

ROUNDING ERROR ANALYSIS OF MIXED-PRECISION HOUSEHOLDER QR ALGORITHMS

L. MINAH YANG, ALYSON FOX, AND GEOFFREY SANDERS

Abstract. Although mixed precision arithmetic has recently garnered interest for training dense neural networks, many other applications could benefit from the speed-ups and lower storage if applied appropriately. The growing interest in employing mixed precision computations motivates the need for rounding error analysis that properly handles behavior from mixed precision arithmetic. We present a framework for mixed precision analysis that builds on the foundations of rounding error analysis presented in [13] and demonstrate its practicality by applying the analysis to various Householder QR Algorithms.

1. Introduction.

2. Background: Build up to rounding error analysis for inner products.

3. Algorithms and existing round-off error analyses. We introduce the Householder QR factorization algorithm (HQR) in subsection 3.1 and two block variants that use HQR within the block in subsections 3.2 and 3.3. The blocked HQR (BQR) in subsection 3.2 partitions the columns of the target matrix and utilizes mainly level-3 BLAS operations and is a well-known algorithm that uses the WY representation of [4]. In contrast, the Tall-and-Skinny QR (TSQR) in subsection 3.3 partitions rows of the matrix and takes a communication-avoiding divide-and-conquer approach that can be easily parallelized (see [7]). We also present the crucial results in standard rounding error analysis of these algorithms that excludes any mixed-precision assumptions. These building steps of round-off error analysis will be easily tweaked for various mixed-precision assumptions in section 4.

3.1. Householder QR (HQR). The HQR algorithm uses Householder transformations to zero out elements below the diagonal of a matrix (see [16]). 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, $\mathbf{P}_{\mathbf{v}}$, such that $\mathbf{P}_{\mathbf{v}}$ zeros out \mathbf{x} below the first element.*

$$(3.1) \quad \begin{aligned} \sigma &= -\text{sign}(x_1)\|\mathbf{x}\|_2, \quad \mathbf{v} = \mathbf{x} - \sigma\mathbf{e}_1, \\ \beta &= \frac{2}{\mathbf{v}^\top \mathbf{v}} = -\frac{1}{\sigma v_1}, \quad \mathbf{P}_{\mathbf{v}} = \mathbf{I}_m - \beta \mathbf{v} \mathbf{v}^\top. \end{aligned}$$

The transformed vector, $\mathbf{P}_{\mathbf{v}}\mathbf{x}$, has the same 2-norm as \mathbf{x} since Householder transformations are orthogonal: $\mathbf{P}_{\mathbf{v}}\mathbf{x} = \sigma\mathbf{e}_1$. In addition, $\mathbf{P}_{\mathbf{v}}$ is symmetric and orthogonal, $\mathbf{P}_{\mathbf{v}} = \mathbf{P}_{\mathbf{v}}^\top = \mathbf{P}_{\mathbf{v}}^{-1}$.

3.1.1. HQR: Algorithm. Given $\mathbf{A} \in \mathbb{R}^{m \times n}$ and Lemma 3.1, HQR is done by repeating the following processes until only an upper triangle matrix remains. For $i = 1, 2, \dots, n$,
Step 1) Compute \mathbf{v} and β that zeros out the i^{th} column of \mathbf{A} beneath a_{ii} (see alg. 1), and
Step 2) Apply $\mathbf{P}_{\mathbf{v}}$ to the bottom right partition, $\mathbf{A}[i : m, i : n]$ (lines 4-6 of alg. 2).

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 and was supported by the LLNL-LDRD Program under Project No. 17-SI-004, LLNL-JRNL-795525-DRAFT.

34 Consider the following 4-by-3 matrix example adapted from [13]. Let \mathbf{P}_i represent the i^{th}
 35 Householder transformation of this algorithm.

$$\begin{aligned}
 36 \quad \mathbf{A} &= \begin{bmatrix} \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \\ \times & \times & \times \end{bmatrix} \xrightarrow{\text{apply } \mathbf{P}_1 \text{ to } \mathbf{A}} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \\ 0 & \times & \times \end{bmatrix} \xrightarrow{\text{apply } \mathbf{P}_2 \text{ to } \mathbf{P}_1 \mathbf{A}} \\
 37 \quad &\begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & \times \end{bmatrix} \xrightarrow{\text{apply } \mathbf{P}_3 \text{ to } \mathbf{P}_2 \mathbf{P}_1 \mathbf{A}} \begin{bmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \end{bmatrix} = \mathbf{P}_3 \mathbf{P}_2 \mathbf{P}_1 \mathbf{A} =: \mathbf{R} \\
 38
 \end{aligned}$$

39 Then, the \mathbf{Q} factor for a full QR factorization is $\mathbf{Q} := \mathbf{P}_1 \mathbf{P}_2 \mathbf{P}_3$ since \mathbf{P}_i 's are symmetric, and the
 40 thin factors for a general matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ are

$$41 \quad (3.2) \quad \mathbf{Q}_{\text{thin}} = \mathbf{P}_1 \cdots \mathbf{P}_n \mathbf{I}_{m \times n} \quad \text{and} \quad \mathbf{R}_{\text{thin}} = \mathbf{I}_{m \times n}^\top \mathbf{P}_n \cdots \mathbf{P}_1 \mathbf{A}.$$

Algorithm 1: $\beta, \mathbf{v}, \sigma = \text{hh_vec}(\mathbf{x})$. Given a vector $\mathbf{x} \in \mathbb{R}^n$, return $\mathbf{v}, \beta, \sigma$ that satisfy
 $(I - \beta \mathbf{v} \mathbf{v}^\top) \mathbf{x} = \sigma \hat{\mathbf{e}}_1$ and $\mathbf{v}_1 = 1$ (see [2, 13]).

Input: $\mathbf{x} \in \mathbb{R}^m$
Output: $\mathbf{v} \in \mathbb{R}^m$, and $\sigma, \beta \in \mathbb{R}$ such that $(I - \beta \mathbf{v} \mathbf{v}^\top) \mathbf{x} = \pm \|\mathbf{x}\|_2 \hat{\mathbf{e}}_1 = \sigma \hat{\mathbf{e}}_1$

42 1 $\mathbf{v} \leftarrow \text{copy}(\mathbf{x})$
 2 $\sigma \leftarrow -\text{sign}(\mathbf{x}_1) \|\mathbf{x}\|_2$
 3 $\mathbf{v}_1 \leftarrow \mathbf{x}_1 - \sigma$
 4 $\beta \leftarrow -\frac{\mathbf{v}_1}{\sigma}$
 5 **return** $\beta, \mathbf{v}/\mathbf{v}_1, \sigma$

Algorithm 2: $\mathbf{V}, \beta, \mathbf{R} = \text{HQR2}(\mathbf{A})$. A Level-2 BLAS implementation of the Householder
 QR algorithm. Given a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ where $m \geq n$, return matrix $\mathbf{V} \in \mathbb{R}^{m \times n}$, vector
 $\beta \in \mathbb{R}^n$, and upper triangular matrix \mathbf{R} . An orthogonal matrix \mathbf{Q} can be generated from
 \mathbf{V} and β , and $\mathbf{Q}\mathbf{R} = \mathbf{A}$.

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

1 $\mathbf{V}, \beta \leftarrow \mathbf{0}_{m \times n}, \mathbf{0}_m$
 2 **for** $i = 1 : n$ **do**
 3 $\mathbf{v}, \beta, \sigma \leftarrow \text{hh_vec}(\mathbf{A}[i : \text{end}, i])$
 4 $\mathbf{V}[i : \text{end}, i], \beta_i, \mathbf{A}[i, i] \leftarrow \mathbf{v}, \beta, \sigma$
 5 $\mathbf{A}[i + 1 : \text{end}, i] \leftarrow \text{zeros}(m - i)$
 6 $\mathbf{A}[i : \text{end}, i + 1 : \text{end}] \leftarrow \mathbf{A}[i : \text{end}, i + 1 : \text{end}] - \beta \mathbf{v} \mathbf{v}^\top \mathbf{A}[i : \text{end}, i + 1 : \text{end}]$
 7 **return** $\mathbf{V}, \beta, \mathbf{A}[1 : n, 1 : n]$

43 **3.1.2. HQR: Rounding Error Analysis.** Now we present an error analysis for [alg. 2](#) by
 44 keeping track of the different operations of [alg. 1](#) and [alg. 2](#).

45 *Calculating the i^{th} Householder vector and constant.* In [alg. 2](#), the i^{th} Householder vector
 46 shares all but the first component with the target column, $\mathbf{A}[i : m, i]$. We first calculate σ as is
 47 implemented in line 2 of [alg. 1](#).

$$48 \quad (3.3) \quad \text{fl}(\sigma) = \hat{\sigma} = \text{fl}(-\text{sign}(\mathbf{A}_{i,i}) \|\mathbf{A}[i : m, i]\|_2) = \sigma + \Delta\sigma, \quad |\Delta\sigma| \leq \gamma_{m-i+1}|\sigma|.$$

49 Note that the backward error incurred here is simply that an inner product of a vector in \mathbb{R}^{m-i+1}
 50 with itself. Let $\tilde{\mathbf{v}}_1 \equiv \mathbf{A}_{i,i} - \sigma$, the penultimate value \mathbf{v}_1 . The subtraction adds a single additional
 51 rounding error via

$$52 \quad \text{fl}(\tilde{\mathbf{v}}_1) = \tilde{\mathbf{v}}_1 + \Delta\tilde{\mathbf{v}}_1 = (1 + \delta)(\mathbf{A}_{i,i} - \sigma - \Delta\sigma) = (1 + \tilde{\theta}_{m-i+2})(\mathbf{A}_{i,i} - \sigma)$$

53 where the last equality is granted because the sign of σ is chosen to prevent cancellation. For
 54 the sake of simplicity, we write $|\Delta\tilde{\mathbf{v}}_1| \leq \tilde{\gamma}_{m-i+1}|\tilde{\mathbf{v}}_1|$ even though a tighter relative upper bound is
 55 θ_{m-i+2} . We sweep that minor difference (in comparison to $\mathcal{O}(m-i)$) under the our use of the $\tilde{\gamma}$
 56 notation defined in [??](#). Since [alg. 1](#) normalizes the Householder vector so that its first component
 57 is 1, the remaining components of \mathbf{v} are divided by $\text{fl}(\tilde{\mathbf{v}}_1)$ incurring another single rounding error.
 58 As a result, the rounding errors in \mathbf{v} are

$$59 \quad (3.4) \quad \text{fl}(\mathbf{v}_j) = \mathbf{v}_j + \Delta\mathbf{v}_j \text{ where } |\Delta\mathbf{v}_j| \leq \begin{cases} 0, & j = 1 \\ \tilde{\gamma}_{m-i+1}|\mathbf{v}_j|, & j = 2 : m - i + 1. \end{cases}$$

60 Next, we consider the Householder constant, β , as is computed in line 4 of [alg. 1](#).

$$61 \quad (3.5) \quad \hat{\beta} = \text{fl}\left(-\frac{\tilde{\mathbf{v}}_1}{\hat{\sigma}}\right) = -(1 + \delta)\frac{\tilde{\mathbf{v}}_1 + \Delta\tilde{\mathbf{v}}_1}{\sigma + \Delta\sigma}$$

$$62 \quad (3.6) \quad = \frac{(1 + \delta)(1 + \theta_{m-i+1})}{(1 + \theta_{m-i+2})}\beta = (1 + \theta_{3(m-i+2)})\beta$$

$$63 \quad (3.7) \quad = \beta + \Delta\beta, \text{ where } |\Delta\beta| \leq \tilde{\gamma}_{m-i+1}\beta.$$

65 We have shown [\(3.5\)](#) to keep our analysis simple in [section 4](#) and [\(3.6\)](#) and [\(3.7\)](#) show that the error
 66 incurred from calculating of $\|\mathbf{A}[i : m, i]\|_2$ accounts for the vast majority of the rounding error so
 67 far.

68 *Applying a Single Householder Transformation.* Now we consider lines 4-6 of [alg. 2](#). Since
 69 the entries in $\mathbf{A}[i + 1 : m, i]$ are simply zeroed out and $\mathbf{A}_{i,i}$ is replaced by σ , we only need to
 70 calculate the errors for applying a Householder transformation with the computed Householder
 71 vector and constant. This is the most crucial building block of the rounding error analysis for any
 72 variant of HQR because the \mathbf{Q} factor is formed by applying the Householder transformations to
 73 the identity and both of the blocked versions in [subsection 3.2](#) and [subsection 3.3](#) require efficient
 74 implementations of this step. In this section, we only consider a level-2 BLAS implementation
 75 of applying the Householder transformation, but in [subsection 3.2](#) we introduce a level-3 BLAS
 76 implementation.

77 A Householder transformation is applied through a series of inner and outer products, since
 78 Householder matrices are rank-1 updates of the identity. That is, computing $\mathbf{P}_\mathbf{v}\mathbf{x}$ for any $\mathbf{x} \in \mathbb{R}^m$
 79 is as simple as computing $\mathbf{y} := \mathbf{x} - (\beta\mathbf{v}^\top\mathbf{x})\mathbf{v}$. Let us assume that \mathbf{x} is an exact vector and there
 80 were errors incurred in forming \mathbf{v} and β . The errors incurred from computing \mathbf{v} and β need to be
 81 included in addition to the new rounding errors accumulating from the action of applying $\mathbf{P}_\mathbf{v}$ to

a column. In practice, \mathbf{x} would be a column in $\mathbf{A}^{(i-1)}[i+1:m, i+1:n]$, where the superscript $(i-1)$ indicates that this submatrix of \mathbf{A} has already been transformed by $i-1$ Householder transformations that zeroed out components below $\mathbf{A}_{j,j}$ for $j = 1:i-1$. We show the error for forming $\text{fl}(\hat{\mathbf{v}}^\top \mathbf{x})$ where we continue to let $\mathbf{v}, \mathbf{x} \in \mathbb{R}^{m-i+1}$ as would be in the i^{th} iteration of the for-loop in [alg. 2](#):

$$\text{fl}(\hat{\mathbf{v}}^\top \mathbf{x}) = (1 + \theta_{m-i+1})(\mathbf{v} + \Delta \mathbf{v})^\top \mathbf{x}.$$

Set $\mathbf{w} := \beta \mathbf{v}^\top \mathbf{x} \mathbf{v}$. Then,

$$\hat{\mathbf{w}} = (1 + \theta_{m-i+1})(1 + \delta)(1 + \tilde{\delta})(\beta + \Delta \beta)(\mathbf{v} + \Delta \mathbf{v})^\top \mathbf{x}(\mathbf{v} + \Delta \mathbf{v}),$$

where θ_{m-i+1} is from computing the inner product $\hat{\mathbf{v}}^\top \mathbf{x}$, and δ and $\tilde{\delta}$ are from multiplying β , $\text{fl}(\hat{\mathbf{v}}^\top \mathbf{x})$, and $\hat{\mathbf{v}}$ together. Finally, we can add in the vector subtraction operation and complete the rounding error analysis of applying a Householder transformation to any vector:

$$(3.8) \quad \text{fl}(\mathbf{x} - \hat{\mathbf{w}}) = (1 + \delta)(\mathbf{x} - \mathbf{w} - \Delta \mathbf{w}) = (1 + \tilde{\theta}_{m-i+1})\mathbf{y}.$$

We can easily switch between forward and errors from [\(3.8\)](#) via

$$\mathbf{y} + \Delta \mathbf{y} = (1 + \tilde{\theta}_{m-i+1})\mathbf{y} = (1 + \tilde{\theta}_{m-i+1})\mathbf{P}_\mathbf{v} \mathbf{x} = (\mathbf{P}_\mathbf{v} + \Delta \mathbf{P}_\mathbf{v})\mathbf{x},$$

where $|\Delta \mathbf{y}| \leq \tilde{\gamma}_{m-i+1}|\mathbf{y}|$ and $|\Delta \mathbf{P}_\mathbf{v}| \leq \tilde{\gamma}_{m-i+1}|\mathbf{P}_\mathbf{v}|$.

Even though we never explicitly form $\mathbf{P}_\mathbf{v}$, forming the normwise error bound for this matrix makes the analysis for HQR simpler. Therefore, we now transition from componentwise error to matrix norm errors: the 2-norm and the Frobenius norm.

First, we transition from componentwise forward error to the 2-norm forward error via

$$(3.9) \quad \|\Delta \mathbf{y}\|_2 = \left(\sum_{i=1}^m \Delta \mathbf{y}_i^2 \right)^{1/2} \leq \left((\tilde{\gamma}_{m-i+1})^2 \sum_{i=1}^m |\mathbf{y}_i|^2 \right)^{1/2} = \tilde{\gamma}_{m-i+1} \|\mathbf{y}\|_2.$$

In exact arithmetic, we are guaranteed $\|\mathbf{y}\|_2 = \|\mathbf{P}_\mathbf{v} \mathbf{x}\|_2 \leq \|\mathbf{P}\|_2 \|\mathbf{x}\|_2 = \|\mathbf{x}\|_2$ since $\mathbf{P}_\mathbf{v}$ is orthogonal and preserves norms. Combining this with [\(3.9\)](#) we find

$$(3.10) \quad \frac{\|\Delta \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \leq \tilde{\gamma}_{m-i+1}.$$

Now we convert this to a normwise backward error. Since $\Delta \mathbf{P}$ is exactly $\frac{1}{\mathbf{x}^\top \mathbf{x}} \Delta \mathbf{y} \mathbf{x}^\top$, we can compute its Frobenius norm by using $\Delta \mathbf{P}_{ij} = \frac{1}{\|\mathbf{x}\|_2^2} \Delta \mathbf{y}_i \mathbf{x}_j$,

$$\|\Delta \mathbf{P}\|_F = \left(\sum_{i=1}^m \sum_{j=1}^m \left(\frac{1}{\|\mathbf{x}\|_2^2} \Delta \mathbf{y}_i \mathbf{x}_j \right)^2 \right)^{1/2} = \frac{\|\Delta \mathbf{y}\|_2}{\|\mathbf{x}\|_2} \leq \tilde{\gamma}_{m-i+1},$$

where the last inequality is a direct application of [\(3.10\)](#). We summarize these results in [Lemma 3.2](#).

LEMMA 3.2. *Let $\mathbf{x} \in \mathbb{R}^m$ and consider the computation of $\hat{\mathbf{y}} = \text{fl}(\mathbf{P}_\mathbf{v} \mathbf{x})$ via*

$$\mathbf{y} + \Delta \mathbf{y} = \text{fl}(\mathbf{P}_\mathbf{v} \mathbf{x}) = \text{fl}(\mathbf{x} - \hat{\beta} \hat{\mathbf{v}} \hat{\mathbf{v}}^\top \mathbf{x})$$

and rounding errors incurred in forming $\hat{\mathbf{v}}$ and $\hat{\beta}$ are expressed componentwise via $\hat{\mathbf{v}} = \mathbf{v} + \Delta\mathbf{v}$ and $\hat{\beta} = \beta + \Delta\beta$. Let us write the componentwise forward error bound as $|\Delta\mathbf{y}| \leq \gamma_y |\mathbf{y}|$. Then, the normwise forward and backward errors are

$$\|\Delta\mathbf{y}\|_2 \leq \gamma_y \|\mathbf{y}\|_2, \quad \|\mathbf{P}_v\|_F \leq \gamma_y.$$

Note that in a uniform precision setting this bound is represented as $\gamma_y = \tilde{\gamma}_m$, where the majority of the round-off errors are attributed to inner product computations for forming $\hat{\beta}$ and \mathbf{v} .

Applying many successive Householder transformations. Consider applying a sequence of transformations in the set $\{\mathbf{P}_i\}_{i=1}^r \subset \mathbb{R}^{m \times m}$ to $\mathbf{x} \in \mathbb{R}^m$, where \mathbf{P}_i 's are all Householder transformations. This is directly applicable to HQR as $\mathbf{Q} = \mathbf{P}_1 \cdots \mathbf{P}_n \mathbf{I}$ and $\mathbf{R} = \mathbf{Q}^\top \mathbf{A} = \mathbf{P}_n \cdots \mathbf{P}_1 \mathbf{A}$. Let us define

$$\mathbf{Q} + \Delta\mathbf{Q}' \equiv \prod_{i=1}^r (\mathbf{P}_i + \Delta\mathbf{P}_i)$$

in the context of applying this matrix to a vector, $\mathbf{x} \in \mathbb{R}^m$, where $\Delta\mathbf{Q}'^\top$ represents the backward error of forming \mathbf{R} , instead of the forward error of the \mathbf{Q} factor. The forward error for \mathbf{Q} is denoted as $\Delta\mathbf{Q} \equiv \text{fl}(\mathbf{Q}) - \mathbf{Q}$ where $\text{fl}(\mathbf{Q})$ is formed via HQR. That is, if $\mathbf{y} = \mathbf{Q}^\top \mathbf{x}$, then $\text{fl}(\mathbf{y}) = \mathbf{y} + \Delta\mathbf{y} = (\mathbf{Q} + \Delta\mathbf{Q}')^\top \mathbf{x}$. Even though an efficient implementation would use that \mathbf{P}_i 's are applied to successively shorter vectors (\mathbf{P}_i is left multiplied to $\mathbf{A}[i : m, i+1 : n]$, which is equivalent to $n-i$ vectors of length $m-i+1$), we assume $\{\mathbf{P}_i\}_{i=1}^r \subset \mathbb{R}^{m \times m}$ to allow for a simpler analysis while forming a looser bound. We will now use Lemma 3.7 from [13] to bound $\Delta\mathbf{Q}'$ with the Frobenius norm.

$$\begin{aligned} \|\Delta\mathbf{Q}'^\top\|_F &= \left\| \prod_{r=1}^1 (\mathbf{P}_i + \Delta\mathbf{P}_i) - \prod_{i=r}^1 \mathbf{P}_i \right\|_F, \\ &\leq \left(\prod_{i=1}^r (1 + \tilde{\gamma}_m) - 1 \right) \prod_{i=r}^1 \|\mathbf{P}_i\|_2 = (1 + \tilde{\gamma}_m)^r - 1. \end{aligned}$$

The last equality results from the orthogonality of Householder matrices, and we further reduce the last term. Generalizing the last rule in Lemma ?? yields

$$(1 + \tilde{\gamma}_m)^r = (1 + \tilde{\gamma}_m)^{r-2} (1 + \tilde{\gamma}_m) (1 + \tilde{\gamma}_m) \leq (1 + \tilde{\gamma}_m)^{r-2} (1 + \tilde{\gamma}_{2m}) \leq \cdots \leq (1 + \tilde{\gamma}_{rm}).$$

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

$$(3.11) \quad 0 < a < b < 1 \Leftrightarrow 1 - a > 1 - b \Leftrightarrow \frac{1}{1-a} < \frac{1}{1-b} \Leftrightarrow \frac{a}{1-a} < \frac{b}{1-b}$$

In addition, we assume $r\tilde{\gamma}_m < \frac{1}{2}$, such that

$$(3.12) \quad (1 + \tilde{\gamma}_m)^r - 1 \leq \gamma_w^{(r\tilde{z})} = \frac{r\tilde{z}u_w}{1 - r\tilde{z}u_w} \quad (\text{by definition})$$

$$(3.13) \quad \leq \frac{r\tilde{\gamma}_m}{1 - r\tilde{\gamma}_m}, \text{ since } r\tilde{z}u_w < r\tilde{\gamma}_m \quad (\text{by Equation 3.11})$$

$$(3.14) \quad \leq 2r\tilde{\gamma}_m \quad (\text{since } r\tilde{\gamma}_m < \frac{1}{2} \text{ implies } \frac{1}{1 - r\tilde{\gamma}_m} < 2)$$

$$(3.15) \quad = r\tilde{\gamma}_m,$$

Therefore, we have $(1 + \tilde{\gamma}_m)^r - 1 \leq r\tilde{\gamma}_m$ and

$$(3.16) \quad \|\Delta \mathbf{Q}'\|_2 \leq \|\Delta \mathbf{Q}'\|_F = \|\Delta \mathbf{Q}'^\top\|_F \leq r\tilde{\gamma}_m$$

In this current uniform precision error analysis, the important quantity $\tilde{\gamma}_m$ is derived from the backward error of applying one Householder transformation. To easily generalize this section for mixed-precision analysis, we benefit from alternatively denoting this quantity as $\tilde{\gamma}_{\mathbf{P}}$ with the understanding that $\tilde{\gamma}_{\mathbf{P}}$ will be some combination of $\tilde{\gamma}$'s of differing precisions. Equation (3.15) would then be

$$(3.17) \quad (1 + \tilde{\gamma}_{\mathbf{P}})^r - 1 \leq r\tilde{\gamma}_{\mathbf{P}}.$$

Next, we apply (3.16) to the i^{th} columns of \mathbf{Q}, \mathbf{R} and set $r = n$ for a full rank matrix, \mathbf{A} . Then,

$$\begin{aligned} \|\Delta \mathbf{R}[:, i]\|_2 &= \|\Delta \mathbf{Q}'^\top \mathbf{A}[:, i]\|_2 \leq \|\Delta \mathbf{Q}'\|_2 \|\mathbf{A}[:, i]\|_2 \leq n\tilde{\gamma}_m \|\mathbf{A}[:, i]\|_2, \\ \|\Delta \mathbf{Q}[:, i]\|_2 &= \|\Delta \mathbf{Q}' \mathbf{I}[:, i]\|_2 \leq \|\Delta \mathbf{Q}'\|_2 \leq n\tilde{\gamma}_m. \end{aligned}$$

These columnwise bounds can now be transformed into matrix norms as follows:

$$\begin{aligned} \|\Delta \mathbf{R}\|_F &= \left(\sum_{i=1}^n \|\Delta \mathbf{R}[:, i]\|_2^2 \right)^{1/2} \leq \left(\sum_{i=1}^n n^2 \tilde{\gamma}_m^2 \|\mathbf{A}[:, i]\|_2^2 \right)^{1/2} = n\tilde{\gamma}_m \|\mathbf{A}\|_F, \\ \|\Delta \mathbf{Q}\|_F &= \left(\sum_{i=1}^n \|\Delta \mathbf{Q}[:, i]\|_2^2 \right)^{1/2} \leq \left(\sum_{i=1}^n \tilde{\gamma}_m^2 \right)^{1/2} = n^{3/2} \tilde{\gamma}_m. \end{aligned}$$

We gather these results into Theorem 3.3.

THEOREM 3.3. *Let $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m \geq n$ have full rank, n . Let $\hat{\mathbf{Q}} \in \mathbb{R}^{m \times n}$ and $\hat{\mathbf{R}} \in \mathbb{R}^{n \times n}$ be the thin QR factors of \mathbf{A} obtained via alg. 2, defined via*

$$\begin{aligned} \hat{\mathbf{R}} &= \mathbf{R} + \Delta \mathbf{R} = \text{fl}(\hat{\mathbf{P}}_n \cdots \hat{\mathbf{P}}_1 \mathbf{A}), \quad n\tilde{\gamma}_m \|\mathbf{A}\|_F \\ \hat{\mathbf{Q}} &= \mathbf{Q} + \Delta \mathbf{Q} = \text{fl}(\hat{\mathbf{P}}_1 \cdots \hat{\mathbf{P}}_n \mathbf{I}), \quad \|\Delta \mathbf{Q}\|_F \leq n^{3/2} \tilde{\gamma}_m. \end{aligned}$$

Let $\mathbf{A} + \Delta \mathbf{A} = \hat{\mathbf{Q}} \hat{\mathbf{R}}$, where $\hat{\mathbf{Q}}$ and $\hat{\mathbf{R}}$ are obtained via Algorithm 2. Then the backward error is

$$(3.18) \quad \|\Delta \mathbf{A}\|_F \leq n^{3/2} \tilde{\gamma}_m \|\mathbf{A}\|_F.$$

The content of this section is largely derived directly from [13], but we kept the analysis general by employing quantities denoted via $\Delta \beta$, $\Delta \mathbf{v}$, $\tilde{\gamma}_y$, and $\tilde{\gamma}_{\mathbf{P}}$. These quantities account for various forward and backward errors formed in computing essential components of HQR, namely the Householder constant and vector, as well as normwise errors of the action of applying Householder transformations. In the next sections, we present blocked variants of HQR that use alg. 2.

3.2. Block HQR with partitioned columns (BQR). We refer to the blocked variant of HQR where the columns are partitioned as BQR. Note that this algorithm relies on the WY representation described in [4] instead of the storage-efficient version of [19], which is widely implemented.

3.2.1. The WY Representation. A convenient matrix representation that accumulates r Householder reflectors is known as the WY representation.

LEMMA 3.4. Suppose $\mathbf{Q} = \mathbf{I}_m - \mathbf{W}\mathbf{Y}^\top \in \mathbb{R}^{m \times m}$ is an orthogonal matrix with $\mathbf{W}, \mathbf{Y} \in \mathbb{R}^{m \times j}$. If $\mathbf{P} = \mathbf{I}_m - \beta \mathbf{v}\mathbf{v}^\top$ with $\mathbf{v} \in \mathbb{R}^m$ and $\mathbf{z} = \beta \mathbf{Q}\mathbf{v}$, then

$$\mathbf{Q}_+ = \mathbf{Q}\mathbf{P} = \mathbf{I} - \mathbf{W}_+ \mathbf{Y}_+^\top,$$

where $\mathbf{W}_+ = [\mathbf{W} | \mathbf{z}]$ and $\mathbf{Y}_+ = [\mathbf{Y} | \mathbf{v}]$ are each m -by- $(j+1)$.

If \mathbf{Q} was already the accumulation of j Householder transformations, then Lemma 3.4 shows us a clever way to build the WY representation of successive Householder transformations. Let us now show the proof for Lemma 3.4.

Proof. A direct right multiplication of $\mathbf{P} := \mathbf{I}_m - \beta \mathbf{v}\mathbf{v}^\top$ onto \mathbf{Q} can be written as

$$\mathbf{Q}\mathbf{P} = \mathbf{Q} - \beta \mathbf{Q}\mathbf{v}\mathbf{v}^\top.$$

Let us use the WY representation of \mathbf{Q} .

$$\mathbf{Q}\mathbf{P} = \mathbf{I}_m - \mathbf{W}\mathbf{Y}^\top - \beta \mathbf{Q}\mathbf{v}\mathbf{v}^\top = \mathbf{I}_m - \mathbf{W}\mathbf{Y}^\top - \mathbf{z}\mathbf{v}^\top$$

Now note that the two subtracted terms are exactly the updated WY factors:

$$\mathbf{W}_+ \mathbf{Y}_+^\top = [\mathbf{W} \quad \mathbf{z}] \begin{bmatrix} \mathbf{Y}^\top \\ \mathbf{v}^\top \end{bmatrix} = \mathbf{W}\mathbf{Y}^\top + \mathbf{z}\mathbf{v}^\top.$$

□

With the correct initialization of \mathbf{W} and \mathbf{Y} , we can build the WY representation of successive Householder transformations as shown in Algorithm 3.

Algorithm 3: $\mathbf{W}, \mathbf{Y} \leftarrow \text{buidlWY}(V, \beta)$: Given a set of householder vectors $\{\mathbf{V}[:, i]\}_{i=1}^r$ and their corresponding constants $\{\beta_i\}_{i=1}^r$, form the final \mathbf{W} and \mathbf{Y} factors of the WY representation of $\mathbf{P}_1 \cdots \mathbf{P}_r$, where $\mathbf{P}_i := \mathbf{I}_m - \beta_i \mathbf{v}_i \mathbf{v}_i^\top$

Input: $\mathbf{V} \in \mathbb{R}^{m \times r}$, $\beta \in \mathbb{R}^r$ where $m > r$.

Output: \mathbf{W}, \mathbf{Y}

1 Initialize: $\mathbf{W} := \beta_1 \mathbf{V}[:, 1]$ and $\mathbf{Y} := \mathbf{V}[:, 1]$.

2 for $j = 2 : r$ do

3 $\mathbf{z} \leftarrow \beta_j [\mathbf{V}[:, j] - \mathbf{W} (\mathbf{Y}^\top \mathbf{V}[:, j])]$

4 $\mathbf{W} \leftarrow [\mathbf{W} \quad \mathbf{z}]$

5 $\mathbf{Y} \leftarrow [\mathbf{Y} \quad \mathbf{V}[:, j]]$

 // \mathbf{W} and \mathbf{Y} are now m -by- j matrices.

/* Update \mathbf{W} . */

/* Update \mathbf{Y} . */

6 return \mathbf{W}, \mathbf{Y}

In the traditional HQR, \mathbf{A} is transformed into an upper triangular matrix \mathbf{R} by first computing the Householder transformation to zero out a column below the diagonal, then applying that Householder transformation to all of the remaining columns to the right. For example, the k^{th} Householder transformation finds an $m - k + 1$ length Householder vector, \mathbf{v}_k , and applies it to an $(m - k + 1)$ -by- $(n - k)$ matrix. The bulk of FLOPs of this step (line 6 in alg. 2) requires two Level-2

193 BLAS operations when computed efficiently, which are $\mathbf{C} := \mathbf{v}_k^\top \mathbf{A}_{k:m, k+1:n} \mathbb{R}^{1 \times (n-k)}$ and \mathbf{vC} , an
 194 outer product.

195 In BQR, the columns of \mathbf{A} are partitioned by groups of r with $\mathbf{A} = [\mathbf{C}_1 \cdots \mathbf{C}_N]$ except for the
 196 last block which is $\mathbf{C}_N = \mathbf{A}[:, (N-1)r+1 : n]$ and $N = \lceil \frac{n}{r} \rceil$. The first block is triangularized
 197 using HQR and the WY representation of $\mathbf{P}_1 \cdots \mathbf{P}_r = \mathbf{I}_m - \mathbf{W}_1 \mathbf{Y}_1^\top$ is built at the end. Both of
 198 these operations are rich in Level-2 BLAS operations. Then, $\mathbf{I}_m - \mathbf{Y}_1 \mathbf{W}_1^\top = \mathbf{P}_r \cdots \mathbf{P}_1$ is applied to
 199 $[\mathbf{C}_2 \cdots \mathbf{C}_N]$ with two Level-3 BLAS operations:

- 200 1. $\mathbf{A} := \mathbf{W}_1^\top [\mathbf{C}_2 \cdots \mathbf{C}_N]$ is a matrix-matrix multiply with m -length inner products.
- 201 2. $[\mathbf{C}_2 \cdots \mathbf{C}_N] - \mathbf{Y}_1 \mathbf{A}$ is a matrix-matrix multiply with subtraction where the product $\mathbf{Y}_1 \mathbf{A}$
 202 computes r -length inner products.

203 We are now ready to triangularize the second block and update rows $r+1 : m$ of $[\mathbf{C}_3 \cdots \mathbf{C}_N]$, and so
 204 on. Algorithm 4 shows the pseudoalgorithm of the described procedure and performs approximately
 $1 - \mathcal{O}(1/N)$ fraction of FLOPs in Level-3 BLAS operations (see section 5.2.3 of [10]).

Algorithm 4: $\mathbf{Q}, \mathbf{R} \leftarrow \text{blockHQR}(\mathbf{A}, r)$: Perform Householder QR factorization of matrix
 \mathbf{A} with column partitions of size r .

Input: $\mathbf{A} \in \mathbb{R}^{m \times n}$, $r \in \mathbb{R}$ where $r < n$.
Output: \mathbf{Q}, \mathbf{R}

```

1  $N = \lceil \frac{n}{r} \rceil$ 
  // Let  $n_i = ri$  for  $i = 1 : N-1$  and  $n_N = n$ .
2 for  $i = 1 : N$  do
3    $\mathbf{V}_i, \beta_i, \mathbf{A}_{n_{i-1}+1:m, n_{i-1}+1:n_i} \leftarrow \text{hhQR}(\mathbf{A}_{n_{i-1}+1:m, n_{i-1}+1:n_i})$            /* Algorithm 2 */
4    $\mathbf{W}_i, \mathbf{Y}_i \leftarrow \text{buildWY}(\mathbf{V}_i, \beta_i)$                                            /* Algorithm 3 */
5   if  $i < N$  then
6      $\mathbf{A}_{n_i+1:m, n_i+1:n} -= \mathbf{Y}_i (\mathbf{W}_i^\top \mathbf{A}_{n_i+1:m, n_i+1:n})$            /* update the rest: BLAS-3 */
  //  $\mathbf{A}$  has been transformed into  $\mathbf{R} = \mathbf{Q}^\top \mathbf{A}$ .
  // Now build  $\mathbf{Q}$ .
7  $\mathbf{Q} \leftarrow \mathbf{I}$                                                                     /*  $\mathbf{I}_m$  if full QR, and  $\mathbf{I}_{m \times n}$  if thin QR. */
8 for  $i = N : -1 : 1$  do
9    $\mathbf{Q}_{n_{i-1}+1:m, n_{i-1}+1:n} = \mathbf{W}_i (\mathbf{Y}_i^\top \mathbf{Q}_{n_{i-1}+1:m, n_{i-1}+1:n})$            /* BLAS-3 */
10 return  $\mathbf{Q}, \mathbf{A}$ 

```

205

206 **3.3. Block HQR with partitioned rows : Tall-and-Skinny QR (TSQR).**

207 **4. Mixed-precision error analysis.**

208 **4.1. Round down at the end of the factorization.**

209 **4.2. Round down at block-level (BLAS-3).**

210 **4.3. Round down at inner-product level (BLAS-2).**

211 **5. Numerical Experiments.**

212 **6. Conclusion.** Though the use of lower precision naturally reduces the bandwidth and stor-
 213 age needs, the development of GPUs to optimize low precision floating point arithmetic have ac-
 214 celerated the interest in half precision and mixed-precision algorithms. Loss in precision, stability,

and representable range offset for those advantages, but these shortcomings may have little to no impact in some applications. It may even be possible to navigate around those drawbacks with algorithmic design.

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 HQR, 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 inner product routine, which can be applied to many other linear algebra tools 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 HQR 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, we experimented with introducing mixed-precision settings into graph partitioning problems. In particular, we applied DBSCAN to the spectral basis of a graph identified via subspace iteration that used our simulated mixed-precision HQR, which yielded clustering results tantamount to results from employing double-precision entirely.

Although this work is focused on QR factorizations and applications in spectral clustering, the mixed precision round-off error analysis can be applied to other tasks and applications that can benefit from employing low precision computations. While the emergence of technology that support low precision floats combats issues dealing with storage, now we need to consider how low precision affects stability of numerical algorithms.

Future work is 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.

REFERENCES

- [1] A. ABDELFAH, S. TOMOV, AND J. DONGARRA, *Fast batched matrix multiplication for small sizes using half-precision arithmetic on GPUs*, in 2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS), May 2019, pp. 111–122, <https://doi.org/10.1109/IPDPS.2019.00022>.
- [2] E. ANDERSON, Z. BAI, C. BISCHOF, L. S. BLACKFORD, J. DEMMEL, J. J. DONGARRA, J. DU CROZ, S. HAMMARLING, A. GREENBAUM, A. MCKENNEY, AND D. SORESENSEN, *LAPACK Users' Guide (Third Ed.)*, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1999; also available online from <http://www.netlib.org>.
- [3] J. APPLEYARD AND S. YOKIM, *Programming Tensor Cores in CUDA 9*, 2017, <https://devblogs.nvidia.com/programming-tensor-cores-cuda-9/> (accessed 2018-07-30).
- [4] C. BISCHOF AND C. VAN LOAN, *The WY Representation for Products of Householder Matrices*, SIAM Journal on Scientific and Statistical Computing, 8 (1987), pp. s2–s13, <https://doi.org/10.1137/0908009>.
- [5] M. COURBARIAUX, Y. BENGIO, AND J.-P. DAVID, *Training deep neural networks with low precision multiplications*, arXiv preprint, arXiv:1412.7024, (2014).
- [6] M. COURBARIAUX, J.-P. DAVID, AND Y. BENGIO, *Low precision storage for deep learning*, arXiv preprint arXiv:1412.7024, (2014).
- [7] J. DEMMEL, I. DUMITRIU, AND O. HOLTZ, *Fast linear algebra is stable*, Numerische Mathematik, 108 (2007), pp. 59–91, <https://doi.org/10.1007/s00211-007-0114-x>, <https://arxiv.org/abs/0612264>.
- [8] J. DEMMEL, L. GRIGORI, M. HOEMMEN, AND J. LANGOU, *Communication-optimal parallel and sequential QR and LU factorizations*, SIAM Journal on Scientific Computing, 34 (2012), <https://doi.org/10.1137/080731992>, <https://arxiv.org/abs/0808.2664>.
- [9] M. FAGAN, J. SCHLACHTER, K. YOSHII, S. LEYFFER, K. PALEM, M. SNIR, S. M. WILD, AND C. ENZ, *Overcoming the power wall by exploiting inexactness and emerging COTS architectural features: Trading precision for improving application quality*, in 2016 29th IEEE International System-on-Chip Conference (SOCC), Sep.

- 264 2016, pp. 241–246, <https://doi.org/10.1109/SOCC.2016.7905477>.
- 265 [10] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, JHU press, 4 ed., 2013.
- 266 [11] A. HAIDAR, A. ABDELFAH, M. ZOUNON, P. WU, S. PRANESH, S. TOMOV, AND J. DONGARRA, *The Design*
267 *of Fast and Energy-Efficient Linear Solvers: On the Potential of Half-Precision Arithmetic and Iterative*
268 *Refinement Techniques*, June 2018, pp. 586–600, https://doi.org/10.1007/978-3-319-93698-7_45.
- 269 [12] A. HAIDAR, S. TOMOV, J. DONGARRA, AND N. J. HIGHAM, *Harnessing GPU tensor cores for fast fp16 arithmetic*
270 *to speed up mixed-precision iterative refinement solvers*, in Proceedings of the International Conference
271 for High Performance Computing, Networking, Storage, and Analysis, SC '18, Piscataway, NJ, USA,
272 2018, IEEE Press, pp. 47:1–47:11, <https://doi.org/10.1109/SC.2018.00050>, [https://doi.org/10.1109/SC.](https://doi.org/10.1109/SC.2018.00050)
273 [2018.00050](https://doi.org/10.1109/SC.2018.00050).
- 274 [13] N. J. HIGHAM, *Accuracy and Stability of Numerical Methods*, 2002, <https://doi.org/10.2307/2669725>.
- 275 [14] N. J. HIGHAM AND T. MARY, *A New Approach to Probabilistic Rounding Error Analysis*, SIAM Journal on
276 Scientific Computing, 41 (2019), pp. A2815–A2835, <https://doi.org/10.1137/18M1226312>, [https://epubs.](https://epubs.siam.org/doi/10.1137/18M1226312)
277 [siam.org/doi/10.1137/18M1226312](https://doi.org/10.1137/18M1226312).
- 278 [15] N. J. HIGHAM AND S. PRANESH, *Simulating Low Precision Floating-Point Arithmetic*, SIAM Journal on Sci-
279 entific Computing, 41 (2019), pp. C585–C602, <https://doi.org/10.1137/19M1251308>, [https://epubs.siam.](https://epubs.siam.org/doi/10.1137/19M1251308)
280 [org/doi/10.1137/19M1251308](https://doi.org/10.1137/19M1251308).
- 281 [16] A. S. HOUSEHOLDER, *Unitary triangularization of a nonsymmetric matrix*, Journal of the ACM (JACM), 5
282 (1958), pp. 339–342.
- 283 [17] I. C. F. IPSEN AND H. ZHOU, *Probabilistic Error Analysis for Inner Products*, (2019), [http://arxiv.org/abs/](http://arxiv.org/abs/1906.10465)
284 [1906.10465](https://arxiv.org/abs/1906.10465), <https://arxiv.org/abs/1906.10465>.
- 285 [18] P. MICIKIEVICIUS, S. NARANG, J. ALBEN, G. DIAMOS, E. ELSSEN, D. GARCIA, B. GINSBURG, M. HOUSTON,
286 O. KUCHAIEV, G. VENKATESH, AND H. WU, *Mixed precision training*, in International Conference on
287 Learning Representations, 2018, <https://openreview.net/forum?id=r1gs9JgRZ>.
- 288 [19] R. SCHREIBER AND C. VAN LOAN, *A Storage-Efficient \$WY\$ Representation for Products of Householder*
289 *Transformations*, SIAM Journal on Scientific and Statistical Computing, 10 (1989), pp. 53–57, [https:](https://doi.org/10.1137/0910005)
290 [//doi.org/10.1137/0910005](https://doi.org/10.1137/0910005).
- 291 [20] G. TAGLIAVINI, S. MACH, D. ROSSI, A. MARONGIU, AND L. BENIN, *A transprecision floating-point platform for*
292 *ultra-low power computing*, in 2018 Design, Automation Test in Europe Conference Exhibition (DATE),
293 March 2018, pp. 1051–1056, <https://doi.org/10.23919/DATE.2018.8342167>.
- 294 [21] U. VON LUXBURG, *A tutorial on spectral clustering*, Statistics and Computing, 17 (2007), pp. 395–416, [https:](https://doi.org/10.1007/s11222-007-9033-z)
295 [//doi.org/10.1007/s11222-007-9033-z">//doi.org/10.1007/s11222-007-9033-z](https://doi.org/10.1007/s11222-007-9033-z), <https://doi.org/10.1007/s11222-007-9033-z>.