



Parallel QR Factorization

Wangshu Li (wangshul)

Jiarui Zhang (jiaruiz)

Description

QR factorization is a decomposition of a non-singular matrix A into a product $A = QR$ of an orthogonal matrix Q and an upper triangular matrix R

$$\begin{array}{c} A \\ \left[\begin{array}{c|c|c} | & | & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \mathbf{a}_3 \\ | & | & | \end{array} \right] \end{array} = \begin{array}{c} Q \\ \left[\begin{array}{c|c|c} | & | & | \\ \mathbf{e}_1 & \mathbf{e}_2 & \mathbf{e}_3 \\ | & | & | \end{array} \right] \end{array} \begin{array}{c} R \\ \left[\begin{array}{ccc} \mathbf{e}_1^T \cdot \mathbf{a}_1 & \mathbf{e}_1^T \cdot \mathbf{a}_2 & \mathbf{e}_1^T \cdot \mathbf{a}_3 \\ 0 & \mathbf{e}_2^T \cdot \mathbf{a}_2 & \mathbf{e}_2^T \cdot \mathbf{a}_3 \\ 0 & 0 & \mathbf{e}_3^T \cdot \mathbf{a}_3 \end{array} \right] \end{array}$$

orthogonal unit vector

Upper Diagonal matrix

Sequential Algorithm

Algorithm 5: QR Decomposition.

input : A non-singular matrix \mathbf{A} of size $n \times n$

output: Orthogonal matrix \mathbf{Q} , upper triangular form \mathbf{U}

Function QRFact(A) is

$N \leftarrow$ order of the square matrix A ;

$\mathbf{u} \leftarrow \mathbf{0}_{N \times N}$, $\mathbf{R} \leftarrow \mathbf{0}_{N \times N}$;

for $j \leftarrow 0$ **to** $n - 1$ **do**

$\mathbf{w}_j \leftarrow A[:, j]$;

 /* Squared norm of \mathbf{w}_j

$\text{sum} \leftarrow 0$;

for $k \leftarrow 0$ **to** $n - 1$ **do**

$\text{sum} \leftarrow \text{sum} + \mathbf{w}_j[k] \times \mathbf{w}_j[k]$;

$|\mathbf{w}_j|^2 \leftarrow \text{sum}$;

*/

 /* Compute the coefficients $\mathbf{R}[i][j]$

$\text{sum} \leftarrow 0$;

for $i \leftarrow 0$ **to** $j - 1$ **do**

$\mathbf{R}[i][j] \leftarrow \mathbf{w}_j \cdot \mathbf{u}[i]$;

$\text{sum} \leftarrow \text{sum} + \mathbf{R}[i][j] \times \mathbf{R}[i][j]$;

$\mathbf{R}[j][j] \leftarrow (|\mathbf{w}_j|^2 - \text{sum})^{1/2}$;

*/

 /* Compute vector $\mathbf{u}[j]$ using fwd substitution

$\text{sumvec} \leftarrow \mathbf{0}_{1 \times N}$;

for $i \leftarrow 0$ **to** $j - 1$ **do**

$\text{sumvec} \leftarrow \text{sumvec} + \mathbf{R}[i][j] \times \mathbf{u}[i]$;

$\mathbf{u}[j] \leftarrow (\mathbf{w}_j - \text{sumvec}) / \mathbf{R}[j][j]$;

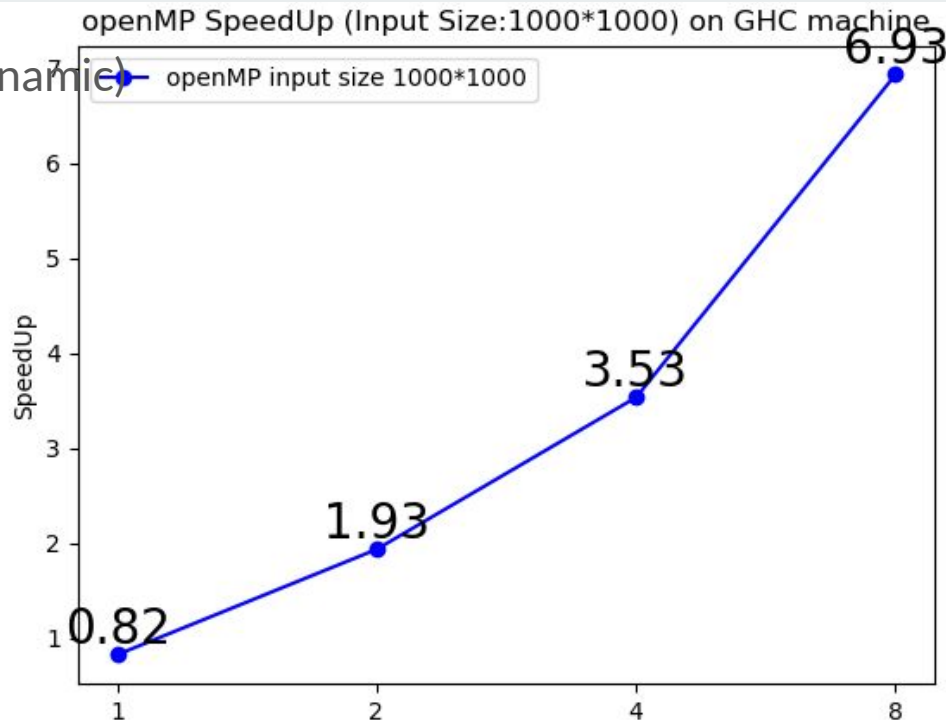
*/

$\mathbf{Q} \leftarrow \mathbf{u}^T$;

return \mathbf{Q}, \mathbf{R} ;

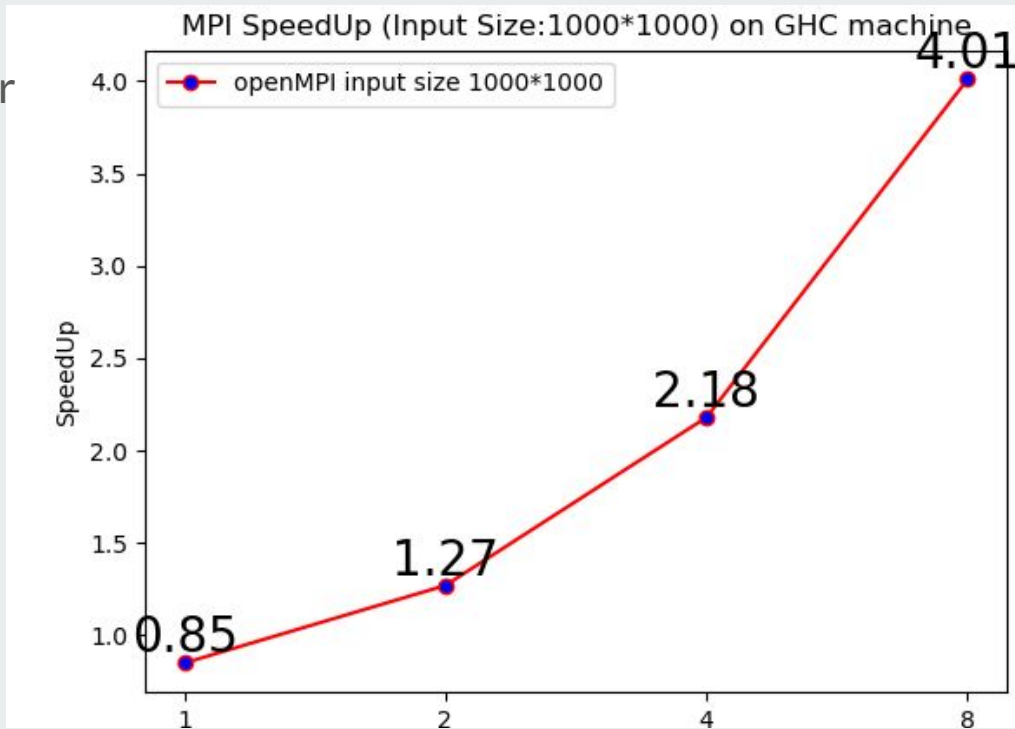
OpenMP

- `#pragma omp parallel for schedule (dynamic)`
 - Parallelize the outer loop
- `#pragma omp barrier`
 - Synchronize between threads
- Almost linear speedup



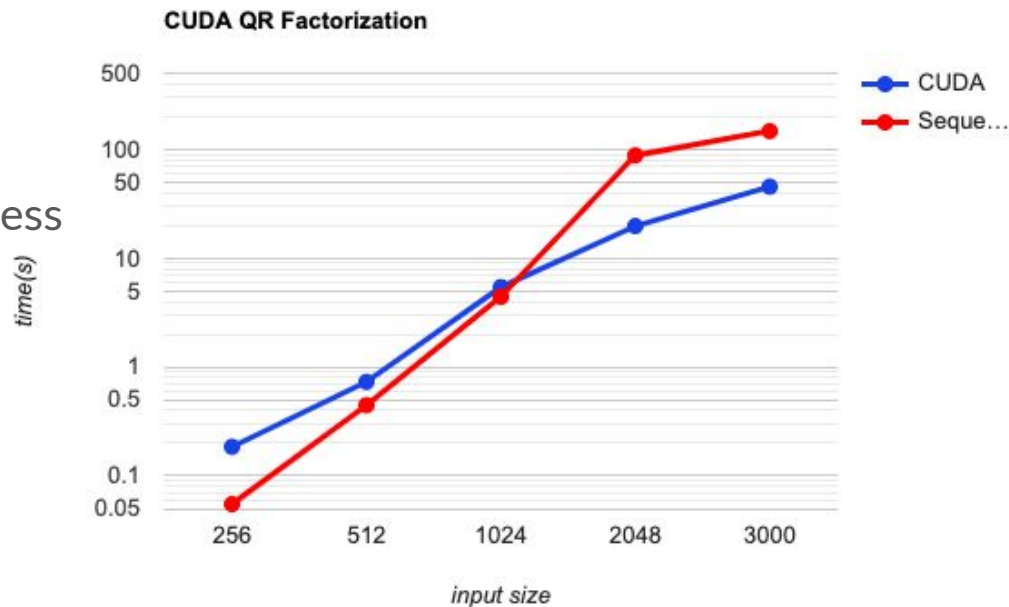
OpenMPI

- Each process works for $1/N$ number of total columns
 - Parallelism over columns
- Processes send its columns to all processes with higher pid
 - High communication cost
- Use a 2-d buffer for receiving
 - Workflow pipeline



CUDA

- Each thread works for a single cell
 - Parallelism over cells
- Each column launches the kernel
 - Guarantee dependency/correctness
- Shared block memory
 - Reduce I/O latency
- Synchronization within kernel
- Bad when input size is small
 - Data dependency



Analytics



- Data dependency due to the nature of QR factorization
 - Each column needs previous columns for computation
- OpenMP has the best performance
 - Shared address space model has minimal communication cost
 - No synchronization is required between threads because there is no data race
- CUDA and MPI have worse performance
 - Too many kernel functions when parallelizing over cells
 - CUDA is not appropriate for dependent workloads
 - OpenMPI has high communication cost