COMPUTER ENGINEERING MASTER DEGREE

COMPUTER ARCHITECTURE

PROJECT DISCUSSION

# INTEGER FACTORIZATION

## PROFESSORS

ANTONIO COSIMO PRETE
ANTONIO DI TECCO
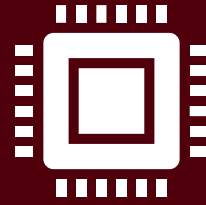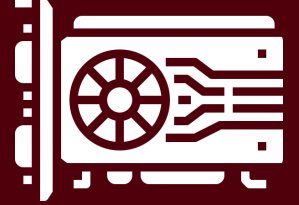
## GROUP MEMBERS

GIOVANNI LIGATO
ALICE ORLANDINI

*IN SUPREMÆ DIGNITATIS*
*1343*

University of Pisa

# SYLLABUS
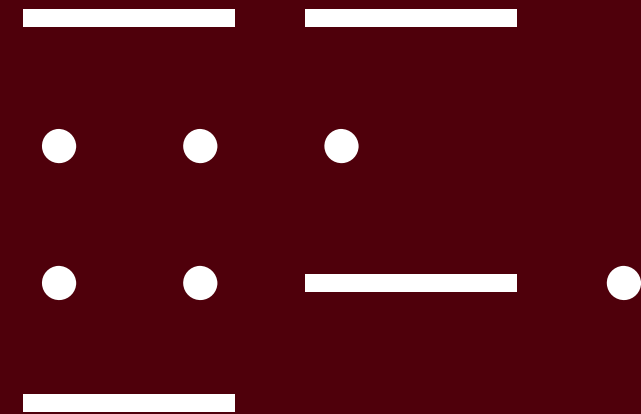
**ALGORITHM** > **GOALS** > **CPU** > **GPU**

# ALGORITHM

# INTEGER
## FACTORIZATION

---

«Every positive integer can be written uniquely as a product of primes»

---

MATHEMATICAL FORMULATION

## TRIAL DIVISION
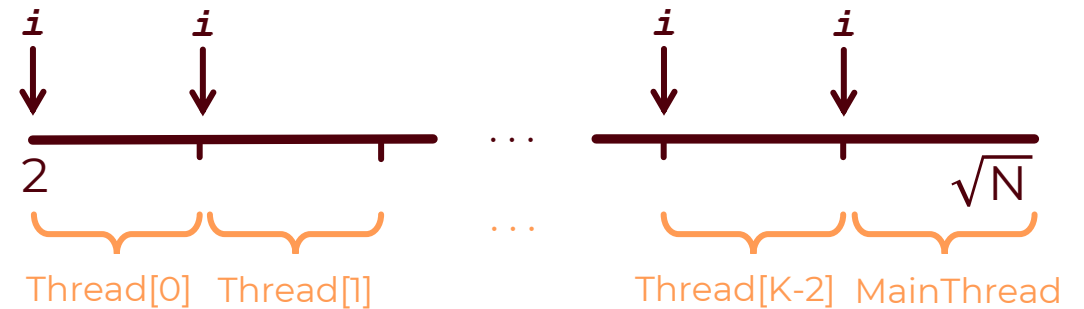### [PARALLEL]

**INPUT**

**N :** Number to be Factorized
**K :** Number of Threads

$i$     $i$        $i$     $i$

2        ...        $\sqrt{N}$

...

Thread[0]   Thread[1]      Thread[K-2]   MainThread

**OUTPUT**

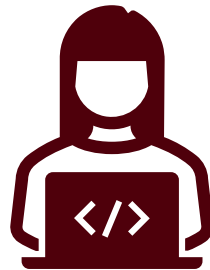**Primes :** Product of primes  🔒 Mutex

# GOALS

# GOALS



**END USER**

Execution Time
≤ **1** second

For numbers up to
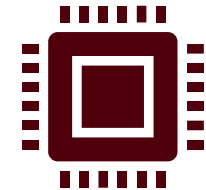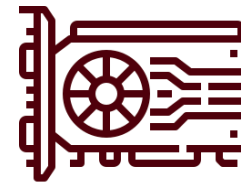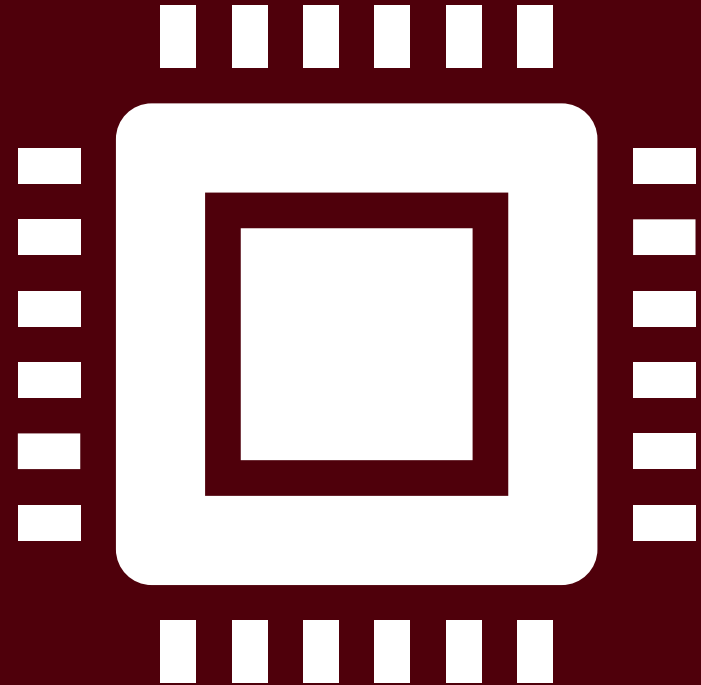**18** DIGITS

e.g. **975734686214396237**

**DEVELOPER**

Develop a **GPU** version that **outperforms** the **CPU** version
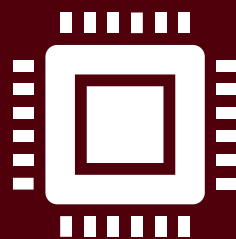
# CPU
IMPLEMENTATION

# 11th Gen **intel** Core(TM) **i5**-11400

**x** Cores per socket: **6**

**=** Threads per core: **2**

TOTAL # LOGICAL CORES: **12**

L1**d** Cache: **288** KiB

L1**i** Cache: **192** KiB

L2 Cache: **3** MiB

L3 Cache: **12** MiB

Base Frequency: **2.60** GHz        Max Turbo Frequency: **4.40** GHz

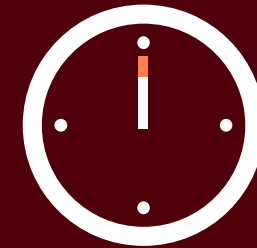# **HARDWARE** DETAILS

# TOOLS ⚒

## 📐 Execution Time

**Wall-Clock Time**

```
chrono::steady_clock::time_point start = chrono::steady_clock::now();

parallelTrialDivision(NUMBER, NUM_THREADS);

chrono::steady_clock::time_point end = chrono::steady_clock::now();


chrono::milliseconds duration = chrono::duration_cast<chrono::milliseconds>(end - start);
```
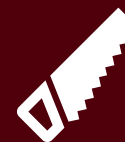
**start**        **end**

🔧 Intel **VTune** Profiler        🪚 Visual Studio Code

# FIRST EXECUTION

**975734686214396237** = 748609 * 1303396948493

**18** DIGITS

EXECUTION TIME
**568340**
MILLISECONDS

≈

**9**
MINUTES

>>

**1**
SECOND

⚠ REQUIREMENT
NOT
SATISFIED

LET'S GO INTO
DETAIL!

# EXECUTION TIME
**FIRST** VERSION

#ITERATIONS : **30**
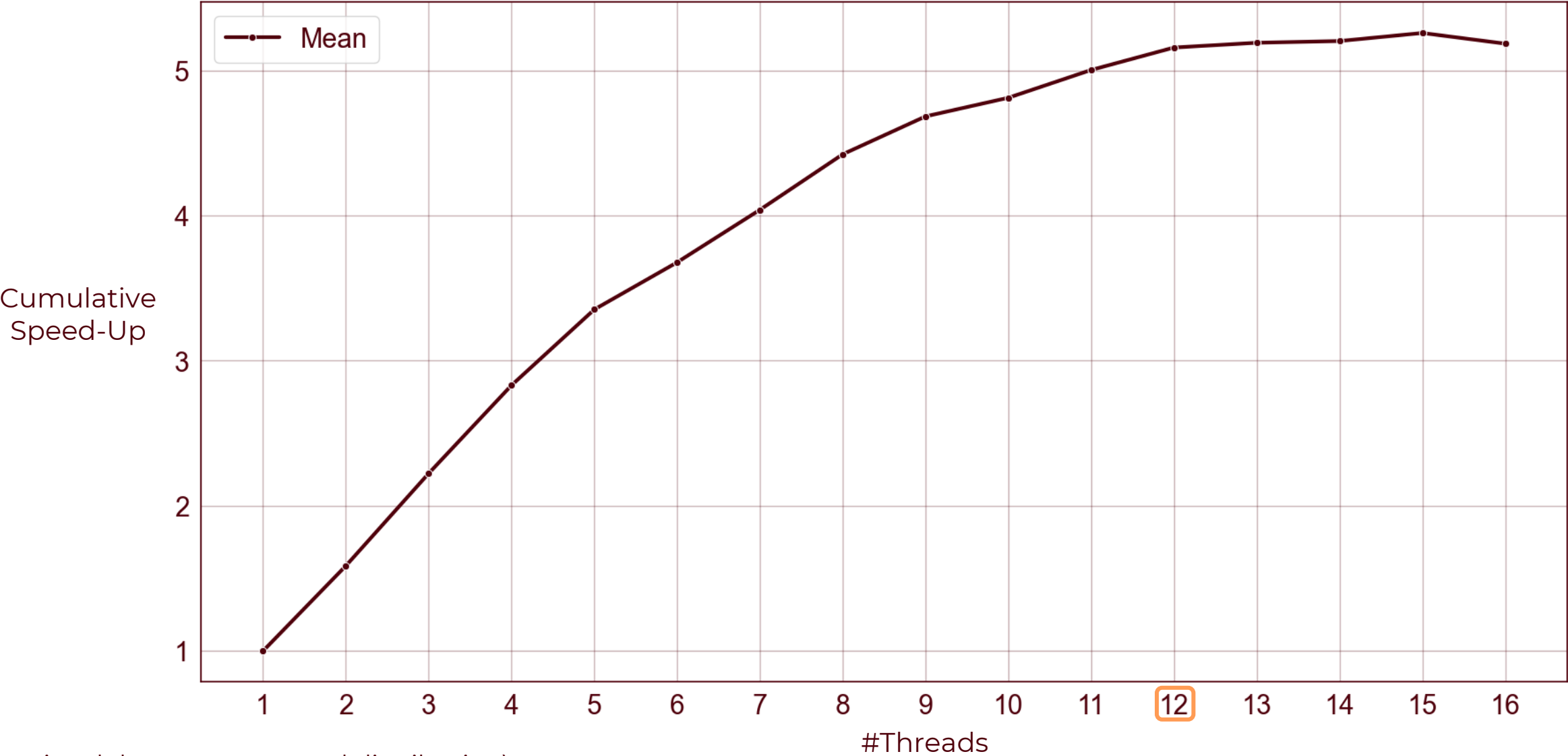
**59999999999991** = 3 * 19999999999997
(**14** DIGITS)

Mean Execution Time [ms]

#Threads

Mean
95% CI

# EFFECTIVE CPU UTILIZATION

Execution Time **[s]**

Logical Cores simultaneously utilized

# TOP **HOTSPOTS**

#THREADS : **12**

**0.588s**
EXECUTION TIME

**4.829s**
CPU TIME

**0s**
PAUSED TIME

**0s** (0.0% of CPU Time)
THREAD OVERSUBSCRIPTION

| Function | CPU Time | % CPU Time |
| --- | --- | --- |
| isPrime | 4.808s | 99.6% |
| findPrimesInRange | 0.013s | 0.3% |
| others | 0.008s | 0.2% |

| FUNCTION | CPU Time | % CPU Time |
|----------|----------|------------|
| ➜ isPrime | 4.808s | 99.6% |

```
bool isPrime(unsigned long long n) {

    …

    for(unsigned long long i = 2; i * i <= n; ++i) {
        if(n % i == 0) {
            is_prime = false;
            break;
        }
    }

    …
}
```

CPU
TIME

Total (%)          Self (s)

92,8%              4,483s

HOTSPOT

LET'S
OPTIMIZE!

$$2 \qquad\qquad \sqrt{N}$$

**MainThread:** N mod 2 == 0

# OPTIMIZATIONS #1

$3$                                              $\sqrt{N}$

Thread[0]     Thread[1]     ...     Thread[K-2]   MainThread

# OPTIMIZATIONS #2

findPrimesInRange

```cpp
    for (unsigned long long i = start; i <= end;  ++i ) {
        if (isPrime(i) && num % i == 0) {

            int exponent = 0;
            while (num % i == 0) {
                exponent++;
                num /= i;
            }

            {
                lock_guard<mutex> lock(mtx);
                primes.push_back({i, exponent});
            }
        }
    }
```
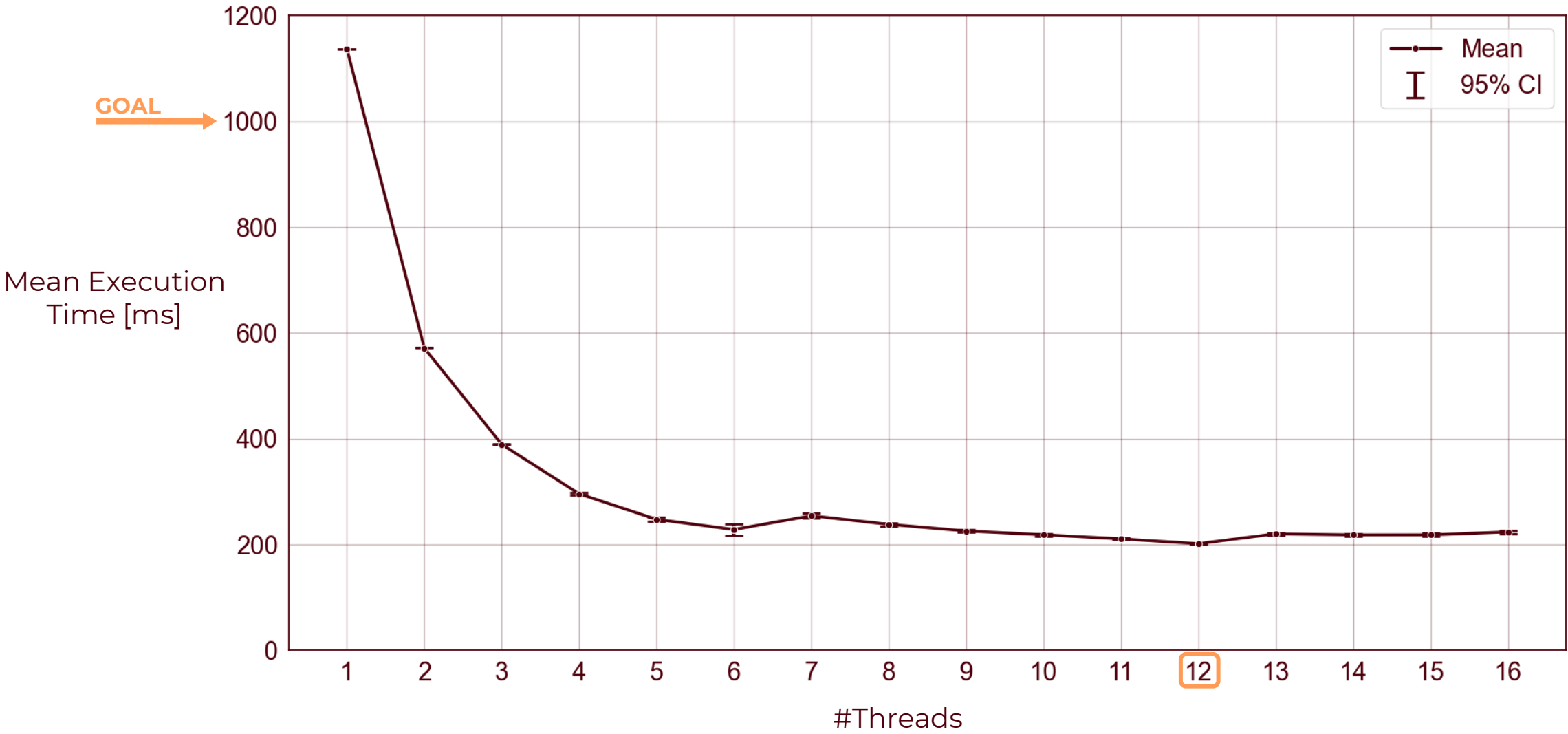
# OPTIMIZATIONS **#2**

## findPrimesInRange

```cpp
        for (unsigned long long i = start; i <= end; i += 2) {
            if (isPrime(i) && num % i == 0) {

                int exponent = 0;
                while (num % i == 0) {
                    exponent++;
                    num /= i;
                }


                {
                    lock_guard<mutex> lock(mtx);
                    primes.push_back({i, exponent});
                }
            }
        }
```

## findPrimesInRange

```cpp
        for (unsigned long long i = start; i <= end; i += 2) {
            if (isPrime(i) && num % i == 0) {

                int exponent = 0;
                while (num % i == 0) {
                        exponent++;
                        num /= i;
                }

                {

                        lock_guard<mutex> lock(mtx);
                        primes.push_back({i, exponent});
                }
            }
        }
```

# OPTIMIZATIONS **#2**

## findPrimesInRange

```
for (unsigned long long i = start; i <= end; i += 2) {
    if (num % i == 0) {

        int exponent = 0;
        while (num % i == 0) {
            exponent++;
            num /= i;
        }


        if(isPrime(i)){
            lock_guard<mutex> lock(mtx);
            primes.push_back({i, exponent});
        }
    }
}
```

# EXECUTION TIME
**OPTIMIZED** VERSION

#ITERATIONS : **30**

**975734686214396237** = 748609 * 1303396948493

(**18** DIGITS)

GOAL →

Mean Execution
Time [ms]

Mean
95% CI

#Threads

**FIRST** VERSION

EXECUTION TIME
**568340**
MILLISECONDS

≈ **9** MINUTES

**OPTIMIZED** VERSION

**MAX** (OVER 30 ITERATIONS)
EXECUTION TIME
**211**
MILLISECONDS

< **1** SECOND

**x2500**
SPEED-UP

#THREADS : **12**

# SPEED-UP
**OPTIMIZED** VERSION

#ITERATIONS : **30**

975734686214396237 = 748609 * 1303396948493

(**18** DIGITS)

Cumulative Speed-Up

#Threads

— · — Mean

(CI omitted due to non-normal distribution)

# SPEED-UP
**OPTIMIZED** VERSION
**+**
# AFFINITY

## CPU
**CORES**
[2 THREADS x CORE]

#1 ██ ██

#2 ▢ ▢

#3 ▢ ▢

#4 ▢ ▢

#5 ▢ ▢

#6 ▢ ▢



Cumulative
Speed-Up

#Threads

THREAD 1

THREAD 2

`if(num % i == 0)`

**+** **−**

DIV DIV

**ALU**

**DIVIDER**
UTILIZATION
~ **56%**

# SPEED-UP
**OPTIMIZED** VERSION

**+** **AFFINITY**

## CPU
**CORES**

[2 THREADS x CORE]

#1

#2

#3

#4

#5

#6



(CI omitted due to non-normal distribution)

# EXECUTION TIME - SCALABILITY
**OPTIMIZED** VERSION

#ITERATIONS : **30**

20 DIGITS **17975734686214396237**

18 DIGITS **975734686214396237**

16 DIGITS **3934686214396237**



Mean Execution Time [ms] vs #Threads

# GPU

IMPLEMENTATION

# OPTIMIZATIONS

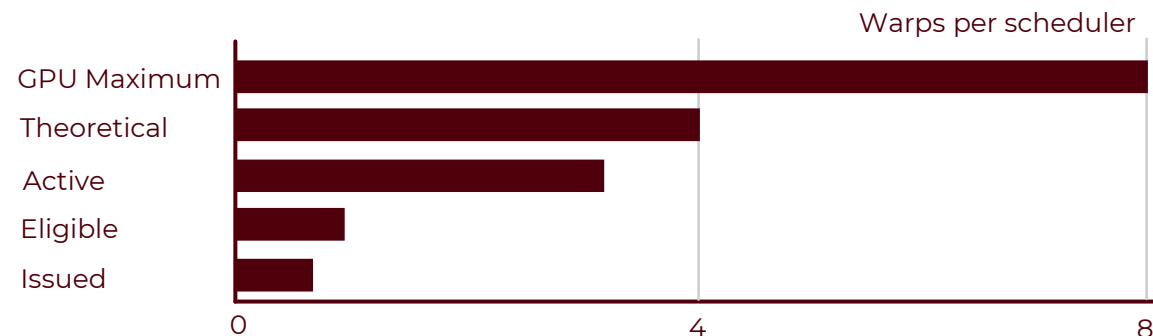## KERNEL

`__global__ void findPrimesInRange(...)`

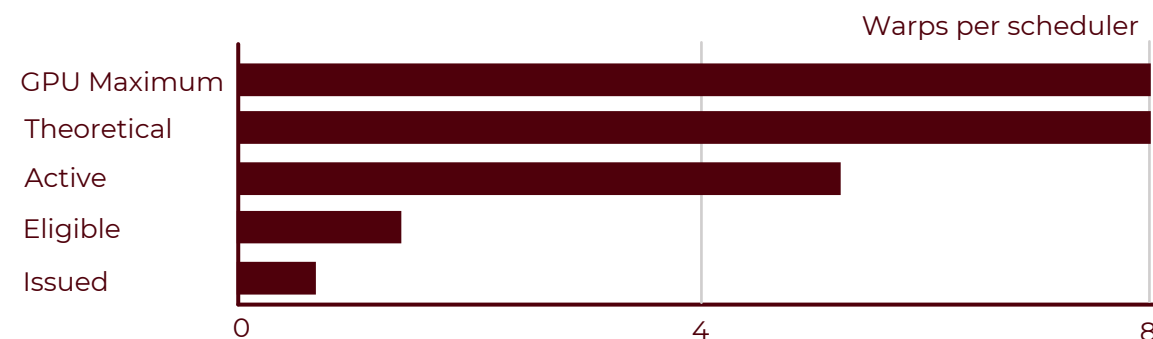JUST **1 THREAD** PER WARP WAS WORKING

∨

**32 THREADS** PER WARP*

*DIVERGENCE CAN HAPPEN BUT ITS EFFECT IS NEGLIGIBLE AS THE AVERAGE NUMBER OF NOT PREDICATED OFF THREADS PER WARP IS **31.64**

## 32 THREADS x BLOCK

Warps per scheduler

| | |
|---|---|
| GPU Maximum | |
| Theoretical | |
| Active | |
| Eligible | |
| Issued | |

0          4          8

## 64 THREADS x BLOCK

Warps per scheduler

| | |
|---|---|
| GPU Maximum | |
| Theoretical | |
| Active | |
| Eligible | |
| Issued | |

0          4          8

max_blocks_per_multiprocessor **16**

max_warps_per_multiprocessor **32**

max_warps_per_scheduler **8**

NVIDIA Nsight™ Compute

# CONCLUSIONS



CPU

GPU

**18** DIGITS

LIMITATION

**64-bit** REGISTER

17975734686214396237

**20** DIGITS

# THANKS
## FOR THE
## ATTENTION