

# université PARIS-SACLAY

M2 E3A - SETI

---

## TP GPU

---

SAŠA RADOSAVLJEVIC



12 FÉVRIER 2023



# Table des matières

<b>Table des matières</b>	<b>1</b>
<b>1 TP A4</b>	<b>2</b>
1.1 TP1 . . . . .	2
1.1.1 Matlab . . . . .	2
1.1.2 Jetson nano CPU . . . . .	4
1.1.3 Jetson Nano CUDA . . . . .	4
1.2 Kmeans Python . . . . .	6
1.3 Conclusion . . . . .	8

# 1 | TP A4

## 1.1 - TP1

---

Lors de ce TP, nous allons programmer un algorithme de seuillage couleur à l'aide de trois langages de programmation. Nous analyserons les résultats relevant des temps d'exécution notamment entre différentes optimisations CPU sous matlab. Nous verrons alors l'exécution de ce même seuillage sur une cible Nvidia Jetson Nano. Nous pourrons alors mesurer les temps d'exécutions entre les versions CPU et CUDA afin de voir les gains que l'on peut tirer d'une utilisation d'un GP-GPU pour le traitement parallèle.

### 1.1.1 - Matlab

La première partie consiste à programmer l'algorithme de seuillage et de changement de couleur sur Matlab. Matlab a la particularité de pouvoir traiter des matrices dans ses opérations. De ce fait, les différentes boucles itératives permettant de parcourir les pixels d'une image peuvent être enlevées afin de tirer un maximum des avantages offerts par Matlab.

L'exécution sera faite sur un PC fixe muni d'un CPU Ryzen 9 3900x de fréquence de base à 3,8 Ghz avec un boost à 4,6 Ghz. Il est composé de 12 coeurs et d'un cache L3 de 64 Mo. La mémoire RAM est composée de 2 barrettes de 16 Go à 3200 MHz C16 en dual channel.

Le code exécuté est donné ci-dessous :

```
1  %tic toc pour mesurer le temps de calcul
2  disp("Seuil :");
3
4  tic;
5  ima_out=ima;
6  nr = ima_r./sqrt(ima_r.^2 + ima_b.^2 + ima_g.^2);
7  ima_seuil = ima.*(nr>0.7);
8  toc;
9  figure('name','RGB out','numbertitle','off');image(ima_seuil);
10
11 disp("Jaune :")
12 tic;
```

```

13 ima_jaune = ima_seuil;
14 ima_jaune(:,:,2) = ima_seuil(:,:,1);
15 toc;
16 figure('name','RGB out','numbertitle','off');image(ima_jaune);
17
18 disp("Reinsertion dans l'image :")
19 tic;
20 ima_out = ima - ima_seuil + ima_jaune;
21 figure('name','RGB out','numbertitle','off');image(ima_out);
22 toc;

```

On obtient les résultats suivants :

- Seuil : Elapsed time is 0.006431 seconds.
- Jaune : Elapsed time is 0.006348 seconds.
- Réinsertion dans l'image : Elapsed time is 0.003195 seconds..
- Exécution totale : 0.015394 seconds.

Si l'on compare ce résultat avec l'exécution du code avec boucles, c'est-à-dire boucle sur i et j avec calcul de l'intensité de rouge dans les boucles, puis en dehors des boucles et pour finir le calcul sans boucles, on obtient les résultats suivants :

- VERSION 1 : Elapsed time is 0.078048 seconds.
- VERSION 2 : Elapsed time is 0.020645 seconds.
- VERSION 3 : Elapsed time is 0.013962 seconds.

On observe bien une amélioration d'un facteur 6 lorsqu'on enlève les boucles. On obtient un résultat similaire entre les deux versions (codée et solution) sans boucles.



(a) Image d'entrée



(b) Image après traitement

FIGURE 1.1 – Exécution du programme de changement de couleur

### 1.1.2 - Jetson nano CPU

La deuxième partie s'effectue sur une carte Nvidia Jetson Nano dotée d'un CPU Quad-core ARM A57 @ 1.43 GHz.

On exécutera sur la cible le code suivant :

```
1 void seuillage_C(float image_out[] [SIZE_J] [SIZE_I], float
  - image_in[] [SIZE_J] [SIZE_I])
2 {
3     float r, g, b;
4     for(int j=0; j<SIZE_J; j++){
5         for(int i=0; i<SIZE_I; i++){
6             r = image_in[0][j][i];
7             g = image_in[1][j][i];
8             b = image_in[2][j][i];
9
10            if(r/sqrt(r*r+g*g+b*b) > 0.7){
11                image_out[0][j][i] = r;
12                image_out[1][j][i] = g;
13                image_out[2][j][i] = b;
14            }
15        }
16    }
17 }
```

En observant l'impact d'inversion des indices i et j sur l'optimisation des calculs, on obtient :

- Boucles bon ordre pour le cache : 95 ms.
- Boucles mauvais ordre pour le cache : 171 ms.

En comparaison avec le temps obtenu avec la version 1 de Matlab on est à un facteur 1,21, ce qui est assez proche au vu du processeur utilisé précédemment. Le même code exécuté sur le ryzen obtient un temps d'exécution de 6,5 ms. On a cette fois un vrai écart d'un facteur 14,6 au vu de la différence de puissance de calcul.

### 1.1.3 - Jetson Nano CUDA

En dernier lieu, on souhaite utiliser le GPU de la Jetson Nano pour observer les bénéfices de l'utilisation de CUDA. On peut tout d'abord regarder les caractéristiques du GPU avec l'exécutable deviceQuery disponible dans le répertoire `/usr/local/cuda/samples/bin`.

deviceQuery :

- Multiprocesseur : 1

- CUDA cores : 128
- Max Fclock = 922Mhz
- Mémoire : 2 Go
- Memory clock : 13 Mhz
- L2 : 262 Ko
- Warp size : 12
- Maximum number of threads per block : 1024

La RTX 3060 ti dispose de 4864 coeurs Cuda @ 1,6Ghz

Disposant de 1024 threads par block, une des solutions serait d'attribuer une ligne par block soit 960 pixels -> 960 threads.

Soit le programme à exécuter,

```

1  __global__ void seuillage_kernel(float d_image_in[][SIZE_J][SIZE_I],
2  - float d_image_out[][SIZE_J][SIZE_I])
3  {
4      // A VOUS DE CODER
5      int i, j;
6      float r, g, b;
7
8      j = blockIdx.x;
9      i = threadIdx.x;
10
11     r = d_image_in[0][j][i];
12     g = d_image_in[1][j][i];
13     b = d_image_in[2][j][i];
14
15     if ((r / sqrt(r * r + g * g + b * b)) > 0.7) {
16         d_image_out[0][j][i] = r;
17         d_image_out[1][j][i] = g;
18         d_image_out[2][j][i] = b;
19     }
20 }

```

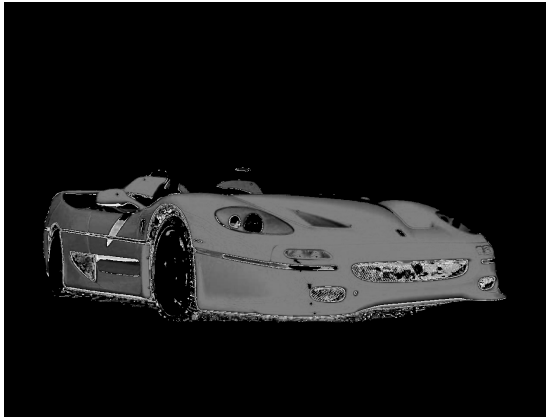
On obtient les résultats suivants :

Cible	CPU (ms)	GPU (ms)	Mem Manag (%)	GPU brut (ms)	R CPU	R Cuda
Nano	84	53	90	5,3	1	1
PC Matlab	13,9				6	
PC C/Cuda	6,5	8,71	98,6	0,12	12,9	44,6

La Nvidia A100 ayant 1,27x plus de coeurs cuda, on peut dire que le temps de traitement brut serait de 0,1 ms et augmentant encore un peu le memory management à 99%. On

peut retomber sur ce résultat en comparant le nombre de coeurs de la Nano, soit un rapport de  $48,35, 5,3/48,35 = 0,109$  ms.

Et les images en sortie :



(a) Image après seuil CPU



(b) Image après seuil GPU et canal rouge

FIGURE 1.2 – Exécution du programme de seuil

## 1.2 - Kmeans Python

---

Le premier TP de l'UE d'outils pour le Machine Learning (D3), a pour objectif d'utiliser l'algorithme des Kmeans pour transformer les pixels d'une image vers le centre du cluster auquel il appartient. Pour ce faire, nous utilisons différents outils disponibles sur Python et scikit-learn. Malheureusement, ce script est assez gourmand et son temps d'exécution est élevé. On propose avec le code suivant d'observer le gain de performance apporté par CUDA.

```
1 import cupy as cp
2 import numpy as np
3 from sklearn.cluster import KMeans
4 from skimage import io
5 import time
6
7 # Load the image using skimage
8 image = io.imread('fruits.jpg')
9
10 # Convert the image to a CuPy array
11 image_cp = cp.asarray(image)
12
13 # Flatten the image into a 2D array of pixels
14 image_flat = image_cp.get().reshape(image_cp.shape[0] *
15     ~ image_cp.shape[1], image_cp.shape[2])
16
17 def Kmeans_cuda(K=1):
```



```

17     # Use KMeans to cluster the pixels into a specified number of
18     ↪ clusters
19     kmeans = KMeans(n_clusters=K, random_state=0).fit(image_flat)
20
21     # Predict the cluster for each pixel
22     clusters = kmeans.predict(image_flat)
23
24     # Create a new CuPy array to hold the modified image
25     new_image_cp = cp.empty_like(image_cp)
26
27     # Iterate over each pixel and assign its value to the corresponding
28     ↪ cluster center
29     for i, cluster in enumerate(clusters):
30         new_image_cp[i // image_cp.shape[1], i % image_cp.shape[1]] =
31         ↪ cp.asarray(kmeans.cluster_centers_[cluster])
32
33     # Convert the CuPy array back to a NumPy array
34     new_image = cp.asnumpy(new_image_cp)
35
36     # Save the modified image using skimage
37     io.imsave("fruits" + "%d" % K + "_cuda.jpg", new_image)
38
39     for K in range(1,256):
40         start_time = time.time()
41         Kmeans_cuda(K=K)
42         end_time = time.time()
43         print(f"It took {end_time-start_time:.2f} seconds to compute for K
44         ↪ =",K)

```

Et le résultat est assez flagrant :

Pour K = 20 clusters, on a un rapport de temps d'exécution de 325 entre CPU et CUDA.

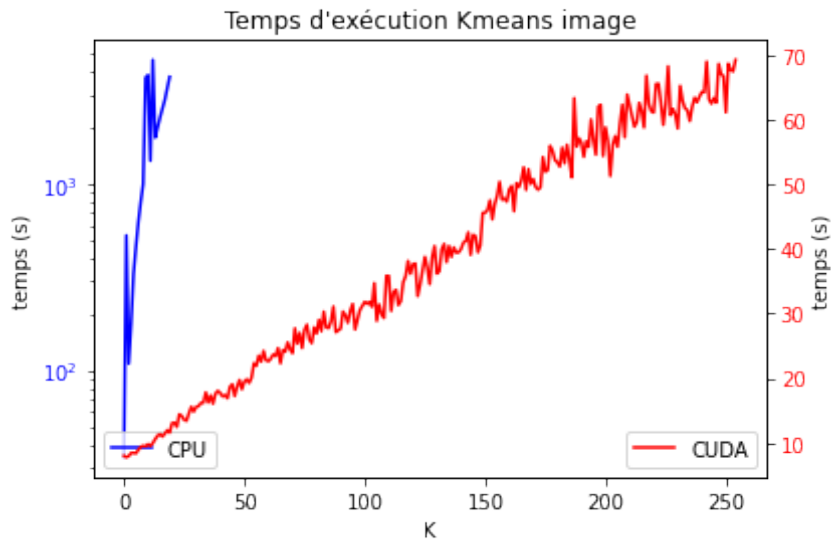


FIGURE 1.3 – Temps d'exécution Kmeans sur CPU et CUDA d'une image de 103 230 pixels

### 1.3 - Conclusion

---

Au cours de ce TP, nous avons appris à utiliser les outils CUDA en passant par une première analyse sous Matlab et en C. Nous avons pu constater la puissance que peut offrir CUDA pour le traitement d'un ensemble très large de données. La majeur problème est l'accès aux données qui peuvent ralentir toute l'exécution du problème, notamment lorsque l'on possède un GPU très puissant créant un bottleneck sur le transport de données.