	0.9407
double-precision	0.9822	0.9405

Table 5. Precision and recall values for 10 trials of DISCAN on graph with 2000 nodes and 10 true clusters.

Did you also redo SI for each trial?
If so, it was lost on me.

For context....

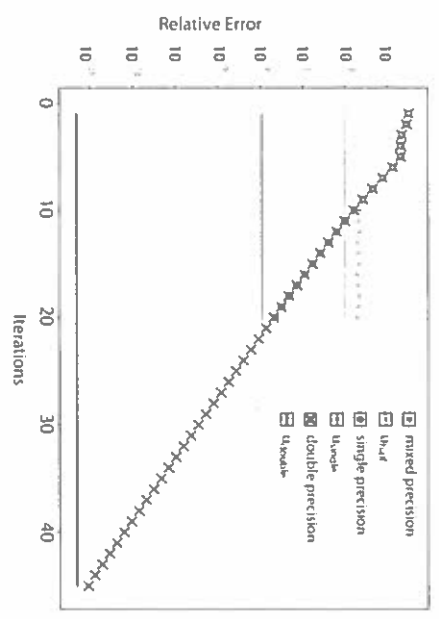


Figure 5. Eigenspace Error for subspace iteration with using double, single, and half precision. Traditional Householder QR factorizations.

6 Conclusion

Development of GPUs that support low precision floating point arithmetic has led to interest in half precision and mixed-precision algorithms that demonstrate speedier times, lower energy consumption, and lower memory usage. Less in precision, stability, and representable range compensate for these advantages, although these shortcomings may have little to no impact in some applications. It may even be possible to mitigate around these drawbacks with algorithmic design. In addition, the existing rounding error analysis cannot accurately bound the behavior of mixed-precision arithmetic.

We have developed a new framework for mixed-precision rounding error analysis and applied it to HQB, a widely used linear algebra routine, and implemented it in an iterative eigensolver in the context of spectral clustering. The mixed-precision error analysis builds from the many product routine, which can be applied to many other linear algebra tasks as well. The new error bounds more accurately describe how rounding errors are accumulated in mixed-precision settings. We also found that TSQR, a communication-avoiding, easily parallelizable QR factorization algorithm for tall-and-skinny matrices, can outperform HQB in mixed-precision settings for ill-conditioned, extremely overdetermined cases, which suggests that some algorithms are more robust against lower precision arithmetic. As QR factorizations of tall-and-skinny matrices are common in spectral clustering,

Given a few more sentences of why your work changes the way we should think about Slab fitting.

we experimented with introducing mixed-precision settings into graph partitioning problems. In particular, we applied DISCAN to the spectral basis of a graph identified via subspace iteration that used our simulated mixed-precision HQR, which yielded interesting results. Confusion to results from employing double-precision entirely. Experiments are needed to test larger, more ill-conditioned problems with different mixed-precision settings, and to explore other divide-and-conquer methods like TSQR that can harness parallel capabilities of GPUs while withstanding lower precisions.

A Proofs

A.1 Lemma 2.4 (Equation 19)

Proof. We wish to round up to the lower precision, μ , since $1 \gg n_p \gg u$. Recall that $d := \lfloor k_p u / n_p \rfloor$ and $r \leq \lfloor n_p / n_s \rfloor$ and note

$$k_p n_p + k_s u_s = (k_p + d)n_p + r u_s \leq (k_p + d + 1)n_p$$

$$\begin{aligned} \gamma_p^{(k_p)} + \gamma_s^{(k_s)} &= \gamma_p^{(k_p)} \gamma_s^{(k_s)} = \frac{k_p n_p}{1 - k_p n_p} + \frac{k_s u_s}{1 - k_p u_s} + \frac{k_p n_p}{k_p n_p + k_s u_s} + \frac{k_p n_p}{1 - k_p u_s} \frac{k_s u_s}{k_p n_p + k_s u_s} \\ &= \frac{1 - (k_p n_p + k_s u_s) + k_p k_s n_p u_s}{1 - (k_p + d + 1)n_p + k_p k_s n_p u_s} \\ &\leq \frac{(k_p + d + 1)n_p}{1 - (k_p + d + 1)n_p} = \gamma_p^{(k_p + d + 1)} \\ &\leq \gamma_p^{(k_p + d + 1)} \end{aligned}$$

□

A.2 Inner Products

A.2.1 Lemma 2.4

Let δ_p and δ_s be rounding error incurred from products and summations, and are bounded by $|\delta_p| < n_p$ and $|\delta_s| < u_s$, following the notation in [5]. Let α_i denote the k^{th} partial sum, and let α_i denote the floating point representation of the calculated α_i

$$\begin{aligned} \alpha_1 &= \mathbf{H}(\mathbf{x}_1 \mathbf{y}) = \mathbf{x}_1 \mathbf{y}(1 + \delta_p, 1) \\ \alpha_2 &= \mathbf{H}(\alpha_1 + \mathbf{x}_2 \mathbf{y}) \\ &= \mathbf{x}_1 \mathbf{y}(1 + \delta_p, 1) + \mathbf{x}_2 \mathbf{y}(1 + \delta_p, 2)(1 + \delta_s, 1) \\ \alpha_3 &= \mathbf{H}(\alpha_2 + \mathbf{x}_3 \mathbf{y}) \\ &= (\mathbf{x}_1 \mathbf{y}(1 + \delta_p, 1) + \mathbf{x}_2 \mathbf{y}(1 + \delta_p, 2)(1 + \delta_s, 1) + \mathbf{x}_3 \mathbf{y}(1 + \delta_p, 3)(1 + \delta_s, 2)) \end{aligned}$$

We can see a pattern emerging. The error for a general length m vector dot product is that:

$$\alpha_m = (\mathbf{x}_1 \mathbf{y} + \mathbf{x}_2 \mathbf{y}) \prod_{i=1}^{m-1} (1 + \delta_{p,i}) + \sum_{i=1}^m \mathbf{x}_i \mathbf{y} (1 + \delta_{p,i}) \left(\prod_{j=i+1}^{m-1} (1 + \delta_{s,j}) \right) \quad (54)$$

27

where each occurrence of δ_p and δ_s are distinct, but are still bound by n_p and u_s .

Using Lemma 2.1, we further simplify

$$\begin{aligned} \mathbf{H}(\mathbf{x}^T \mathbf{y}) &= \alpha_m = (1 + \theta_p^{(1)}) (1 + \theta_p^{(m-1)}) \mathbf{x}^T \mathbf{y} \\ &= (\mathbf{x} + \Delta \mathbf{x})^T \mathbf{y} = \mathbf{x}^T (\mathbf{y} + \Delta \mathbf{y}) \end{aligned}$$

Here $\Delta \mathbf{x}$ and $\Delta \mathbf{y}$ are vector perturbations.

By using Lemma 2.3 equation 19, we can bound the perturbations componentwise. Let $d = \lfloor \frac{(m-1)n_p}{m} \rfloor$ such that $(m-1)n_p = dn_p + ru_s$,

$$\begin{aligned} |\Delta \mathbf{x}| &\leq \gamma_p^{(d+12)} |\mathbf{x}| \\ |\Delta \mathbf{y}| &\leq \gamma_p^{(d+12)} |\mathbf{y}| \end{aligned}$$

Furthermore, these bounds lead to a forward error result as shown in Equation 55.

$$\|\mathbf{x}^T \mathbf{y} - \mathbf{H}(\mathbf{x}^T \mathbf{y})\| \leq \gamma_p^{(d+12)} \|\mathbf{x}\|^T \|\mathbf{y}\| \quad (55)$$

A.2.2 Corollary 2.5

This proof follows similarly to the proof for Lemma 2.4. Since no error is incurred in the multiplication portion of the inner products, δ_p and δ_s are rounding error incurred from summations and storage. The error for a general m -length vector dot product is then,

$$\alpha_m = (\mathbf{x}_1 \mathbf{y} + \mathbf{x}_2 \mathbf{y}) \prod_{i=1}^{m-1} (1 + \delta_{s,i}) + \sum_{i=1}^m \mathbf{x}_i \mathbf{y} \left(\prod_{j=i+1}^{m-1} (1 + \delta_{s,j}) \right) \quad (56)$$

Using Lemma 2.1, we further simplify

$$\begin{aligned} \mathbf{H}(\mathbf{x}^T \mathbf{y}) &= \alpha_m = (1 + \theta_s^{(m-1)}) \mathbf{x}^T \mathbf{y} \\ &= (\mathbf{x} + \Delta \mathbf{x})^T \mathbf{y} = \mathbf{x}^T (\mathbf{y} + \Delta \mathbf{y}) \end{aligned}$$

Here $\Delta \mathbf{x}$ and $\Delta \mathbf{y}$ are vector perturbations.

By using Lemma 2.4 equation 19, we can bound the perturbations componentwise. Let $d = \lfloor \frac{(m-1)u_s}{m} \rfloor$ such that $(m-1)u_s = du_s + ru_p$,

$$\begin{aligned} |\Delta \mathbf{x}| &\leq \gamma_p^{(d+11)} |\mathbf{x}| \\ |\Delta \mathbf{y}| &\leq \gamma_p^{(d+11)} |\mathbf{y}| \end{aligned}$$

Furthermore, these bounds lead to a forward error result as shown in Equation 57.

$$\|\mathbf{x}^T \mathbf{y} - \mathbf{H}(\mathbf{x}^T \mathbf{y})\| \leq \gamma_p^{(d+11)} \|\mathbf{x}\|^T \|\mathbf{y}\| \quad (57)$$

A.3 Proof for Mixed-Precision HQR result

Here, we show a few results that are necessary for the proof for Theorem 3.5. Lemma A.1 shows normwise results for a single mixed-precision Householder transformation performed on a vector, and Lemma A.2 builds on Lemma A.1 to show normwise results for multiple mixed-precision Householder transformations on a vector. We build column-wise results for HQR based on these lemmas, then compute the matrix norms at the end.

28

Lemma A.1. Let $\mathbf{x} \in \mathbb{R}^m$ and consider the computation of $\mathbf{y} = \mathbf{P}^T \mathbf{x} = \mathbf{x} - \beta \mathbf{v} \mathbf{v}^T \mathbf{x}$, where \mathbf{v} has accumulated error shown in Lemma 3.4. Then, the computed \mathbf{y} satisfies

$$\mathbf{y} = (\mathbf{P} + \Delta \mathbf{P}) \mathbf{x}, \quad \|\Delta \mathbf{P}\|_F \leq \gamma_{\text{acc}}^{(6d+13)} \quad (36)$$

where $\mathbf{P} = \mathbf{I} - \beta \mathbf{v} \mathbf{v}^T$ is a Householder transformation.

Proof. Recall that the computed \mathbf{y} accumulates component-wise error shown in Equation 32. Even though we don't explicitly form \mathbf{P} , forming the normwise error bound for this matrix makes the analysis simple. First, recall that any matrix \mathbf{A} with rank r has the following relations between its 2-norm and Frobenius norm,

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F \leq \sqrt{r} \|\mathbf{A}\|_2. \quad (35)$$

Then, we have

$$\|\mathbf{y}\|_2 = \|\mathbf{P} \mathbf{x}\|_2 \leq \|\mathbf{P}\|_2 \|\mathbf{x}\|_2 = \|\mathbf{x}\|_2 \quad (34)$$

since \mathbf{P} is orthogonal and $\|\mathbf{P}\|_2 = 1$. We now transition from the componentwise error to normwise error for $\Delta \mathbf{y}$ and write $\hat{\mathbf{z}} = \mathbf{e}_d + \mathbf{e}_2 + \mathbf{e}_3$.

$$\|\Delta \mathbf{y}\|_2 = \left(\sum_{i=1}^m \Delta y_i^2 \right)^{1/2} \leq \gamma_{\text{acc}}^{(13)} \left(\sum_{i=1}^m y_i^2 \right)^{1/2} = \gamma_{\text{acc}}^{(13)} \|\mathbf{y}\|_2 \quad (41)$$

Combining Equations (40) and (41), we result in

$$\frac{\|\Delta \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \leq \gamma_{\text{acc}}^{(13)}. \quad (42)$$

Now, notice that $\Delta \mathbf{P}$ is exactly $\frac{1}{\mathbf{x}^T \mathbf{x}} \Delta \mathbf{y} \mathbf{x}^T$.

$$\begin{aligned} (\mathbf{P} + \Delta \mathbf{P}) \mathbf{x} &= (\mathbf{P} + \frac{1}{\mathbf{x}^T \mathbf{x}} \Delta \mathbf{y} \mathbf{x}^T) \mathbf{x} \\ &= \mathbf{P} \mathbf{x} + \frac{\mathbf{x}^T \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \Delta \mathbf{y} = \mathbf{y} + \Delta \mathbf{y} \end{aligned}$$

We can compute the Frobenius norm of $\Delta \mathbf{P}$ by using $\Delta \mathbf{P} \mathbf{x} = \frac{1}{\mathbf{x}^T \mathbf{x}} \Delta \mathbf{y} \mathbf{x}$,

$$\begin{aligned} \|\Delta \mathbf{P}\|_F &= \left(\sum_{i=1}^m \sum_{j=1}^m \Delta P_{ij}^2 \right)^{1/2} = \left(\sum_{i=1}^m \sum_{j=1}^m \left(\frac{1}{\|\mathbf{x}\|_2^2} \Delta y_i x_j \right)^2 \right)^{1/2} \\ &= \left(\frac{1}{\|\mathbf{x}\|_2^2} \sum_{i=1}^m \sum_{j=1}^m \Delta y_i^2 x_j^2 \right)^{1/2} = \frac{1}{\|\mathbf{x}\|_2} \left(\sum_{i=1}^m \Delta y_i^2 \left(\sum_{j=1}^m x_j^2 \right) \right)^{1/2} \\ &= \frac{1}{\|\mathbf{x}\|_2} \left(\|\mathbf{x}\|_2^2 \sum_{i=1}^m \Delta y_i^2 \right)^{1/2} = \frac{\|\mathbf{x}\|_2 \|\Delta \mathbf{y}\|_2}{\|\mathbf{x}\|_2} = \frac{\|\Delta \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \end{aligned}$$

Finally, using Equation 42, we result in

$$\|\Delta \mathbf{P}\|_F \leq \gamma_{\text{acc}}^{(13)}. \quad (43)$$

□

Lemma A.2. Consider applying a sequence of transformations on the set $\{\mathbf{P}_i\}_{i=1}^d \subseteq \mathbb{R}^{m \times m}$ to $\mathbf{x} \in \mathbb{R}^m$, where $\mathbf{P}_i, \hat{\mathbf{y}}$ are all Householder transformations, and we will assume that $\gamma_{\text{acc}}^{(4)} < \frac{1}{2}$

$$\mathbf{y} = \mathbf{P}_1 \mathbf{P}_2 \dots \mathbf{P}_d \mathbf{x} = \mathbf{Q}^T \mathbf{x}$$

Then, $\mathbf{y} = (\mathbf{Q} + \Delta \mathbf{Q})^T \mathbf{x}$, where

$$\|\Delta \mathbf{Q}\|_F \leq \gamma_{\text{acc}}^{(4)} \quad (44)$$

$$\|\Delta \mathbf{y}\|_2 \leq \gamma_{\text{acc}}^{(4)} \|\mathbf{y}\|_2. \quad (45)$$

In addition, if we let $\mathbf{y} = \mathbf{Q}^T (\mathbf{x} + \Delta \mathbf{x})$, then

$$\|\Delta \mathbf{x}\|_2 \leq \gamma_{\text{acc}}^{(4)} \|\mathbf{x}\|_2 \quad (46)$$

Proof. As was for the proof for Lemma A.1, we know $\Delta \mathbf{Q}^T = \frac{1}{\mathbf{x}^T \mathbf{x}} \Delta \mathbf{y} \mathbf{x}^T$. Recall that the HQR factorization applies a series of Householder transformations on \mathbf{A} to form \mathbf{R} , and the same series of Householder transformations in reverse order to \mathbf{I} to form \mathbf{Q} . Therefore, it is appropriate to assume that \mathbf{x} is exact in this proof, and we form a forward bound on \mathbf{y} . However, we can still easily switch between forward and backward errors in the following way.

$$\begin{aligned} \mathbf{y} &= \mathbf{y} + \Delta \mathbf{y} = \mathbf{Q}^T (\mathbf{x} + \Delta \mathbf{x}) \\ &= (\mathbf{Q} + \Delta \mathbf{Q})^T \mathbf{x} \\ \Delta \mathbf{y} &= \Delta \mathbf{Q}^T \mathbf{x} = \mathbf{Q}^T \Delta \mathbf{x} \end{aligned}$$

In addition, we can switch between $\Delta \mathbf{y}$ and $\Delta \mathbf{x}$ by using the fact that $\|\mathbf{Q}\|_2 = 1$

Error bound for $\|\Delta \mathbf{y}\|_2$: We will first find $\|\Delta \mathbf{Q}\|_2$, where this is NOT the forward error from forming \mathbf{Q} with Householder transformations, but rather a backward error in accumulating Householder transformations. From Lemma A.1, we have $\|\Delta \mathbf{P}\|_F \leq \gamma_{\text{acc}}^{(4)} = \gamma_{\text{acc}}^{(4)} \|\mathbf{P}\|_2$ for any Householder transformation $\mathbf{P} \in \mathbb{R}^{m \times m}$, where $\hat{\mathbf{z}} = \mathbf{e}_d + \mathbf{e}_2 + \mathbf{e}_3$ and $d = \lfloor \frac{(m-1)\alpha}{2} \rfloor$, $\hat{\mathbf{z}} \in \{1, 2\}$. Therefore, this applies to the sequence of \mathbf{P}_i 's that form \mathbf{Q} as well.

We will now use Lemma 3.7 from [3] to bound $\Delta \mathbf{Q}$

$$\begin{aligned} \Delta \mathbf{Q}^T &= (\mathbf{Q} - \mathbf{Q})^T = \prod_{i=1}^d (\mathbf{P}_i + \Delta \mathbf{P}_i) - \prod_{i=1}^d \mathbf{P}_i \\ \|\Delta \mathbf{Q}\|_F &= \|\Delta \mathbf{Q}^T\|_F = \left\| \prod_{i=1}^d (\mathbf{P}_i + \Delta \mathbf{P}_i) - \prod_{i=1}^d \mathbf{P}_i \right\|_F \\ &\leq \left(\prod_{i=1}^d (1 + \gamma_{\text{acc}}^{(4)}) - 1 \right) \prod_{i=1}^d \|\mathbf{P}_i\|_2 = \prod_{i=1}^d (1 + \gamma_{\text{acc}}^{(4)}) - 1 \end{aligned}$$

The last equality results from the orthogonality of Householder matrices.

Let's take a look at the constant, $(1 + \gamma_{\text{acc}}^{(4)})^d - 1$. From the very last rule in Lemma 2.2, we can generalize the following:

$$(1 + \gamma_{\text{acc}}^{(4)})^d \leq (1 + \gamma_{\text{acc}}^{(4)})^{d-2} (1 + \gamma_{\text{acc}}^{(4)}) (1 + \gamma_{\text{acc}}^{(4)}) \quad (48)$$

$$\leq (1 + \gamma_{\text{acc}}^{(4)})^{d-2} (1 + \gamma_{\text{acc}}^{(4)}) \leq \dots \leq (1 + \gamma_{\text{acc}}^{(4)}) \quad (49)$$

□

So our quantity of interest can be bound by $(1 + \gamma_w^{(i)})^r - 1 \leq \gamma_w^{(i)r}$.

Now we will use the following equivalent algebraic inequalities to get the final result

$$0 < a < b < 1 \Leftrightarrow 1 - a > 1 - b \Leftrightarrow \frac{1}{1-a} < \frac{1}{1-b} \Leftrightarrow \frac{a}{1-a} < \frac{b}{1-b} \quad (68)$$

In addition, we assume $\gamma_w^{(i)} < \frac{1}{2}$.

$$\begin{aligned} (1 + \gamma_w^{(i)})^r - 1 &\leq \gamma_w^{(i)r} = \frac{r^2 n_w}{1 - r^2 n_w}, \text{ by definition.} \\ &\leq \frac{\gamma_w^{(i)r}}{1 - r^2 n_w}, \text{ since } r^2 n_w < r^2 n_w^{(i)}. \text{ By Equation 68} \\ &\leq 2r^2 n_w^{(i)}, \text{ since } r^2 n_w^{(i)} < \frac{1}{2} \text{ implies } \frac{1}{1 - r^2 n_w^{(i)}} < 2. \\ &= r_w^{(i)r}. \end{aligned} \quad (69)$$

where $\gamma_w^{(i)r} = \frac{r^2 n_w}{1 - r^2 n_w}$ for some small integer r . If we had started with $(1 + \gamma_w^{(i)})^r - 1$, we can still result in $r_w^{(i)r}$ assuming that $2r$ is still a small integer. In conclusion, we have

$$(1 + \gamma_w^{(i)})^r - 1 \leq r_w^{(i)r},$$

which results in the bound for ΔQ as shown in Equation 63.

$$\|\Delta Q\|_F \leq \|\Delta Q\|_F \leq r_w^{(i)r}.$$

Now, we can get the norm bound for ΔY .

$$\|\Delta Y\|_2 = \|\Delta Q X\|_2 \leq \|\Delta Q\|_2 \|X\|_2 \leq r_w^{(i)r} \|X\|_2 \quad (69)$$

Bound for $\|\Delta X\|_2$: We use the above result.

$$\|\Delta X\|_2 = \|Q \Delta Y\|_2 \leq \|Q\|_2 \|\Delta Y\|_2 = \|\Delta Y\|_2 \leq r_w^{(i)r} \|X\|_2 \quad (70)$$

While $r_w^{(i)r} = r \frac{r^2 n_w}{1 - r^2 n_w} < r \frac{r^2 n_w}{1 - r^2 n_w}$ holds true when $r > 0$ and $r k w < 1$ are satisfied, the strict inequality implies that $r_w^{(i)r}$ is a tighter bound than $\gamma_w^{(i)r}$. However, $\gamma_w^{(i)r}$ is easier to work with using the rules in Lemma 2.1. \square

A.3.1 Proof for Theorem 3.5

First, we use Lemma A.2 directly on columns of A and I_{row} to get a result for columns of R and Q . Figure 1 shows that each element of Q and R each go through different numbers of Householder transformations and Householder transformations for different lengths of vectors as well. We will use the maximum number of transformations and the length of the longest vector we perform a Householder transformation onto—that is n transformations of vectors of length m . For j in $\{1, \dots, n\}$, the j^{th} column of R and Q are the results of j Householder transformations on A and I .

$$\|\Delta Q\|_2 \|j\|_2 \leq j^{\gamma_w^{(i)r}} \|j\|_2 < \gamma_w^{(i)r} \quad (71)$$

$$\|\Delta R\|_2 \|j\|_2 \leq j^{\gamma_w^{(i)r}} \|A\|_2 < \gamma_w^{(i)r} \|A\|_2 \quad (72)$$

31

Finally, we relate columnwise 2-norms to matrix Frobenius norms. It is straightforward to see the result for the Q factor.

$$\|\Delta Q\|_F = \left(\sum_{j=1}^n \|\Delta Q\|_2 \|j\|_2^2 \right)^{1/2} \leq \left(\sum_{j=1}^n (j^{\gamma_w^{(i)r}})^2 \|j\|_2^2 \right)^{1/2} \leq n^{1/2} \gamma_w^{(i)r} \quad (73)$$

Note that we bound $\sum_{j=1}^n j^2$ by n^3 but the summation is actually exactly $\frac{n(n+1)(2n+1)}{6}$. Therefore, a tighter bound would replace $n^{1/2}$ with $\left(\frac{n(n+1)(2n+1)}{6}\right)^{1/2}$.

We can bound the R factor in a similar way.

$$\|\Delta R\|_F = \left(\sum_{j=1}^n \|\Delta R\|_2 \|j\|_2^2 \right)^{1/2} \leq \left(\sum_{j=1}^n (j^{\gamma_w^{(i)r}})^2 \|A\|_2^2 \|j\|_2^2 \right)^{1/2} \leq n^{1/2} \gamma_w^{(i)r} \|A\|_F \quad (74)$$

Now let us get the backward error from the UQR factorization.

$$\begin{aligned} \Delta A &= A - QR = A - QR + QR - QR \\ &= Q(R - R) + (Q - Q)R = Q\Delta R + \Delta Q R \end{aligned}$$

A columnwise result for A is shown by:

$$\begin{aligned} \|\Delta A\|_F \|j\|_2 &= \|(Q\Delta R + \Delta Q R)\|_F \|j\|_2 \\ &\leq \|Q\Delta R\|_F \|j\|_2 + \|\Delta Q R\|_F \|j\|_2 \\ &\leq \|\Delta R\|_2 \|j\|_2 + \|\Delta Q\|_2 \|R\|_2 \|j\|_2 \\ &\leq \|\Delta R\|_2 \|j\|_2 + \|\Delta Q\|_F \|R\|_F \|j\|_2 \\ &\leq j^{\gamma_w^{(i)r}} \|A\|_2 \|j\|_2 + n^{1/2} \gamma_w^{(i)r} \|Q\|_2^T \|A + \Delta R\|_2 \|j\|_2 \\ &\leq j^{\gamma_w^{(i)r}} \|A\|_2 \|j\|_2 + n^{1/2} \gamma_w^{(i)r} (\|A\|_2 \|j\|_2 + \|\Delta R\|_2 \|j\|_2) \\ &\leq (j^{\gamma_w^{(i)r}} + n^{1/2} \gamma_w^{(i)r}) (\|A\|_2 \|j\|_2 + \|\Delta R\|_2 \|j\|_2) \\ &= n^{1/2} \gamma_w^{(i)r} \|A\|_2 \|j\|_2, \end{aligned}$$

where we assume $n^{1/2} \gamma_w^{(i)r} < 1$ and the last equality swaps all non-leading order terms under the arbitrary constant ϵ within the definition of $\tilde{\gamma}$.

$$\|\Delta A\|_F = \left(\sum_{j=1}^n \|\Delta A\|_F \|j\|_2^2 \right)^{1/2} \leq \left(\sum_{j=1}^n (n^{1/2} \gamma_w^{(i)r})^2 \|A\|_2^2 \|j\|_2^2 \right)^{1/2} \leq n^{1/2} \gamma_w^{(i)r} \|A\|_F \quad (75)$$

References

- [1] M. Goulou, Y. Deng, and L.-P. David, "Training deep neural networks with low precision multiplications," *arXiv preprint arXiv:1412.7024*, 2014.

- [3] M. Conranton, J.-P. David, and Y. Bergho, "Low precision storage for deep learning," *arXiv preprint arXiv:1412.0024*, 2014.
- [4] N. J. Higham, *Accuracy and Stability of Numerical Algorithms*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2nd ed., 2002.
- [5] A. S. Hestenes, "Unitary triangulization of a nonsymmetric matrix," *Journal of the ACM (JACM)*, vol. 5, no. 4, pp. 339–342, 1958.
- [6] D. Mori, Y. Yamamoto, and S.-L. Zhang, "Backward error analysis of the allative algorithm for householder qr decomposition," *Japan Journal of Industrial and Applied Mathematics*, vol. 29, pp. 111–130, Feb 2012.
- [7] J. Demmel, L. Grigori, M. Hoemmen, and J. Langou, "Communication-optimal parallel and sequential qr and lu factorizations," *SIAM Journal on Scientific Computing*, vol. 34, no. 1, pp. A206–A239, 2012.
- [8] G. Ballard, J. W. Demmel, L. Grigori, M. Juregan, H. Dipp, Xinyu, and E. Solomonik, "Reconstructing householder vectors from tall-skinny qr," vol. 85, pp. 1159–1170, 05 2014.
- [9] Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, and van der Vorst H., *Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 1st ed., 2000.
- [10] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, et al., "A density-based algorithm for discovering clusters in large spatial databases with noise," in *Kdd*, vol. 96, pp. 226–231, 1996.
- [11] E. Kim, Y. Guletskii, M. Hurley, M. Jones, J. Kepner, S. Mahapatra, P. Montarola, A. Reuther, S. Sami, W. Song, D. Stahle, and S. Smith, "Streaming graph challenge: Stochastic block partition," in *2017 IEEE High Performance Extreme Computing Conference (HPEC)*, pp. 1–12, Sep. 2017.