

MIXED - PRECISION BLOCK QR DECOMPOSITION ON GPU

STUDENTS: Jaidon Lybbert, Fulin Li, Mike Pao, Alice Lin, Shashank Shivashankar

Motivation

General QR Decomposition

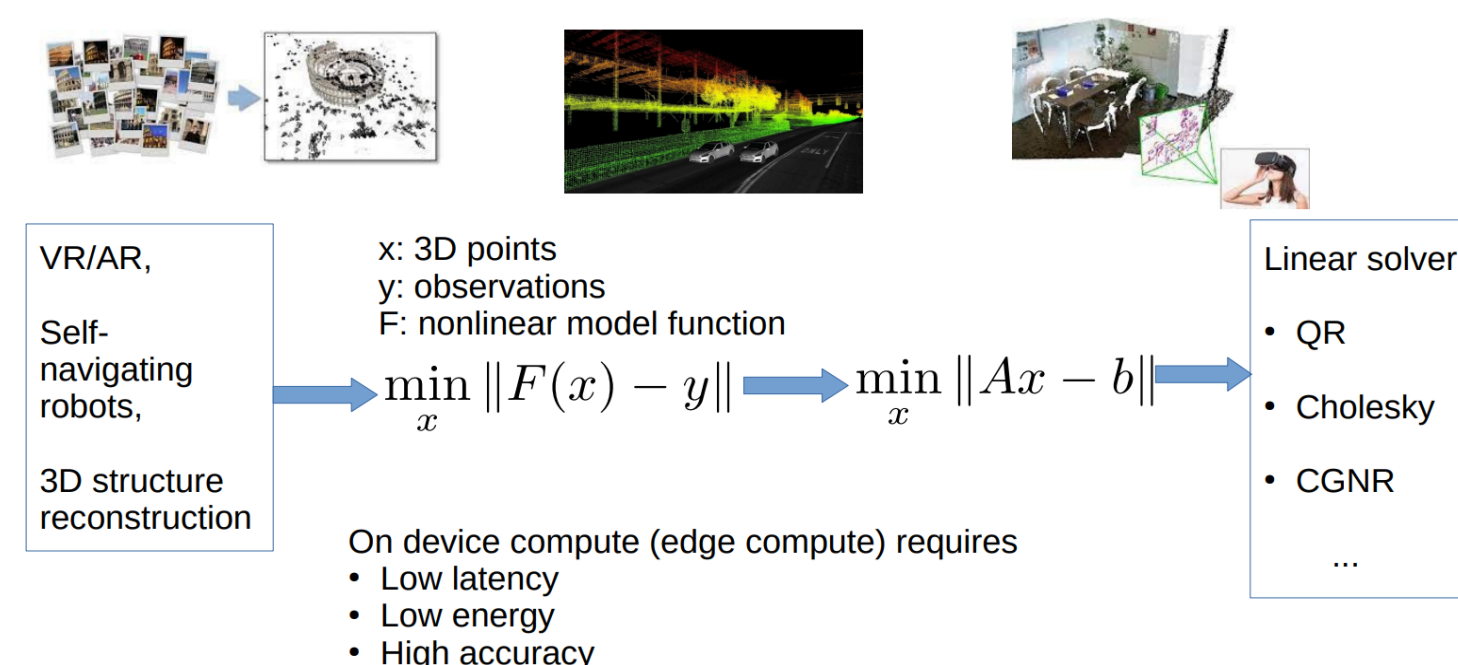
$$A = QR$$



Our project aims to optimize QR decomposition, which factorizes a matrix into an orthogonal matrix (Q) and an upper triangular matrix (R), by developing efficient algorithms and techniques to improve its computational efficiency and precision.

Application

QR decomposition is crucial for various applications in linear algebra and numerical computation.

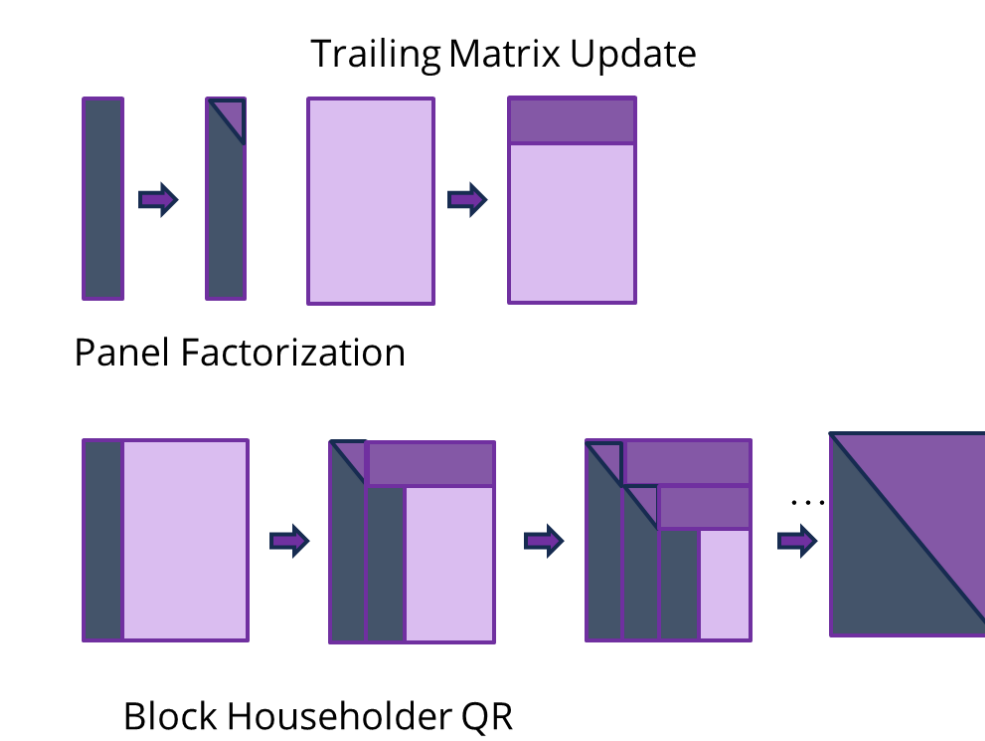


Software Design

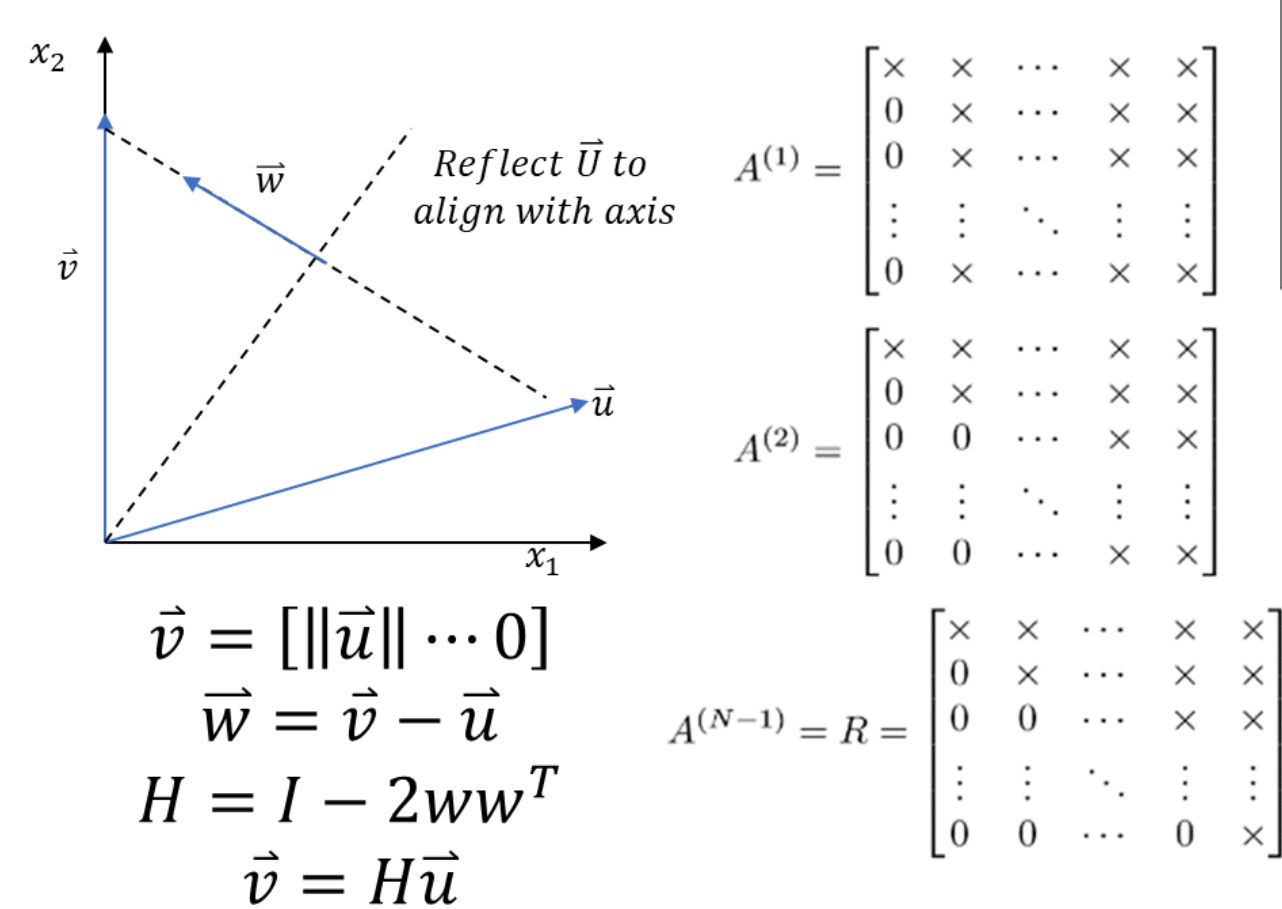
To address this challenge, our team implements several techniques to accelerate the QR decomposition process.

Block QR

Split the input matrix into smaller block matrices and perform QR decomposition on each block.

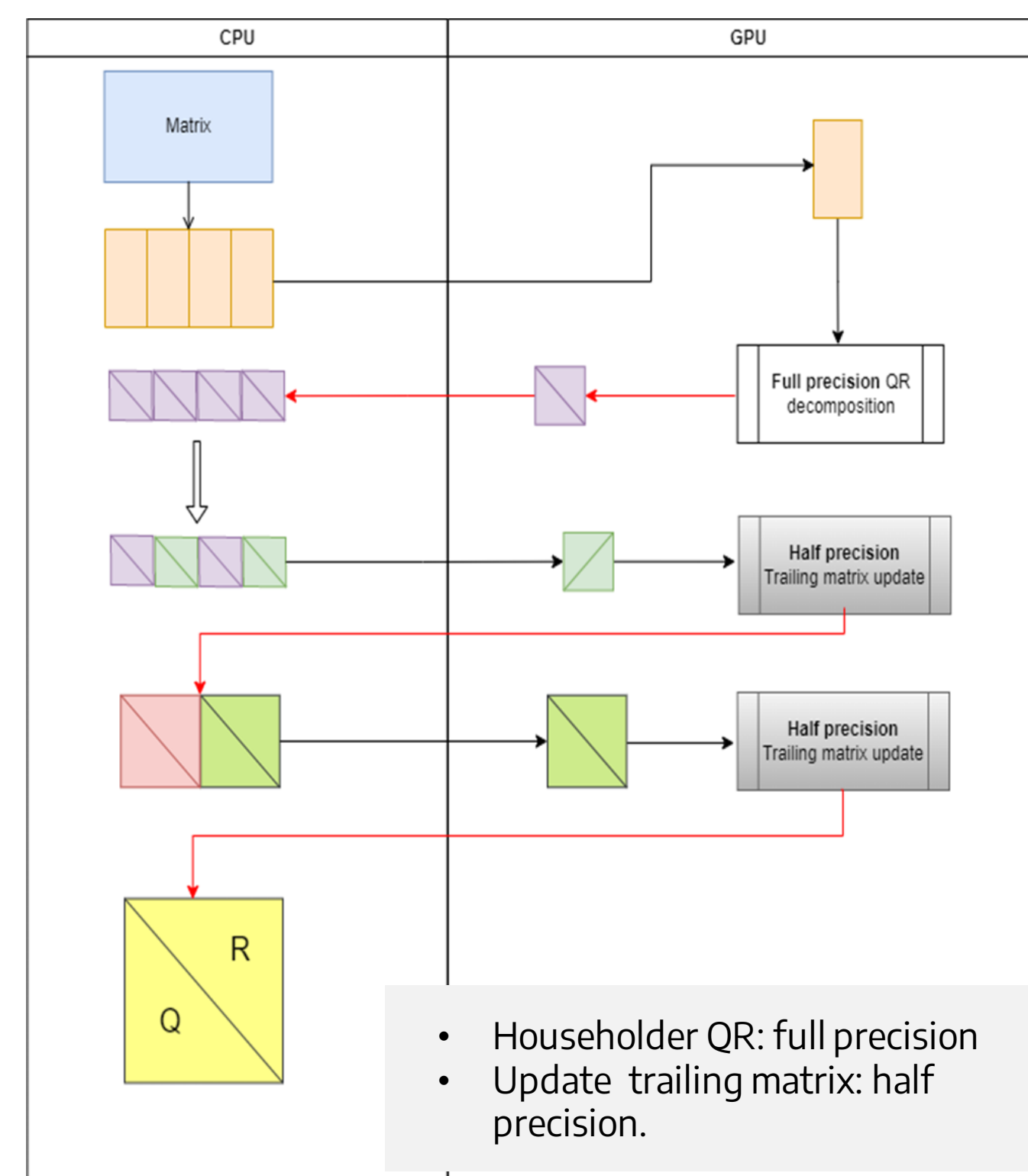


Householder QR



- Iterate over columns
 - Compute householder matrix to zero column below diagonal
 - Update matrix A by $A = HA$
- After iterating over all columns $A = R$

Mixed Precision



WY Transformation

Combine multiple householder transformations into a single matrix via the WY-representation of matrix products before doing the matrix update.

$Q \in R^{M \times M} = Q_1 Q_2 \dots Q_j \dots Q_n$
 where $Q_j = I_m - \beta_j w_j w_j^T$
 and the factors β_j, w_j are stored as

$V \in R^{M \times n} = [w_1 w_2 \dots w_j \dots w_n]$
 $B \in R^n = [\beta_1 \beta_2 \dots \beta_j \dots \beta_n]$

the W and Y factors such that $Q = I_m - WY^T$ can be calculated from V, and B.

Requirements

Targets

- Implement a fast and correct version of mixed precision QR using GPU.
- Implementation should use the GPU's half precision data type for the matrix multiplications and single precision for remaining operations.
- The computing accuracy and speed to be higher than that of a naive implementation on an x86 CPU architecture.
- Integration into an Open-Source Package.

Hardware Constraints

NVIDIA RTX 2080	
# of SMs	46
Threads/SM	1024
Blocks/SM	16
TC FMA dimension	16x16x16
TC / SM	8
SMEM	64KB
RF	4 * 64kB

Software Dependencies

- CUDA Toolkit:**
 - Version: CUDA Toolkit 9.1 or higher
 - Provides development tools and libraries for GPU programming
- NVIDIA GPU Driver:**
 - Required for NVIDIA RTX 2080 hardware compatibility

CUDA programming

Procedures

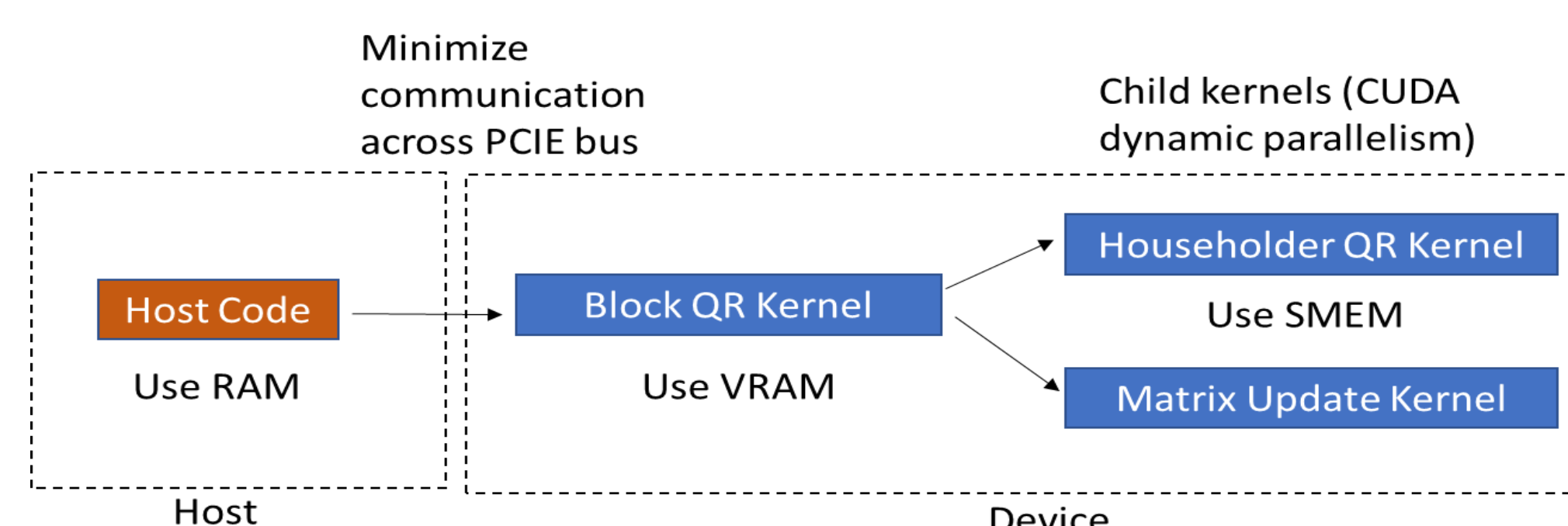
Divide the matrix to N blocks. N is a power of 2.

Send every block to the GPU core, apply full precision **Householder QR** in every block.

Combine the panel Q matrices from each block into a single matrix Q using the **WY transformation**.

Get all QR deposition of every block, then recursively **update the trailing matrix**.

(Requires large amounts of computation)



Parallel computing in multiple GPU cores.

To enhance efficiency, half precision is used to accelerate the process.

Algorithm

Algorithm 1 Calculate $A = QR$ using Householder reflections

```

1: for k = 1 to n do
2:    $u = A_{k:m,k}$ 
3:    $v_k = \text{sign}(u_1) \|u\|_2 e_1 + u$ 
4:    $v_k = v_k / \|v_k\|_2$ 
5:    $A_{k:m,k:n} = A_{k:m,k:n} - 2v_k v_k^T A_{k:m,k:n}$ 
6: end for
  
```

Algorithm 2 Calculate W, Y from the factored form of Q: V and B

```

1:  $Y = w_1$ 
2:  $W = \beta_1 w_1$ 
3: for j = 2 to r do
4:    $z = \beta_j (I_m - WY^T) w_j$ 
5:    $W = [W | z]$ 
6:    $Y = [Y | w_j]$ 
7: end for
  
```

Algorithm 3 Block Householder QR Decomposition

```

1:  $Q = I_m$ 
2:  $\lambda = 1$ 
3:  $k = 0$ 
4: while  $\lambda \leq n$  do
5:    $\tau \leftarrow \min(\lambda + r - 1, n)$ 
6:    $k = k + 1$ 
7:    $A_{\lambda:m,\lambda:r} \leftarrow \text{Householder\_qr}(A_{\lambda:m,\lambda:r})$ 
8:    $W_k, Y_k \leftarrow \text{WY\_transform}(Y_k)$ 
9:    $A_{\lambda:m,\tau+1:n} = (I - W_k Y_k^T) A_{\lambda:m,\tau+1:n}$ 
10:   $Q_{:, \lambda:m} = Q_{:, \lambda:m} (I - W_k Y_k^T)$ 
11:   $\lambda = \tau + 1$ 
12: end while
  
```

Performance results

Error criteria

$$\begin{aligned} \|QR - A\| &\leq m \cdot 2^{-23} \|A\| \\ \|Q^T Q - I\| &\leq m \cdot 2^{-23} \\ \|L\| &\leq m \cdot 2^{-23} \end{aligned}$$

* L: the trapezoidal submatrix below the main diagonal of R

By utilizing GPU parallel computing, we can perform several tasks at once, leading to significant speedup compared to traditional sequential CPU processing.

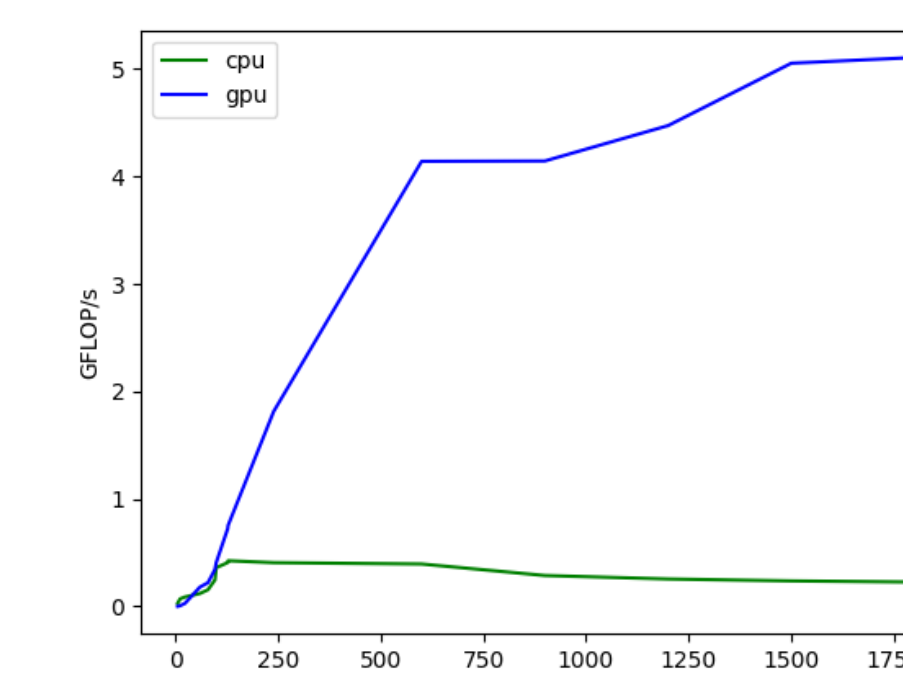


Figure 2:
Sustained performance of QR decomposition

<https://projects.asl.ethz.ch/datasets/doku.php?id=kmavvisualinertialdatasets>

Input Matrix A: Random Float
Matrix Sizes: Various Sizes

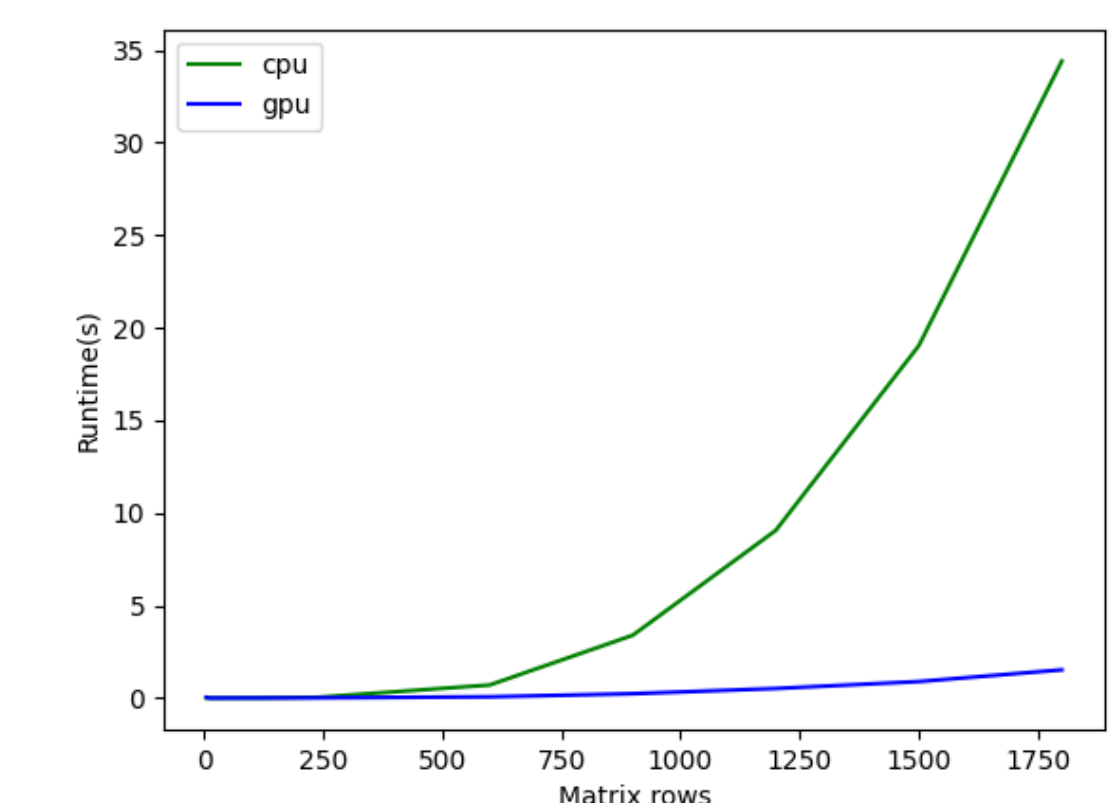


Figure 1:
Runtimes of QR decomposition

Conclusion

- By implementing the block QR decomposition in parallel on a GPU, we saw a speedup of **over 10x** compared to the sequential CPU implementation for large matrices, meeting our success criteria for execution time
- Our performance bottleneck is in the construction of matrix Q, which takes **about 80%** of the execution time, this *can be accelerated* on the GPU

Future Work, References, and Acknowledgments

- Implement tiled QR to enhance the parallelism and performance of the QR decomposition process.
- Investigate and incorporate any other relevant advancements or optimizations in the field to enhance the overall algorithm.

References

- Zhang, S., Baharlouei, E., & Wu, P. (2020). High Accuracy Matrix Computations on Neural Engines: A Study of QR Factorization and its Applications. Proceedings of the 29th International Symposium on High-Performance Parallel and Distributed Computing, 17–28. <https://doi.org/10.1145/3369583.3392685> | PDF
- Bouwmeester, H., Mathias Jacquelin, Langou, J., & Yves, R. (2011). Tiled QR factorization algorithms. arXiv.org. <https://arxiv.org/pdf/1104.4475.pdf>