

université  
PARIS-SACLAY

M2 E3A - SETI

---

## TP GPU

---

SAŠA RADOSAVLJEVIC



27 FÉVRIER 2023



# Table des matières

<b>Table des matières</b>	<b>1</b>
<b>1 TP1 Seuillage</b>	<b>2</b>
1.1 Contexte . . . . .	2
1.2 Matlab . . . . .	2
1.3 Jetson nano CPU . . . . .	4
1.4 Jetson Nano CUDA . . . . .	4
1.5 Kmeans Python . . . . .	6
1.6 Conclusion . . . . .	8
<b>2 TP2 Réduction sommation</b>	<b>9</b>
2.1 Contexte . . . . .	9
2.2 Analyse du code . . . . .	9
2.3 Réduction GPU . . . . .	10
2.3.1 Méthode non parallélisée . . . . .	11
2.3.2 Méthode par réduction . . . . .	11
2.4 Réduction optimisée . . . . .	12
2.4.1 Méthode non parallélisée . . . . .	12
2.4.2 Méthode optimisée . . . . .	13
2.4.3 Méthode optimisée sans mysum+= . . . . .	15
2.5 Résultats . . . . .	16
2.6 Conclusion . . . . .	16

# 1 | TP1 Seuillage

## 1.1 - Contexte

---

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.

## 1.2 - 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

## 1.3 - 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.

## 1.4 - 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 }

```

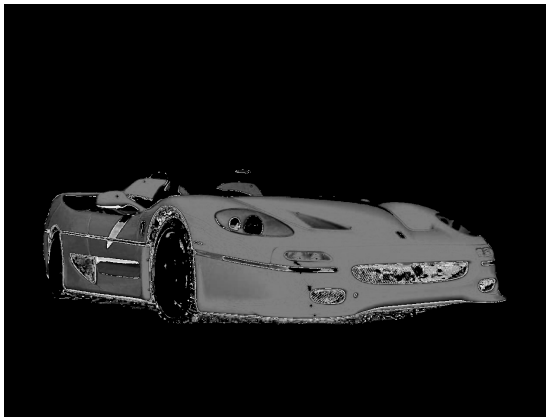
On obtient les résultats suivants :

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

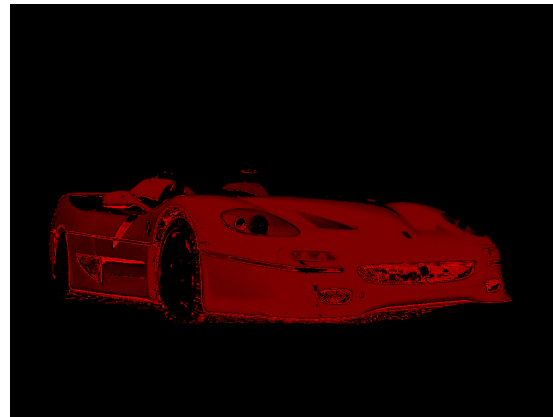
La Nvidia A100 ayant 1,27x plus de coeurs cuda, on peut dire que le temps de traitement brut serait de 0,1 ms et augmentant encore un peu le memory management à 99%. On

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

## 1.5 - 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.

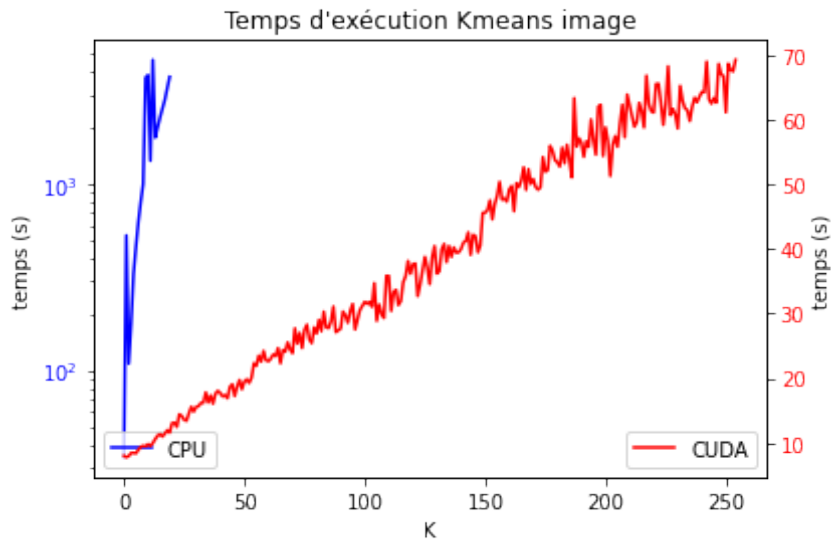


FIGURE 1.3 – Temps d'exécution Kmeans sur CPU et CUDA d'une image de 103 230 pixels

## 1.6 - 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.

## 2 | TP2 Réduction sommation

### 2.1 - Contexte

---

Au cours de ce TP, nous travaillerons soit sur Nvidia Jetson Nano, soit une carte Nvidia RTX 3060 Ti personnelle. (Voir TP Seuillage pour les specs de la Nano)

RTX 3060 ti :

- 38 MP
- 4864 Cuda Cores
- Max clock rate = 1665 MHz - Max threads/block = 1024 - Max threads/multiprocessor = 1536 - Memory 8192 MB

On souhaite mesurer les temps exécutions que peuvent proposer les différentes solutions proposées à notre problème de sommation sur des variables dépendantes. L'objectif est d'essayer de profiter de toutes les capacités offerte par un GPU pour la parallélisation de ce type d'algorithme.

### 2.2 - Analyse du code

---

Nous avons 5 voire 6 (solution) méthodes pour essayer de réduire le temps de calcul au maximum.

Calcul sur CPU :

```
1 int sum = 0;
2 for( int i_test = 0 ; i_test < NUM_TESTS ; ++i_test )
3 {
4     sum = 0;
5     for( int i = 0 ; i < N ; ++i )
6         sum += a_host[i];
7 }
```

Calcul sur CPU + OpenMP :

```

1 omp_sum = 0;
2 #pragma omp parallel shared(omp_sum){
3     #pragma omp for reduction(+:omp_sum)
4     for(int i=0; i<N; i++){
5         omp_sum += a_host[i];
6     }
7 }

```

Calcul sur GPU avec Thrust :

```

1 int thrust_sum = 0;
2 for( int i_test = 0 ; i_test < NUM_TESTS ; ++i_test )
3 {
4     thrust_sum = thrust::reduce( thrust::device_ptr<int>(a_device),
5     thrust::device_ptr<int>(a_device+N) );
6 }

```

Calcul sur GPU :

```

1 int gpu_sum = 0;
2 for( int i_test = 0 ; i_test < NUM_TESTS ; ++i_test )
3 {
4     gpu_sum = reduce_on_gpu( N, a_device );
5 }

```

Calcul sur GPU optimisé :

```

1 int optim_gpu_sum = 0;
2 for( int i_test = 0 ; i_test < NUM_TESTS ; ++i_test )
3 {
4     optim_gpu_sum = reduce_on_gpu_optimized<BLOCK_DIM>( N, a_device );
5 }

```

## 2.3 - Réduction GPU

---

La première solution non optimisée consiste à utiliser les blocs de threads disponibles sur le GPU.

### 2.3.1 - Méthode non parallélisée

```
1 // Compute the sum of my elements.
2 for (int i = threadIdx.x + range.x; i < range.y; i += blockDim.x) {
3     my_sum += in_buffer[i];
4 }
5
6 // Copy my sum in shared memory.
7 s_mem[threadIdx.x] = my_sum;
8
9 // Make sure all the threads have copied their value in shared memory.
10 __syncthreads();
11
12 // Compute the sum inside the block.
13 // The first thread of the block stores its result.
14 if (threadIdx.x == 0) {
15     float sum2 = 0;
16     for (int i = 0; i < blockDim.x; i++) {
17         sum2 += s_mem[i];
18     }
19     out_buffer[blockIdx.x] = sum2;
20 }
```

### 2.3.2 - Méthode par réduction

```
1 // Compute the sum of my elements.
2 int my_sum = 0;
3 for (int idx = range.x + threadIdx.x; idx < range.y; idx += blockDim.x)
4     my_sum += in_buffer[idx];
5
6 // Copy my sum in shared memory.
7 s_mem[threadIdx.x] = my_sum;
8
9 // Make sure all the threads have copied their value in shared memory.
10 __syncthreads();
11
12 int offset;
13 // Compute the sum inside the block.
14 for (offset = blockDim.x / 2; offset > 16; offset /= 2)
15 {
16     if (threadIdx.x < offset)
17         s_mem[threadIdx.x] += s_mem[threadIdx.x + offset];
18     __syncthreads();
19 }
20
```

```

21 //INSIDE WARP 0 SYNC NOT NECESSARY
22 for (; offset > 0; offset /= 2)
23 {
24     if (threadIdx.x < offset)
25         s_mem[threadIdx.x] += s_mem[threadIdx.x + offset];
26     //__syncthreads( );
27 }
28
29
30 // The first thread of the block stores its result.
31 if (threadIdx.x == 0)
32     out_buffer[blockIdx.x] = s_mem[0];

```

## 2.4 - Réduction optimisée

Cette fois-ci, on va chercher à utiliser les plus petits blocs de threads pouvant être utilisés, aussi appelés des warps. Ils contiennent chacun 32 threads.

### 2.4.1 - Méthode non parallélisée

```

1 // Compute the sum of my elements.
2 for (int i = threadIdx.x + range.x; i < range.y; i += blockDim.x) {
3     my_sum += in_buffer[i];
4 }
5 // Copy my sum in shared memory.
6 s_mem[threadIdx.x] = my_sum;
7
8 // Compute the sum inside each warp.
9 __syncthreads();
10
11 // Each warp leader stores the result for the warp.
12 if (lane_id == 0) {
13     float sum2 = 0;
14     for (int i = threadIdx.x; i < threadIdx.x + WARP_SIZE; i++) {
15         sum2 += s_mem[i];
16     }
17     s_mem[threadIdx.x] = sum2;
18 }
19 __syncthreads();
20
21 if (warp_id == 0)
22 {
23     // Read my value from shared memory and store it in a register.
24     // Sum the results of the warps.

```

```

25     if (threadIdx.x == 0) {
26         float sum3 = 0;
27         for (int i = threadIdx.x; i < threadIdx.x + blockDim.x; i +=
← WARP_SIZE) {
28             sum3 += s_mem[i];
29         }
30
31         out_buffer[blockIdx.x] = sum3;
32     }
33 }

```

## 2.4.2 - Méthode optimisée

```

1  // Compute the sum of my elements.
2  int my_sum = 0;
3  for (int idx = range.x + threadIdx.x; idx < range.y; idx += BLOCK_DIM)
4      my_sum += in_buffer[idx];
5
6  // Copy my sum in shared memory.
7  s_mem[threadIdx.x] = my_sum;
8
9  // Compute the sum inside each warp.
10 #pragma unroll
11 for (int offset = 16; offset > 1; offset >>= 1)
12     if (lane_id < offset)
13         s_mem[threadIdx.x] = my_sum += s_mem[threadIdx.x + offset];
14
15 __syncthreads();
16
17 // Each warp leader stores the result for the warp.
18 if (lane_id == 0)
19     s_mem[warp_id] = my_sum += s_mem[threadIdx.x + 1];
20
21 __syncthreads();
22
23 if (warp_id == 0)
24 {
25     // Read my value from shared memory and store it in a register.
26     my_sum = s_mem[lane_id];
27
28     // Sum the results of the warps.
29 #pragma unroll
30     for (int offset = NUM_WARPS / 2; offset > 1; offset >>= 1)
31         if (threadIdx.x < offset)
32             s_mem[threadIdx.x] = my_sum += s_mem[threadIdx.x + offset];
33 }

```

```
34  
35 // The 1st thread stores the result of the block.  
36 if (threadIdx.x == 0)  
37     out_buffer[blockIdx.x] = my_sum += s_mem[1];
```



### 2.4.3 - Méthode optimisée sans mysum+=

```
1 // Compute the sum of my elements.
2 int my_sum = 0;
3 for (int idx = range.x + threadIdx.x; idx < range.y; idx += BLOCK_DIM)
4     my_sum += in_buffer[idx];
5
6 // Copy my sum in shared memory.
7 s_mem[threadIdx.x] = my_sum;
8
9 // Compute the sum inside each warp.
10 #pragma unroll
11 for (int offset = 16; offset > 0; offset >>= 1)
12     if (lane_id < offset)
13         s_mem[threadIdx.x] += s_mem[threadIdx.x + offset];
14
15 __syncthreads();
16
17 // Each warp leader stores the result for the warp.
18 if (lane_id == 0)
19     s_mem[warp_id] = s_mem[threadIdx.x];
20
21 __syncthreads();
22
23 if (warp_id == 0)
24 {
25     // Sum the results of the warps.
26     #pragma unroll
27     for (int offset = NUM_WARPS / 2; offset > 0; offset >>= 1)
28         if (threadIdx.x < offset)
29             s_mem[threadIdx.x] += s_mem[threadIdx.x + offset];
30 }
31
32 // The 1st thread stores the result of the block.
33 if (threadIdx.x == 0)
34     out_buffer[blockIdx.x] = s_mem[0];
```

## 2.5 - Résultats

---

Les codes présentés précédemment sont exécutés sur un Desktop composé d'un AMD Ryzen 9 3900x et d'une Nvidia 3060ti.

	Non parallélisé (ms)	Parallélisé (ms)
CPU 1 Thread	103	103
CPU 24 Threads	8.7	8.7
GPU Thrust	124	124
GPU	2.32	1.99
GPU opti	2.12	1.98
GPU mysum		1.93

TABLEAU 2.1 – Temps d'exécution des différentes versions de la sommation

Il semblerait que le PC ai du mal avec Thrust, mais il n'y a pas d'erreur lors de la compilation sur visual studio.

## 2.6 - Conclusion

---

Nous avons pu apprendre à utiliser les threads en blocs et en warps pour paralléliser au maximum un algorithme de sommation sur une variable dépendante à l'aide de la méthode par réduction. La méthode de réduction est disponible par défaut sur openMP et avec Thrust nous permettant en une simple instruction de réaliser le code de réduction qui, cependant, obtient de meilleure performances lorsqu'il est réécrit pour l'application en ayant connaissance de l'architecture de l'algorithme et d'un GPU.