# Mean shift clustering algorithm
## comparison between sequential, openMP and CUDA implementations

**Parallel Computing (9 CFU)**

*Francesca Del Lungo and Matteo Petrone*

Prof. Marco Bertini

January 2021

Introduction
000

Implementation
00000000

Experiments
0000000

Conclusions
000

## Contents overview

# Introduction

## Mean Shift

- Non-parametric **clustering** procedure
- Computational Complexity: $\mathcal{O}(n^2)$.
- **Embarrassingly parallel** structure

# Mean Shift

- It builds upon the concept of **kernel density estimation** (KDE).
- At each step of the procedure a kernel function is applied to each point that causes the points to shift in the direction of the nearest peak of the KDE surface:

$$K(x) = e^{-\frac{x^2}{2\sigma^2}} \tag{1}$$

where x: distance, $\sigma$ : *bandwidth*.

Introduction
000

Implementation
●0000000

Experiments
0000000

Conclusions
000

# Implementation

## Sequential Version

---

**Algorithm 1** Mean shift algorithm

---

**procedure** $\textsc{MeanShift}(\textit{original\_points})$
  $\textit{shifted\_points} \leftarrow \textit{original\_points}$
  **while** $\textit{iteration} < \textit{N\_ITERATIONS}$ **do**
    **for** $\textit{point}$ **in** $\textit{shifted\_points}$ **do**
      $\textit{point} \leftarrow \textsc{ShiftPoint}(\textit{point}, \textit{original\_points})$

---

Implemented in C++

Introduction
000

Implementation
00●00000

Experiments
0000000

Conclusions
000

## Sequential Version

---

**Algorithm 2** Mean shift subroutine

---

**procedure** SHIFTPOINT(*point*, *original_points*)

    $num \leftarrow 0$

    $den \leftarrow 0$

    **for** *op* **in** *original_points* **do**

        $dist \leftarrow$ COMPUTEDISTANCE(*point*, *op*)

        $w \leftarrow$ KERNEL(*dist*, *BANDWIDTH*)

        $num \leftarrow num + op * w$

        $den \leftarrow den + w$

    **return** $num/den$

---

## Parallel Version: OpenMP

- OpenMP (Open Multiprocessing) is an API for shared-memory parallel programming.
- It follows the the **fork-join** programming model.
- It does not require program restructuring, but only the addition of **compiler directives**

## Parallel Version: OpenMP

---

**Algorithm 3** Mean shift openMP algorithm

---

**procedure** $\mathrm{OMPMEANSHIFT}(\textit{original\_points})$

   $\textit{shifted\_points} \leftarrow \textit{original\_points}$

   **while** $\textit{iteration} < \textit{N\_ITERATIONS}$ **do**

      #pragma omp parallel for schedule(static)

      **for** $\textit{point}$ **in** $\textit{shifted\_points}$ **do**

         $\textit{point} \leftarrow \mathrm{SHIFTPOINT}(\textit{point}, \textit{original\_points})$

---

## Parallel Version: CUDA naive

- Each Point to be shifted assigned to a thread.
- **Coalesced access** to global memory: points are stored as *Structure of Arrays* $[x_1, ..., x_n, y_1, ...y_n, z_1, ..., z_n]$.
- Access to the array is done: *BlockDim.x * BlockIdx.x + threadIdx.x*.

Introduction
000

Implementation
00000000

Experiments
0000000

Conclusions
000

## Parallel Versione: CUDA naive

---

**Algorithm 4** CUDA naive kernel

---

**procedure** NAIVEKERNEL(*shifted_pts*, *orig_pts*)

  $idx \leftarrow threadIdx.x + blockIdx.x * blockDim.x$

  **if** $idx < |orig\_pts|$ **then**

    $sp \leftarrow shifted\_pts[idx]$

    $tot\_weight \leftarrow 0$

    $new\_spt \leftarrow 0$

    **for** *op* **in** *orig_pts* **do**

      $dist \leftarrow$ COMPUTEDISTANCE(*op*, *sp*)

      $weight \leftarrow$ KERNEL(*dist*, *BW*)

      $new\_spt \leftarrow new\_spt + op * weight$

      $tot\_weight \leftarrow tot\_weight + weight$

    $shifted\_pts[idx] \leftarrow new\_spt/tot\_weight$

Introduction
000

**Implementation**
0000000●

Experiments
0000000

Conclusions
000

## Parallel Version: CUDA Tiling

- Limit global memory access making use of shared memory

---

**Algorithm 5** CUDA kernel tiling - data loading

---

$tile \leftarrow shared\_mem\_array[TILE\_WIDTH]$

$n\_tiles \leftarrow (|orig\_pts| - 1)/(TILE\_WIDTH + 1)$

**for** $tile\_i$ **in** $n\_tiles$ **do**

  $tile\_idx \leftarrow tile\_i * TILE\_WIDTH + threadIdx.x$

  **if** $tile\_idx < |orig\_pts|$ **then**

    $tile[tx] \leftarrow orig\_pts[tile\_idx]$

  **else**

    $tile[tx] \leftarrow null\_pt$

  $\_\_syncthreads()$

---

# Experiments

Introduction
000

Implementation
00000000

Experiments
0●00000

Conclusions
000

Dataset

- We generated sets of 3d points using *sklearn.dataset.make_blobs()*.
- It generates **gaussian distribution** with 3 centers and 1.5 standard deviation.
- Datasets dimensions: 100, 1.000, 10.000, 100.000, 500.000, 1.000.000.

Introduction
000

Implementation
00000000

Experiments
0000000

Conclusions
000

## Experiments setting

All the tests have been performed on a machine equipped with:

- CPU: Intel Core i7-860 @ 2.80GHz, with 4 cores/ 8 threads
- GPU: NVidia GeForce GTX 980, 4 GB (with CUDA 10.1)
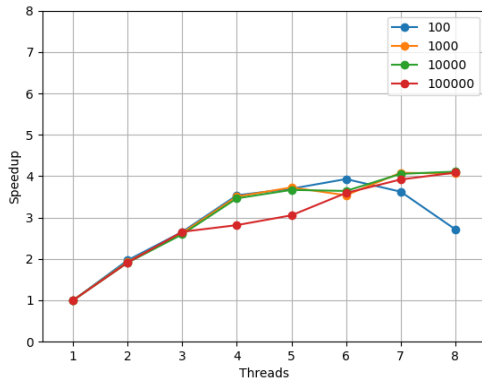- Average over 15 runs

# Speedup

To compare the performance of a sequential respect to a parallel
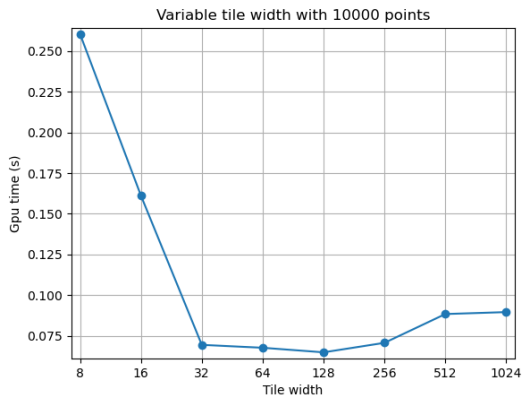implementation: **speedup** metric.

$$S = \frac{t_S}{t_P}$$

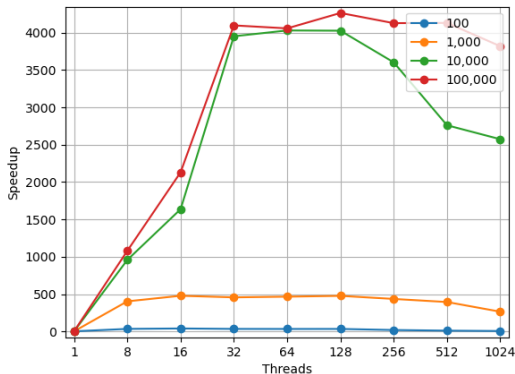where $t_s$: sequential time execution, $t_p$: parallel time execution.

Introduction
ooo

Implementation
oooooooo

Experiments
oooooooo

Conclusions
ooo

# Speedup: Sequential over OpenMP

Introduction
000

Implementation
00000000

Experiments
0000000

Conclusions
000

# GPU Time varying tile width



Variable tile width with 10000 points

Introduction
ooo

Implementation
oooooooo

Experiments
ooooooo●

Conclusions
ooo

# Speedup: Sequential over CUDA

Introduction
000

Implementation
00000000

Experiments
0000000

Conclusions
●00

# Conclusions

Introduction
000

Implementation
00000000

Experiments
0000000

Conclusions
0●0

## Conclusions

- Mean shift $\mathcal{O}(n^2)$, but **embarrassingly parallel** structure

### OpenMP

- Implementation difficulty: very low
- Speedup: sub-linear

### CUDA

- Implementation difficulty: higher
- Speedup: up to 4400x time faster.

Introduction
ooo

Implementation
oooooooo

Experiments
ooooooo

Conclusions
ooo●

# Thanks for the attention