

Article

# Parallel K-Means Clustering for Brain Cancer Detection Using Hyperspectral Images

Emanuele Torti <sup>1,\*</sup>, Giordana Florimbi <sup>1</sup>, Francesca Castelli <sup>1</sup>, Samuel Ortega <sup>2</sup>,  
Himar Fabelo <sup>2</sup>, Gustavo Marrero Callicó <sup>2</sup>, Margarita Marrero-Martin <sup>2</sup> and  
Francesco Loporati <sup>1</sup>

<sup>1</sup> Department of Electrical, Computer and Biomedical Engineering, University of Pavia, I-27100 Pavia, Italy; giordana.florimbi01@ateneopv.it (G.F.); francesca.castelli02@ateneopv.it (F.C.); francesco.leporati@unipv.it (F.L.)

<sup>2</sup> Institute for Applied Microelectronics (IUMA), University of Las Palmas de Gran Canaria (ULPGC), 35017 Las Palmas de Gran Canaria, Spain; sortega@iuma.ulpgc.es (S.O.); hfabelo@iuma.ulpgc.es (H.F.); gustavo@iuma.ulpgc.es (G.M.C.); margarita@iuma.ulpgc.es (M.M.-M.)

\* Correspondence: emanuele.torti@unipv.it; Tel.: +39-0382-985678

Received: 9 October 2018; Accepted: 26 October 2018; Published: 30 October 2018



**Abstract:** The precise delineation of brain cancer is a crucial task during surgery. There are several techniques employed during surgical procedures to guide neurosurgeons in the tumor resection. However, hyperspectral imaging (HSI) is a promising non-invasive and non-ionizing imaging technique that could improve and complement the currently used methods. The HypErspectraL Imaging Cancer Detection (HELICoiD) European project has addressed the development of a methodology for tumor tissue detection and delineation exploiting HSI techniques. In this approach, the K-means algorithm emerged in the delimitation of tumor borders, which is of crucial importance. The main drawback is the computational complexity of this algorithm. This paper describes the development of the K-means clustering algorithm on different parallel architectures, in order to provide real-time processing during surgical procedures. This algorithm will generate an unsupervised segmentation map that, combined with a supervised classification map, will offer guidance to the neurosurgeon during the tumor resection task. We present parallel K-means clustering based on OpenMP, CUDA and OpenCL paradigms. These algorithms have been validated through an in-vivo hyperspectral human brain image database. Experimental results show that the CUDA version can achieve a speed-up of ~150× with respect to a sequential processing. The remarkable result obtained in this paper makes possible the development of a real-time classification system.

**Keywords:** Graphics Processing Units (GPUs); CUDA; OpenMP; OpenCL; K-means; brain cancer detection; hyperspectral imaging; unsupervised clustering

## 1. Introduction

One of the most diffused types of cancer is the brain tumor, which has an estimated incidence of 3.4 per 100,000 subjects [1]. There are different types of brain tumors; the most common one concerns the *glial* cells of the brain and is called *glioma*. It accounts from the 30% to the 50% of the cases. In particular, in the 85% of these cases, it is a malignant tumor called *glioblastoma*. Moreover, these kind of *gliomas* are characterized by fast-growing invasiveness, which is locally very aggressive, in most cases unicentric and rarely metastasizing [2].

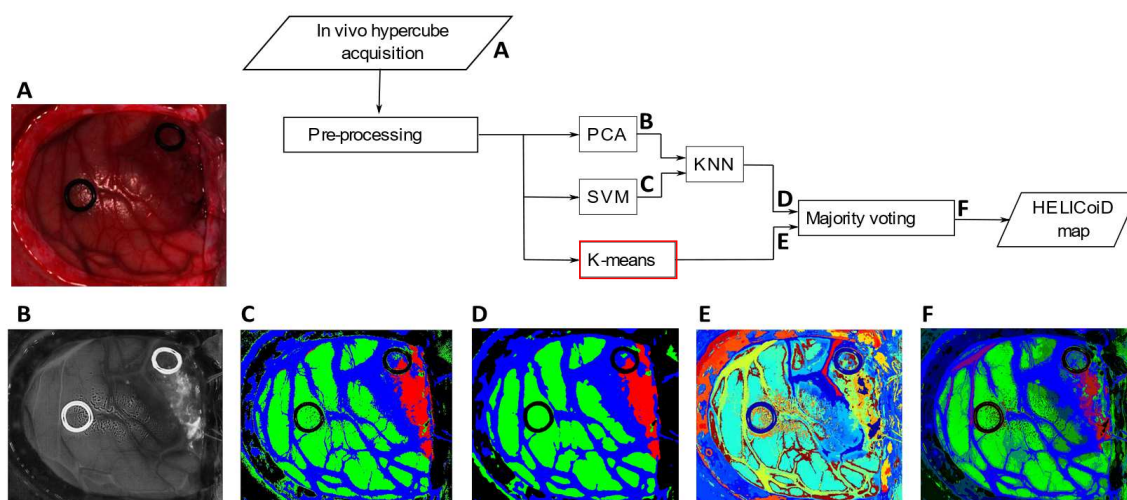
Typically, the first diagnosis is performed through the Magnetic Resonance Imaging (MRI) and the Computed Tomography (CT). Those techniques are capable to highlight possible lesions. However, it is not always possible to use them, since they can, for example, make interference with

pacemakers or other implantable devices. Moreover, the certainty of the diagnosis only comes from the histological and pathological analyses, which require samples of the tissue. In order to obtain this tissue, an *excisional biopsy* is necessary, which consists in the removal of tissue from the living body through surgical cutting. It is important to notice that all those approaches have some disadvantages; in particular they are not capable of providing a real-time response and, most important, they are invasive and/or ionizing.

The clinical practice for brain cancers is the tumor resection, which can cure the lowest grade tumors and prolongs the life of the patient in the most aggressive cases. The main issue about this approach is the inaccuracy of the human eye in distinguishing between healthy tissue and cancer. This is because the cancer often infiltrates and diffuses into the surrounding healthy tissue and this is particularly critical for brain cancers. As a consequence, the surgeon can unintentionally leave behind tumor tissue during a surgery routine potentially causing tumor recurrence. On the other hand, if the surgeon removes too much healthy tissue, a permanent disability to the patient can be provoked [3].

The HELICoiD European project aims at providing to the surgeon a system which can accurately discriminate between tumor and healthy tissue in real-time during surgery routines [4,5]. Traditional imaging techniques feature a low grade of sensitivity and often cannot clearly determine the tumor region and its boundaries. Therefore, the HELICoiD project exploits Hyperspectral Imaging (HSI) techniques in order to solve this critical issue. Hyperspectral images (HS) can be acquired over a wide range of the electromagnetic spectrum, from visible to near-infrared frequencies and beyond. Hyperspectral sensors acquire the so-called *HS cube* where the spatial information is in the  $x$ -axis and in the  $y$ -axis, while the spectral information is in the  $z$ -axis. Thus, a single hyperspectral pixel can be seen as a mono-dimensional vector, which contains the spectral response across the different wavelengths. Moreover, it is important to notice that the spectral information is strictly correlated with the chemical composition of the specific material. It is possible to say that each hyperspectral pixel contains the so-called *spectral signature* of a certain substance. Thus, different substances can be distinguished by properly analyzing those images [6].

A previous study [4,7] proposed a processing chain for hyperspectral image analysis acquired during brain surgery. The framework developed in this work is depicted in Figure 1.



**Figure 1.** Hyperspectral brain cancer detection algorithm proposed in [7]. After acquiring (A) and pre-processing the image, the system performs a supervised classification through Principal Component Analysis (PCA) (B), Support Vector Machine (SVM) (C) and K-Nearest Neighbor (KNN) (D). Moreover, it generates a segmentation map through the K-means (E). The (D,E) maps are merged using the majority voting (F).

First, the acquired HS cube (Figure 1A) is pre-processed in order to perform a radiometric calibration, reduce the noise and the dimensionality of the HS image and normalize it. After this

preparatory step, the image is analyzed using the supervised and unsupervised classification. The former is performed exploiting the Principal Component Analysis (PCA), the Support Vector Machine (SVM) and the K-Nearest Neighbor (KNN). The KNN filters the spatial information given by the PCA (Figure 1B) and the classification map generated by the SVM (Figure 1C). Its output (Figure 1D) is a map where tissues are displayed with different colors representing the associated classes. The unsupervised classification is based on the K-means algorithm. Despite the supervised classification output, the unsupervised result is a segmentation map (Figure 1E), whose clusters are semantically meaningless. However, the K-means provides a good delimitation of the different areas present in the scene. Since the goal of the system is to accurately delineate the tumor borders, the K-means plays a crucial role for its ability to clearly separate different areas. For these reasons, it is important to merge the two outputs in order to exploit the benefits of the two approaches. The majority voting provides the final output combining the supervised and unsupervised classifications (Figure 1F).

While the PCA, the SVM and the KNN filter are executed through a fixed number of steps, the K-means algorithm iterates until a certain condition is satisfied. In order to provide the real-time classification during surgery, parallel computing is required, since the computational load of the algorithms is extremely high. The other algorithms of this framework have been already developed in parallel, in particular the SVM [8] and the KNN filtering [9] have been recently proposed in the literature. Those works target Graphics Processing Units (GPUs) technology since the considered algorithms have an intrinsically parallel structure. Previously, other parallel technologies have been evaluated in order to provide faster implementations of the PCA [10], SVM [11] and KNN [12] compared to the serial ones. Despite this, in our work we choose to exploit GPUs since they assure higher performance. Moreover, GPUs are going to be increasingly used for real-time image processing [13–15], together with other scientific applications related to simulation and modeling or machine learning in biomedical applications [16,17].

In this paper, we present the parallelization of the K-means algorithm on different parallel architectures in order to evaluate which one is more suitable for real-time processing. In particular, we consider multi-core CPUs through the OpenMP API and the GPU technology using NVIDIA CUDA framework. We also propose OpenCL-based implementations in order to address code portability.

In other words, the work performed allows identifying the best suitable parallel approach between one that could be more appealing since it requires low programming effort and another one more efficient but also more demanding in terms of optimization and tuning. A tool that allows intra-architectures portability (OpenCL) was also considered but due to its lower performance it is not competitive with the other two approaches.

The paper is organized as follows: Section 2 describes the K-means algorithm for hyperspectral images, while Section 3 details the different parallel versions. Section 4 contains the experimental results and their discussion, making comparisons between the different approaches described in this paper. Section 5 concludes the papers and addresses some possible future research lines.

## 2. K-Means Algorithm for Hyperspectral Images

As already said, the K-means algorithm, unlike the other ones of the hyperspectral brain cancer detection algorithm, is not performed through a fixed number of steps. It performs an unsupervised learning since no previous knowledge of the data is needed. The algorithm separates the input data into  $K$  different clusters with a  $K$  value fixed a priori. Data are grouped together on the basis of feature similarity. The first step of the algorithm is the definition of  $K$  centroids, one for each cluster. Using those centroids, a first grouping is performed on the basis of the distance of each point to the centroids. A point is associate to the cluster represented by the nearest centroid. At this moment, each  $k$  centroid is updated as the baricenter of the group it represents. This process iterates until the difference between the centroids of two consecutive iterations are smaller than a fixed threshold or if the maximum number of iterations is reached.

The pseudo-code of the K-means algorithm is shown in Algorithm 1, where  $Y$  indicates a hyperspectral image made up of  $N$  pixels and  $L$  bands. Therefore, the hyperspectral image can be seen as an  $N \times L$  matrix. The number of clusters to produce is determined by  $K$ , the threshold error by  $\text{min\_error}$  and the maximum number of iterations by  $\text{max\_iter}$ . The K-means algorithm produces as a result a  $K \times L$  array containing the centroids, which will be referred as  $\text{cluster\_centroids}$  in Algorithm 1 and an  $N$ -dimensional array containing the label of the cluster assigned to each pixel. This array is denoted by  $\text{assigned\_cluster}$ .

---

**Algorithm 1** K-means
 

---

**Input:**  $Y, K, \text{min\_error}, \text{max\_iter}$

```

1: Pseudo-random initialization of cluster_centroids
2: Initialize previous_centroids at 0                                ▷ previous_centroids is an  $K \times L$  array
3: n_iter ← 0                                                       ▷ initialize the iteration counter to 0
4: Initialize actual_error with a huge value
5: while actual_error > min_error and n_iter < max_iter do
6:   for i:=1 to  $N$  do
7:     Initialize centroid_distances to 0                            ▷ centroid_distances is a  $K$ -dimensional array
8:     for j:=1 to  $K$  do
9:       centroid_distancej ← distance between the  $j$ -th centroid and the  $i$ -th pixel
10:    end for
11:    assigned_clusteri ← index of min centroid_distance
12:  end for
13:  previous_centroids ← cluster_centroids
14:  update cluster_centroids
15:  actual_error ←  $\frac{\sum_{i=1}^K \sum_{j=1}^L |\text{previous\_centroids}_{i,j} - \text{cluster\_centroids}_{i,j}|}{K \cdot L}$ 
16:  n_iter ← n_iter + 1
17: end while

```

**Output:**  $\text{assigned\_cluster}, \text{cluster\_centroids}$

---

In Algorithm 1, lines 1 and 2 contain the initialization of the variables. In particular,  $\text{cluster\_centroids}$  is initialized with  $K$  different hyperspectral pixels pseudo-randomly chosen from the input image  $Y$ . The variable  $\text{actual\_error}$  is initialized with a huge value in order to ensure that the main loop of the algorithm (from line 5 to 17) is performed at least one time. Inside this main loop there are two *for* loops that iterate over the number of pixels  $N$  and the number of clusters  $K$  (lines from 6 to 12). For each pixel, a temporary array  $\text{centroid\_distances}$  is set to 0, used for storing the distances between the considered hyperspectral pixel and the centroids. The distance metric used for hyperspectral pixels is usually the *Spectral Angle* (SA) which is defined as:

$$SA = \theta(x, y) = \cos^{-1} \left( \frac{\sum_{h=1}^L x_h y_h}{\left( \sum_{h=1}^L x_h^2 \right)^{1/2} \left( \sum_{h=1}^L y_h^2 \right)^{1/2}} \right) \quad (1)$$

where  $x$  and  $y$  are the spectral vectors and  $x_h$  and  $y_h$  represent the response of the  $h$ -th band of  $x$  and  $y$  respectively, being  $L$  the number of bands.

The label assigned to the  $i$ -th pixel corresponds to the group represented by the centroid with the minimum SA value, as shown in line 11. This phase is repeated for each pixel.

After these steps, the centroids used for the SAs computation are stored in the  $\text{previous\_centroids}$  array. Successively, the centroids are updated by computing the barycenter of each group that is

computing the mean value, for each band, of the pixels belonging to the group. Using the updated centroids and the previous ones, it is possible to evaluate the variation from the previous iteration. It represents how much the centroids have changed and it can be used as a stopping criterion when these variations become small (line 15). The last step of the *while* loop is the increment of the *n\_iter* variable, used for controlling the maximum number of iterations performed by the algorithm.

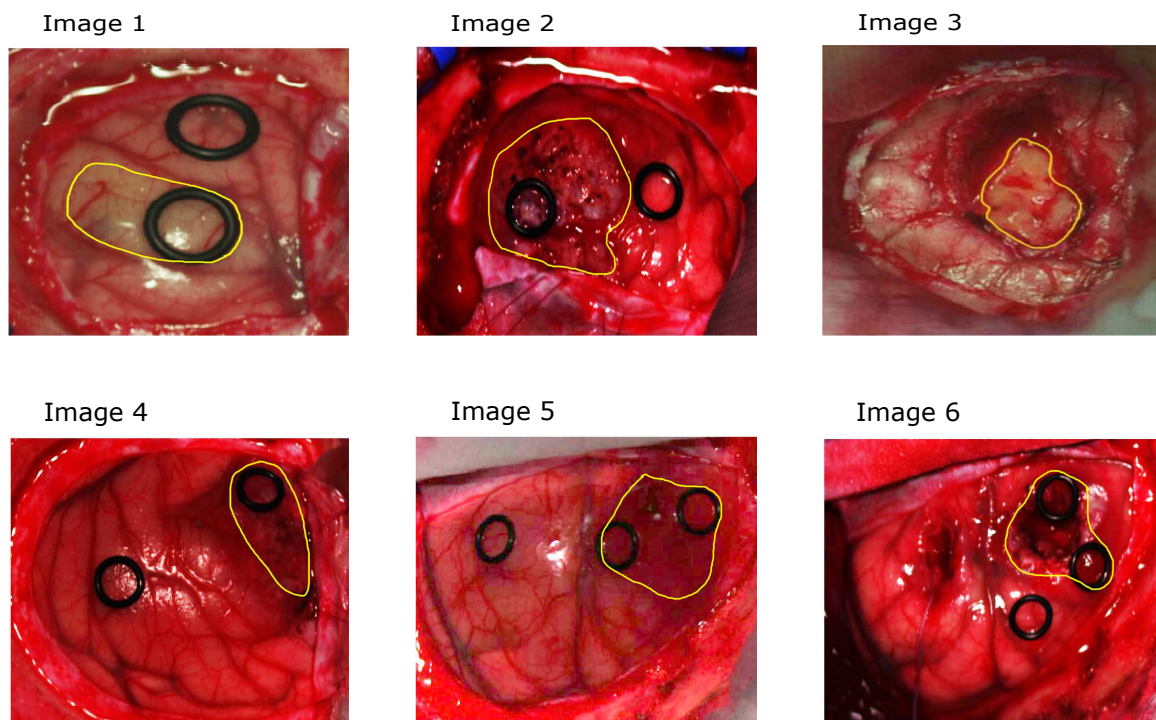
The next section describes the serial and the parallel versions of this algorithm that we developed using different parallel approaches, together with a code profiling carried out in order to identify the heaviest code parts from the computational point of view.

### 3. Parallel K-Means Implementations

First, we developed a serial version of the K-mean algorithm written in C code. It serves both as reference for validating the results of the parallel implementations and for performing a careful code profiling needed to identify the most complex code parts. The numerical representation used is the IEEE-754 floating-point single precision.

#### 3.1. Serial Code Profiling

The code profiling was performed using a dataset formed by real HS images and assuming  $K = 24$ ,  $\text{min\_error} = 10^{-3}$  and  $\text{max\_iter} = 50$ . This  $K$  value was established during the development of the HS brain cancer algorithm presented in [7]. Using this configuration, the execution of the algorithm never reached the maximum number of iterations. The characteristics of the dataset are shown in Table 1, while Figure 2 shows the RGB representation of the images.



**Figure 2.** RGB representations of each hyperspectral cube of the brain cancer dataset. The yellow line represents the tumor location identified by the neurosurgeon.

**Table 1.** Dataset characteristics.

Image ID	# of Rows	# of Columns	Total # of Pixels	# of Bands	Size (MB)
Image 1	329	379	124,691	128	60.88
Image 2	493	376	185,368	128	90.51
Image 3	402	472	189,744	128	92.65
Image 4	496	442	219,232	128	107.05
Image 5	548	459	251,532	128	122.82
Image 6	552	479	264,408	128	129.11

The profiling highlighted that the heaviest code parts are the computation of distances, which are evaluated between each hyperspectral pixel and each centroid. In the considered cases, the *for* loops of lines 6–12 (Algorithm 1) take from 94 to 98% of the time for the smallest and the biggest image, respectively. Notice that these computations can be performed in parallel, since there is no dependency between the evaluations needed by a single pixel and the others.

### 3.2. OpenMP Algorithms

OpenMP (<https://www.openmp.org/>) is a parallel programming framework capable of exploiting multi-core architectures. It is based on a set of simple *#pragma* statements used for code annotations that indicates to the compiler which parts should be parallelized. An example is the *#pragma omp parallel for* statement, which generates a set of parallel threads and assigns to each one a group of iterations. It is also possible to indicate to the compiler which variables should be shared among the threads and which ones are private through the *shared* and *private* clauses, respectively. Finally, it is possible to choose the scheduling algorithm to use through the *schedule* option. The supported scheduling algorithms are *static*, *dynamic* and *guided*. In the first case, the number of iterations are equally or as equal as possible subdivided among the threads. Thus, each thread performs the same number of iterations. The *dynamic* scheduling uses the internal work queue to give a chunk-sized block of loop iterations to each thread. When a thread finishes, it retrieves the next block of loop iterations from the top of the work queue. The default value of the chunk size is 1, but it is possible to change it by a proper command. Finally, the *guided* scheduling is similar to the *dynamic* one, but the chunk size starts from a big value and then decreases in order to manage load imbalance between different iterations.

We developed two different OpenMP versions of the K-means algorithm. The first one parallelizes the *for* loop which iterates over the hyperspectral pixels (line 6, Algorithm 1). In this way, at each iteration of the main *while* loop, a set of parallel threads are generated and each one computes the SA between a certain group of pixels and the centroids. All the other operations are performed in a serial way. The shared arrays are the *cluster\_centroids* and the input image *Y*, while all the other variables are private. In this version, the parallel region is created and destroyed at each iteration of the main *while* loop.

Concerning the second implementation, the majority of the operations are performed in parallel. The operations that continue to be performed sequentially are the *actual\_error* computation and the increment of *n\_iter*, at lines 15 and 16 of Algorithm 1, respectively. A barrier must be placed after the *actual\_error* computation, in order to prevent the other threads to evaluate the *while* condition with an inconsistent old value. In this case, also the *centroid\_distance* is declared as shared. Notice that, in this version, the parallel region is created and destroyed only once, at the beginning and at the end of the main *while* loop. However, in this case, it is necessary to introduce a barrier in order to ensure the correct execution of the program.

### 3.3. CUDA Algorithms

CUDA (<https://developer.nvidia.com/cuda-zone>) is a parallel programming framework developed by NVIDIA to exploit GPU computing power. In this framework, the GPU, also called *device*, is seen as a parallel co-processor, with separated address space with respect to the CPU, also called *host*. The execution of a CUDA program always begins from the host, using a serial thread. When it is necessary to perform a parallel operation, the host allocates memory on the GPU and transfers the data to that memory. Those two operations are performed through the *cudaFree* and *cudaMemcpy* routines. At this point, the GPU generates thousands of parallel threads, which cooperate in order to perform the desired computation. The function performed by the GPU is called *kernel*. When the kernel execution ends, the CPU retrieves the results from the GPU memory through memory transfer (*cudaMemcpy* routine). The GPU memory is then deallocated by the *cudaFree* routine. The execution proceeds then in a serial way.

The threads generated by the GPU are grouped into *blocks*, which form the *grid*. The blocks can be mono-dimensional, bi-dimensional, or three-dimensional and the number of threads within a block can be chosen by the programmer.

The typical bottleneck of GPU applications is represented by memory transfers. Therefore, it is necessary to properly manage them in order to achieve the best performance.

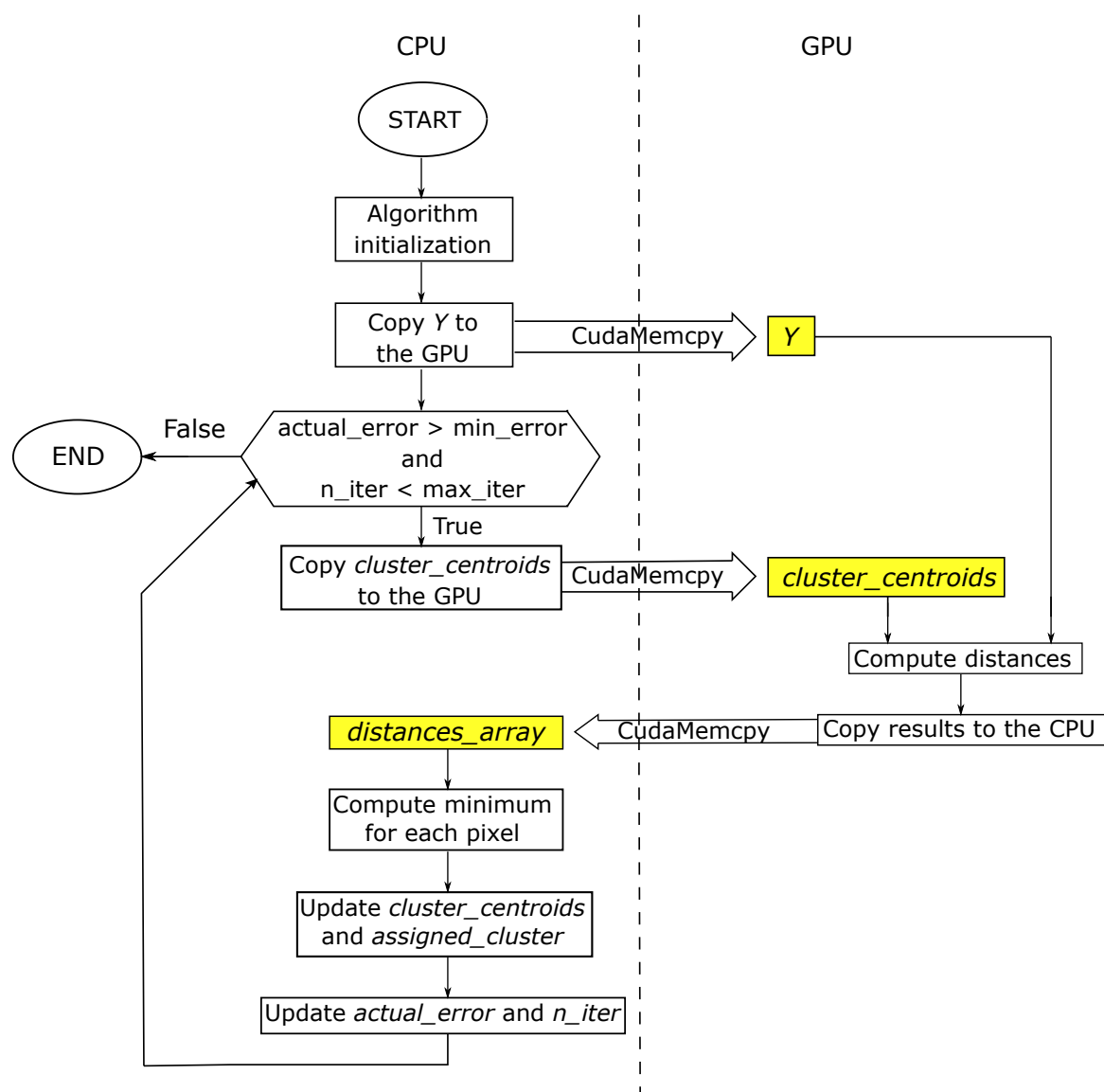
In this work, we present three different parallel versions of the K-means algorithm. The first one is based on the parallelization of the distance computation. In this case, the thread performs the computation of the distance between the assigned pixel and the  $K$  centroids. This kernel takes as inputs the hyperspectral image  $Y$  and the  $K$  centroids stored in the *cluster\_centroids* variable. It produces as output a  $N \times K$  array which contains the distances between each pixel and each centroid. In particular, the  $i$ -th row and the  $j$ -th column of this array store the distance between the  $i$ -th pixel and the  $j$ -th centroid. Therefore, it is necessary to add a supplementary temporary array ( $N \times K$ ) with respect to the serial implementation. The schematization of this implementation is shown in Figure 3.

The hyperspectral image  $Y$  is copied to the GPU memory only once, before the beginning of the main *while* loop. This has been done since the image is not modified by the algorithm. At every iteration, the only data transferred to the GPU is the matrix containing the centroids, that is used, together with the image, to compute the distances, stored in a temporary matrix (*distances\_array* in Figure 3). Data are sent back to the host, which computes the minimum distance for each pixel (i.e., each row of this matrix) and then updates the centroids and computes the error in order to evaluate convergence.

The second CUDA version has been developed in order to avoid the limit of the amount of data transferred during each iteration of the main *while* loop. Therefore, the minimum distance computations, the centroids update, and the error evaluation have been performed on the GPU side. Since the distances are stored in an  $N \times K$  array, the kernel used to find the index of the minimum distance is executed by  $N$  threads. The  $i$ -th thread performs a *for* loop in order to evaluate the minimum distance of the  $i$ -th centroid. In other words, this task has been parallelized by assigning to each thread the computation of the minimum distance for one pixel. The index of the minimum distance is stored in the *assigned\_cluster* array, which contains the classification obtained at the current iteration of the main loop. The update of the centroids has been performed by a simple kernel where the  $i$ -th thread computes the update of the  $i$ -th centroid. Concerning the error evaluation, it is possible to use the highly optimized routines offered by the CUBLAS library. In particular, it is possible to use the *cublasSasum* routine, which calculates the sum of the absolute values stored in the input array. Before activating this kernel the element-wise difference between the values stored in the *actual\_centroids* and *previous\_centroids* arrays must be computed. This is done by a kernel in which a single thread computes the difference between two elements. The sum computed by the *cublasSasum* routine is returned to the host, which performs the final division needed for error evaluation and increments the number of iterations. The schematization of this CUDA version is shown in Figure 4.

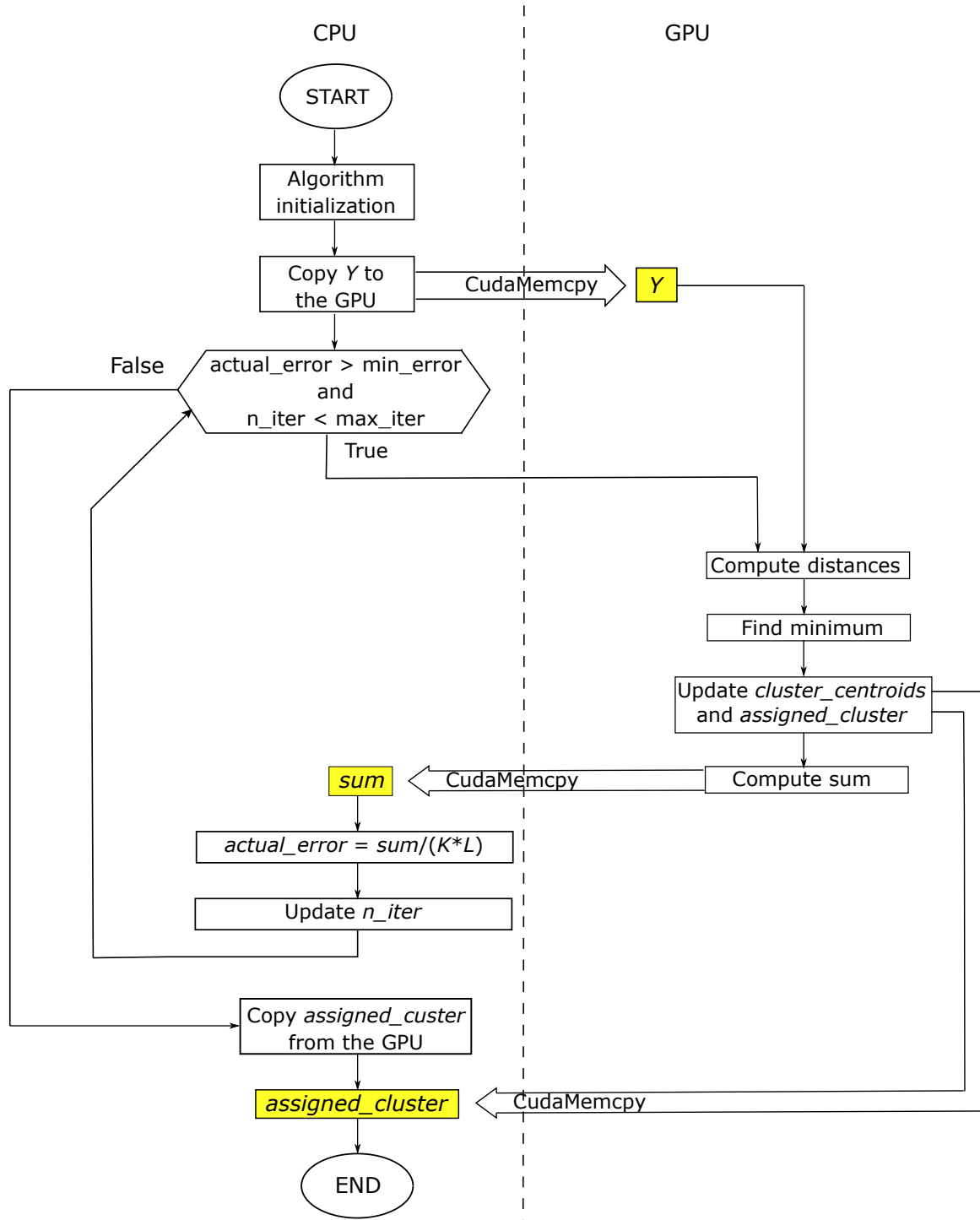
It is important to notice that, in this version, only a single precision floating-point value is transferred at each iteration of the main loop. However, an additional data transfer at the end of the main loop must be performed, since it is necessary to retrieve the *assigned\_cluster* that contains the hyperspectral pixel classification.

The last CUDA version developed in this work exploits the *dynamic parallelism* introduced by CUDA 6.0. This allows to use a thread inside a kernel in order to generate a grid which executes another kernel. In the proposed case, it is possible to take advantage of dynamic parallelism by moving the main *while* loop inside the kernel. In other words, this version is made up of a single kernel executed on the GPU by a single thread, which manages the activation of the kernels already described for the second CUDA version. In this case, the only memory transfers are performed before (the hyperspectral image *Y*) and after the main loop (the classified pixels *assigned\_cluster*). However, it is worth noting that the activation of a kernel from another kernel requires a launching overhead, which will be discussed in Section 4.



**Figure 3.** Schematization of the first GPU implementation. The operations are in white boxes, while data are in yellow boxes.





**Figure 4.** Schematization of the second GPU implementation. The operations are in white boxes, while data are in yellow boxes.

### 3.4. OpenCL Algorithms

OpenCL (<https://www.khronos.org/opencl/>) is a parallel programming framework maintained by the Khronos Group which addresses the issue of portability between devices from different vendors. It assumes a model similar to CUDA, with the difference that the blocks are called *working groups* and the threads are called *working items*. The computing platforms that can be programmed using OpenCL range from multi-core CPUs to manycore GPU and finally to Field Programmable Gate Arrays (FPGAs). Similarly to CUDA, this paradigm assumes that the computing platform is made

up of a serial processor, called *host*, and one or more parallel *devices*. At the beginning of an OpenCL application, it is important to correctly initialize the execution context. This has the effect of pointing out to the OpenCL environment which devices will be used in the case that there are more than one OpenCL compatible boards installed on the same machine. Data that must be processed by the devices are stored into *buffers*, which can be of different types, depending on the targeting devices of the implementation. In particular, considering a generic GPU as a device, a buffer is a portion of the device memory where data are copied from the host memory. On the other hand, if we consider as a device a multi-core CPU or an integrated GPU which shared the RAM memory with the host, the buffer is only a reference to the RAM portion where data are stored. Notice that, in this last case, there will be no data transfers between host and device since the RAM address space is shared. An important difference when compared with CUDA is the absence of dynamic parallelism, thus, it is not possible to activate a kernel from another one. Therefore, we only developed two versions based on OpenCL. The first one performs the parallel computation of the distances on the device while the other operations are performed on the host. The second version performs all the operations inside the main *while* loop on the device, as we implemented in the second CUDA version already described. However, it is not possible to exploit the CUBLAS library, since it is strictly correlated with the adoption of NVIDIA devices. Therefore, we exploit the *clBLAS* library which is very similar to CUBLAS. In particular, we use the *clblasSasum* routine which performs the summation of all the elements stored in a given array. All the other operations have been performed as described for the second CUDA version.

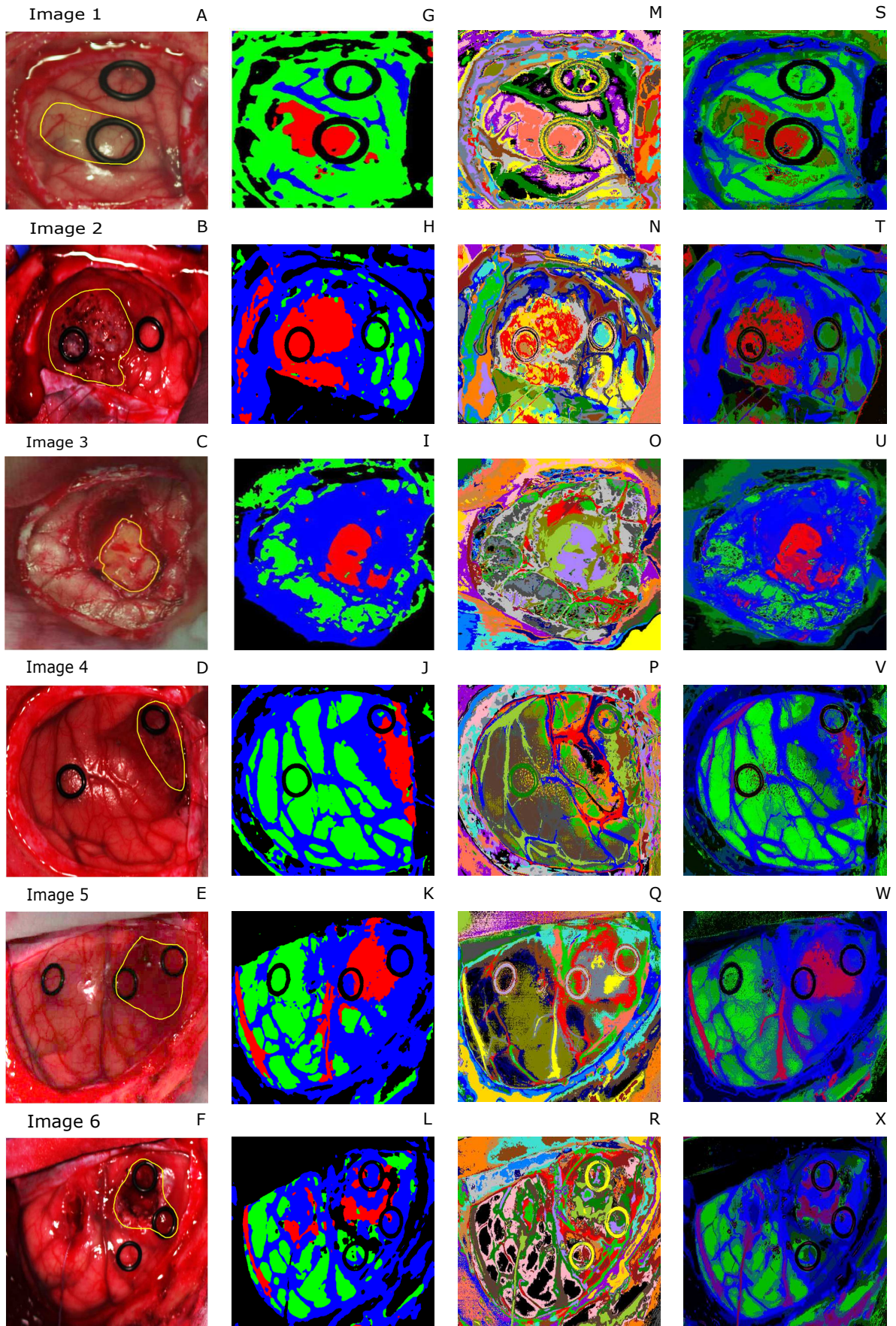
#### 4. Experimental Results and Discussion

All the parallel implementations have been tested with the hyperspectral dataset already used for serial code profiling and shown in Table 1. The images have been employed both to evaluate the processing times of the different versions and to validate the results. All the parallel versions have obtained the same results than the serial one when removing the random initialization. In other words, if the serial and the parallel versions have the same initialization, they perform the same number of iterations and produce the same outputs.

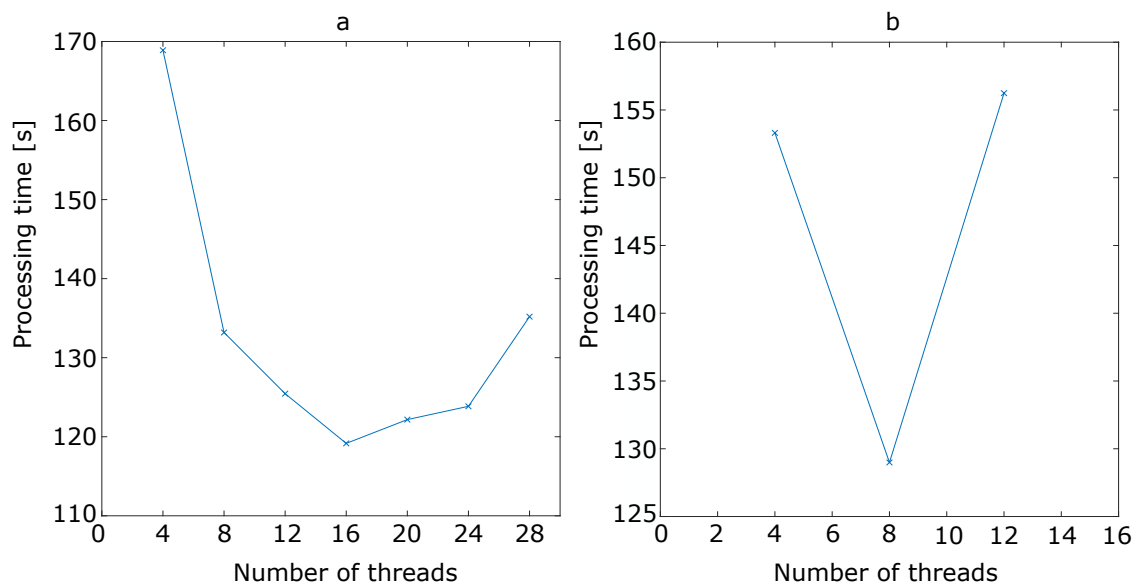
To graphically analyze the classification results, Figure 5A–F shows the RGB images where the tumor is highlighted. Moreover, Figure 5G–L depicts the supervised classification maps where the tumor is indicated with the red color, the healthy tissue in green, the hypervascularized in blue and the background in black. Images in Figure 5M–R are the segmentation maps produced by the K-means algorithm. As can be seen from the images, this algorithm is capable of distinguishing blood vessels, different tissue regions and the ring markers (used by neurosurgeons to label the image for the supervised classification). As said before, the algorithm defines with high accuracy the boundaries of each area, but the clusters do not correspond to specific classes. Moreover, colors are randomly assigned to each cluster. For this reason, it is crucial to combine this segmentation map with the supervised classification result in obtain the final output, shown in Figure 5S–X.

##### 4.1. OpenMP Performance Evaluation

First, for each OpenMP implementation, several tests have been conducted in order to establish the optimal number of threads to be generated, using the biggest image (Image 6). These tests have been performed on an Intel i7 6700 processor working at 3.40 GHz equipped with 32 GB of RAM. The codes have been compiled with the *vc140* compiler, using compilation options to indicate the target architecture (i.e., x64 processor) and to maximize the processing speed. The processing times have been measured through the *omp\_get\_wtime* routine. The obtained results are shown in Figure 6.



**Figure 5.** RGB representations of the dataset (A–F), supervised classification maps (G–L), segmentation maps (M–R), final output (S–X).



**Figure 6.** Processing time for Image 6 with respect to the number of threads for the first (a) and for the second (b) OpenMP version.

It is important to highlight that the experiments have been conducted with the same initialization, and they provide the same number of iterations and the same classification results. For the first OpenMP implementation, we measured the processing time from 4 to 28 threads. We did not test the application with more threads since the processing time begins to significantly grow after 28 threads. By analyzing Figure 6a, it is possible to see that the optimal number of threads for the first OpenMP version is 16. These measures have also been performed for the second OpenMP version, but in this case the maximum number of threads tested was 12 since the processing times begins to grow. In this case, as highlighted by Figure 6b, the optimal number of threads is 8.

After establishing the optimal number of threads, we tested the two OpenMP versions on the entire dataset presented in Table 1. In order to allow a direct comparison between the serial and the parallel versions, we initialized the algorithm with the same values for a given image. The obtained results, together with the speed-up values, are reported in Table 2, where OpenMPv1 indicates the first version (the parallelization of the distance metric computation), while OpenMPv2 indicates the second one (the whole main loop parallelized). Those results, together with the others obtained by the CUDA and OpenCL versions, will be discussed in Section 4.4.

**Table 2.** Comparison between the serial and the OpenMP versions of the algorithm. The speed-up is reported between brackets.

Image ID	# of Iterations	Serial [s]	OpenMPv1 [s]	OpenMPv2 [s]
Image 1	32	272.20	73.18 (3.72×)	75.88 (3.59×)
Image 2	13	162.37	44.42 (3.65×)	45.80 (3.55×)
Image 3	23	289.64	77.81 (3.72×)	81.35 (3.56×)
Image 4	10	151.50	40.98 (3.70×)	43.42 (3.49×)
Image 5	13	214.52	59.00 (3.64×)	63.09 (3.40×)
Image 6	25	465.51	118.31 (3.93×)	136.99 (3.40×)

#### 4.2. CUDA Performance Evaluation

The three CUDA versions have been compiled using the NVIDIA *nvcc* compiler, which is part of the CUDA 9.0 environment. Compilation options have been chosen in order to maximize the execution speed. The tests have been conducted using two different GPUs. The first one is a NVIDIA Tesla K40 GPU equipped with 2880 CUDA cores working at 750 MHz and with 12 GB of DDR5 RAM.

It is based on the Kepler architecture that does not have a graphical output port since it is optimized for scientific computations. The second GPU is a NVIDIA GTX 1060 equipped with 1152 CUDA cores working at 1.75 GHz and with 3 GB of DDR5 RAM. This GPU is more recent than the first one and it is based on the Pascal architecture, having a graphical output port. In order to take full advantage of the specific architecture of each GPU, we indicate to the compiler which is the target micro-architecture. Specifically, we used the options *sm\_35* and *compute\_35* for the Tesla K40 GPU and the options *sm\_60* and *compute\_60* for the GTX 1060, where the values 35 and 60 represent the Kepler and the Pascal architecture, respectively. The results obtained using the Tesla K40 GPU are reported in Table 3, while the results obtained by the GTX 1060 GPU are reported in Table 4.

**Table 3.** Comparison between the serial and the CUDA versions of the algorithm on a Tesla K40 GPU. The speed-up is reported between brackets.

Image ID	# of Iterations	Serial [s]	CUDA <sub>v1</sub> [s]	CUDA <sub>v2</sub> [s]	CUDA <sub>v3</sub> [s]
Image 1	32	272.20	87.48 (3.11×)	4.52 (60.22×)	4.87 (55.89×)
Image 2	13	162.37	54.51 (2.98×)	2.97 (54.67×)	3.34 (48.61×)
Image 3	23	289.64	100.67 (2.88×)	4.99 (58.04×)	5.12 (56.57×)
Image 4	10	151.50	56.41 (2.69×)	2.96 (51.18×)	3.47 (43.66×)
Image 5	13	214.52	79.74 (2.69×)	3.99 (53.76×)	4.14 (51.82×)
Image 6	25	465.51	159.77 (2.91×)	7.45 (62.48×)	7.83 (59.45×)

**Table 4.** Comparison between the serial and the CUDA versions of the algorithm on a GTX 1060 GPU. The speed-up is reported between brackets.

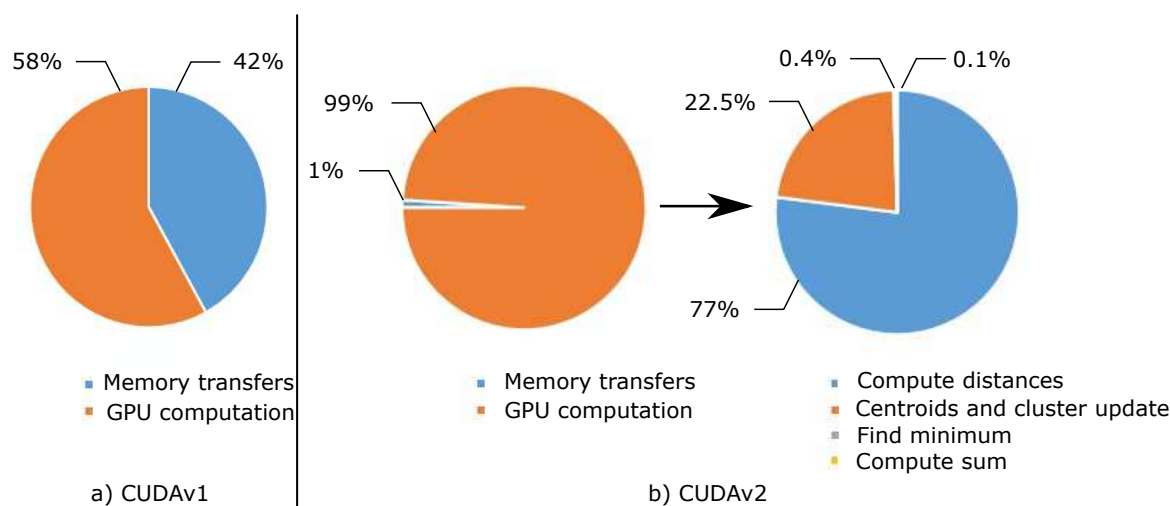
Image ID	# of Iterations	Serial [s]	CUDA <sub>v1</sub> [s]	CUDA <sub>v2</sub> [s]	CUDA <sub>v3</sub> [s]
Image 1	32	272.20	80.84 (3.37×)	2.37 (114.85×)	4.21 (64.66×)
Image 2	13	162.37	44.54 (3.65×)	1.94 (83.70×)	2.93 (55.42×)
Image 3	23	289.64	98.61 (2.94×)	2.66 (108.89×)	2.96 (97.85×)
Image 4	10	151.50	52.47 (2.89×)	2.02 (75.00×)	3.25 (46.62×)
Image 5	13	214.52	75.00 (2.86×)	2.41 (89.01×)	3.48 (61.64×)
Image 6	25	465.61	147.50 (3.16×)	3.16 (147.34×)	3.69 (126.18×)

In both tables, CUDA<sub>v1</sub> indicates the version where only the distance computation is computed on the GPU, CUDA<sub>v2</sub> indicates the version where all the operations are performed in parallel and, finally, CUDA<sub>v3</sub> indicates the version exploiting dynamic parallelism.

Concerning the GPU implementation, we also conducted a profiling using the NVIDIA Visual Profiler. This tool allows to profile the code execution on GPU together with memory transfers, in order to evaluate the efficiency of the implementation. Figure 7 shows the results obtained by profiling the Image 1 (the most demanding one) processing on the GTX 1060 with the CUDA<sub>v1</sub> (a) and CUDA<sub>v2</sub> (b) codes. Concerning the CUDA<sub>v2</sub>, in Figure 7b the different kernels executions percentage on the GPU are detailed. Profiling of CUDA<sub>v3</sub> is not shown since it is very similar to CUDA<sub>v2</sub> as well as the code profiling on the NVIDIA Tesla K40 GPU.

#### 4.3. OpenCL Performance Evaluation

OpenCL codes have been compiled using vendor-specific compilers. In particular, the OpenCL version without memory transfers have been tested on an Intel i7 6700 processor working at 3.40 GHz, equipped with 32 GB of RAM and on an Intel HD Graphics 530 integrated GPU with 16 cores working at 350 MHz. The integrated board shares the RAM with the CPU. Concerning the OpenCL version which performs memory transfers, it has been tested on the NVIDIA GTX 1060 GPU. Results obtained by the OpenCL versions are reported in Table 5.



**Figure 7.** Profiling of GPU versions on the NVIDIA GTX 1060 board for the CUDA v1 (a) and CUDA v2 (b) versions.

**Table 5.** Comparison between the serial and the OpenCL versions of the algorithm. The speed-up is reported between brackets.

Image ID	# of Iterations	Serial [s]	Intel i7 [s]	Intel HD 530 [s]	GTX 1060 [s]
Image 1	32	272.20	74.13 (3.67×)	183.96 (1.48×)	57.61 (4.72×)
Image 2	13	162.37	44.32 (3.66×)	114.24 (1.42×)	35.61 (4.56×)
Image 3	23	289.64	79.52 (3.64×)	203.62 (1.42×)	63.75 (4.54×)
Image 4	10	151.50	40.52 (3.74×)	113.82 (1.33×)	35.36 (4.28×)
Image 5	13	214.52	59.92 (3.58×)	156.56 (1.37×)	47.03 (4.56×)
Image 6	25	465.51	121.59 (3.83×)	312.68 (1.49×)	93.65 (4.97×)

#### 4.4. Comparisons and Discussion

The OpenMP version that offers the better results is the one where only the distance evaluations are processed in parallel. In this version, a parallel region is created and then destroyed at every iteration of the main loop, while in the second version the parallel region is created only once before the beginning of the main loop and is destroyed after the end of the main loop. However, the second version requires synchronization barriers between the threads, since there are operations that should be performed sequentially to obtain correct results. As an example, the increment of the number of iterations and the check of the conditions for repeating or not the main loop should be performed by a single thread. This is a critical issue since the advantage of creating and destroying the parallel region only once is thwarted by the synchronization bottleneck. The result analysis shows that the processing times are very similar, but the processor manages better the first version (OpenMPv1). Moreover, the speed-up values of the two versions are similar. Finally, it is important to highlight that the considered processor is equipped with 4 physical cores and the obtained speed-up is always greater than 3.5×. This means that the parallelization efficiency is close to the theoretical value.

Concerning the CUDA versions, by analyzing Tables 3 and 4, it is possible to observe that, for both GPU boards, the first CUDA version (CUDA v1) performs worse than the OpenMP ones. The reason is highlighted in Figure 7a, where the profiling results show that the memory transfers take about 42% of the time and only the remaining 58% is used for the computation. This is the typical bottleneck of GPU computing since the memory transfers are performed by the PCI-express external bus. The second CUDA version (CUDA v2) is not affected by this issue since the amount of data transferred at each main iteration is significantly lower than in the previous case. It is possible to parameterize the amount of transferred data at each iteration for these two versions. In the first case (CUDA v1), the data transferred is the distance\_array matrix, which is made up of  $N \times K$  elements represented in single

precision floating-point arithmetic, while in the second case (CUDA<sub>v2</sub>) only one single precision floating-point value is copied back to the host. In the third case (CUDA<sub>v3</sub>), when dynamic parallelism is used, there are no data transfers inside the main loop. Therefore the third CUDA version is the one which transfers the minimum possible amount of data, but it does not perform better than the second version. This is because the dynamic parallelism produces an overhead due to the GPU switch between the main kernel and the subroutine kernel. This overhead affects every sub-kernel activation. In this specific case, four different sub-kernels are activated at every iteration. This overhead is not negligible and, as it can be seen from the results of Tables 3 and 4, it takes longer than the copy of a single float value from device to host. In other words, the time needed by the GPU to manage the generation of four sub-kernels (CUDA<sub>v3</sub>) is comparable with the time taken by a single value copy from the host to the device memory (CUDA<sub>v2</sub>). Finally, for all the CUDA versions, the GTX 1060 board performs better than the Tesla K40, even if the last board is optimized for scientific computations. This is because the first board is equipped with a more recent architecture which has better CUDA cores working at a higher frequency than the Tesla GPU.

The analysis of the OpenCL versions highlight that, considering the Intel i7, the processing times are close to the OpenMP ones. For what concerns the Intel HD 530 integrated GPU, the performance is very poor, and the speed-ups are negligible. This is probably due to the low-end integrated GPU with a working frequency of 350 MHz and only 16 parallel processing elements. Therefore, it is not possible to obtain a significant speed-up compared to the serial version. Comparing the OpenCL version and the CUDA versions running on the GTX 1060 GPU, it is possible to notice that the OpenCL version performs better than the CUDA<sub>v1</sub>, but is significantly slower than the other two CUDA versions. These CUDA versions (CUDA<sub>v2</sub> and CUDA<sub>v3</sub>) employ highly optimized routines, which exploits all the hardware features of a GPU. Moreover, the CUDA versions have been compiled using compilation options in order to produce an executable code which fully exploits the specific target architecture. This is not possible in OpenCL since it targets portability between different devices as main feature.

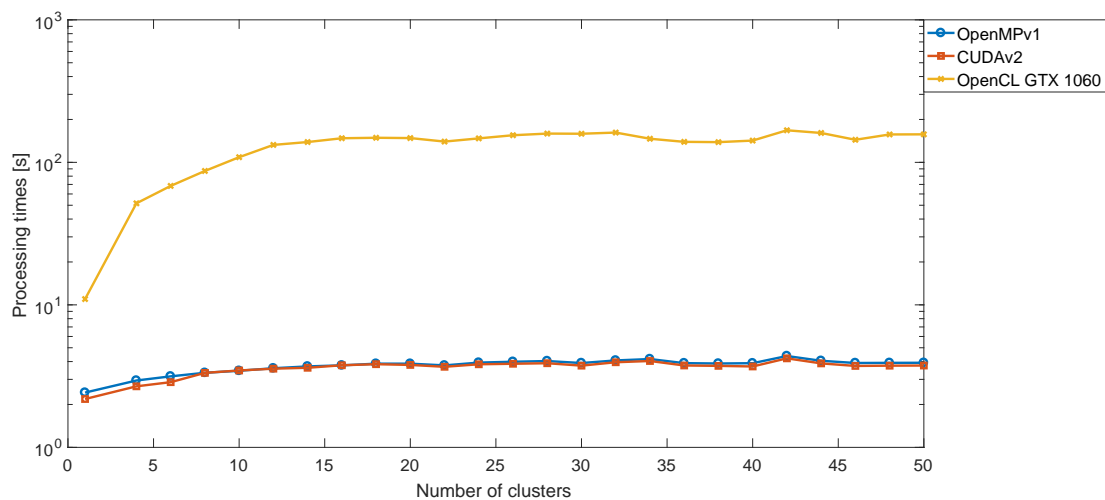
We also performed a comparative study between the three best performing versions of the three considered technologies in order to characterize how the speed-up varies with respect to the number of clusters. In particular, we performed experiments using Image 6 and  $K$  values varying from 2 to 50. The speed-ups of the OpenMP, CUDA and OpenCL best versions with respect to the serial implementation are shown in Figure 8 using a semi-logarithmic scale. It is possible to see that the CUDA version has a speed-up that ranges from  $10\times$  to  $\sim 150\times$  and from 12 clusters on it becomes nearly constant. On the other hand, the solutions based on a multi-core processor have speed-ups that are close to  $4\times$ .

In the literature, there are different works about parallel K-means.

Baramkar et al. [18] performed a review of different parallel GPU-based K-means, but the considered works were focused only on general classification, without considering high data dimensionality, which is the case explored in our work. Therefore it is hard to perform direct comparisons with these works, which achieve very different speed-ups ranging from  $11\times$  to  $220\times$ .

Zechner et al. [19] proposed a parallel implementation of this algorithm using both CPU and GPU. In particular, the GPU was only employed for distance computation, while centroids update was left to the CPU. They classified an artificial dataset with two-dimensional elements ranging from 500 to 500,000. The maximum speed-up achieved was  $14\times$ , lower than the one obtained in our work. This is because the optimization proposed in [19] is only valid for low-dimensional data and cannot be employed for classifying high-dimensional data such as hyperspectral images.

A similar approach is shown in [20,21], with the difference that also the clusters update has been performed on the GPU. However, between the distance computation and centroids update they performed host computation for updating each pixel label. This choice leads to a maximum speed-up of 60 in both works, lower than our one, since we moved all the computation on the device side.



**Figure 8.** Speed-ups achieved by the three best parallel implementations (one for each evaluated parallel technology) with respect to the number of clusters.

In [22], a GPU-based K-means algorithm is proposed, with a distance computation that is evaluated through a simple Cartesian distance. Under this assumption, they classify 1,000,000 pixels with 32 features in 1.15 s. In our case, the bigger image has 264,408 pixels, the features (i.e., the bands) are 128 and it is processed in ~3.56 s. Moreover, the distance metric that we adopt (the spectral angle) is more complex than the one proposed in [22].

Baydoun et al. [23] developed a parallel K-means for RGB images classification. They adopted as metric a simple Cartesian distance and they parallelize only this computation, achieving a maximum speed-up of ~25 $\times$ . In this case, the metric and the data dimensionality are very different compared to this work.

In [24], the K-means algorithm was modified to further reduce the distance computation. The speed-up varies from 4 $\times$  to 386 $\times$ , but also, in this case, it is not possible to perform a direct comparison since there are not enough details about the dataset composition. Finally, in [25], the K-means algorithm was developed on GPU with the Cartesian distance. They adopt a modern GPU with 1536 CUDA cores obtaining a maximum speed-up of 88 $\times$ , which is very similar to our results.

Lutz et al. [26] proposed a parallel K-means implementation using an NVIDIA GTX 1080 GPU. They performed only experiments producing four groups and no further details are given in the paper about the dataset. They achieved a maximum speed-up of 18.5 $\times$  which is nearly an order of magnitude smaller than one of our implementations.

A comparative analysis similar to the one we conducted is reported in [27]. Authors exploited GPUs, OpenMP, Message Passing Interface (MPI) and FPGAs. However, also in this case, they considered only a dataset made up of 10-dimensional points, therefore the computational complexity of the distance computation is lower than our one. On the other hand, the results were not as good as our ones, since the maximum GPU speed-up achieved is ~60 $\times$ . They also demonstrated that the speed-up could reach a value up to 200 $\times$  if the number of clusters to produce was significantly increased (i.e., more than 2000 clusters), but a study of how this speed-up varied also with respect to the data dimensionality were not carried out. Concerning OpenMP, the classification of 20,000 10-dimensionality points took ~3 s. Our smallest image is 6 times bigger than this one and with a dimensionality 18 times greater than the one considered. Keeping this in mind, the performance of our best OpenMP version is quite similar to this one. Finally, concerning the FPGA implementation, experiments were reported only with a 17,692 9-dimensionality dataset. Classification time is ~100 ms, but, as stated in [27], the FPGA resources, especially memory banks, were not enough to process bigger datasets. The authors of this work do not use external DDR memory, therefore, the FPGA performance is limited due to this design choice.



The comparison between our work and the literature is summarized in Table 6.

**Table 6.** Comparison between the proposed work and the literature.

Paper	Maximum Image Size	Data Dimensionality	Technology	Speed-Up
[18]	2,000,000	8	GPU NVIDIA GTX 280	220
[19]	500,000	2	GPU NVIDIA 9600 GT	14
[20]	1,000,000	2	GPU NVIDIA 8800 GTX	60
[21]	15,052,800	3	4 × GPU NVIDIA GTX 750Ti	60
[22]	1,000,000	32	GPU NVIDIA GTX 280	N. A.
[23]	16,777,216	3	GPU NVIDIA Tesla C2050	25
[24]	245,057	4	GPU NVIDIA GeForce 210	386
[25]	500,000	16	GPU NVIDIA Quadro K5000	88
[26]	N. A.	N. A.	GPU NVIDIA GTX 1080	18.5
[27]	20,000	10	2 × AMD Opteron quad-core	8
[27]	65,536	10	GPU NVIDIA Tesla 2050	60
[27]	17,692	9	Mittrion MVP FPGA Simulator	N. A.
Our work	264,408	128	GPU NVIDIA GTX 1060	126

## 5. Conclusions

In this paper, we presented different parallel implementations of the K-means algorithm for hyperspectral medical image clustering. In particular, we evaluated multi-core CPUs and manycore GPUs through the OpenMP and CUDA frameworks, respectively. Moreover, we also addressed the problem of code portability by developing OpenCL-based versions. We performed experiments with a dataset made up of *in-vivo* hyperspectral human brain images. Those experiments validated the results of all the proposed parallel implementations. Among them, CUDA achieved the better performance, outperforming OpenMP implementations. The cost of the better performance is the parallelization effort, which is significantly greater when working with CUDA. In fact, the development of the CUDA versions required the development of custom kernels and ad-hoc memory transfer management, while OpenMP only required code annotations with suitable pragmas. Code portability has also been addressed with OpenCL. However, this technology is not yet competitive with OpenMP or CUDA, achieving the worst results among the developed parallel applications. Moreover, OpenCL guarantees portability among different devices, but, for obtaining the best performance from a given device, it is necessary to tune the code with respect to specific hardware features. The comparison of the proposed implementations shows that the best one is based on CUDA and executed on the GTX 1060 board, achieving a maximum speed-up of  $\sim 125\times$ . In particular, the best CUDA version performs all the computations on the GPU without exploiting dynamic parallelism.

We also made comparisons with other recent works in the literature, that only in one case achieved results comparable but not better than ours, except for the FPGA solution proposed in [27]. However, FPGA memory constraint does not allow to process images with more than 17,692 pixels. This limits the use of this technology and, in particular, it makes FPGAs not suitable for our target application.

Summarizing, the proposed work confirms that the GPU technology is the best solution for these class of problems, even when considering a data dimensionality bigger than the ones considered before. It also highlights that the GPU algorithm has a good scalability with respect to the number of clusters ( $K$ ). Moreover, when considering high data dimensionality, the parallelization of the distance computation is not enough, since also the centroids update, and the error computation can be parallelized. This ensures a supplementary speed-up. Finally, the technological evolution of GPUs offers increasing computing power at relatively low cost. In our case, a consumer GPU sold at about \$200 outperforms a more expensive Tesla K40 GPU ( $\sim$ \$5000) of a previous generation, but optimized for scientific computations.

In conclusion, this work provides an efficient GPU implementation of the K-means algorithm, that will be included in the parallel version of the complete system already shown in Figure 1. Since the K-means is one of the most computational demanding algorithms in the system, this remarkable result is essential to satisfy the real-time constraint.

Future research will be focused on integrating this parallel algorithm in more complicated classification frameworks, such as the one proposed in [5,7].

**Author Contributions:** E.T. performed the GPU implementations, the algorithms optimizations, designed and performed experiments and wrote the manuscript. G.F. designed experiments and edited the manuscript. F.C. performed the GPU implementations and experiments. S.O., H.F. performed the serial algorithm implementation and edited the manuscript. G.M.C., M.M.-M., F.L. supervised the project and edited the manuscript.

**Funding:** This work has been supported in part by the Canary Islands Government through the ACIISI (Canarian Agency for Research, Innovation and the Information Society), ITHACA project “Hyperspectral Identification of Brain Tumors” under Grant Agreement ProID2017010164 and it has been partially supported also by the Spanish Government and European Union (FEDER funds) as part of support program in the context of Distributed HW/SW Platform for Intelligent Processing of Heterogeneous Sensor Data in Large Open Areas Surveillance Applications (PLATINO) project, under contract TEC2017-86722-C4-1-R. Additionally, this work was completed while Samuel Ortega was beneficiary of a pre-doctoral grant given by the “Agencia Canaria de Investigación, Innovación y Sociedad de la Información (ACIISI)” of the “Conserjería de Economía, Industria, Comercio y Conocimiento” of the “Gobierno de Canarias”, which is part-financed by the European Social Fund (FSE) (POC 2014–2020, Eje 3 Tema Prioritario 74(85%)). Finally, this work has been also supported in part by the 2016 PhD Training Program for Research Staff of the University of Las Palmas de Gran Canaria.

**Acknowledgments:** The authors would like to thank NVIDIA Corporation for the donation of the NVIDIA Tesla K40 GPU used for this research.

**Conflicts of Interest:** The authors declare no conflict of interest.

## References

1. Ferlay, J.; Soerjomataram, I.; Ervik, M.; Dikshit, R.; Eser, S.; Mathers, C.; Rebelo, M.; Parkin, D.M.; Forman, D.; Bray, F. *Cancer Incidence and Mortality Worldwide: IARC CancerBase No. 11*; International Agency for Research on Cancer: Lyon, France, 2013.
2. Louis, D.N.; Perry, A.; Reifenberger, G.; von Deimling, A.; Figarella-Branger, D.C.; Webster, K.; Ohgaki, H.; Wiestler, O.D.; Kleihues, P.; Ellison, D.W. The 2016 World Health Organization Classification of Tumors of the Central Nervous System: A summary. *Acta Neuropathol.* **2016**, *131*, 803–820. [[CrossRef](#)] [[PubMed](#)]
3. Sanai, M.; Berger, M.S. Operative techniques for gliomas and the value of extent of resection. *Neurotherapeutics* **2009**, *6*, 478–486. [[CrossRef](#)] [[PubMed](#)]
4. Fabelo, H.; Ortega, S.; Kabwama, S.; Callicó, G.M.; Bulters, D.; Szolna, A.; Pineiro, J.F.; Sarmiento, R. HELICoiD project: A new use of hyperspectral imaging for brain cancer detection in real-time during neurosurgical operations. *Proc. SPIE Int. Soc. Opt. Eng.* **2016**, *12*, 9860. [[CrossRef](#)]
5. Fabelo, H.; Ortega, S.; Lazcano, R.; Madronal, D.; Callicó, G.M.; Juárez, E.; Salvador, R.; Bulters, D.; Bulstrode, H.; Szolna, A.; et al. An intraoperative visualization system using hyperspectral imaging to aid in brain tumor delineation. *Sensors* **2018**, *18*, 430. [[CrossRef](#)] [[PubMed](#)]
6. Chang, C.-I. *Hyperspectral Data Processing: Algorithm Design and Analysis*; John Wiley & Sons: Hoboken, NJ, USA, 2013; ISBN 978-0-471-69056-6.
7. Fabelo, H.; Ortega, S.; Ravi, D.; Kiran, B.R.; Sosa, C.; Bulters, D.; Callicó, G.M.; Bulstrode, H.; Szolna, A.; Pineiro, J.F.; Kabwama, S.; et al. Spatio-spectral classification of hyperspectral images for brain cancer detection during surgical operations. *PLoS ONE* **2018**, *13*, e0193721. [[CrossRef](#)] [[PubMed](#)]
8. Torti, E.; Fontanella, A.; Florimbi, G.; Leporati, F.; Fabelo, H.; Ortega, S.; Callicó, G.M. Acceleration of brain cancer detection algorithms during surgery procedures using GPUs. *Microprocess. Microsyst.* **2018**, *61*, 171–178. [[CrossRef](#)]
9. Florimbi, G.; Fabelo, H.; Torti, E.; Lazcano, R.; Madronal, D.; Ortega, S.; Salvador, R.; Leporati, F.; Danese, G.; Báez-Quevedo, A.; et al. Accelerating the K-Nearest Neighbors Filtering Algorithm to Optimize the Real-Time Classification of Human Brain Tumor in Hyperspectral Images. *Sensors* **2018**, *18*, 2314. [[CrossRef](#)] [[PubMed](#)]
10. Lazcano, R.; Madronal, D.; Salvador, R.; Desnos, K.; Pelcat, M.; Guerra, R.; Fabelo, H.; Ortega, S.; Lopez, S.; Callico, G.M.; et al. Porting a PCA-based hyperspectral image dimensionality reduction algorithm for brain cancer detection on a manycore architecture. *J. Syst. Archit.* **2017**, *77*, 101–111. [[CrossRef](#)]
11. Madronal, D.; Lazcano, R.; Salvador, R.; Fabelo, H.; Ortega, S.; Callico, G.M.; Juarez, E.; Sanz, C. SVM-based real-time hyperspectral image classifier on a manycore architecture. *J. Syst. Archit.* **2017**, *80*, 30–40. [[CrossRef](#)]

12. Domingo, R.; Salvador, R.; Fabelo, H.; Madronal, D.; Ortega, S.; Lazcano, R.; Juarez, E.; Callico, G.M.; Sanz, C. High-level design using Intel FPGA OpenCL: A hyperspectral imaging spatial-spectral classifier. In Proceedings of the 2017 12th International Symposium on Reconfigurable Communication-centric Systems-on-Chip (ReCoSoC), Madrid, Spain, 12–14 July 2017; pp. 1–8. [[CrossRef](#)]
13. Fontanella, A.; Marenzi, E.; Torti, E.; Danese, G.; Plaza, A.; Leporati, F. A suite of parallel algorithms for efficient band selection from hyperspectral images. *J. Real-Time Image Process.* **2018**, 1–17. [[CrossRef](#)]
14. Marenzi, E.; Carrus, A.; Danese, G.; Leporati, F.; Callicó, G.M. Efficient Parallelization of Motion Estimation for Super-Resolution. In Proceedings of the 2017 25th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP), St. Petersburg, Russia, 6–8 March 2017; pp. 274–277. [[CrossRef](#)]
15. Lopez-Fandino, J.; Heras, D.B.; Arguello, F.; Dalla Mura, M. GPU Framework for Change Detection in Multitemporal Hyperspectral Images. *Int. J. Parallel Program.* **2017**, 1–21. [[CrossRef](#)]
16. Florimbi, G.; Torti, E.; Danese, G.; Leporati, F. High Performant Simulations of Cerebellar Golgi Cells Activity. In Proceedings of the 2017 25th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP), St. Petersburg, Russia, 6–8 March 2017; pp. 527–534. [[CrossRef](#)]
17. Feng, X.; Jin, H.; Zheng, R.; Zhu, L.; Dai, W. Accelerating Smith-Waterman Alignment of Species-Based Protein Sequences on GPU. *Int. J. Parallel Program.* **2015**, *43*, 359–380. [[CrossRef](#)]
18. Baramkar, P.P.; Kulkarni, D.B. Review for K-Means On Graphics Processing Units (GPU). *Int. J. Eng. Res. Technol.* **2014**, *3*, 1911–1914.
19. Zechner, M.; Granitzer, M. K-Means on the Graphics Processor: Design and Experimental Analysis. *Int. J. Adv. Syst. Meas.* **2009**, *2*, 224–235. [[CrossRef](#)]
20. Hong-tao, B.; Li-li, H.; Dan-tong, O.; Zhan-shan, L.; He, L. K-Means on Commodity GPUs with CUDA. In Proceedings of the WRI World Congress on Computer Science and Information Engineering, Los Angeles, CA, USA, 31 March–2 April 2009; pp. 651–655. [[CrossRef](#)]
21. Fakhi, H.; Bouattane, O.; Youssfi, M.; Hassan, O. New optimized GPU version of the k-means algorithm for large-sized image segmentation. In Proceedings of the Intelligent Systems and Computer Vision, Fez, Morocco, 17–19 April 2017; pp. 1–6. [[CrossRef](#)]
22. Li, Y.; Zhao, K.; Chu, X.; Liu, J. Speeding up k-Means algorithm by GPUs. *J. Comput. Syst. Sci.* **2013**, *79*, 216–229. [[CrossRef](#)]
23. Baydoun, M.; Dawi, M.; Ghaziri, H. Enhanced parallel implementation of the K-Means clustering algorithm. In Proceedings of the 3rd International Conference on Advances in Computational Tools for Engineering Applications (ACTEA), Beirut, Lebanon, 13–15 July 2016; pp. 7–11. [[CrossRef](#)]
24. Saveetha, V.; Sophia, S. Optimal Tabu K-Means Clustering Using Massively Parallel Architecture. *J. Circuits Syst. Comput.* **2018**, In press. [[CrossRef](#)]
25. Cuomo, S.; De Angelis, V.; Farina, G.; Marcellino, L.; Toraldo, G. A GPU-accelerated parallel K-means algorithm. *Comput. Electr. Eng.* **2017**, 1–13. [[CrossRef](#)]
26. Lutz, C.; Bress, S.; Rabl, T.; Zeuch, S.; Markl, V. Efficient k-means on GPUs. In Proceedings of the 14th International Workshop on Data Management on New Hardware, Huston, ID, USA, 11 June 2018. [[CrossRef](#)]
27. Yang, L.; Chiu, S.C.; Liao, W.K.; Thomas, M.A. High Performance Data Clustering: A Comparative Analysis of Performance for GPU, RASC, MPI, and OpenMP Implementations. *J. Supercomput.* **2014**, *70*, 284–300. [[CrossRef](#)] [[PubMed](#)]

