

Second assignment

Cholesky factorization

Benatti Sebastiano, Carpi Samuele,
Pasquali Mattia

Cholesky decomposition factors a symmetric, positive-definite matrix **A** into

$$\mathbf{A} = \mathbf{L}\mathbf{L}^T$$

where **L** is a **lower triangular matrix** and **L^T** is the **transpose of L**. It provides a numerically efficient way to solve linear systems and is widely used in optimization, statistics, and numerical analysis.

Diagonal elements

$$l_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} l_{ik}^2}$$

Off-diagonal elements

$$l_{ij} = \frac{1}{l_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} \right)$$

```
int i, j, k;
DATA_TYPE x;
for (i = 0; i < _PB_N; ++i)
{
    x = A[i][i];
    for (j = 0; j <= i - 1; ++j){
        x -= A[i][j] * A[i][j];
    }
    p[i] = 1.0 / sqrt(x);

    for (j = i + 1; j < _PB_N; ++j)
    {
        x = A[i][j];
        for (k = 0; k <= i - 1; ++k)
            x -= A[j][k] * A[i][k];
        A[j][i] = x * p[i];
    }
}
```

OPTIMIZATION APPROACHES

- UVM
- Pinned memory
- Streams
- **Shared memory**
- **Stride**
- **Tiling**

OUR BEST SOLUTION CODE

Sequential

```
int i, j, k;  
DATA_TYPE x;  
for (i = 0; i < _PB_N; ++i)  
{  
    x = A[i][i];  
    for (j = 0; j <= i - 1; ++j){  
        x -= A[i][j] * A[i][j];  
    }  
    p[i] = 1.0 / sqrt(x);  
}
```

```
for (j = i + 1; j < _PB_N; ++j)  
{  
    x = A[i][j];  
    for (k = 0; k <= i - 1; ++k)  
        x -= A[j][k] * A[i][k];  
    A[j][i] = x * p[i];  
}
```

compute_diagonal

Kernel 1

compute_column

Kernel 2

Sequential

```
int i, j, k;
DATA_TYPE x;
for (i = 0; i < _PB_N; ++i)
{
    x = A[i][i];
    for (j = 0; j <= i - 1; ++j){
        x -= A[i][j] * A[i][j];
    }
    p[i] = 1.0 / sqrt(x);

    for (j = i + 1; j < _PB_N; ++j)
    {
        x = A[i][j];
        for (k = 0; k <= i - 1; ++k)
            x -= A[j][k] * A[i][k];
        A[j][i] = x * p[i];
    }
}
```

Optimized

```
1  __global__ void compute_diagonal(DATA_TYPE* __restrict__ p, DATA_TYPE* __restrict__ A, int n, int i){
2      __shared__ DATA_TYPE sharedSum[BLOCK_SIZE];
3      int tid = threadIdx.x;
4
5      // 1) ogni thread calcola una somma parziale in parallelo con stride (per i > BLOCK_SIZE)
6      DATA_TYPE localSum = 0.0;
7      for (int k=tid; k<i; k+=blockDim.x){
8          localSum += A[i*n+k]*A[i*n+k];
9      }
10     sharedSum[tid] = localSum;
11     __syncthreads();
12
13     // 2) somma parziale parallela riducendo da BLOCK_SIZE a 32 valori finali (warp size per sincronizzazione implicita)
14     if (tid<32){
15         DATA_TYPE partialSum = 0.0;
16         for (int t=tid; t<blockDim.x; t+=32){
17             partialSum += sharedSum[t];
18         }
19         sharedSum[tid] = partialSum;
20     }
21
22     // 3) somma finale dei 32 valori rimanenti sequenzialmente nel thread 0
23     if (tid == 0){
24         DATA_TYPE sum = 0.0;
25         for (int t=0; t<32; t++){
26             sum += sharedSum[t];
27         }
28         p[i] = 1.0/sqrt(A[i*n+i] - sum);
29     }
30 }
```


Optimized

Sequential

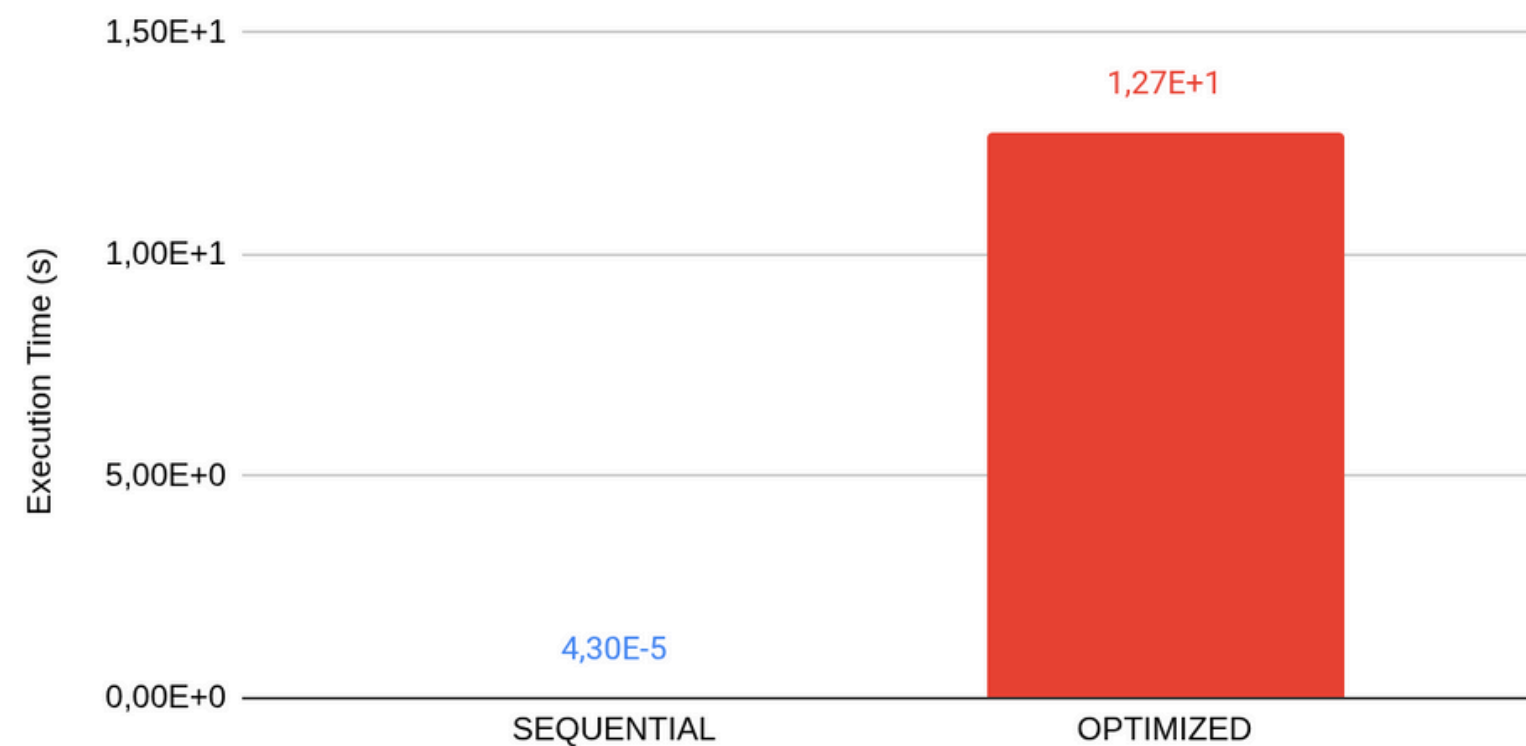
```
int i, j, k;
DATA_TYPE x;
for (i = 0; i < _PB_N; ++i)
{
    x = A[i][i];
    for (j = 0; j <= i - 1; ++j){
        x -= A[i][j] * A[i][j];
    }
    p[i] = 1.0 / sqrt(x);

    for (j = i + 1; j < _PB_N; ++j)
    {
        x = A[i][j];
        for (k = 0; k <= i - 1; ++k)
            x -= A[j][k] * A[i][k];
        A[j][i] = x * p[i];
    }
}
```

```
1 __global__ void compute_column(DATA_TYPE* __restrict__ p, DATA_TYPE* __restrict__ A, int n, int i){
2     int j = blockIdx.x*blockDim.x+threadIdx.x; // Ogni thread elabora la riga j
3     int tid = threadIdx.x;
4     __shared__ DATA_TYPE sharedPivotRow[BLOCK_SIZE]; // Shared memory per il tile della riga pivot A[i][0..i-1]
5     __shared__ DATA_TYPE sharedP; // Shared memory per p[i]
6
7     // 1) Thread 0 carica p[i] in shared memory
8     if (tid == 0){
9         sharedP = p[i];
10    }
11
12    // 2) x = A[i][j] inizializzato per il triangolo inferiore
13    DATA_TYPE x = 0.0;
14    if (j>i && j<n){
15        x = A[i*n+j];
16    }
17
18    // Tiling della riga pivot di BLOCK_SIZE elementi.
19    for (int tile=0; tile<i; tile+=BLOCK_SIZE){
20
21        // 3) Caricamento del tile in shared memory
22        int k_curr = tile+tid; // Indice globale dell'elemento da caricare
23        if (k_curr<i){
24            sharedPivotRow[tid] = A[i*n+k_curr];
25        }else{
26            sharedPivotRow[tid] = 0.0; // Padding per evitare accessi fuori limite
27        }
28        __syncthreads(); // Aspetta che tutti i thread abbiano caricato il loro elemento
29
30        // 4) Calcolo del prodotto scalare parziale
31        if (j>i && j<n){
32            int pivotRowSize = min(BLOCK_SIZE,i-tile); // Elementi validi nel tile corrente
33            for (int k=0; k<pivotRowSize; ++k){
34                x -= A[j*n+(tile+k)]*sharedPivotRow[k]; // A[j][k] in global memory, A[i][k] in shared memory
35            }
36        }
37        __syncthreads(); // Aspetta prima di sovrascrivere sharedPivotRow nel prossimo tile
38    }
39
40    // 5) Scrittura del risultato finale
41    if (j>i && j<n){
42        A[j*n+i] = x*sharedP;
43    }
44 }
```

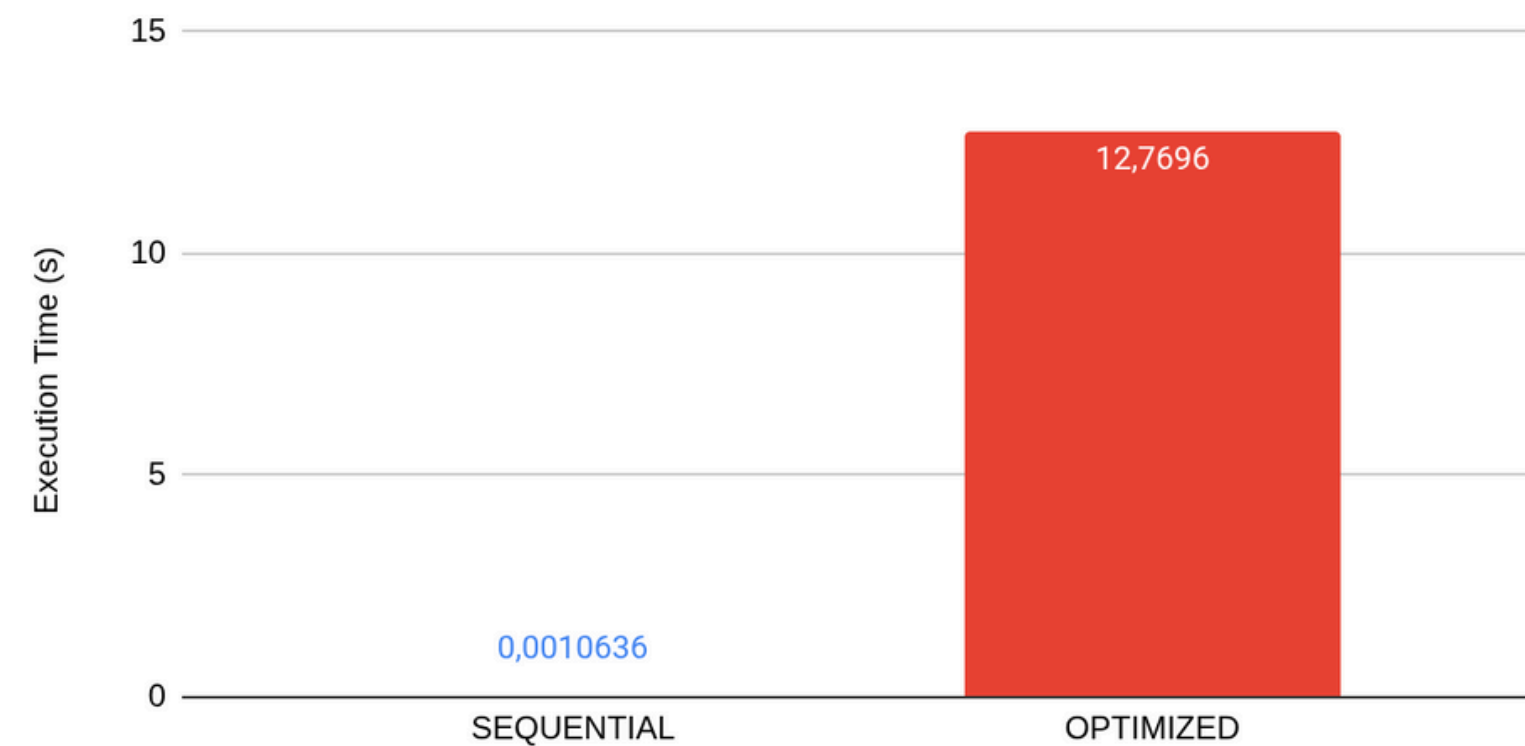

EXPERIMENTAL RESULTS

Mini



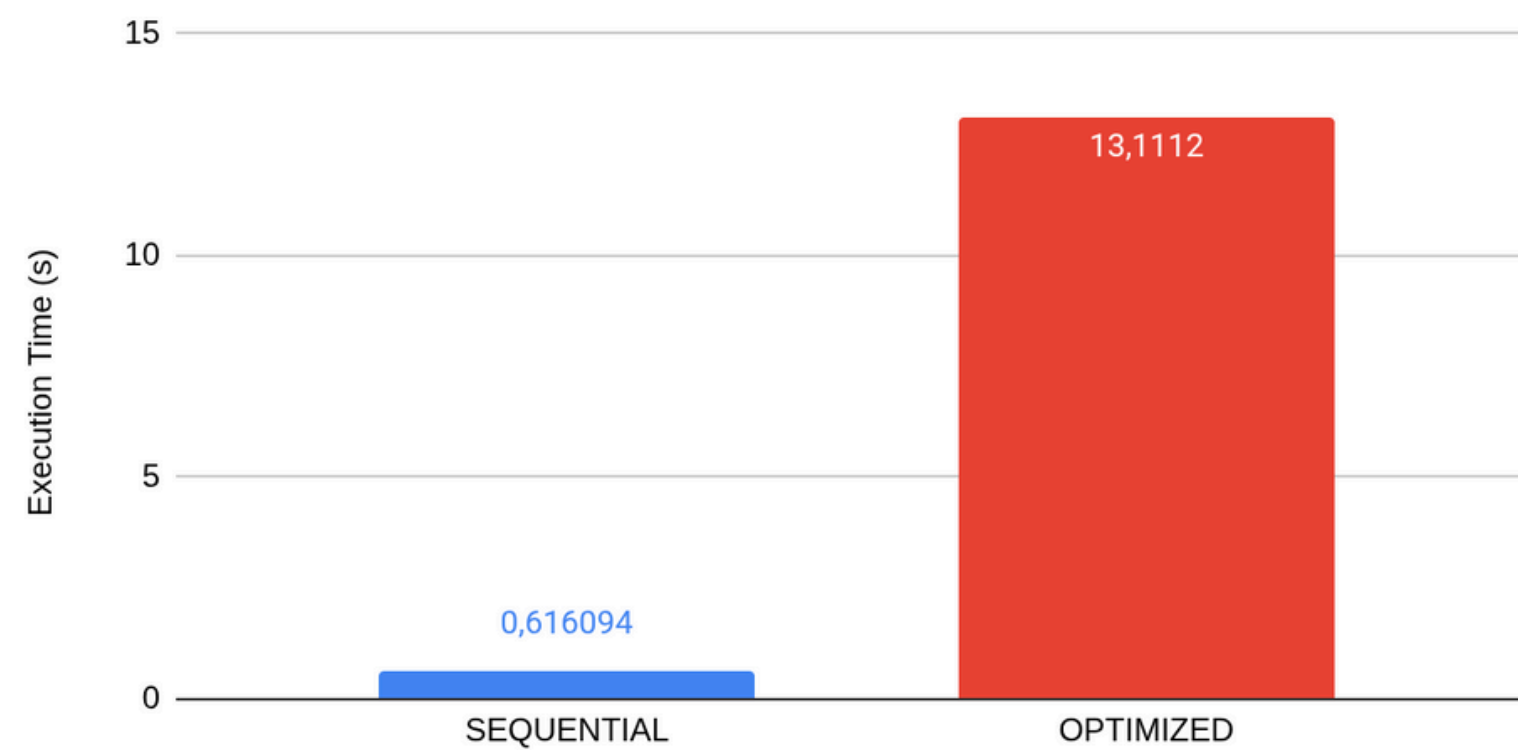
Optimization

Small



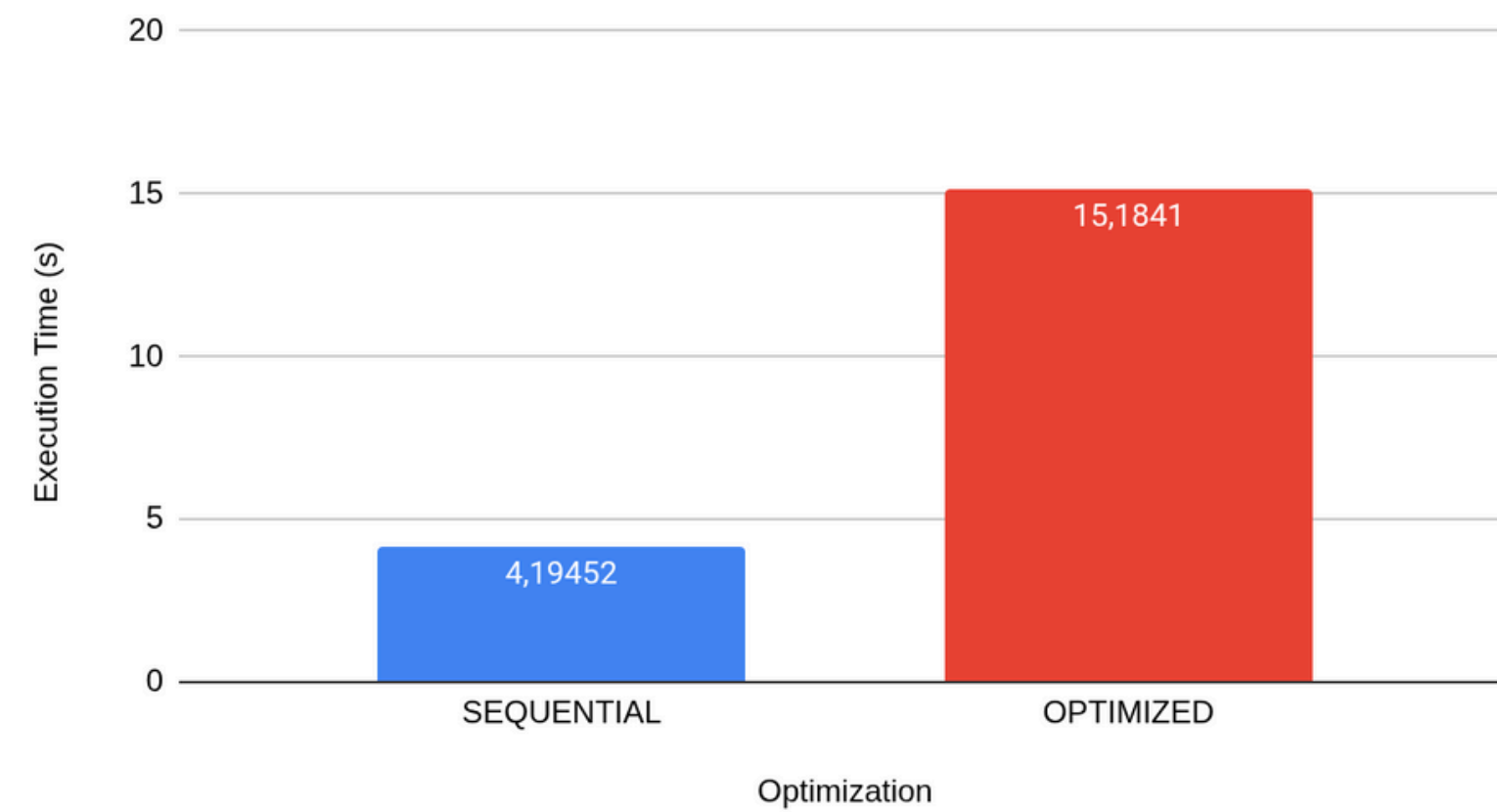
Optimization

Standard

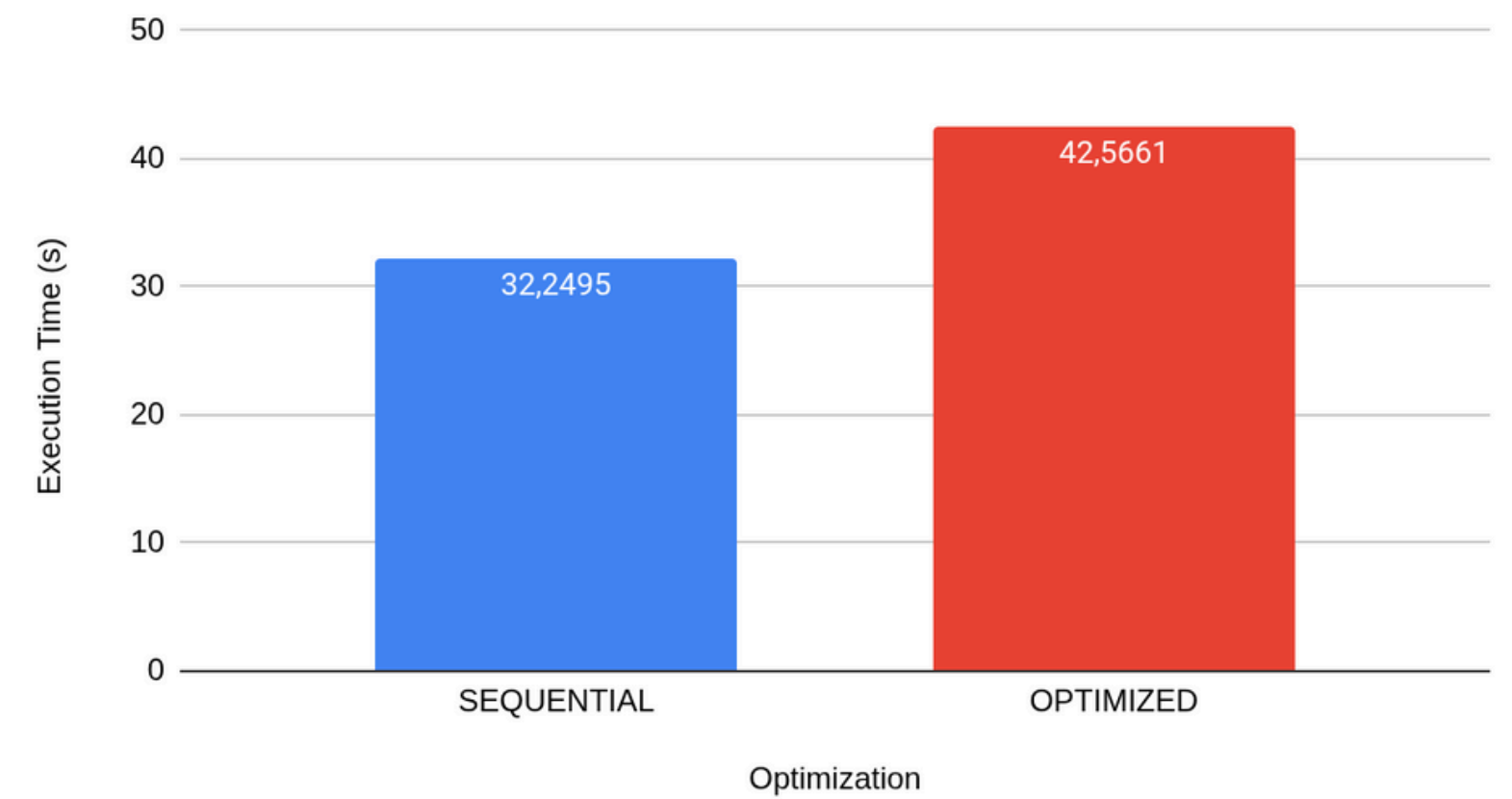


Optimization

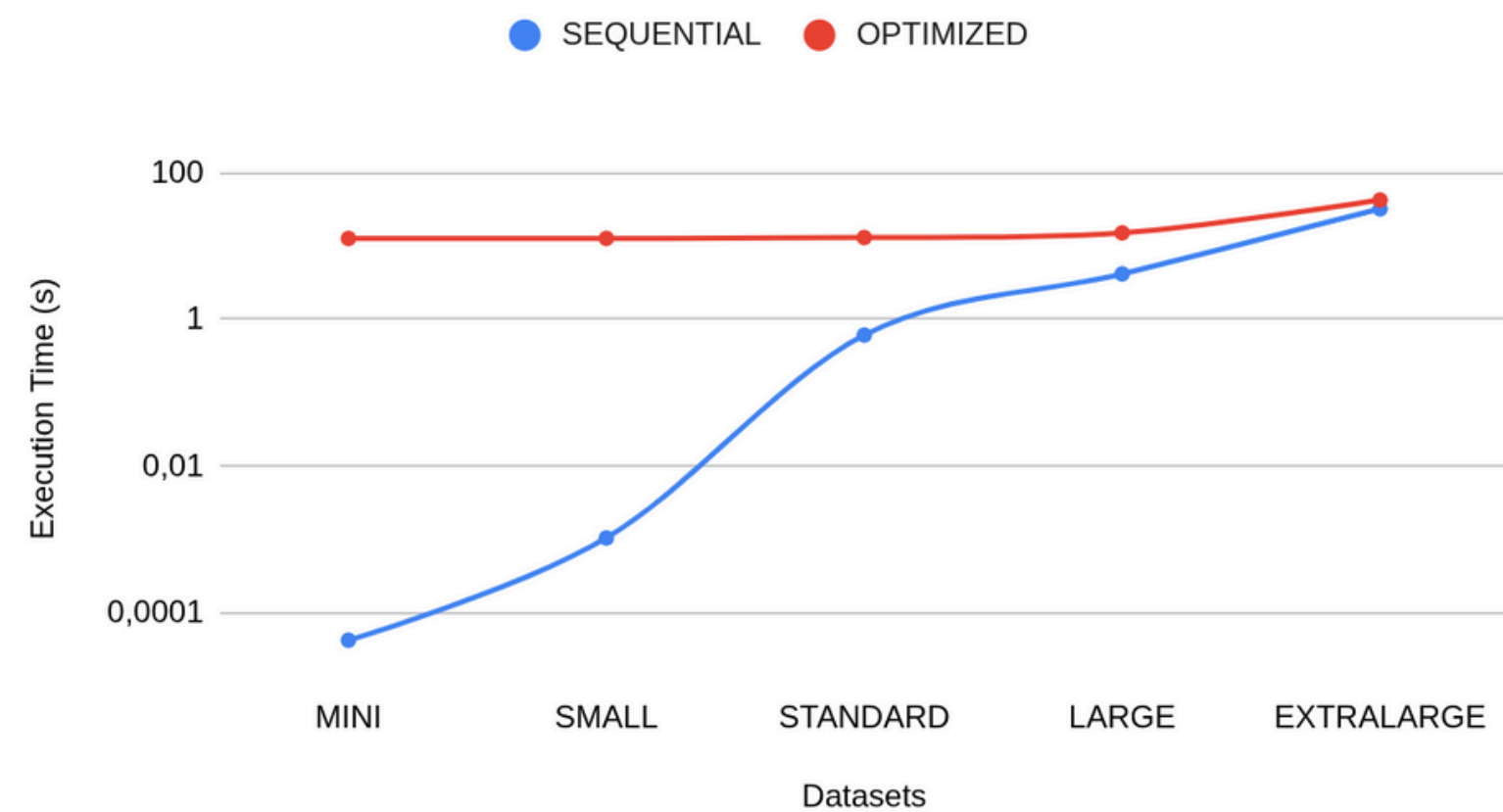
Large



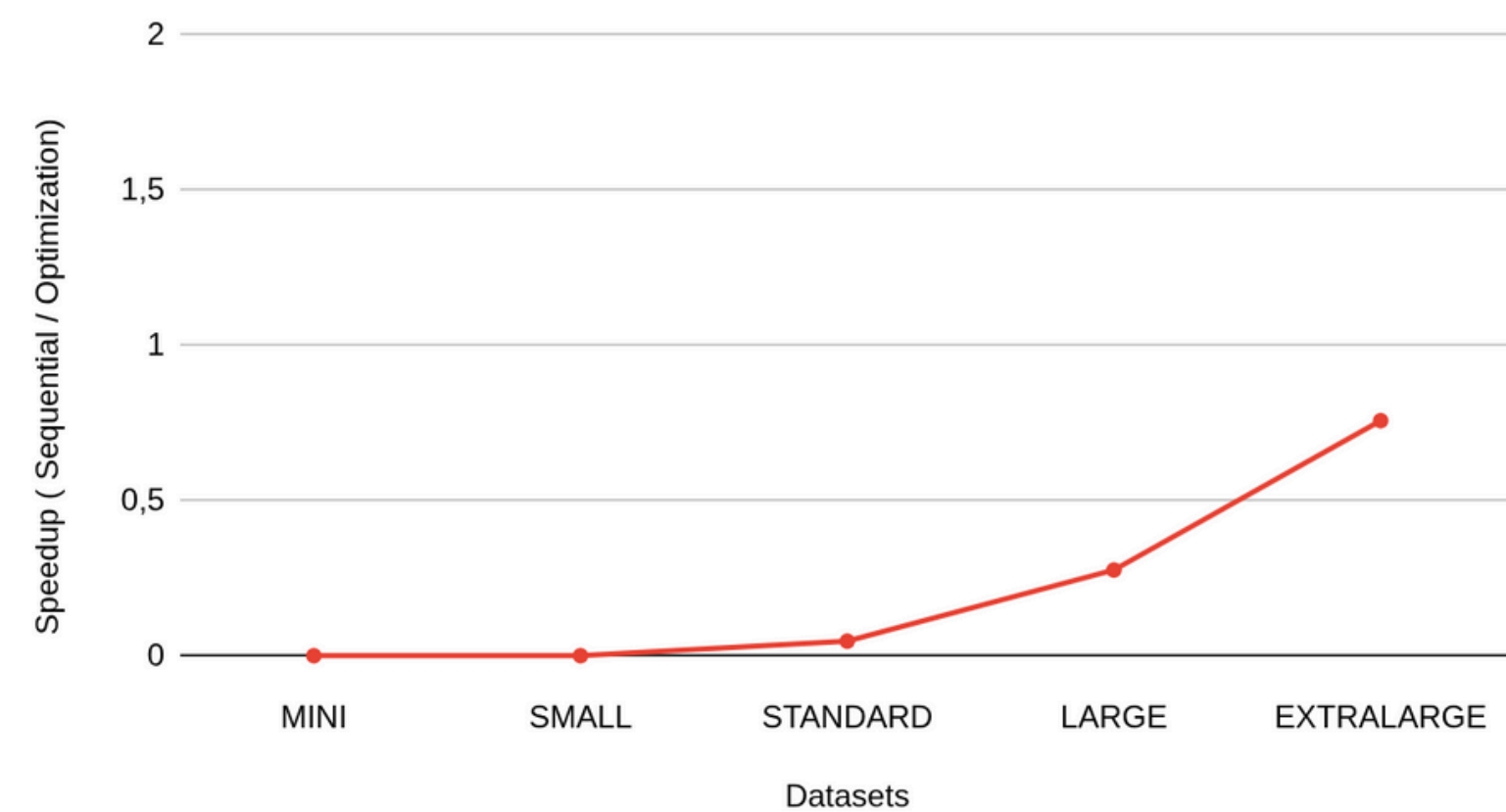
Extralarge



Execution Times



Speedup



CONCLUSIONS

- It **does not** make sense use **CUDA** for Cholesky Algorithm because of **data dependencies** between for loops
- But if the matrix is **large enough**, the execution time would be **similar** or better than the **sequential** version