

GPU implementation of a region based algorithm for large images segmentation

Gilles Perrot¹, Stéphane Domas¹, Raphaël Couturier¹, Nicolas Bertaux²

¹*Distributed Numerical Algorithmics team (AND), Laboratoire d'Informatique de Franche-comté
Rue Engel Gros, 90000 Belfort, France.
forename.name@univ-fcomte.fr*

²*Institut Fresnel, CNRS, Aix-Marseille Université, Ecole Centrale Marseille,
Campus de Saint-Jérôme, 13013 Marseille, France.
nicolas.bertaux@ec-marseille.fr*

Abstract—Image segmentation is one of the most challenging issues in image computing. In this work, we focus on region-based active contour techniques (snakes) as they seem to achieve a high level of robustness and fit with a large range of applications. Some algorithmic optimizations provide significant speedups, but even so, execution times are still non-neglectable with the continuing increase of image sizes. Moreover, these algorithms are not well suited for running on multi-core CPU's. At the same time, recent developments of Graphical Processing Units (GPU) suggest that higher speedups could be obtained by use of their specific design. We have managed to adapt a specially efficient snake algorithm that fits recent Nvidia GPU architecture and takes advantage of its massive multithreaded execution capabilities. The speedup obtained is most often around 7.

Keywords-GPU; segmentation; snake;

I. INTRODUCTION

Segmentation and shape detection are still key issues in image computing. These techniques are used in numerous fields ranging from medical imaging to video tracking, shape recognition or localization. Since 1988, the active contours (snakes) introduced by Kass et al. [1], have proved to be efficient and robust, especially against noise, for a wide range of image types.

The main shortcoming of these algorithms is often their high dependence on the initial contour, though several contributions have lowered this dependency and also brought more accurate segmentation of non convex shape [2].

The information that drives a contour model comes either from the contour itself or from the characteristics of the regions it defines. For noisy images, the second option is often more suitable as it takes into account the statistical fluctuations of the pixels. One approach [2] proposes a geometric (polygonal) region-based snake driven by the minimization of the likelihood (ML).

An important issue of image processing, especially segmentation, has always been the computation time of most algorithms. Over the years, the increase of CPU computing capabilities, although quite impressive, has not been able to fulfill the combined needs of growing resolution and

real-time computation. Since having been introduced in the early 1980's, the capabilities and speed of graphics accelerators have always been increasing. So much so that the recent GPGPU (General Purpose Graphic Processing Units) currently benefit by a massively parallel architecture for general purpose programming, especially when dealing with large matrices or vectors. On the other hand, their specific design obviously imposes a number of limitations and constraints.

A. Related work

Since the main issue with most of the segmentation methods is the computational effort it implies, researchers have recently tried to benefit from the GPGPU architecture to reduce time costs. In [3] authors achieve impressive speed-ups on a range of contour detection algorithms on the condition that they are applied to images with good contrast and SNR (Signal-to-Noise Ratio). Others, like [4], have focused on the related issues of segmentation and tracking: this proves efficient in processing low-contrast images, but imposes limits as to the size of images that cannot be of big dimensions (several million pixels). One third option, parametric snakes, has also been investigated with some success, as in [5], although the principle of computation per small tile is not suited to the algorithm we have implemented. In the medical imaging field, some researchers have implemented efficient parallel segmentation algorithms. For example in [6], authors implement a GPU watershed-based algorithm, but source images are considered with a good contrast.

B. Contribution

The geometric snake we have focused on has proved to be efficient processing real-life images, with poor SNR and contrast. However, the required computational effort still imposes limits to the size of a wide range of images to be processed, such as the output of SAR (Synthetic Aperture Radar). Our goal was then to propose a way to fit such a region-based snake algorithm to the Nvidia Tesla[©] GPU

architecture. The remainder of this paper exposes the principles of the algorithm and notations in section II. Section III deals with the details of sequential CPU implementation. Section IV summarizes the main features of the Nvidia[©] GPU and explains how to deal with them efficiently. Then sections V and VI detail our GPU implementation and timing results. In our conclusion VII we attempt to evaluate the pros and cons of this implementation, and suggest further tracks to be investigated in future research.

II. SEQUENTIAL ALGORITHM: OUTLINES

The goal of the active contour segmentation method (snake) we studied [2] is to distinguish, inside an image I , a target region T from the background region B . The size of I is $L \times H$ pixels of coordinates (i, j) and gray level $z(i, j)$. Z represents the gray levels data of I . We assume that the gray levels of T and B are vectors of independent and identically distributed values, each with a probability density function (PDF) p^Ω ($\Omega \in \{T; B\}$). The present implementation uses a Gaussian PDF, but another one can easily be used as Gamma or Poisson (Cf. [2]).

The *active contour* S , which defines the shape of T is chosen as polygonal. The purpose of the segmentation is then to determine the shape that optimizes a generalized log-likelihood-based criterion (GL). This is done by an iterative process which is initialized with an arbitrary shape, then at each step:

- 1) it modifies the shape
- 2) it estimates the parameters of the Gaussian functions for the two regions and evaluates the criterion.
- 3) it validates the new shape if the criterion has a better value.

A simplified description of it is given in *Algorithm 1* which features two nested loops: the main one, on iteration level, is responsible for tuning the number of nodes; the inner one, on step level, takes care of finding the best shape for a given number of nodes. *Figure 1* shows intermediate results at iteration level. Sub-figure *1a* shows the initial rectangular shape, *1b* shows the best four-node shape that ends the first iteration. Sub-figures *1c* and *1d* show the best shape for an eight-node contour (resp. 29-node) which occurs at the end of the second iteration (resp. fourth).

III. SEQUENTIAL ALGORITHM: DETAILS

A. Criterion

Let p^Ω be a Gaussian PDF. Its vector of parameters Θ_Ω ($\Omega \in \{T; B\}$) has two components, the average value μ and the standard deviation σ . The likelihood for the regions Ω ($\Omega \in \{T; B\}$) is given by

$$P[Z|T, B, \Theta_T, \Theta_B] = P(Z|T, \Theta_T)P(Z|B, \Theta_B)$$

where

$$P(Z|\Omega, \Theta_\Omega) = \prod_{(i,j) \in \Omega} p^\Omega[z(i, j), \Theta_\Omega] \quad (\Omega \in \{T; B\})$$

Algorithm 1: Sequential algorithm: outlines

```

1: begin with a rectangular 4 nodes contour;
2: repeat /* iteration level */
3:   repeat /* step level */
4:     Test some other positions for each node, near
       its current position;
5:     Find the best GL and adjust the node's position;
6:   until no more node can be moved;
7:   Add a node in the middle of each long enough
       segment;
8: until no more node can be added;

```

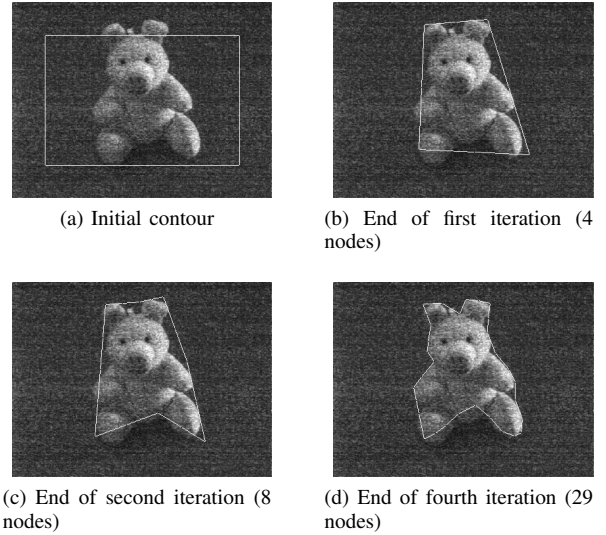


Figure 1. segmentation of a noisy image

The log-likelihood of the region Ω is then

$$-N_\Omega \log(\sqrt{2\pi}) - N_\Omega \cdot \log(\sigma) - \frac{1}{2\sigma^2} \sum_{(i,j) \in \Omega} (z(i, j) - \mu)^2$$

inside which the vectors of parameters Θ_Ω are determined by ML estimation

$$\widehat{\Theta}_\Omega \begin{cases} \widehat{\mu}_\Omega = \frac{1}{N_\Omega} \sum_{(i,j) \in \Omega} z(i, j) \\ \widehat{\sigma}_\Omega^2 = \frac{1}{N_\Omega} \sum_{(i,j) \in \Omega} (z(i, j) - \widehat{\mu}_\Omega)^2 \end{cases}$$

Considering the two regions, the criterion to be optimized is then, up to a constant, the Generalized Likelihood (GL):

$$GL = \frac{1}{2} \left(N_B \log(\widehat{\sigma}_B^2) + N_T \log(\widehat{\sigma}_T^2) \right)$$

B. CPU implementation

Let $S_{n,l}$ be the polygonal contour state at step l of iteration n , and $S_{n,l}^i$ the node i of $S_{n,l}$ ($i \in [0; N_n]$).

$S_{n,l}^{i,w}$ is the neighbor of index w of the node $S_{n,l}^i$ in a 8-connectivity meaning with d pixels scope. Each segment of $S_{n,l}$ is considered as an oriented list of discrete points. Chesnaud & Réfrégier, based on the Green-Ostogradski theorem, have shown how to replace the 2 dimensions (2D) sums needed to estimate Θ_Ω by 1 dimension sums along $S_{n,l}$ [2]. This approach leads to compute a pair of transformed images, called cumulated images, at the very beginning of the process, which are then used as lookup tables. It also involves weighting coefficients for pixels and segments of the contour. See [2] for details. Therefore, beyond this point, we will talk about the *contribution* of each point to the 1D sums. By extension, we also talk about the *contribution* of each segment to the 1D sums.

A more detailed description of the sequential algorithm is given by *Algorithm 2*. The process starts with the computation of cumulated images; an initialization stage takes place from line 3 to line 9. Then we recognize the two nested loops (line 10 and line 11) and finally the heart of the algorithm stands on line 15 which represents the main part of the calculations to be done:

- 1) compute the various sums without the contributions of both segments connected to current node $S_{n,l}^i$.
- 2) compute the contributions of both segments, which requires:
 - To determine the coordinates of every discrete pixel of both segments connected to $S_{n,l}^i$.
 - To compute every pixel contribution.
 - To sum pixel contributions to obtain segment contributions.
- 3) compute the GL given the contribution of each segment of the tested contour.

The profiling results of the CPU implementation shown in *Figure 2* display the relative costs of the most time-consuming functions. It appears that more than 80% of the total execution time is always spent by only three functions:

- `compute_segment_contribution()` which is responsible for point 2 above,
- `compute_cumulated_images()` which computes the 3 lookup tables at the very beginning,
- `compute_pixels_coordinate()` which is called by `compute_segment_contribution()`.

Measurements have been performed for several image sizes from 15 MPixels (about 3900 x 3900) to 144 MPixels (about 12000 x 12000). On the one hand, we can notice that function `compute_segment_contribution()` always lasts more than 45% of the total running time, and even more when the image gets larger. On the other hand, the function `compute_cumulated_images()` costs more than 23%, decreasing with image size, while function `compute_pixels_coordinate()` always takes around 6%. It confirms that the need for parallelization

Algorithm 2: Sequential simplified algorithm

```

1: read image from source (Hard Disk Drive);
2: compute_cumulated_images();
3: iteration  $n \leftarrow 0$ ;
4:  $N_0 \leftarrow 4$ ;
5:  $S_{n,l} \leftarrow S_{0,0}$ ;
6: step  $d \leftarrow d_{max}$  an arbitrary power of 2 value;
7: current node  $S_{0,0}^i \leftarrow S_{0,0}^0$ ;
8:  $l \leftarrow 0$ ;
9: compute  $GL_{ref}$ , the GL of  $S_{n,0}$ ;
10: repeat /* iteration level, n index */
11:   repeat /* step level, l index */
12:     for  $i = 0$  to  $N_n$  do
13:        $S_{n,l}^{i,w}$  ( $w \in [0;7]$ ) are the neighbors of  $S_{n,l}^i$ 
         by  $d$  pixels;
14:       for  $w = 0$  to 7 do
15:         compute  $GL_w$  for  $S_{n,l}$  when  $S_{n,l}^{i,w}$ 
           replaces  $S_{n,l}^i$ ;
16:         if  $GL_w$  is better than  $GL_{ref}$  then
17:            $GL_{ref} \leftarrow GL_w$ ;
18:           move node  $S_{n,l}^i \leftarrow S_{n,l}^{i,w}$ ;
19:         end
20:       end
21:     end
22:      $l \leftarrow l + 1$ ;
23:   until no node move occured;
24:   add new nodes,  $N_n \leftarrow N_n + N_{newnodes}$ ;
25:   if  $d > 1$  then  $d \leftarrow d/2$  else  $d = 1$ ;
26:    $n \leftarrow n + 1$ ;
27:   compute  $GL_{ref}$ , the GL of  $S_{n,0}$ ;
28: until no new node added;
```

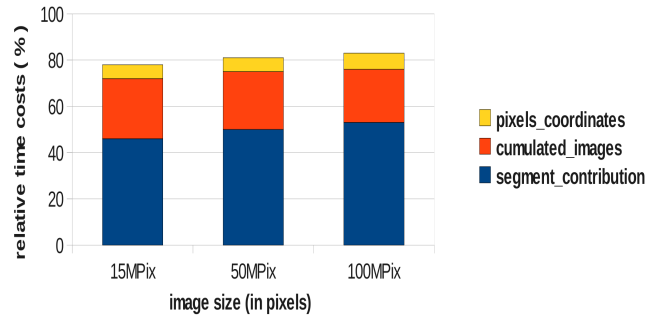


Figure 2. the three most-consuming functions for various image sizes

resides in line 15 and line 2 of Algorithm 2 as they contain every call to those three functions.

The following sections detail how we managed to implement these time-consuming functions in parallel, but a brief reminder on GPU's recent architecture is presented first.

IV. NVIDIA'S GPU ARCHITECTURE

GPUs are multi-core, multi-threaded processors, optimized for highly parallel computation. Their design focuses on a Single Instruction Multiple Threads (SIMT) model by devoting more transistors to data processing rather than data-caching and flow control [7]. For example, a C2050 card features 3GB of global memory and a total of 448 cores bundled in several Streaming Multiprocessors (SM). An amount of shared memory, much faster than the global memory, is available on each SM (from 16 KB to 487 KB)

Writing efficient code for such architectures is not obvious, as re-serialization must be avoided as much as possible. Thus, when designing, one must keep a few key points in mind:

- the CUDA model organizes threads by a) threads blocks in which synchronization is possible, b) a grid of blocks with no possible synchronization between blocks.
- there is no way to know in what order the blocks are to be scheduled during one single kernel execution.
- data must be kept in GPU memory, to reduce the overhead due to copying between CPU and GPU memories.
- the total amount of threads running the same computation must be maximized.
- the number of execution branches inside a block should be reduced as much as possible.
- global memory accesses should be coalescent, *ie.* memory accesses done by physically parallel threads (16 at a time) must be consecutive and contained in a 128 Bytes range.
- shared memory is organized by 16 x 32 bits wide banks. To avoid bank conflicts, each parallel thread (16 at a time) must access a different bank.

All the above characteristics make it always a quite constrained problem to solve when designing a GPU code. Moreover, a non suited code would probably run even slower on GPU than on CPU due to the automatic serialization which would be done at run time.

V. GPU IMPLEMENTATION

In the implementation described below, pre-computations and proper segmentation are discussed separately. To keep data in GPU memory, the whole computation is assigned to the GPU. CPU still hosts:

- data reading from HDD
- data writing on HDD if needed
- main loops control (corresponding to lines 10 and 11 of Algorithm 2)

It must be noticed that controlling these loops is achieved with only a very small amount of data being transferred between host (CPU) and device (GPU), which does not produce high overhead.

Moreover, the structures described below need 20 Bytes per pixel of the image to process (plus an offset of about

50 MByte). It defines the maximum image size we can accept: approximately 150 M Pixels.

A. Pre-computations

To replace 2D sums by 1D sums, Chesnaud *et al.* [2] have shown that the three matrices below should be computed:

$$C_1(i, j) = \sum_{k=0}^{k=j} (1 + k)$$

$$C_z(i, j) = \sum_{k=0}^{k=j} z(i, k)$$

and

$$C_{z^2}(i, j) = \sum_{k=0}^{k=j} z^2(i, k)$$

Where $z(i, k)$ is the gray level of pixel of coordinate (i, j) , so that C_1 , C_z and C_{z^2} are the same size as image I .

First, we chose not to generate $C_1(i, j)$, which requires that values should be computed when needed, but saves global memory and does not lead to any overhead. The computation of C_z and C_{z^2} easily decomposes into series of *inclusive prefixsums* [8]. However, by keeping the *1 thread per pixel* rule, as the total number of threads that can be run in a grid cannot exceed 2^{25} (Cf. [7]), slicing is necessary for images exceeding a size threshold which can vary according to the GPU model (e.g. 33 MPix for sm13 GPU, eg. C1060). It's quite easy to do, but it leads to a small overhead as the process requires multiple calls to one kernel. Slicing can be done in two ways:

- all slices are of the same size (balanced)
- slices fit the maximum size allowed by the GPU, leaving one smaller slice at the end of the process (full-sized).

The balanced slice option has proved to run faster.

For example: if a given image has 9000 lines and the GPU can process up to 4000 lines at a time, it's faster to run 3 times with 3000 lines rather than twice with 4000 and once with 1000.

As the sums in C_z and C_{z^2} are row-wide, it is easy to see that every block-wide sum will be needed before being able to use it in the global sum. But as mentioned earlier, the scheduling of blocks must be considered as random. So, in order to ensure synchronizations, each row of the original image is then treated by three different kernels:

- `compute_blocks_prefixes()`.
- `scan_blocksums()`.
- `add_sums2prefixes()`.

Figures 3, 4 and 5 show relevant data structures for a given row i of I . We assume that each thread block runs bs threads in parallel and each row of C_z needs n blocks to cover its L pixels.

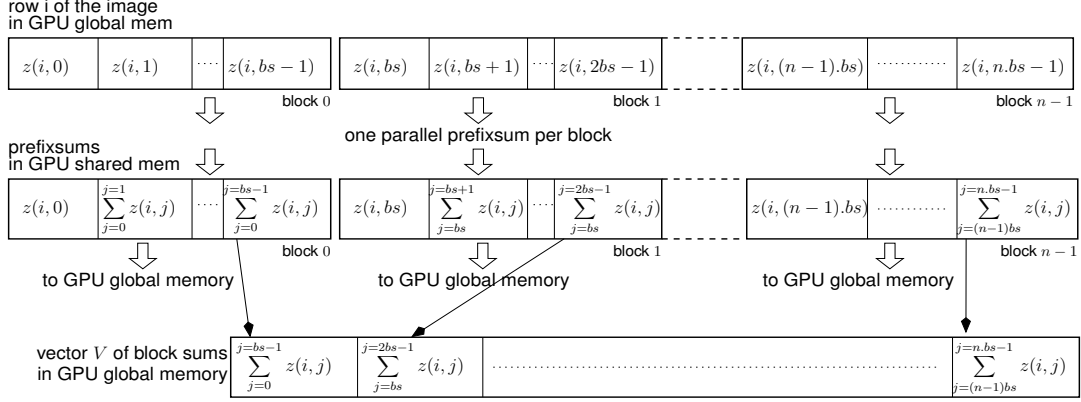


Figure 3. `compute_blocks_prefixes()` details.

Figure 3 shows the details of the process for row i of the original image I , already stored in GPU global memory. Operands are first copied into GPU shared memory for efficiency reasons. An inclusive prefixsum is then performed inside each independent thread block. At this point, only the first shared memory block contains the final values. Its last element contains the sum of all elements in the corresponding block of I . In order to obtain the right values for the row i of C_z , every element value in the other blocks must then be summed with an offset value. This offset value is the sum of all element values in every corresponding previous block of row i .

As the scheduling of blocks is fully unpredictable, the necessary intermediate results have to be stored in GPU global memory before exiting from kernel. Each element of the prefixsums in GPU shared memory has been stored in its corresponding position in C_z (GPU global mem), along with the vector of block sums which will be passed later to the next kernel `scan_blocksums()`.

The kernel `scan_blocksums()` (Figure 4) only makes an exclusive prefixsum on the vector of block sums described above. The result is a vector containing, at index x , the value to be added to every element of block x in each line of C_z .

This summing is done in shared memory by kernel `add_sums2prefixes()` as described by Figure 5.

The values of C_{z^2} are obtained together with those of C_z and in exactly the same way. For publishing reasons, figures do not show the C_{z^2} part of structures.

With this implementation, speedups are quite significant (Table I). Moreover, the larger the image, the higher the speedup is, as the step-complexity of the sequential algorithm is of $O(N^2)$ and $O(N \log(N))$ for the parallel version. Even higher speedups are achieved by adapting the code to specific-size images, especially when the number of columns is a power of 2. This avoids inactive threads in the grid, and thus improves efficiency. However, since the use of 64-bit sums is imposed by image sizes (up to 12000 pixel wide) and 16-bit pixel coding, computations are made with a 2-

way bank conflict as sums are based on 64-bit words, thus creating overhead.

B. Segment contributions

The choice made for this implementation has been to keep the *1 thread per pixel* rule for the main kernels. Of course, some reduction stages need to override this principle and will be pointed out.

As each of the N_n nodes of the contour $S_{n,l}$ may move to one of the eight neighbor positions as shown in Figure 6, there is $16N_n$ segments whose contribution has to be estimated. The best combination is then chosen to obtain $S_{n,l+1}$ (Figure 6). Segment contributions are computed in parallel by kernel `GPU_compute_segments_contrib()`.

The grid parameters for this kernel are determined according to the size of the longest segment $npix_{max}$. If bs_{max} is the maximum theoretical blocksize that a GPU can accept,

- the block size bs is taken as
 - $npix_{max}$'s next power of two if $npix_{max} \in [33; bs_{max}]$
 - 32 if $npix_{max} < 32$
 - bs_{max} if $npix_{max} > 256$
- the number of threads blocks assigned to each segment, $N_{TB} = \frac{npix_{max} + bs - 1}{bs}$

Our implementation makes intensive use of shared memory and does not allow the use of the maximum theoretical blocksizes (512 for sm13, 1024 for sm20, see [9] and [7]). Instead we set $bs_{max}^{sm13} = 256$ and $bs_{max}^{sm20} = 512$. Anyway, testing has shown that most often, the best value is 256 for both *sm13* and *sm20* GPU's.

Then `GPU_compute_segments_contrib()` computes in parallel:

- each pixel coordinates for all $16N_n$ segments. Since the contour is only read in one direction, we have been able to use a very simple parallel algorithm instead of Bresenham's. It is based on the slope k of each segment: one pixel per row if $|k| > 1$, one pixel per column otherwise.

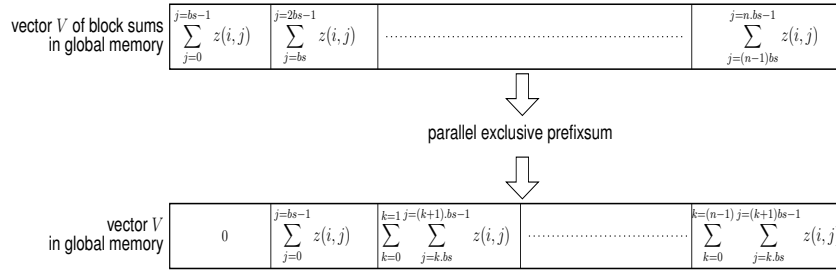


Figure 4. scan_blocksums() details.

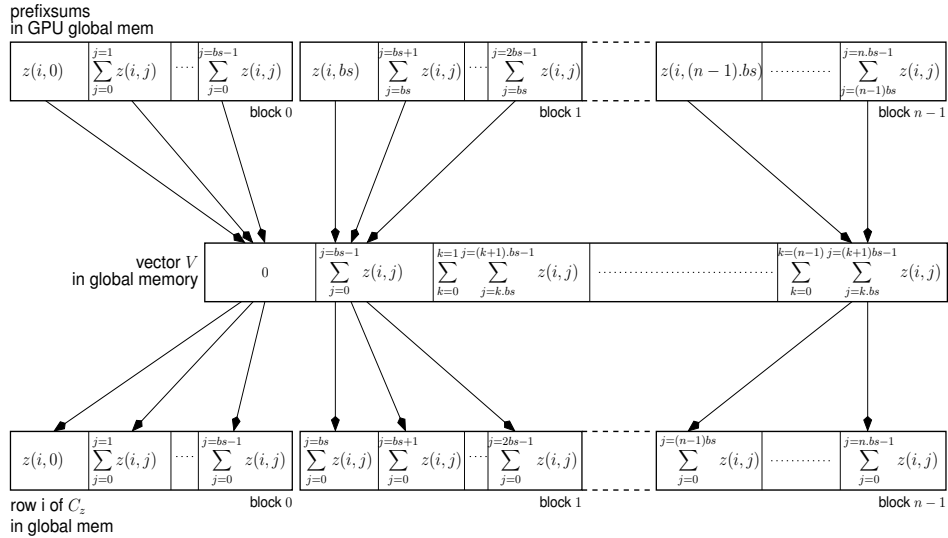


Figure 5. add_sums2prefixes() details.

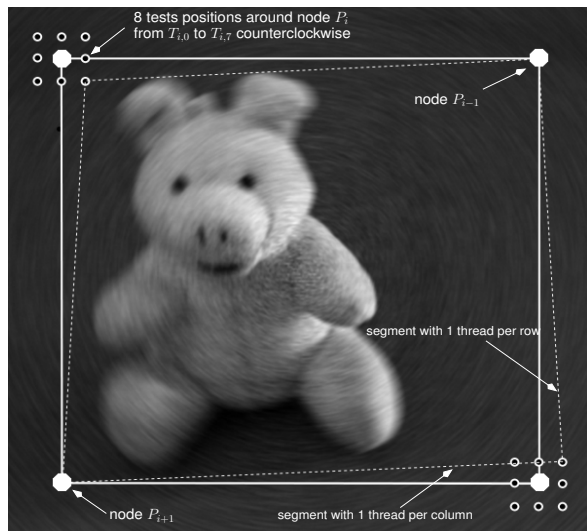


Figure 6. Optimization of node locations using 8 position tests around each node.

- each pixel contribution by reading the corresponding values in the lookup tables.

- each thread-block sum of individual pixel contributions by running a *reduction* stage for each block.

The top line of Figure 7 shows the base data structure in GPU shared memory which is relative to one segment. We concatenate the single segment structure as much as necessary to create a large vector representing every pixel of every test segment. As each segment has a different size (most often different from any power of two), there is a non-neglectable number of inactive threads scattered in the whole structure. Two stages are processed separately: one for all even nodes and another one for odd nodes, as shown in the two bottom lines of Figure 7.

The process is entirely done in shared memory; only a small amount of data needs to be stored in global memory for each segment:

- the coordinates of its middle point, in order to be able to add nodes easily if needed.
- the coordinates of its first and last two points, to compute the slope at each end of the segment.

The five values above are part of the weighting coefficients determination for each segment and node.

The GPU_sum_contribs() takes the blocks sums ob-

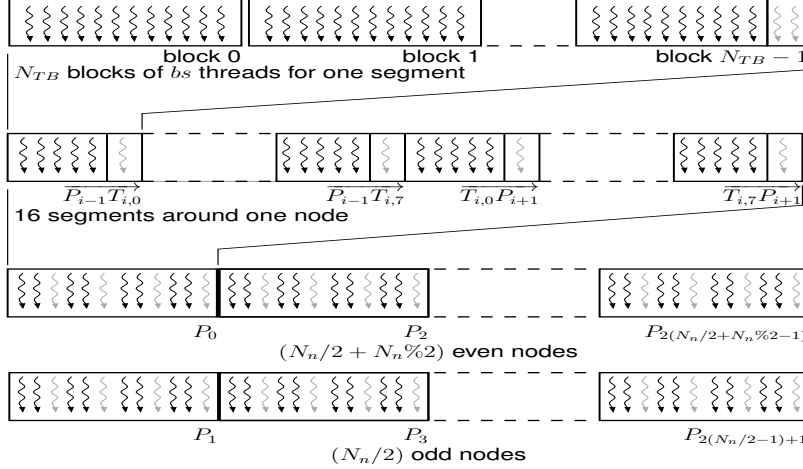


Figure 7. structure for segments contributions computation. Gray symbols help to locate inactive threads as opposed to black ones that figure active threads.

tained by `GPU_compute_segments_contrib()` and computes a second stage parallel summing to provide the $16N_n$ segment contributions.

C. Segments with a slope k such as $|k| \leq 1$

Such a segment is treated with 1 thread per column and consequently, it often has more than one pixel per row as shown by Figure 8. In an image row, consecutive pixels which belong to the target define an interval which can only have one low and one high ends. That's why, on each row, we choose to consider only the contributions of the innermost pixels. This selection is also done inside `GPU_compute_segments_contrib()` when reading the lookup tables for each pixel contribution. We simply set a null contribution for pixels that need to be ignored.

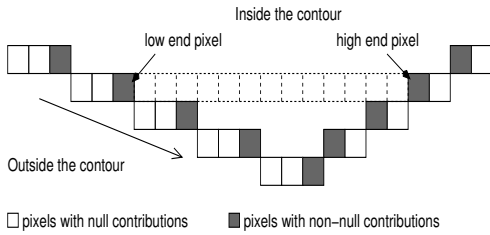


Figure 8. zoom on part a of segment with $|k| < 1$, at pixel level.

D. Parameters estimation

A `GPU_compute_GL()` kernel computes in parallel:

- every $8N_n$ vector of parameters values corresponding to each possible next state of the contour. Summing is done in shared memory but relevant data for these operations are stored in global memory.
- every associated pseudo likelihood value.
- every node substitution when better GL have been found and if it does not lead to segments crossing.

E. End of segmentation

Segmentation is considered achieved out when no other node can be added to the contour (Algorithm 3). A very simple GPU kernel adds every possible node and returns the number it added.

Algorithm 3: Parralel GPU algorithm: outlines.

<<<. . .>>> indicates a GPU kernel parallel process.

- 1: load images;
 - 2: transfer image from CPU to GPU;
 - 3: <<<compute the 2 cumulated images>>>;
 - 4: <<<initialize the contour>>>;
 - 5: **repeat** /* iteration level */
 - 6: **repeat** /* step level */
 - 7: <<<find best neighbor contour>>>;
 - 8: <<<adjust node's positions>>>;
 - 9: transfer the number of moves achieved from GPU memory to CPU memory.
 - 10: **until** no more node can be moved;
 - 11: <<<Add nodes>>>;
 - 12: transfert the number of nodes added from GPU memory to CPU memory.
 - 13: **until** no more node can be added;
-

VI. SPEEDUPS

Results are given in Table I. CPU timings were measured on an Intel Xeon E5530-2.4GHz with 12Go RAM (LIFC cluster). GPU timings were obtained on a C2050 GPU with 3GB RAM (adonis-11.grenoble.grid5000.fr).

Execution times reported are means on ten executions. The image of figure 1a (scaled down for printing reasons) is based on a real noisy image (800 x 800), 16-bit gray level. Contrast has been enhanced for better viewing; its various

		CPU	GPU	Speedup
Image 15MP	total	0.51 s	0.06 s	x8.5
	pre-comp.	0.13 s	0.02 s	x6.5
	segment.	0.46 s	0.04 s	x11.5
Image 100MP	total	4.08 s	0.59 s	x6.9
	pre-comp.	0.91 s	0.13 s	x6.9
	segment.	3.17 s	0.46 s	x6.9
Image 150Mp	total	5.7 s	0.79 s	x7.2
	pre-comp.	1.4 s	0.20 s	x7.0
	segment.	4.3 s	0.59 s	x7.3

Table I
GPU (C2050) vs CPU TIMINGS.

sizes have been obtained by interpolation and addition of gaussian noise.

We separately give the timings of pre-computations as they are a very general purpose piece of code. Segmentations have been performed with strictly the same parameters (initial shape, threshold length). The neighborhood distance for the first iteration is 32 pixels. It has a slight influence on the time process, but it leads to similar speedups values of approximately 7 times faster than CPU.

Though it does not appear in Table I, we observed that during segmentation stage, higher speedups are obtained in the very first iterations, when segments are made of a lot of pixels, leading to a higher parallelism ratio.

Several parameters prevent from achieving higher speedups:

- accesses in the lookup tables in global memory cannot be coalescent. It would imply that the pixel contributions of a segment are stored in consecutive spaces in C_z and C_{z^2} . This is only the case for horizontal segments.
- the use of 64-bit words for computations in shared memory often leads to 2-way bank conflicts.
- the level of parallelism is not so high, ie. the total number of pixel is not large enough to achieve impressive speedups. For example, on C2050 GPU, a grid can run about 66 million of threads, but a contour in a 10000 x 10000 image would be less than 0.1 million pixel long.

VII. CONCLUSION

The algorithm we have focused on is not easy to adapt for high speedups on GPGPU, though we managed to make it work quite faster than on CPU. The main drawback is clearly its relative low level of parallelism. Nevertheless, we proposed different kernels that allowed us to take advantage of the computation power of GPUs. In future works, we plan to try and manage to benefit from larger computing grids of thread blocks. Among the possible solutions, we plan to work on:

- slicing the image and processing the parts in parallel.

This is made possible since sm20 GPU provide multi kernel capabilities.

- slicing the image and processing the parts on two different GPUs, hosted by the same CPU.

To extend the scope of this work beyond our present hypothesis (based on *single* target segmentation), we are also going to investigate achieving speedups in *multiple* target segmentation of large images. This might be useful in a wide range of applications.

REFERENCES

- [1] M. Kass, A. P. Witkin, and D. Terzopoulos, "Snakes: Active contour models," *International Journal of Computer Vision*, vol. 1, no. 4, pp. 321–331, 1988.
- [2] C. Chesnaud, P. Réfrégier, and V. Boulet, "Statistical region snake-based segmentation adapted to different physical noise models," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 21, no. 11, pp. 1145–1157, 1999.
- [3] C. Bryan, S. Bor-Yiing, S. Narayanan, Lee, Yunsup, M. Mark, and K. Kurt, "Efficient, high-quality image contour detection," *International Conference on Computer Vision*, pp. 2381–2388, 2009.
- [4] T. Schoenemann and D. Cremers, "A combinatorial solution for model-based image segmentation and real-time tracking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, pp. 1153–1164, 2010.
- [5] E. Dipl.-Inf. Kienel and G. Prof. Dr. Brunnett, "Gpu-accelerated contour extraction on large images using snakes," 2009. [Online]. Available: <http://archiv.tu-chemnitz.de/pub/2009/0035>
- [6] C. Kauffmann and N. Piche, "Cellular automaton for ultra-fast watershed transform on gpu," in *ICPR*, 2008, pp. 1–4.
- [7] *NVIDIA CUDA C Programming Guide v3.1.1*, NVIDIA Corporation, 7 2010.
- [8] M. Harris, S. Sengupta, and J. D. Owens, *Gpu gems 3*, 1st ed. Addison-Wesley Professional, 2007, ch. 39 - Parallel Prefix Sum with CUDA.
- [9] *NVIDIA Fermi Tuning Guide*, NVIDIA Corporation, 7 2010.