

# PC-2016/17 Progetto mid-term del corso

## Implementazione in CUDA dell'algoritmo KMean

Tommaso Ceccarini  
6242250

tommaso.ceccarini1@stud.unifi.it

Federico Schipani  
6185896

federico.schipani@stud.unifi.it

### Abstract

*L'algoritmo di KMeans è uno dei più popolari metodi per la clusterizzazione. Nel nostro lavoro forniamo una implementazione in CUDA che fa uso di GPU Nvidia per l'esecuzione dell'algoritmo. Inoltre forniamo un'analisi delle performance con lo scopo di comparare la nostra implementazione parallela con una implementazione sequenziale scritta in C.*

## 1. Introduzione

L'algoritmo KMeans è uno degli algoritmi di clustering più famosi. Lo scopo del clustering è di dividere dati in gruppi significativi chiamati cluster. Dato un insieme di dati  $(x_1, x_2, \dots, x_N)$  dove ogni dato è un vettore reale di dimensione  $P$ , lo scopo di KMeans è di partizionare  $N$  dati in  $K (\leq N)$  insiemi  $S = \{S_1, S_2, \dots, S_K\}$  per minimizzare la distanza dentro ai singoli cluster. In altre parole viene definita una funzione obiettivo del tipo:

$$\arg \min_S \sum_{i=1}^k \sum_{x \in S_i} \|x - \mu_i\|^2 \quad (1)$$

dove  $\mu_i$  è la media dei punti in  $S_i$ . [1]

L'Algoritmo 1 rappresenta lo pseudocodice di un'implementazione iterativa dell'algoritmo KMeans.

I principali passi dell'algoritmo sono:

1. Selezionare casualmente  $K$  dati per inizializzare la media dei  $K$  cluster. Questo passo è usualmente chiamato *inizializzazione*.
2. Assegnare ogni dato al cluster più vicino, in accordo ad certa funzione di distanza. Questo passo è chiamato *assegnamento*.
3. Calcolare la nuova media dei cluster. Questo passo è chiamato *aggiornamento*.
4. Ripetere il passo due e tre finché nessun dato cambia cluster, oppure se solo pochi lo fanno.

---

### Algoritmo 1 KMeans

---

```
1: procedure KMEANS(data[n][p])
2:   mean[k][p], oldMean[k][p]
3:   assignment[n]
4:   (mean, assignment)  $\leftarrow$  initAss(data)
5:   while !stop do
6:     assignment  $\leftarrow$  calcMin(mean, data)
7:     oldMean  $\leftarrow$  mean
8:     mean  $\leftarrow$  calcMean(assignment, data)
9:     stop  $\leftarrow$  stopCrit(mean, oldMean)
10:  end while
11:  return mean
12: end procedure
```

---

#### 1.1. KMeans sequenziale

Prima del ciclo *while* c'è un passo di assegnazione iniziale. In questo step i primi  $K$  dati sono assegnati ai primi  $K$  cluster. Uno pseudocodice di questo assegnamento iniziale è mostrato in Algoritmo 2

---

### Algoritmo 2 Initial Assignment

---

```
1: procedure INITASS(data[n][p])
2:   mean[k][p]
3:   assignment[n]
4:   for i = 0; i < k; i ++ do
5:     assignment[i] = i
6:     for j = 0; j < p : j ++ do
7:       mean[i][j] = data[i][j]
8:     end for
9:   end for
10:  return (mean, assignment)
11: end procedure
```

---

Dopo questo step iniziale di assegnamento parte il vero e proprio algoritmo. Nella prima parte del ciclo *while* viene calcolata la distanza euclidea minima tra ognuno dei dati e tutti i cluster. Successivamente si assegna ogni dato al cluster più vicino. Il passo successivo consiste nel ricalcolare

le nuove medie dei cluster, secondo la (2).

$$m_i^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j \quad (2)$$

L'ultimo step è il criterio d'arresto. In questo semplice algoritmo, mostrato in 3, si iterano l'attuale matrice dei cluster e la matrice dei cluster calcolata al passo  $i - 1$ . Quando  $ActualValue - OldValue$  eccede una prefissata tolleranza TOL il metodo ritorna il valore *false*. Questo valore di ritorno porterà ad un'altra iterazione del ciclo *while* esterno.

---

#### Algoritmo 3 Stop Criterion

---

```

1: procedure STOPCRIT(mean[k][p], oldMean[k][p])
2:   for all value, oldValue in mean, oldMean do
3:     if abs(value - oldValue) > TOL then
4:       return false
5:     end if
6:   end for
7:   return true
8: end procedure

```

---

## 2. Implementazione

### 2.1. Come vengono rappresentati i dati ed i cluster

Le due differenti implementazioni dell'algoritmo di KMeans, riportate in sezione 2.2 e 2.3, risolvono il problema del clustering per un valore arbitrario del parametro  $P$ . Tuttavia le analisi delle performance, mostrate in Sezione 3, sono state effettuate nel caso in cui  $P = 2$ , ovvero il caso in cui i dati sono vettori bidimensionali.

Le due maggiori strutture su cui il codice lavora sono la matrice dei dati e la matrice dei centroidi, rispettivamente di dimensione  $N * P$  e  $K * P$ . In entrambe le implementazioni viene usato anche un vettore chiamato *assignment* che, per tutti gli  $N$  dati, memorizza l'indice del cluster al quale appartengono.

#### 2.1.1 Come viene generato il dataset

Dopo che sono state dichiarate, e successivamente istanziate, le strutture che sono necessarie per rappresentare i dati e i centroidi viene effettuata una generazione pseudocasuale dei dati. Questa generazione viene effettuata in maniera parallela attraverso le API CUDA, in particolare è stata usata la libreria cuRAND. Il metodo che è stato sviluppato per questo scopo è `generateRandomDataOnDevice()`. Questo metodo usa un seed prefissato e, con l'aiuto di un altro metodo che fa parte della libreria cuRAND, genera i valori del dataset. Nella Sezione 3 sono mostrate le analisi delle performance nel caso in cui i dati sono generati con una distribuzione uniforme.

## 2.2. Una implementazione sequenziale in C

### Lo step di inizializzazione

Dopo aver generato casualmente il dataset nella memoria globale del device è necessario istanziare la memoria richiesta per memorizzare la matrice dei dati e dei centroidi nella memoria dell'host. Inoltre, per come funziona l'algoritmo, è necessaria un'altra matrice che memorizza il valore della media durante la precedente iterazione. Questa matrice è usata per implementare il criterio d'arresto. L'inizializzazione è effettuata da un semplice ciclo *for* che assegna `data[i][j]` a `mean[i][j]` per  $i$  tra 0 e  $K - 1$ . Questa strategia d'assegnamento è motivata dal fatto che i dati sono generati casualmente.

### Lo step di assegnamento

Dopo aver effettuato l'assegnazione iniziale l'algoritmo sequenziale implementa il passo di assegnamento attraverso il Codice 1. Nel Codice 1 vengono utilizzati due cicli *for*: un ciclo esterno che itera su  $N$  dati e un ciclo annidato che per ogni dato itera sui  $K$  cluster per cercare il più vicino. Al termine del ciclo annidato l'indice del cluster più vicino al dato  $i$ -esimo verrà memorizzato alla cella  $i$  del vettore `assignment[]`.

#### Codice 1: Passo di assegnamento sequenziale

```

1  for (int i = 0; i < N; i++) {
2    float minDistance = 999999.9;
3    short minIndex = -1;
4    float distance = 0;
5    for (int z = 0; z < K; z++) {
6      distance = 0;
7      for (int j = 0; j < COMPONENTS; j++) {
8        distance += pow(
9          (hostData[i * COMPONENTS + j]
10           - mean[z * COMPONENTS + j]), 2);
11      }
12      if (distance < minDistance) {
13        minIndex = z;
14        minDistance = distance;
15      }
16    }
17    assignment[i] = minIndex;
18  }

```

---

### Lo step di aggiornamento

Successivamente all'assegnazione dei dati al cluster più vicino è necessario ricalcolare nuovamente i valori dei centroidi per i cluster. Questo passo dell'algoritmo è implementato mediante il Codice 2. Come si può notare dal Codice 2 questo passo è implementato mediante due cicli *for*: un primo ciclo esterno che itera sui  $K$  cluster e un ciclo annidato in cui, utilizzando il vettore `assignment[]` viene calcolata la somma dei valori delle singole componenti dei

dati che appartengono al cluster e il numero di dati che appartengono al cluster. Quindi alla fine del ciclo interno, per ogni componente, viene diviso il valore della somma per il numero di dati che appartengono al cluster per calcolare il valore effettivo dei nuovi centroidi.

#### Codice 2: Passo di aggiornamento della media

```

1  for (int i = 0; i < K; i++) {
2      int numberOfData = 0;
3      float *arraySum = (float*) malloc(COMPONENTS * sizeof(
4          ↪ float));
5      for (int x = 0; x < COMPONENTS; x++) {
6          arraySum[x] = 0.0;
7      }
8      for (int j = 0; j < N; j++) {
9          if (assignment[j] == i) {
10             numberOfData++;
11             for (int x = 0; x < COMPONENTS; x++) {
12                 arraySum[x] += hostData[j * COMPONENTS + x];
13             }
14         }
15     }
16     for (int j = 0; j < COMPONENTS; j++) {
17         oldMean[i * COMPONENTS + j] = mean[i * COMPONENTS + j
18             ↪ ];
19         mean[i * COMPONENTS + j] = arraySum[j] / numberOfData
20             ↪ ;
21     }
22 }
```

### Il criterio d'arresto

Il criterio d'arresto che viene implementato verifica se il vettore dei centroidi è cambiato nel corso dell'ultima iterazione rispetto a una prefissata tolleranza TOL. Per implementare questa strategia, alla riga 16 del Codice 2, prima di aggiornare il valore della media, vengono copiati i valori di `mean[]` nel vettore di appoggio `oldMean[]`. Il Codice 3 verifica quindi se i valori dei centroidi sono cambiati iterando mediante un ciclo for sulle componenti dei vettori `oldMean[]` e `mean[]` verificando se i valori sono uguali a meno di una tolleranza TOL.

#### Codice 3: Criterio di arresto sequenziale

```

1  stopCriterion = true;
2  for (int i = 0; i < K * COMPONENTS; i++) {
3      if (abs(mean[i] - oldMean[i]) > TOL) {
4          stopCriterion = false;
5          break;
6      }
7  }
```

## 2.3. Implementazione parallela tramite CUDA

### 2.3.1 Benefici di un'implementazione parallela

In riferimento al codice sequenziale si può notare come gran parte dei calcoli che vengono effettuati siano indipendenti

gli uni dagli altri. Questo fatto porta a pensare che un'implementazione parallela avrà dei miglioramenti in termini di prestazioni. In particolare le operazioni parallelizzabili sono:

- Assegnazione iniziale. Infatti è possibile assegnare ad un thread il compito di assegnare il valore di una componente del centroide.
- Calcolo delle distanze. Ciascuna delle  $N * K$  distanze può essere calcolata indipendentemente dalle altre.
- Determinazione del cluster più vicino. In questo contesto è possibile assegnare ad un thread il compito di determinare quale sia il cluster più vicino ad un dato. In ogni caso è comunque necessario effettuare una ricerca sequenziale scorrendo  $K$  medie.
- Calcolo della media. Utilizzando il vettore degli assegnamenti ogni thread può: determinare a quale cluster il dato corrisponde appartiene, e incrementare il contatore del numero di elementi per cluster che servirà poi per il calcolo effettivo della media.

Seguendo queste idee quindi è stato implementato il metodo `kMeans(float *devData, short* devAssignment, float* devMean)` che riceve in input la matrice dei dati, il vettore degli assegnamenti e la matrice dei centroidi, tutte queste strutture sono allocate nella memoria del device.

### Lo step di inizializzazione

Per realizzare il passo di inizializzazione vengono utilizzati i due metodi Kernel riportati in Codice 4 e Codice 5.

#### Codice 4: Assegnazione iniziale dei dati ai cluster

```

1  __global__ void initializeAssignment(short *assignment, int n
2      ↪ ) {
3      short index = (blockIdx.x * blockDim.x) + threadIdx.x;
4      if (index < K) {
5          assignment[index] = index;
6      }
7  }
```

Nel Codice 4 viene effettuata l'assegnazione iniziale. Questo viene eseguito con un totale di  $K$  thread. Ogni thread che esegue questa funzione Kernel si occupa di inizializzare il vettore degli assegnamenti scrivendo in posizione `index` il proprio indice `index`. A questo scopo la dimensione dei blocchi viene impostata a `dim3 DimBlockInizPart(1024, 1, 1)` e di conseguenza la griglia sarà inizializzata a `dim3 DimGridInizPart((k / 1024) + 1, 1, 1)`.

#### Codice 5: Inizializzazione della matrice dei centroidi

```

1  __global__ void initializeMean(float *mean, float* dataset,
2      ↪ int n,
```

```

2   int P) {
3   int tx = threadIdx.x;
4   int ty = threadIdx.y;
5   int bdx = blockDim.x;
6   int bdy = blockDim.y;
7   int bx = blockIdx.x;
8   int by = blockIdx.y;
9   int row = by * bdy + ty;
10  int col = bx * bdx + tx;
11  if (row < K && col < P) {
12      mean[row * P + col] = dataset[row * P + col];
13  }
14 }

```

Dopo aver assegnato i primi  $K$  dati ai relativi cluster il Codice 5 si occupa di assegnare sfruttando questo fatto: infatti è sufficiente assegnare i valori delle componenti del dato  $(i, j)$  al centroide  $(i, j)$ . Quindi in questo caso vengono usati  $K * P$  thread e la dimensione dei blocchi viene impostata a `dim3 DimBlockInizMean(P, 16, 1)`, e di conseguenza la griglia sarà inizializzata a `dim3 DimGridInizMean(1, k / 16 + 1, 1)`.

### Lo step di assegnamento

Per parallelizzare il passo di assegnamento vengono usate due funzioni, riportati in Codice 6 e Codice 7. L'implementazione parallela del passo di assegnamento è realizzata attraverso i due metodi riportati in Codice 6 e Codice 7. Nel Codice 6 ogni thread  $(i, j)$  è responsabile del calcolo della distanza del dato  $i$  dal cluster  $j$ . Quindi la griglia dei blocchi è organizzata in `dim3 DimGridDistances((k / 4) + 1, (n / 256) + 1, 1)` e di conseguenza i blocchi sono stati impostati a `dim3 DimBlockDistances(4, 256, 1)`.

#### Codice 6: Kernel per il calcolo delle distanze

```

1 __global__ void computeDistances(float* dataset, float*
    ↪ centroids,
2   float* distances, int P, int k) {
3   int bx = blockIdx.x;
4   int by = blockIdx.y;
5   int tx = threadIdx.x;
6   int ty = threadIdx.y;
7   int row = by * blockDim.y + ty;
8   int col = bx * blockDim.x + tx;
9   float dist = 0.0;
10  if (row < N && col < K) {
11      for (int i = 0; i < P; i++) {
12          dist += (dataset[row * P + i]
13                  - centroids[col * P + i])
14                  * (dataset[row * P + i]
15                    - centroids[col * P + i]);
16      }
17      distances[row * k + col] = dist;
18  }
19 }

```

L'altro Kernel riportato nel Codice 7 effettua la ricerca del cluster più vicino al dato utilizzando la matrice delle distanze appena calcolata dal Codice 6. In questo contesto ven-

gono utilizzati  $N$  thread in cui ciascuno effettua una ricerca sequenziale della distanza minima, memorizzando allo stesso tempo l'indice del cluster più vicino. La griglia è organizzata come `dim3 DimGridMin(1, (n / 256) + 1, 1)`, e i blocchi come `dim3 DimBlockAss(1, 256, 1)`.

#### Codice 7: Kernel per il calcolo del minimo

```

1 __global__ void computeMin(float* distances, short*
    ↪ devAssignment, int k,
2   int n) {
3   unsigned int ty = threadIdx.y;
4   int by = blockIdx.y;
5   unsigned int row = by * blockDim.y + ty;
6   if (row < n) {
7       float min = distances[row * k];
8       short minInd = 0;
9       for (int i = 1; i < k; i++) {
10          bool conf = (min - distances[row * k + i] <= 0);
11          min = conf * min + (1 - conf) * distances[row * k + i];
12          minInd = conf * minInd + (1 - conf) * i;
13      }
14      devAssignment[row] = minInd;
15  }
16 }

```

### Lo step di aggiornamento

L'aggiornamento della media è divisa in tre Kernel, riportati in Codice 8, Codice 9 e Codice 10. Il metodo Kernel, riportato in Codice 8, si occupa di reinizializzare il valore dei centroidi a 0. Questo per realizzare il calcolo delle somme delle componenti dei dati in modo atomico, riportato nel Codice 9.

#### Codice 8: Kernel per reinizializzare un vettore

```

1 __global__ void initializeVectorToValue(float* vector, float
    ↪ value, int bound) {
2   int tx = threadIdx.x;
3   int bdx = blockDim.x;
4   int bx = blockIdx.x;
5   int index = bx * bdx + tx;
6   if (index < bound) {
7       vector[index] = value;
8   }
9 }

```

Il Codice 9 utilizza  $N * P$  thread in cui il thread  $(i, j)$  si occupa di verificare, attraverso il vettore degli assegnamenti a quale cluster appartiene la componente  $j$  del dato  $i$  e di incrementare il conteggio degli elementi appartenenti al cluster del dato.

#### Codice 9: Calcolo della somma e del numero di elementi di un cluster

```

1 __global__ void computeSum37(float* dataset, short*
    ↪ devAssignment,
2   float* centroids, int* counter) {
3   int tx = threadIdx.x;

```

```

4  int ty = threadIdx.y;
5  int bdx = blockDim.x;
6  int bdy = blockDim.y;
7  int bx = blockIdx.x;
8  int by = blockIdx.y;
9  int row = by * bdy + ty;
10 int col = bx * bdx + tx;
11 if (row < N) {
12     int clusterIndex = devAssignment[row];
13     atomicAdd(&(centroids[clusterIndex * P + col]),
14             dataset[row * P + col]);
15     atomicAdd(&(counter[clusterIndex]), 1);
16 }
17
18 }

```

Infine nel Codice 10,  $K * P$  thread si occupano di fare il calcolo effettivo della media dividendo il risultato delle somme calcolate nel Codice 9 per il numero di elementi che appartengono al relativo cluster. I  $K * P$  thread sono divisi in blocchi `dim3 DimBlockMean(P, 256, 1)` e la griglia è formata da `dim3 DimGridMean(1, (K / 256) + 1, 1)` blocchi.

#### Codice 10: Kernel per il calcolo della media

```

1  __global__ void computeMean37(float* centroids, int* counter)
2  {
3      int tx = threadIdx.x;
4      int ty = threadIdx.y;
5      int bdx = blockDim.x;
6      int bdy = blockDim.y;
7      int bx = blockIdx.x;
8      int by = blockIdx.y;
9      int row = by * bdy + ty;
10     int col = bx * bdx + tx;
11
12     if (row < K) {
13         centroids[row * P + col] = centroids[row * P + col]
14         / (counter[row] / 2);
15     }
16 }

```

### Il criterio d'arresto

Per implementare il criterio d'arresto sono state prese in considerazione due possibili strade:

1. Utilizzare una variabile globale che conta il numero di dati che cambiano il cluster di appartenenza nel corso di una singola iterazione.
2. Riutilizzare il criterio di arresto implementato per la versione sequenziale, mostrato in Codice 3

L'implementazione della possibilità 1 è stata realizzata incrementando atomicamente la variabile globale in questione tutte le volte in cui il Codice 7 determina un cambiamento di cluster. Questa soluzione porta ad un incremento della divergenza a causa del controllo necessario per verificare se

il cluster del dato cambia. Quindi è stata valutata la possibilità di riutilizzare il codice che effettuava il criterio d'arresto nell'implementazione dell'algoritmo sequenziale. Sono state condotte una serie di verifiche per determinare se ci fossero delle differenze in termini di prestazioni tra le due diverse possibilità. A seguito di queste verifiche è stato riscontrato che la versione sequenziale del criterio d'arresto è generalmente, anche se di poco, migliore della controparte parallela.

Per realizzare una versione più efficiente della possibilità 2 è stata considerata la possibilità di effettuare una verifica parallela sul vettore delle medie utilizzando la tecnica di riduzione. Ad ogni modo verifiche condotte in merito hanno mostrato che tale approccio non portava alcun tipo di beneficio in termini di prestazioni. Questo tipo di tecnica ha il merito di sfruttare al meglio le potenzialità della memoria shared utilizzandole in maniera efficiente, e senza generare conflitti, per ridurre un vettore ad un singolo elemento. D'altra parte questo approccio richiede di iterare sempre e comunque su tutto il vettore per decidere il risultato finale del controllo, mentre la versione sequenziale appena verifica che una media nel corso dell'ultima iterazione è cambiata a meno di una tolleranza  $TOL$  si ferma restituendo il fallimento del controllo. Ad ogni modo dovendo operare su strutture allocate nella memoria del host sono necessarie due `cudaMemcpy()` del tipo `DeviceToHost` per il vettore delle medie aggiornato e quello vecchio.

## 3. Analisi delle performance e risultati sperimentali

Una volta che le due diverse versioni dell'algoritmo di KMeans sono state sviluppate sono stati condotti una serie di esperimenti per verificare quale fosse il vantaggio effettivo ottenuto dall'implementazione parallela. Le Figure 2 3 4 riportano il risultato della clusterizzazione ottenuta attraverso l'implementazione in CUDA con diversi valori di  $N$  e  $K$ .

### 3.1. Tempi di esecuzione

In Tabella 1 sono riportati i tempi di esecuzione per i due algoritmi implementati. Come si può osservare consultando la tabella la versione parallela è migliore di quella sequenziale. I risultati riportati sono stati ottenuti impostando un valore della tolleranza di  $5 \cdot 10^{-6}$ . In Tabella 2 è riportato lo speedup che la versione parallela raggiunge rispetto a quella sequenziale per gli esperimenti riportati nella Tabella 1. In Figura 1 sono riportati le percentuali di tempo dei metodi kernel eseguiti dall'implementazione parallela. Tali dati sono stati ottenuti con  $N = 80000$  e  $K = 200$ , eseguendo il programma con lo strumento di profiling di Nvidia.

Tabella 1: Tempi di esecuzione

| Punti  | Cluster | Sequenziale  | Parallelo   | Iterazioni |
|--------|---------|--------------|-------------|------------|
| 2000   | 100     | 0.142902 s   | 0.009067 s  | 20         |
| 10000  | 200     | 2.45001 s    | 0.067496 s  | 42         |
| 25000  | 350     | 12.564984 s  | 0.533056 s  | 52         |
| 50000  | 500     | 56.518483 s  | 2.519471 s  | 86         |
| 100000 | 750     | 182.013450 s | 8.965721 s  | 93         |
| 200000 | 1000    | 758.144052 s | 32.357209 s | 143        |

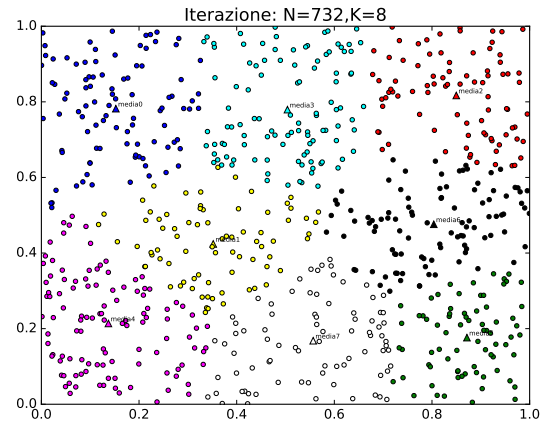


Figura 2: Plot per  $N = 732$  e  $K = 8$

Tabella 2: Speedup

| Punti  | Cluster | Speedup |
|--------|---------|---------|
| 2000   | 100     | 15.760  |
| 10000  | 200     | 36.300  |
| 25000  | 350     | 23.571  |
| 50000  | 500     | 22.433  |
| 100000 | 750     | 20.301  |
| 200000 | 1000    | 23.430  |

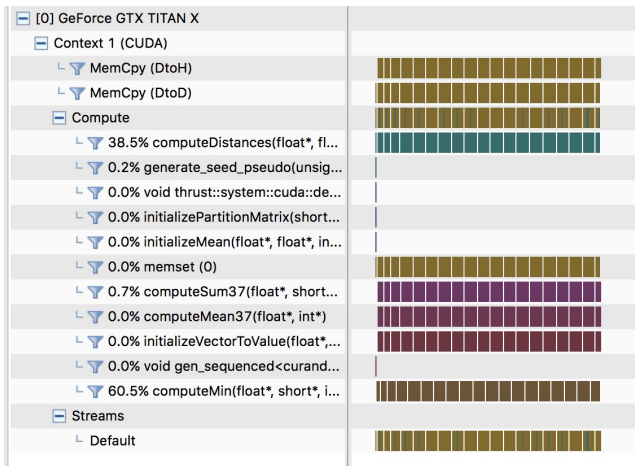


Figura 1: Percentuale del tempo di esecuzione dei diversi kernel

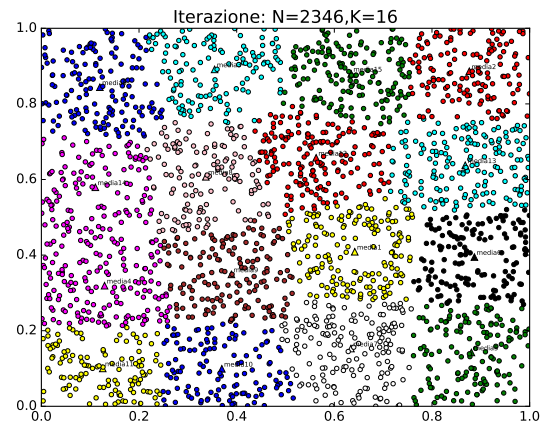


Figura 3: Plot per  $N = 2346$  e  $K = 16$



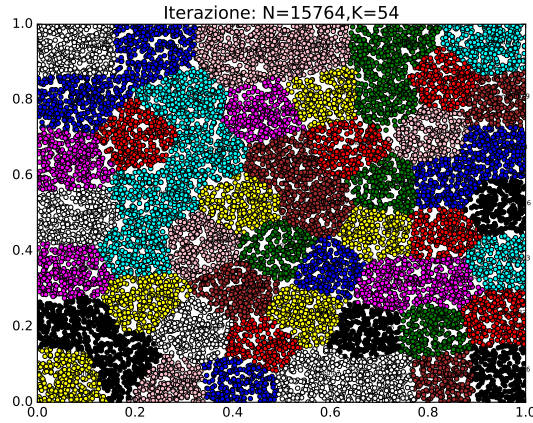


Figura 4: Plot per  $N = 15764$  e  $K = 54$

## 4. Conclusioni

È quindi possibile affermare che l'algoritmo di KMeans trae notevoli benefici da un'implementazione parallela che fa uso di una GPU. Infatti per gran parte dei calcoli è stato possibile realizzare una versione parallela in quanto indipendenti gli uni dagli altri. Alla luce degli esperimenti riportati non è possibile affermare che la percentuale di miglioramento della versione parallela rispetto alla sequenziale (speedup) sia legata alle due dimensioni principali del problema che sono state prese in considerazione, ovvero  $N$  e  $K$ . Ad ogni modo la versione implementata in CUDA si è rivelata essere notevolmente più efficiente rispetto alla versione sequenziale sviluppata in C.

## A. Codice

### Codice 11: Codice completo

```

1
2 #include <iostream>
3 #include <numeric>
4 #include <stdlib.h>
5 #include <stdio.h>
6 #include <curand.h>
7 #include <curand_kernel.h>
8 #include <math.h>
9 #include <cuda.h>
10 #include <stdbool.h>
11 #include <thrust/device_ptr.h>
12 #include <thrust/fill.h>

```

```

13 #include <float>
14 #include <fstream>
15 #include <sstream>
16 #include <thrust/transform_reduce.h>
17 #include <thrust/functional.h>
18 #include <thrust/device_vector.h>
19 #include <thrust/host_vector.h>
20 #include <cmath>
21 #include <time.h>
22
23 static void CheckCudaErrorAux(const char *, unsigned, const
    ↪ char *,
24                               cudaError_t);
25 #define CUDA_CHECK_RETURN(value) CheckCudaErrorAux(__FILE__,
    ↪ __LINE__, #value, value)
26 #define P 2
27 #define K 233
28 #define N 35067
29 #define TOL 0.000005
30 #define SEED 1234ULL
31 #define TD 1024
32
33 #define CURAND_CALL(x) do { \
34     if ((x) != CURAND_STATUS_SUCCESS) { \
35         printf("ERROR AT %s:%d\r\n", __FILE__, __LINE__); \
36         return EXIT_FAILURE; \
37     } \
38 } while(0)
39
40 using namespace std;
41
42 __global__ void initializeVectorToValue(float *vector, float
    ↪ value, int bound) {
43     int tx = threadIdx.x;
44     int bdx = blockDim.x;
45     int bx = blockIdx.x;
46     int index = bx * bdx + tx;
47     if (index < bound) {
48         vector[index] = value;
49     }
50 }
51
52 __global__ void initializeAssignment(short *assignment, int n
    ↪ ) {
53     short index = (blockIdx.x * blockDim.x) + threadIdx.x;
54     if (index < K) {
55         assignment[index] = index;
56     }
57 }
58 __global__ void initializeMean(float *mean, float* dataset,
    ↪ int n,
59                               int P) {
60     int tx = threadIdx.x;
61     int ty = threadIdx.y;
62     int bdx = blockDim.x;
63     int bdy = blockDim.y;
64     int bx = blockIdx.x;
65     int by = blockIdx.y;
66     int row = by * bdy + ty;
67     int col = bx * bdx + tx;
68     if (row < K && col < P) {
69         mean[row * P + col] = dataset[row * P + col];
70     }
71 }
72
73 __global__ void computeSum37(float* dataset, short*
    ↪ devAssignment,

```

```

74     float* centroids, int* counter) {
75     int tx = threadIdx.x;
76     int ty = threadIdx.y;
77     int bdx = blockDim.x;
78     int bdy = blockDim.y;
79     int bx = blockIdx.x;
80     int by = blockIdx.y;
81     int row = by * bdy + ty;
82     int col = bx * bdx + tx;
83     if (row < N) {
84         int clusterIndex = devAssignment[row];
85         atomicAdd(&(centroids[clusterIndex * P + col]),
86                 dataset[row * P + col]);
87         atomicAdd(&(counter[clusterIndex]), 1);
88     }
89 }
90 }
91
92 __global__ void computeMean37(float* centroids, int* counter)
93     ↪ {
94     int tx = threadIdx.x;
95     int ty = threadIdx.y;
96     int bdx = blockDim.x;
97     int bdy = blockDim.y;
98     int bx = blockIdx.x;
99     int by = blockIdx.y;
100    int row = by * bdy + ty;
101    int col = bx * bdx + tx;
102
103    if (row < K) {
104        centroids[row * P + col] = centroids[row * P + col]
105        / (counter[row] / 2);
106    }
107 }
108
109 __global__ void computeMean(float* centroids, float* dataset,
110                             short* devAssignment, int n, int P) {
111     int tx = threadIdx.x;
112     int ty = threadIdx.y;
113     int bdx = blockDim.x;
114     int bdy = blockDim.y;
115     int bx = blockIdx.x;
116     int by = blockIdx.y;
117     int row = by * bdy + ty;
118     int col = bx * bdx + tx;
119     float mean = 0.0;
120     int counter = 0;
121     if (row < K) {
122         for (int i = 0; i < n; i++) {
123             mean += (devAssignment[i] == row) * dataset[i * P + col]
124             ↪ ;
125             counter += (devAssignment[i] == row);
126         }
127         centroids[row * P + col] = mean / counter;
128     }
129 }
130
131 __global__ void computeDistances(float* dataset, float*
132     ↪ centroids,
133     float* distances, int P, int k) {
134     int bx = blockIdx.x;
135     int by = blockIdx.y;
136     int tx = threadIdx.x;
137     int ty = threadIdx.y;
138     int row = by * blockDim.y + ty;
139     int col = bx * blockDim.x + tx;

```

```

138     float dist = 0.0;
139     if (row < N && col < K) {
140         for (int i = 0; i < P; i++) {
141             dist += (dataset[row * P + i]
142                     - centroids[col * P + i])
143                     * (dataset[row * P + i]
144                       - centroids[col * P + i]);
145         }
146         distances[row * k + col] = dist;
147     }
148 }
149
150 __global__ void computeMin(float* distances, short*
151     ↪ devAssignment, int k,
152     int n) {
153     unsigned int ty = threadIdx.y;
154     int by = blockIdx.y;
155     unsigned int row = by * blockDim.y + ty;
156     if (row < n) {
157         float min = distances[row * k];
158         short minInd = 0;
159         for (int i = 1; i < k; i++) {
160             bool conf = (min - distances[row * k + i] <= 0);
161             min = conf * min + (1 - conf) * distances[row * k + i];
162             minInd = conf * minInd + (1 - conf) * i;
163         }
164         devAssignment[row] = minInd;
165     }
166 }
167
168 __global__ void assignPartition(short* S, int* minIndex, int
169     ↪ n) {
170     int by = blockIdx.y;
171     int ty = threadIdx.y;
172     int row = by * blockDim.y + ty;
173
174     if (row < n) {
175         int pos = minIndex[row];
176         S[pos * n + row] = true;
177     }
178 }
179
180 int generateRandomDataOnDevice(float* devData, int n, int P)
181     ↪ {
182     curandGenerator_t generator;
183     CURAND_CALL(curandCreateGenerator(&generator,
184     ↪ CURAND_RNG_PSEUDO_DEFAULT));
185     CURAND_CALL(curandSetPseudoRandomGeneratorSeed(generator,
186     ↪ SEED));
187     CURAND_CALL(curandGenerateUniform(generator, devData, n * P
188     ↪ ));
189     CURAND_CALL(curandDestroyGenerator(generator));
190     return 0;
191 }
192
193 void printFile(std::string path, std::string filename, int n,
194     ↪ int k,
195     int P, float* devDataset, short* devAssignment,
196     float* devMean) {
197     float* mean, *dataset;
198     short* assignment;
199     mean = (float*) malloc(k * P * sizeof(float));
200     dataset = (float*) malloc(n * P * sizeof(float));
201     assignment = (short*) malloc(n * sizeof(short));

```



```

198 CUDA_CHECK_RETURN(
199     cudaMemcpy(mean, devMean, P * k * sizeof(float),
200     cudaMemcpyDeviceToHost));
201 CUDA_CHECK_RETURN(
202     cudaMemcpy(dataset, devDataset, n * P * sizeof(float),
203     cudaMemcpyDeviceToHost));
204 CUDA_CHECK_RETURN(
205     cudaMemcpy(assignment, devAssignment, n * sizeof(short)
206     ↪ ,
207     cudaMemcpyDeviceToHost));
208 cudaDeviceSynchronize();
209 ofstream outputfile;
210 outputfile.open(path + filename, ios::out);
211 for (int i = 0; i < K; i++) {
212     std::string toFile = "media" + to_string(i);
213     for (int j = 0; j < P; j++) {
214         toFile = toFile + " " + to_string(mean[i * P + j]);
215     }
216     toFile = toFile + "|";
217     outputfile << toFile;
218 }
219 for (int i = 0; i < N; i++) {
220     std::string toFile = "";
221     toFile = to_string(assignment[i]);
222     for (int j = 0; j < P; j++) {
223         toFile = toFile + " " + to_string(dataset[i * P + j]);
224     }
225     toFile = toFile + "|";
226     outputfile << toFile;
227 }
228 outputfile.close();
229 free(mean);
230 free(dataset);
231 free(assignment);
232 }
233 bool stopCriterionSequential(short *assignment, short *
234     ↪ oldAssignment, int len) {
235     for (int i = 0; i < len; i++) {
236         if (assignment[i] != oldAssignment[i]) {
237             return false;
238         }
239     }
240     return true;
241 }
242 bool stopCriterionSequential(float *vector1, float *vector2,
243     ↪ int len) {
244     for (int i = 0; i < len; i++) {
245         if (abs(vector1[i] - vector2[i]) > TOL) {
246             return false;
247         }
248     }
249     return true;
250 }
251 int kMeans(float *devData, short *devAssignment, float *
252     ↪ devMean, int n, int k,
253     int P) {
254     printf("KMEANS nuovo \r\n");
255     dim3 DimBlockInizToZero(256, 1, 1);
256     dim3 DimGridInizToZero((2 * K) / 256 + 1, 1, 1);
257     dim3 DimBlock2(n, k, 1);
258     dim3 DimBlockAssignment(1, n, 1);
259     dim3 DimBlockUpdate(P, k, 1);
260     dim3 DimGridSum(1, (N / 256) + 1, 1);

```

```

261     dim3 DimGridSum(1, (N / 256) + 1, 1);
262     dim3 DimBlockMean(P, 256, 1);
263     dim3 DimGridMean(1, (K / 256) + 1, 1);
264     dim3 DimBlockDistances(4, 256, 1);
265     dim3 DimGridDistances((k / 4) + 1, (n / 256) + 1, 1);
266     dim3 DimBlockMin(1, 256, 1);
267     dim3 DimGridMin(1, (n / 256) + 1, 1);
268     dim3 DimBlockAss(1, 256, 1);
269     dim3 DimGridAss(1, (n / 256) + 1, 1);
270     dim3 DimGridInizPart((k / 1024) + 1, 1, 1);
271     dim3 DimBlockInizPart(1024, 1, 1);
272     dim3 DimGridInizMean(1, k / 16 + 1, 1);
273     dim3 DimBlockInizMean(P, 16, 1);
274     cudaMemset(devAssignment, 0, n * sizeof(short));
275     initializeAssignment<<<DimGridInizPart, DimBlockInizPart
276     ↪ >>>(
277         devAssignment, n);
278     cudaDeviceSynchronize();
279     initializeMean<<<DimGridInizMean, DimBlockInizMean>>>(
280     ↪ devMean, devData, n,
281     P);
282     cudaDeviceSynchronize();
283     short *devOldAssignment;
284     CUDA_CHECK_RETURN(
285         cudaMalloc((void**) &devOldAssignment, n * sizeof(short)
286         ↪ ));
287     float *hostData, *hostMean, *hostOldMean;
288     hostData = (float*) malloc(n * P * sizeof(float));
289     short *hostAssignment, *hostOldAssignment;
290     hostAssignment = (short*) malloc(n * sizeof(short));
291     hostOldAssignment = (short*) malloc(n * sizeof(short));
292     int *deviceOutputVector;
293     hostMean = (float*) malloc(k * P * sizeof(float));
294     hostOldMean = (float*) malloc(k * P * sizeof(float));
295     CUDA_CHECK_RETURN(
296         cudaMalloc((void**) &deviceOutputVector, k * sizeof(int)
297         ↪ ));
298     bool stopCriterion = false;
299     float *devDistances;
300     CUDA_CHECK_RETURN(
301         cudaMalloc((void**) &devDistances, n * k * sizeof(float)
302         ↪ ));
303     int* devCounter, *hostCounter;
304     CUDA_CHECK_RETURN(cudaMalloc((void**) &devCounter, k *
305     ↪ sizeof(int)));
306     hostCounter = (int*) malloc(k * sizeof(int));
307     while (!stopCriterion) {
308         CUDA_CHECK_RETURN(cudaMemset(devOldAssignment, 0, n *
309         ↪ sizeof(short)));
310         CUDA_CHECK_RETURN(cudaMemset(devCounter, 0, k * sizeof(
311         ↪ int)));
312         cudaDeviceSynchronize();
313         CUDA_CHECK_RETURN(
314             cudaMemcpy(devOldAssignment, devAssignment, n *
315             ↪ sizeof(short),
316             cudaMemcpyDeviceToDevice));
317         cudaDeviceSynchronize();
318         computeDistances<<<DimGridDistances, DimBlockDistances
319         ↪ >>>(devData,
320             devMean, devDistances, P, k);
321         cudaDeviceSynchronize();
322         computeMin<<<DimGridMin, DimBlockMin>>>(devDistances,
323         ↪ devAssignment, k,
324             n);
325         cudaDeviceSynchronize();
326         CUDA_CHECK_RETURN(

```

```

316     cudaMemcpy(hostOldMean, devMean, P * k * sizeof(float)
317         ↪ ),
318     cudaMemcpyDeviceToHost));
319 initializeVectorToValue<<<DimGridInizToZero,
320     ↪ DimBlockInizToZero>>>(
321     devMean, 0.0, P * K);
322 cudaDeviceSynchronize();
323 computeSum37<<<DimGridSum, DimBlockSum>>>(devData,
324     ↪ devAssignment,
325     devMean, devCounter);
326 cudaDeviceSynchronize();
327 computeMean37<<<DimGridMean, DimBlockMean>>>(devMean,
328     ↪ devCounter);
329 cudaDeviceSynchronize();
330 CUDA_CHECK_RETURN(
331     cudaMemcpy(hostMean, devMean, P * k * sizeof(float),
332     cudaMemcpyDeviceToHost));
333 stopCriterion = stopCriterionSequential(hostOldMean,
334     ↪ hostMean,
335     K * P);
336 /*
337     printf("/home/cecca/Desktop/RemoteCompilation/TXTDATA
338     ↪ /",
339     to_string(iter) + ".txt", n, k, P, devData,
340     devAssignment, devMean);
341 */
342 }
343 return 0;
344 }
345
346 void printFileForSeq(std::string path, std::string filename,
347     ↪ float *dataset,
348     short *assignment, float *mean) {
349     ofstream outputfile;
350     outputfile.open(path + filename, ios::out);
351     for (int i = 0; i < K; i++) {
352         std::string toFile = "media" + to_string(i);
353         for (int j = 0; j < P; j++) {
354             toFile = toFile + " " + to_string(mean[i * P + j]);
355         }
356         toFile = toFile + "|";
357         outputfile << toFile;
358     }
359     for (int i = 0; i < N; i++) {
360         std::string toFile = "";
361         toFile = to_string(assignment[i]);
362         for (int j = 0; j < P; j++) {
363             toFile = toFile + " " + to_string(dataset[i * P + j]);
364         }
365         toFile = toFile + "|";
366         outputfile << toFile;
367     }
368     outputfile.close();
369 }
370
371 int main(int argc, char *argv[]) {
372     system("rm -rf /home/cecca/Desktop/RemoteCompilation/
373     ↪ TXTDATA");
374     system("mkdir /home/cecca/Desktop/RemoteCompilation/TXTDATA
375     ↪ ");
376     size_t n;
377     size_t i;
378     int k;
379     int P = 2;
380
381     if (argc > 1) {
382         n = atoi(argv[1]);
383     } else {
384         n = N;
385     }
386
387     if (argc > 2) {
388         k = atoi(argv[2]);
389     } else {
390         k = K;
391     }
392
393     float *devCentroids, *hostCentroids;
394     int sizeOfData = P * n * sizeof(float);
395     short *hostAssignment, *devAssignment;
396     float *devData, *hostData;
397     float *devMean, *hostMean;
398     hostData = (float*) malloc(sizeOfData);
399     hostCentroids = (float*) malloc((sizeOfData * k) / n);
400     hostMean = (float*) malloc((sizeOfData * k) / n);
401     hostAssignment = (short*) malloc(n * sizeof(short));
402     CUDA_CHECK_RETURN(cudaMalloc((void**) &devData, sizeOfData)
403     ↪ );
404     CUDA_CHECK_RETURN(cudaMalloc((void**) &devMean, (sizeOfData
405     ↪ * k) / n));
406     CUDA_CHECK_RETURN(cudaMalloc((void**) &devAssignment, n *
407     ↪ sizeof(short)));
408     CUDA_CHECK_RETURN(cudaMalloc((void**) &devCentroids, (
409     ↪ sizeOfData * k) / n));
410     generateRandomDataOnDevice(devData, n, P);
411     //#####
412     clock_t start, end;
413     double cpu_time_used;
414     start = clock();
415     KMeans(devData, devAssignment, devCentroids, n, k, P);
416     end = clock();
417     cpu_time_used = ((double) (end - start)) / CLOCKS_PER_SEC;
418     printf("TEMPO: %f SECONDI \r\n", cpu_time_used);
419     //#####
420     CUDA_CHECK_RETURN(cudaFree(devAssignment));
421     CUDA_CHECK_RETURN(cudaFree(devMean));
422     CUDA_CHECK_RETURN(cudaFree(devData));
423     free(hostData);
424     free(hostMean);
425     free(hostAssignment);
426
427     return 0;
428 }
429
430 static void CheckCudaErrorAux(const char *file, unsigned line
431     ↪ ,
432     const char *statement, cudaError_t err) {
433     if (err == cudaSuccess)
434         return;
435     std::cerr << statement << " returned " <<
436     ↪ cudaGetErrorString(err) << "("
437     << err << ") at " << file << ":" << line << std::endl;
438     exit(1);
439 }

```

## Riferimenti bibliografici

- [1] Wikipedia. K-means clustering — wikipedia, the free encyclopedia, 2017. [Online; accessed 23-January-2017]. 1