



Fast parallel vessel segmentation

Nitin Satpute^{a,*}, Rabia Naseem^b, Rafael Palomar^c, Orestis Zachariadis^a, Juan Gómez-Luna^d, Faouzi Alaya Cheikh^b, Joaquín Olivares^a

^a Department of Electronic and Computer Engineering, Universidad de Córdoba, Spain

^b Norwegian Colour and Visual Computing Lab, Norwegian University of Science and Technology, Norway

^c The Intervention Centre, Oslo University Hospital, Norway

^d Department of Computer Science, ETH Zurich, Switzerland

ARTICLE INFO

Article history:

Received 17 December 2019

Revised 17 February 2020

Accepted 2 March 2020

Keywords:

Seeded region growing

GPU

Kernel termination and relaunch (KTRL)

Persistent

Grid-stride loop

ABSTRACT

Background and Objective: Accurate and fast vessel segmentation from liver slices remain challenging and important tasks for clinicians. The algorithms from the literature are slow and less accurate. We propose fast parallel gradient based seeded region growing for vessel segmentation. Seeded region growing is tedious when the inter connectivity between the elements is unavoidable. Parallelizing region growing algorithms are essential towards achieving real time performance for the overall process of accurate vessel segmentation.

Methods: The parallel implementation of seeded region growing for vessel segmentation is iterative and hence time consuming process. Seeded region growing is implemented as kernel termination and relaunch on GPU due to its iterative mechanism. The iterative or recursive process in region growing is time consuming due to intermediate memory transfers between CPU and GPU. We propose persistent and grid-stride loop based parallel approach for region growing on GPU. We analyze static region of interest of tiles on GPU for the acceleration of seeded region growing.

Results: We aim fast parallel gradient based seeded region growing for vessel segmentation from CT liver slices. The proposed parallel approach is 1.9x faster compared to the state-of-the-art.

Conclusion: We discuss gradient based seeded region growing and its parallel implementation on GPU. The proposed parallel seeded region growing is fast compared to kernel termination and relaunch and accurate in comparison to Chan-Vese and Snake model for vessel segmentation.

© 2020 Elsevier B.V. All rights reserved.

1. Introduction

In medical imaging, vessel segmentation from liver slices is one of the challenging tasks. Seeded region growing (SRG) is a widely used approach for semi automatic vessel segmentation [1,2]. Delibasis et. al. [3] have proposed a tool based on a modified version of SRG algorithm, combined with a priori knowledge of the required shape. SRG starts with a set of pixels called seeds and grows a uniform, connected region from each seed. Key steps to SRG are to define seed(s) and a classifying criterion that relies on the image properties and user interaction [4]. SRG starts from a seed and finds the similar neighboring points based on the threshold criteria using 4 or 8 connectivity. Region is grown if the thresh-

old criteria is satisfied. Similar neighbors are new seed points for the next iteration. This process is repeated until the region can not be grown further. In practice, it demands high computational cost to the large amount of dependent data to be processed in SRG especially in the medical image analysis and still requires efficient solutions [5].

SRG is an iterative process. SRG is invoked continuously until region can not be grown further. Iterative process in SRG, when implemented on GPU requires terminating kernel and relaunching from CPU (Kernel Termination and Relaunch (KTRL)) and data transfers between CPU and GPU [1,4]. So our main objective is to reduce these data transfers using different inter block GPU synchronization (IBS) methods resulting in an efficient parallel implementation of SRG. IBS provides flexibility to move all the computations on GPU by providing visibility to updated intermediate data without any intervention from CPU.

* Corresponding author.

E-mail address: el2sasan@uco.es (N. Satpute).

Table 1
List of abbreviations with full forms.

List	Full Forms
SRG	Seeded Region Growing
GPU	Graphics Processing Unit
CPU	Central Processing Unit
RoI	Region of Interest
KTRL	Kernel Termination and Relaunch
IBS	Inter Block GPU Synchronization
CT	Computed Tomography
PT	Persistent Threads
SM	Streaming Multiprocessor
CUDA	Compute Unified Device Architecture
DS	Dice Score

In this paper, we propose persistent, grid-stride loop and IBS based GPU approach for SRG to avoid intermediate memory transfers between CPU and GPU. This also reduces processing over unnecessary image voxels providing significant speedup. Persistent thread block (PT) approach is basically dependent on number of active thread blocks and grid-stride loop becomes essential when the number of threads in the grid are not enough to process the image voxels independently [6,7].

We implement parallel image gradient using grid-stride loop and propose gradient and shared memory based fast parallel SRG implemented entirely on GPU without any intermediate transfers between CPU and GPU. This is inspired by parallel processing on static region of interest (RoI) of tiles on GPU. We compare the proposed persistent based parallel SRG with KTRL for accurate vessel segmentation. The gradient based fast parallel SRG for 2D vessel segmentation is $1.9 \times$ faster compared to the state-of-the-art.

The rest of the paper is structured as follows. Section 2 briefs relevant works and state-of-the-art with respect to SRG. Section 3 explains GPU approaches (KTRL and Static) for SRG implementation using persistence and grid-stride loop. The application of parallel SRG to vessel segmentation is discussed in the Section 4. Performance results and comparison of persistent and grid-stride loop based parallel SRG for vessel segmentation are mentioned in the Section 5. Section 6 concludes summarizing the main conclusions of this paper and indicating future directions. List of abbreviations with explanations are mentioned in Table 1.

2. Background and motivation

There are many works done on image segmentation recently which are based on snake based model [8], gradient vector flow [9,10], and level set based Chan-Vese model [11]. Scientists have explored the snake model for segmentation. Snakes are defined as a set of points around a contour [8]. But the problem with the snake model is that the contour never sees the strong edges that are far away and the snake gets hung up due to many small noises in the image [8]. Hence researchers came up with the solution called gradient vector flow (GVF). In GVF, instead of using image gradient, a new vector field is created over the image plane [9,10]. Cost of GVF includes smoothness and edge map but it requires keeping track of the number of points and point distribution. Hence researchers came up with another solution called as level sets based Chan-Vese model for image segmentation [11]. In the absence of strong edges, a region based formulation for image segmentation is proposed by Chan-Vese model. Chan-Vese model for active contours is a powerful and flexible method which is able to segment many types of images. But amongst all, SRG is the simplest algorithm and plays a vital role in medical image segmentation [1,12].

Smistad et al. [13,14] have discussed parallel SRG for image segmentation. The reference implementation is shown in Fig. 1. Medi-

cal image dataset is cropped before processing. Then the CPU allocates the memory equivalent to the cropped size to copy the data to the cropped image on the GPU. Further SRG is performed for image segmentation. This is the simplistic representation of the work by Smistad et al. [14]. We have not considered pre-processing stage in this work assuming the images are pre-processed. Smistad et al. [14] have proposed non persistent thread (non-PT) approach for SRG based vessel segmentation.

Smistad et al. [4] have proposed parallel region growing with double buffering algorithm based on the parallel breadth first search algorithm by Harish and Narayanan [15]. They have suggested a dynamic queue for SRG and mentioned that changing the number of threads (due to border expansion of the region) typically involves restarting the kernel, and this requires reading all the values from global memory again. But they have not recommended probable solution for this problem. Smistad et al. [14] have presented a data parallel version of the SRG based Inverse Gradient Flow Tracking Segmentation algorithm using KTRL. Zhang et al. [16] have implemented bidirectional region growing where they have used a dynamic queue (stack). Jiang et al. [17] have proposed improved branch based region growing vessel segmentation algorithm using stack.

GPU based implementation of SRG needs a dynamic queue (stack). CPUs provide hardware support for stacks but GPUs do not [7]. Any queuing system has a large number of pieces of work to do and a fixed number of workers corresponding to the fixed number of computing units. Pieces are then assigned dynamically to the workers. The problem is deciding the maximum number of pieces of work in the queuing system. If decided, persistent blocks iterate through these pieces of work in the queuing system.

GPU implementation of a stack requires continuous changes in memory allocations which in turn requires iterative GPU kernel invocation from CPU in other words kernel termination and relaunch as discussed in the algorithms IVM backtracking and work stealing phase by Pessoa et al. [18]. Task-parallel run-time system, called TREES, that is designed for high performance on CPU/GPU platforms by Hechtman et al. [19] have shown the invocation of GPU kernels from CPU iteratively for updating task mask stack (TMS) in TREES execution. The loop involved while implementing data flow through the stream kernels of the rendering system (involving stack) on GPU controlled by CPU (that is KTRL) is proposed by Ernst et al. [20].

Nevertheless, there is an alternate GPU implementation of queuing system (stack) using dynamic kernel launching. Chen et al. [7] have proposed free launch based dynamic kernel launches through thread reuse technique [7]. This technique requires no hardware extensions, immediately deployable on existing GPUs. By turning subkernel launch into a programming feature independent of hardware support, free launch provides alternate approach for subkernel launch which can be used beneficially on GPUs.

KTRL includes terminating a GPU kernel and invoking it from the CPU if the region can be grown further [4,14]. GPU kernel SRG is called from CPU. Region grows from a seed based on the threshold criteria. SRG kernel is terminated and relaunched from CPU if region is not grown completely. This process continues until region can not be grown further. The process involves transfer of data to and fro from CPU and GPU. In KTRL, SRG kernel operates on each voxel of whole image data in all the iterations. It includes redundant memory transfers and unnecessary computations over complete image. Hence the main contributions of this paper are the implementation of persistence based approaches to improve the performance of SRG by reducing unwanted computations and avoiding intermediate memory transfers between CPU and GPU. Memory on the GPU is limited and may not be enough for processing large medical datasets. However, most medical datasets contain a lot of data that is not part of the RoI.

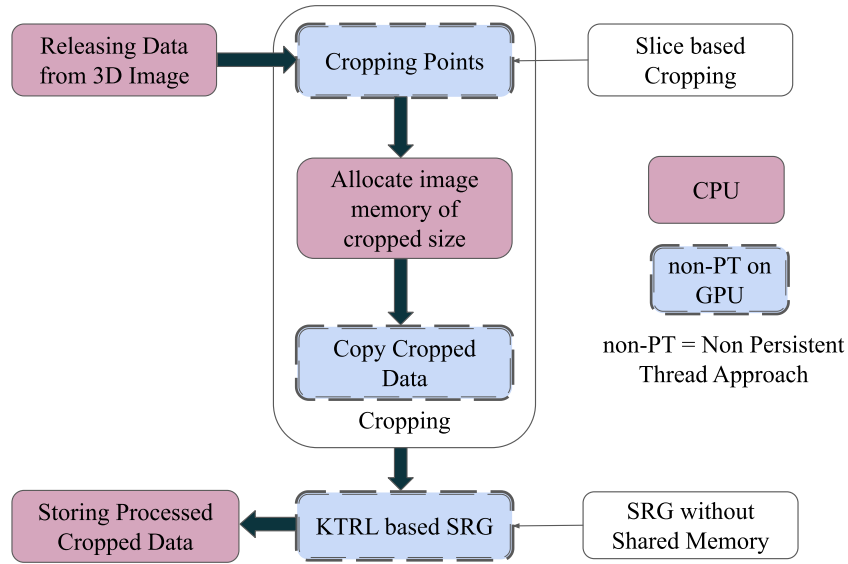


Fig. 1. Reference approach derived from Smistad et. al. [14]

The process of KTRL which involves iterative calling of the kernel is not efficient when implemented on GPU. Hence, as an optimized solution to KTRL, we propose persistent and grid-stride loop based GPU approaches. These approaches are based on processing over static RoI of tiles and dynamic RoI of tiles. We discuss the further details in the upcoming sections.

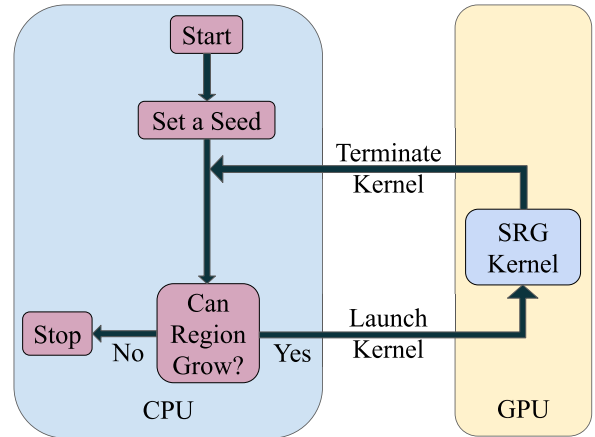
3. Parallel SRG

GPU is a grid of block of threads. Thread is the smallest computational unit mapped on the cores and block of threads are mapped on the streaming multiprocessors (SMs). Each SM can occupy more than one block. The threads from independent blocks can access data via shared memory in the SM [21]. In order to communicate valid data between the blocks, these persistent blocks need to be synchronized via IBS through device memory. Persistence implies maximum number thread blocks that can be active at the time of computation depending upon the GPU resources available [6,22].

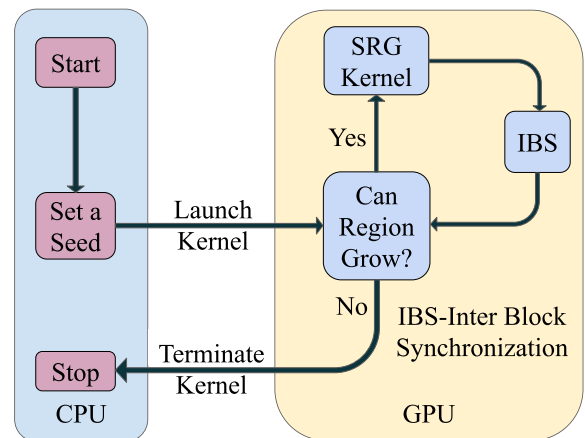
We use PT and shared memory based approaches for SRG implementations. Shared memory and grid-stride loop based SRG reduces total memory transfers and computations. Grid-stride is inspired when the grid is not large enough to occupy all the data elements [23,24]. Rather than assuming that the thread grid is large enough to cover the entire image elements, the kernel loops over the image one grid-size at a time. The stride of the loop is the total number of threads on the grid [23]. These threads (or block of threads) iterate over the image until the process of SRG terminates.

For each thread in parallel on GPU, SRG starts from the seed thread and finds similar neighbours surrounding it Region is grown by making similar neighbouring elements as new seeds. The process of SRG is repeated until similar neighbours can not be found. Normally SRG can be implemented on the GPU as a recursive or iterative kernel calling (KTRL) as shown in Fig. 2a. Kernel calling involves invocation of a grid of block of threads. The blocks are executed on streaming multiprocessors and threads are executed on cores. Park et al. [25] and Smistad et al. [4] have given brief introduction about CUDA (Compute Unified Device Architecture) architecture and GPU computing. They have detailed the information on grid, blocks, threads and memory hierarchy of CUDA architecture.

SRG can be recursive or iterative process. Recursive kernel calling can not utilize GPU cores efficiently due to hardware limita-



(a) using Kernel Termination and Relaunch [13, 14]



(b) using Proposed Grid-Stride Loop

Fig. 2. GPU implementations of seeded region growing (SRG).

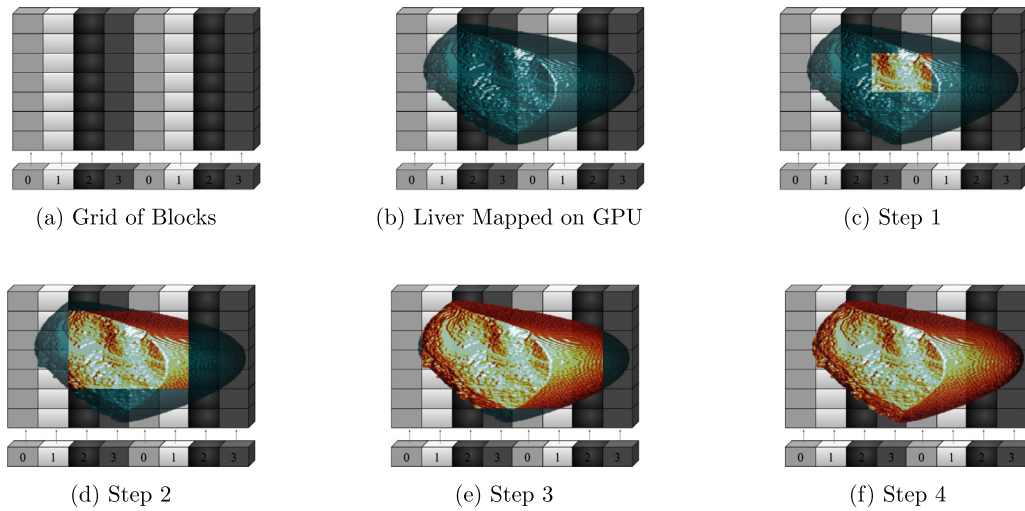


Fig. 3. SRG using persistence and grid-stride loop through complete image.

tions [26]. Iterative GPU kernel call from CPU is costlier due to memory transfers between CPU and GPU and it involves all the image elements to be considered in each step of SRG. GPU implementation of SRG using KTRL is shown in the Fig. 2a. It shows that, the kernel SRG is called on GPU continuously from the host CPU until the region can not be grown further. It starts from the seed, finds similar neighbours and grows the region. This process continues until the region can not be grown further. The process of the KTRL causes unnecessary image elements to be part of computations and intermediate memory transfers between CPU and GPU.

Hence in order to avoid these problems, we propose grid-stride loop through complete image based GPU approach as shown in the Fig. 2b. SRG starts from the seed and the control goes to GPU. The SRG kernel is launched if the region is not grown completely. IBS is needed in order to transfer valid data in between the active thread blocks. The number of active thread blocks on SMs are limited due to resource constraints. These maximum number of active blocks are persistent blocks [6,21,22]. The looping i.e. grid-stride loop terminate when the region can not be grown further and control returns to the host CPU as shown in the Fig. 2b. We have discussed KTRL based GPU approach for SRG implementation and its disadvantages. Now, we are going to analyze PT based GPU approaches for high performance SRG implementation. Proposed approaches exploit parallelism using persistence and IBS as detailed in the static and dynamic approaches.

3.1. Static approach

In the proposed approach, we apply grid-stride loop through static RoI (complete image) using persistence and IBS [6,22]. The complete liver image is mapped on the GPU as grid of block of threads as shown in the Fig. 3b. CPU invokes SRG kernel on GPU. Persistent blocks iterate through complete image and grow region from the seed in each and every iteration on GPU. This iteration of persistent blocks over the tiles of the image and the grid-stride loop based SRG is shown in Fig. 2b. Steps of SRG in Figs. 3c–f show the grown region of the liver. SRG kernel terminates when the region is grown completely. We copy the data from the device memory to the shared memory. This data is shared by all the threads inside the blocks. This is necessary to share the neighbouring elements between different voxels of the image. For each parallel thread in the block, if seed is found and is not the boundary element of the block, we calculate similar neighbouring elements.

Region is grown by making similar neighbouring elements as new seeds.

There are four persistent blocks shown in Fig. 3. These four persistent blocks are iterated through liver elements. Tiles with the same color are iterated by same persistent block. In KTRL, these tiles are processed by the thread blocks randomly. Grid-stride loop by persistent blocks is applied on the tiles over complete liver image (Step 1 in Fig. 3c). Region grows around the seed containing similar elements. IBS is applied to communicate valid data in between the blocks for the next step of SRG as shown in Fig. 2b.

Persistent blocks iterate over the liver image and the region is grown again in step 2 as shown in Fig. 3d. IBS is applied and valid data is communicated in between the blocks so that the region can be grown further as shown in Fig. 3e and f. After step 4 in Fig. 3f, SRG stops as region can not be grown further. Each step contain many iterations where region starts growing when persistent blocks iterate through tiles of the image. This iterative process continues until region can not be grown further. Code snapshot of the complete process is provided in the Algorithm 1.

Algorithm 1: Grid-stride loop through complete image.

```

1: unfinished=1;
2: while unfinished==1 do
3:   unfinished=0;
4:   for int i=blockIdx.x;i <=
      width/(blockDim.x - 2);i+=gridDim.x do
5:     for int j=blockIdx.y;j <=
      height/(blockDim.y - 2);j+=gridDim.y do
6:       for int k=blockIdx.z;k <=
      depth/(blockDim.z - 2);k+=gridDim.z do
7:         Region_Growing(arguments, unfinished);
8:       end for
9:     end for
10:  end for
11:  Inter_Block_GPU_Sync();
12: end while

```

Global variable “unfinished” is 1 if region has to be grown further else it is 0. Persistent blocks in x , y and z directions iterate through complete image. Two is subtracted from block dimensions to avoid computations around boundary voxels (from left and right in each dimensions) from shared memory as region can not be grown further in the block. After each step of SRG, when the processing on complete image is done then all the persistent blocks are globally synchronized via “Inter_Block_GPU_Sync()” barrier. This ensures that valid data is communicated for the next step of SRG. This barrier can be Atomic(), Quasi(), LockFree() or can be implemented using NVIDIA CUDA API Cooperative-groups [22,27,28]. We use quasi based IBS because of its efficient implementation [27].

The difference between static and dynamic approach by Nitin et. al. [1] can be explained in terms of static and dynamic RoI of tiles. In static approach, RoI remains constant and SRG happens within the constant RoI until the region can not be grown further. Whereas in dynamic approach, SRG starts within the initial RoI. RoI increases and includes more elements uniformly in all the directions for the next step of SRG. SRG takes place, RoI increases and the region is grown further. Hence RoI changes in each step of SRG until the region can not be grown further in a dynamic approach. In the next section, we present 2D vessel segmentation as an application to static RoI based SRG.

4. Application to 2D vessel segmentation

The 2D segmentation algorithm is inspired by the gradient based SRG algorithm developed by Rai and Nair [29]. We proposed the fast parallel SRG based segmentation algorithm on GPU for vessel segmentation. We discuss the two important modules i.e. image gradient and SRG for the fast parallel 2D segmentation of vessels from CT liver images.

4.1. Parallel image gradient

Rai and Nair [29] have presented homogeneity criterion selection and its impact on the quality of segmentation using SRG. In general, the threshold criteria include object contrast, region boundary, homogeneity of the region, intensities values and texture features like shape and color. But we include cost functions mainly based on intensity values and their gradient direction and magnitude.

The cost function exploits certain features of the image around the seed. Gradient based cost function requires gradient of the image, largest gradient magnitude (\max_g) and minimum gradient (\min_g) present in the image. The cost functions are:

$$\text{cost1} = g/(k * \max_g) \quad 0 < \text{cost1} < 1 \quad (1)$$

$$\text{cost2} = (\max_g - g)/(\max_g - \min_g) \quad 0 < \text{cost2} < 1 \quad (2)$$

where g is gradient magnitude of the pixel under consideration and k is the constant parameter which controls the region growth. The pixel under consideration is added in the growing region if it matches with the seed elements i.e. cost functions specified by Eqs. (1) and (2) are satisfied otherwise it is excluded from consideration.

We propose grid-stride loop based parallel image gradient method in Algorithm 2. For each pixel in parallel, we calculate its gradient magnitude (g) with respect to neighbouring element. Horizontal and vertical gradient components are given by g_x and g_y . The magnitude of maximum and minimum gradients are updated simultaneously. The gradient of the image is desired input for SRG based segmentation along with the seed. This is discussed in the next section.

Algorithm 2: Parallel image gradient using grid-stride loop.

```

1: voxel.x = blockIdx.x * blockDim.x + threadIdx.x;
2: voxel.y = blockIdx.y * blockDim.y + threadIdx.y;
3: stridex = blockDim.x * gridDim.x;
4: stridey = blockDim.y * gridDim.y;
5: for int k=voxel.x; k < rows; k=k+stridex do
6:   for int l=voxel.y; l < cols; l=l+stridey do
7:     candidate.x = k + 1; candidate.y = l + 1;
8:     check if neighbour candidate is within image dimensions;
9:     gx = 0.5*(data[candidate.x*cols + l] - data[k*cols + l]);
10:    gy = 0.5*(data[k*cols + candidate.y] - data[k*cols + l]);
11:    g = sqrt(gx*gx + gy*gy);
12:    data_g[k*cols + l]=g;
13:    if(max_g < g) atomicMax(&max_g, g);
14:    if(min_g > g) atomicMin(&min_g, g);
15:   end for
16: end for

```

4.2. Parallel vessel segmentation

We propose fast parallel vessel segmentation as shown in Fig. 4. The algorithm is inspired from gradient based segmentation algorithm by Rai and Nair [29]. Fig. 4 shows parallel implementation of vessel segmentation where the user selects seed(s). These seed(s) along with the image are transferred to the GPU. Device kernel calculates the image gradient in parallel as discussed in the earlier section. The IBS is necessary to reflect the updated image gradients in the device memory.

Further we apply SRG algorithm. The cost functions based on gradient are shown in Eqs. (1) and (2). For each pixel in parallel, the pixel under consideration invokes SRG kernel if it satisfies the cost functions. The seeds are updated after IBS and the gradient based cost functions are verified again for new pixels. This process continues until no new seeds are formed i.e. no new pixels are added to the growing region.

The kernel is terminated and the control returns to the CPU when the region is grown completely. The segmented image is transferred to the CPU. The process of segmentation stops. This GPU implementation avoids iterative call of SRG kernel from CPU. We use gradient and persistent based parallel SRG for vessel segmentation.

5. Performance evaluation

We propose persistent and grid-stride based GPU approaches for fast parallel 2D vessel segmentation. The performance results are obtained from KTRL and proposed persistent based GPU approach. We compare proposed approaches with KTRL. We use Intel(R) Core(TM) i7-7700HQ CPU @ 2.80GHz RAM 24 GB, NVIDIA GPU 1050 (RAM 4GB), OpenCL 1.2 (ref. [30]) and CUDA Toolkit 10.1 for the implementation.

5.1. Liver dataset and ground truth

Liver data for the research work has been acquired from The Intervention Center, University of Oslo, Norway [31]. The ground truths for vessel segmentation are provided by the clinician. The modality used is Computed tomography (CT). For the ground truth, images are pre-processed through locally developed application with 3D Slicer to enhance vessels [32]. In some cases, the same application is used for vessel segmentation and separation of portal and hepatic vessels although another possibility is to use active contour tool using ITK-SNAP and manual correction [1,31]. Table 2 shows information about images of different sizes including total

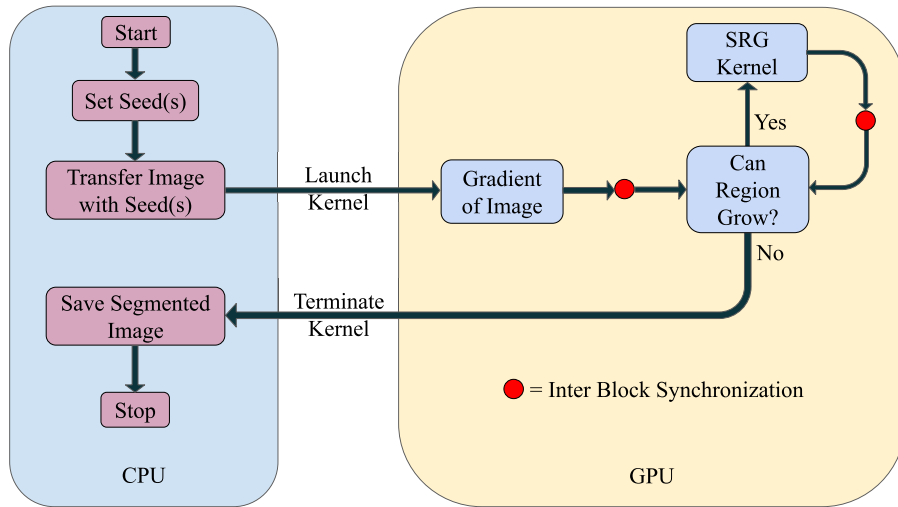


Fig. 4. Proposed parallel vessel segmentation.

Table 2
Liver dataset with vessels.

Volume #	Total # of Slices	Image Size ($w \times h$)	# of Slices with Vessels
10504	59	460 × 306	7
18152	139	512 × 512	5
23186	87	405 × 346	6
28059	59	462 × 321	6

number of vessel slices used for experimentation from a particular volume.

5.2. Parallel 2D vessel segmentation

We propose persistent and grid-stride based GPU approaches for fast parallel 2D vessel segmentation. Variations in vessel segmentation with parameter ' k ' using parallel SRG is shown in Fig. 5. Input to parallel SRG is CT liver slice as shown in Fig. 5a. Gradient of input CT image and the ground truth for the segmentation are shown in Fig. 5b and c respectively. Dice similarity coefficient (DS) and Precision [33,34] are used to assess the quality of vessel segmentation. Dice similarity coefficient measures the similarity between ground truth and the segmented output. If they are identical (i.e. they contain the same elements), the coefficient is equal to 1.0, while if they have no elements in common, it is equal to 0.0. Otherwise it is somewhere in between 0 to 1. Precision describes the number of positive detections with respect to the ground truth. Of all of the elements that are segmented in a given liver vessel image, the number of these elements actually had a matching ground truth annotation can be called as precision.

We show the two segmented vessels with change in parameter k i.e. 0.04, 0.05, and 0.06. The first segmented vessel as shown in Fig. 5d, e, f is accurate at 0.04 with high dice score i.e. $DS = 0.77$. Similarly we show the segmentation of second vessel from the same slice. The variations in the quality segmentation due to k are shown in Fig. 5g, h, i. The more accurate segmentation is obtained at 0.05 as the dice similarity coefficient value is higher i.e. $DS = 0.60$.

Further we show the quality of segmentation on another CT Slice as shown in Fig. 6a and calculate the gradient (Fig. 6b) of the input CT image. GPU computes parallel SRG using gradient based thresholding criteria giving more accurate results with high dice score at $k = 0.05$ for two vessels inside the CT slice as shown in Fig. 6c and d. The ground truth for the segmentation is shown in Fig. 6e. We analyze that the vessels are more accurately segmented

due to better value of dice similarity coefficient when the parameter k takes the value 0.05 as shown in Fig. 6.

The speedup obtained by proposed parallel static approach over KTRL on first two CT liver slices are shown in the Table 3. The maximum speedup for vessel segmentation by proposed parallel static SRG is $1.67 \times$ in comparison to KTRL on the first liver slice. But the average speedup obtained by proposed parallel static approach for all the vessels (in 6 slices tested) is $1.9 \times$ compared to KTRL. We evaluate the speedup of the vessel segmentation when the vessel segmentation is more accurate with better Dice score value (Fig. 5).

Further we analyze the effect of parallel SRG on different slices for multiple vessel segmentation using multiple seeds as shown in Figs. 7–10. Segmentation of the long vessel as shown in Fig. 8d is slightly extended compared to the ground truth shown in Fig. 8e. It can be seen from input CT image and gradient image (Figs. 8a and b), long vessel has extension which is not shown in the ground truth. We show the thick vessel segmentation in Figs. 7d, 9 c, 8 c and thin vessel segmentation in Figs. 7c, 8 c and d. The results of segmentation in terms of Dice Score and Precision are provided in Table 4. The highest and lowest value of precision for these slices are 0.94 and 0.79 respectively. This implies the number of positive detections in the segmented images are higher. It is possible to use multiple seeds for the same vessel in the proposed vessel segmentation approach. We have the flexibility to provide two seeds on the same vessel and then the proposed approach can create a curve or a line of initial seeds (if needed) as an input for SRG. It is useful in order to increase the quality of vessel segmentation.

There are many works done on image segmentation recently which are based on snake based model [8], gradient vector flow [9,10], and level set based Chan-Vese model [11]. We validate the performance on 72 vessels from 24 vessel slices obtained from 4 different volumes in Table 5. Table shows the comparison of the vessel segmentation accuracy between three different models i.e. Snake model [8], Chan-Vese [11], and proposed SRG in terms of dice score and precision. Average dice score value and precision obtained by proposed SRG is outperforming the Chan-Vese and Snake based vessel segmentation.

Our proposed parallel implementations of SRG is not only fast but also accurate for vessel segmentation. The accuracy of the segmentation depends on the parameter ' k '. Clinicians get the flexibility to decide which segmentation is more accurate. The process takes very less time (few ms). Hence this reduces the overall time for segmentation for various values of parameter ' k ' if the

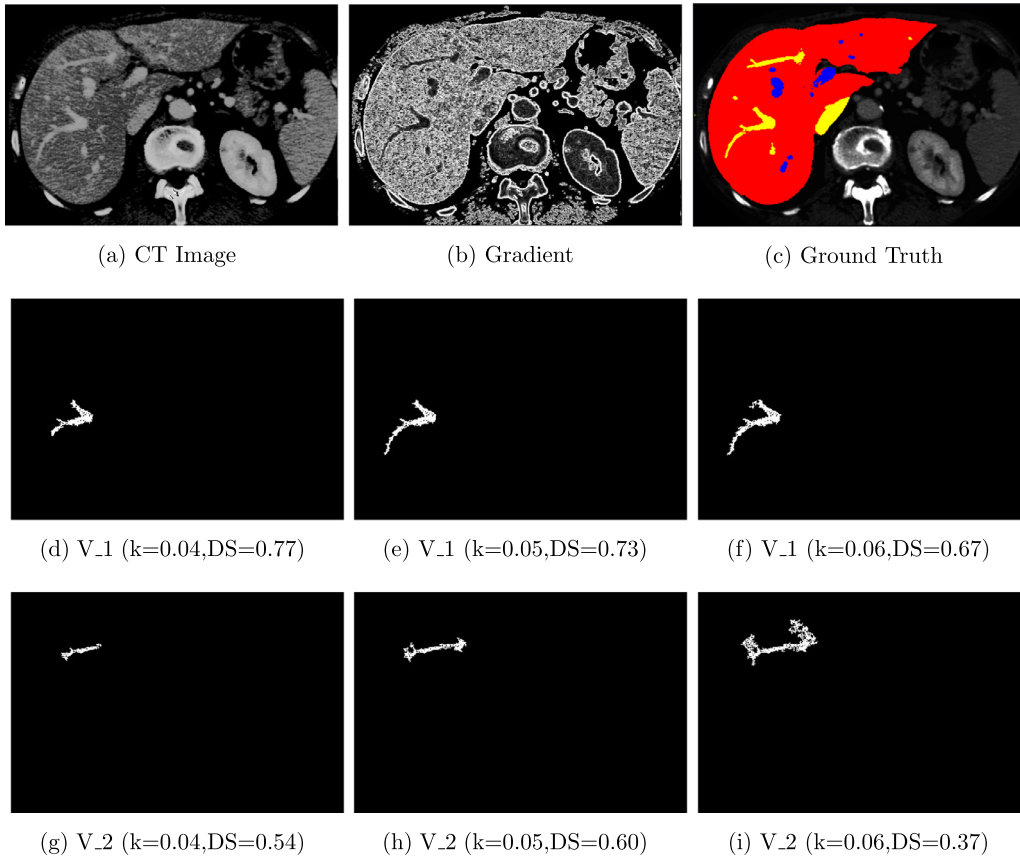


Fig. 5. Variations in fast parallel vessel segmentation with constant parameter 'k' using parallel SRG on first liver slice.

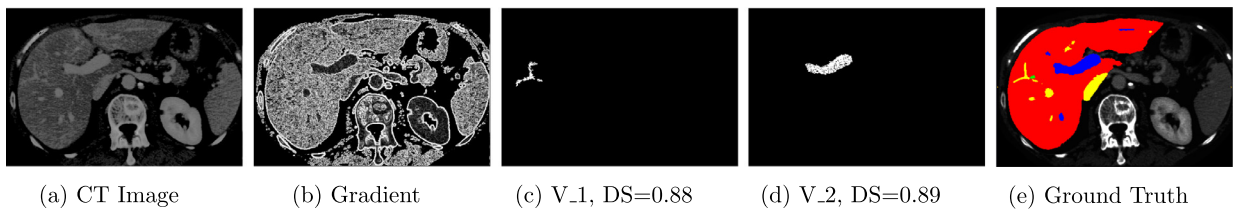


Fig. 6. Vessel segmentation (for $k = 0.05$) using parallel SRG on second liver slice.

Table 3
Time and speedup for vessel segmentation.

Data → GPU Approaches → Metrics ↓	Vessel Segmentation	
	KTRL	Static (Speedup)
Time in ms for kernel SRG - 1st Slice ($k = 0.05$) 1st vessel	5.7	3.4 (1.67 ×)
Time in ms for kernel SRG - 1st Slice ($k = 0.05$) 2nd vessel	2.1	1.5 (1.4 ×)
Time in ms for kernel SRG - 2nd Slice ($k = 0.05$) 1st vessel	1.5	1 (1.5 ×)
Time in ms for kernel SRG - 2nd Slice ($k = 0.05$) 2nd vessel	3.5	2.4 (1.45 ×)

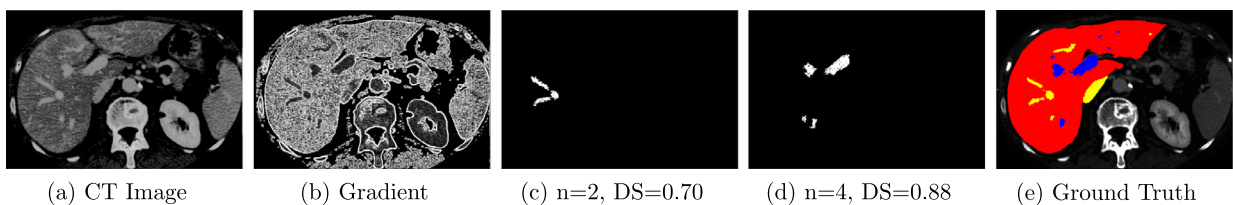


Fig. 7. Fast parallel vessel segmentation using parallel SRG on 3rd liver slice using multiple seeds (n).

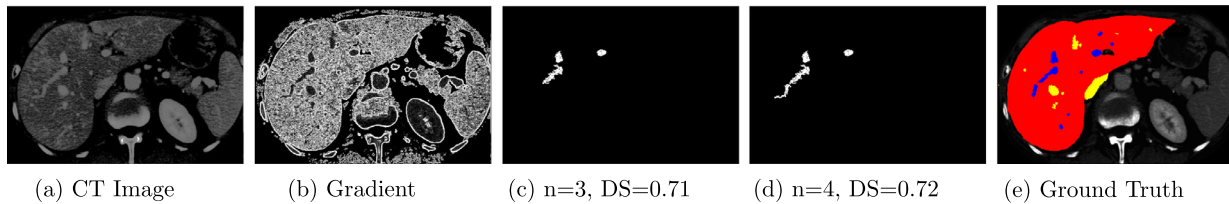


Fig. 8. Fast parallel vessel segmentation using parallel SRG on 4th liver slice using multiple seeds (n).

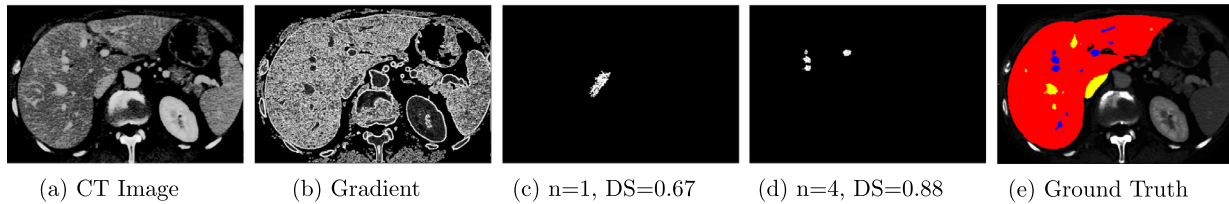


Fig. 9. Fast parallel vessel segmentation using parallel SRG on 5th liver slice using multiple seeds (n).

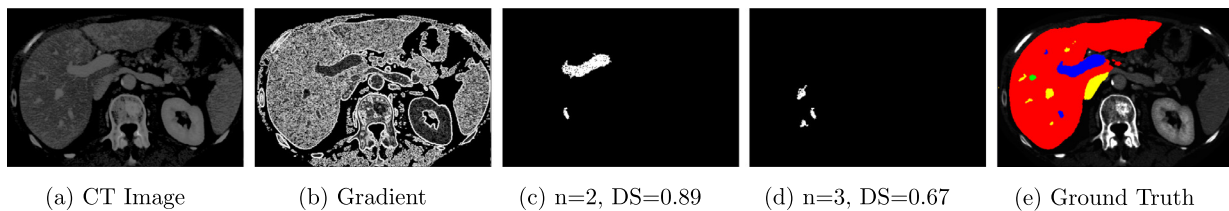


Fig. 10. Fast parallel vessel segmentation using parallel SRG on 6th liver slice using multiple seeds (n).

Table 4

Quality of vessel segmentation in terms of dice score and precision for Figs. 5–10.

Sr. No.	Ground Truth Figure	Segmented Image Figure	Value of 'k'	Dice Score	Precision
1	Fig. 5c	Fig. 5d	0.04	0.77	0.86
2		Fig. 5e	0.05	0.73	0.94
3		Fig. 5f	0.06	0.67	0.94
4	Fig. 5c	Fig. 5g	0.04	0.54	0.86
5		Fig. 5h	0.05	0.60	0.83
6		Fig. 5i	0.06	0.37	0.86
7	Fig. 6e	Fig. 6c	0.05	0.88	0.86
8		Fig. 6d	0.05	0.89	0.85
9	Fig. 7e	Fig. 7c	0.05	0.70	0.91
10		Fig. 7d	0.05	0.88	0.91
11	Fig. 8e	Fig. 8c	0.05	0.71	0.79
12		Fig. 8d	0.05	0.72	0.90
13	Fig. 9e	Fig. 9c	0.05	0.67	0.82
14		Fig. 9d	0.05	0.88	0.88
15	Fig. 10e	Fig. 10c	0.05	0.89	0.94
16		Fig. 10d	0.05	0.67	0.91

Table 5

Segmentation accuracy comparison for 4 volumes, 24 vessel slices and 72 vessels.

Volume #	Image Size ($w \times h$)	# of Vessel Slices	Total # of Vessels	Chan-Vese [11]		Snake model [8]		Proposed SRG	
				Average Dice	Average Precision	Average Dice	Average Precision	Average Dice	Average Precision
10504	460 × 306	7	21	0.78	0.82	0.77	0.81	0.85	0.83
18152	512 × 512	5	15	0.75	0.85	0.78	0.83	0.82	0.87
23186	405 × 346	6	18	0.77	0.83	0.81	0.82	0.84	0.86
28059	462 × 321	6	18	0.72	0.81	0.76	0.78	0.81	0.84

clinician wants to have more accuracy. The advantage of k is, by adjusting the value of k from the threshold criteria, clinicians have the flexibility to fine tune the accuracy of the vessel segmentation. But, vessels vary in shapes, sizes, texture features etc not only in different CT slices but also in the same CT slice. Even if

the authors propose $k = 0.05$ provides better accuracy for the provided CT images, sometimes finding the same value of k for different vessels in the same CT slice becomes difficult. The range of k lies between 0.03 to 0.12 for the better accuracy of vessel segmentation.

5.3. Discussion

In this paper, we propose persistence and grid-stride loop based SRG implementation. In order to obtain significant speedup, we need to exploit parallelism by using persistence and IBS. It involves change in the large body of SRG algorithm. We want algorithms that require as less synchronization as possible. In general if algorithm requires IBS, it is probably not going to be particularly fast. The fastest algorithms on GPUs are ones that fit nicely into the GPU programming model, where blocks are independent from each other and do not require synchronization [35].

But the problem arises when iterative calling of the kernel can not be avoided. It incurs memory transfers from CPU to GPU when KTRL is used for global synchronizations. Hence it has to go through synchronizations as the next step of SRG which is dependent on the current step. Terminating a kernel and relaunching incurs data transfers from CPU to GPU and vice versa. It is time consuming.

If we use IBS method along with persistence, then we can map whole algorithm on GPU with synchronization. Control comes back to CPU only if the kernel task is over. CPU launches a kernel on GPU, GPU executes it and final results are copied to CPU. No intermediate data communication occurs in the proposed approach (unlikely in KTRL).

6. Conclusion

In this paper, we discuss SRG based vessel segmentation and its parallel implementation on GPU. We propose persistence and grid-stride loop based GPU approach for SRG providing significant speedup. Normally recursion/iterative calling of a kernel is generally a bad idea on GPUs. We use persistence and grid-stride approach as an alternate implementation for KTRL. We compare proposed GPU optimization strategy for SRG implementation. The proposed persistent and gradient based parallel SRG for 2D vessel segmentation is accurate with high dice scores and $1.9 \times$ faster compared to the KTRL.

Declaration of Competing Interest

We wish to confirm that there are no known conflicts of interest associated with this publication and there has been no significant financial support for this work that could have influenced its outcome.

Acknowledgements

The work is supported by the project High Performance soft tissue Navigation (HiPerNav). This project has received funding from the European Union Horizon 2020 research and innovation program under grant agreement No. 722068. We thank The Intervention Centre, Oslo University Hospital, Oslo, Norway for providing the CT images with ground truths for the clinical validation of vessel segmentation.

Supplementary material

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.cmpb.2020.105430.

References

- [1] N. Satpute, R. Naseem, E. Pelanis, J. Gomez-Luna, F. Alaya Cheikh, O.J. Elle, J. Olivares, GPU acceleration of liver enhancement for tumor segmentation, *Comput. Methods Programs Biomed.* 184 (2020) 105285, doi:10.1016/j.cmpb.2019.105285.
- [2] R. Palomar, F.A. Cheikh, B. Edwin, Å. Fretland, A. Beghdadi, O.J. Elle, A novel method for planning liver resections using deformable Bézier surfaces and distance maps, *Comput. Methods Programs Biomed.* 144 (2017) 135–145.
- [3] K.K. Delibasis, A. Kechriniotis, I. Maglogiannis, A novel tool for segmenting 3d medical images based on generalized cylinders and active surfaces, *Comput. Methods Programs Biomed.* 111 (1) (2013) 148–165, doi:10.1016/j.cmpb.2013.03.009.
- [4] E. Smistad, T.L. Falch, M. Bozorgi, A.C. Elster, F. Lindseth, Medical image segmentation on GPUs—a comprehensive review, *Med. Image Anal.* 20 (1) (2015) 1–18.
- [5] J. Wassenberg, W. Middelmann, P. Sanders, An efficient parallel algorithm for graph-based image segmentation, in: *International Conference on Computer Analysis of Images and Patterns*, Springer, 2009, pp. 1003–1010.
- [6] K. Gupta, J.A. Stuart, J.D. Owens, A study of persistent threads style GPU programming for GPGPU workloads, in: *Innovative Parallel Computing—Foundations & Applications of GPU, Manycore, and Heterogeneous Systems (INPAR 2012)*, IEEE, 2012, pp. 1–14.
- [7] G. Chen, X. Shen, Free launch: optimizing GPU dynamic kernel launches through thread reuse, in: *Proceedings of the 48th International Symposium on Microarchitecture*, ACM, 2015, pp. 407–419.
- [8] S. Roy, S. Mukhopadhyay, M.K. Mishra, Enhancement of morphological snake based segmentation by imparting image attachment through scale-space continuity, *Pattern Recognit.* 48 (7) (2015) 2254–2268.
- [9] H. Zhou, X. Li, G. Schaefer, M.E. Celebi, P. Miller, Mean shift based gradient vector flow for image segmentation, *Comput. Vis. Image Underst.* 117 (9) (2013) 1004–1016.
- [10] E. Smistad, A.C. Elster, F. Lindseth, Real-time gradient vector flow on GPUs using OpenCL, *J. Real-Time Image Process.* 10 (1) (2015) 67–74.
- [11] S.K. Siri, M.V. Latte, Combined endeavor of neutrosophic set and Chan-Vese model to extract accurate liver image from ct scan, *Comput. Methods Programs Biomed.* 151 (2017) 101–109.
- [12] R.P. Kumar, F. Albrechtsen, M. Reimers, B. Edwin, T. Langø, O.J. Elle, Three-dimensional blood vessel segmentation and centerline extraction based on two-dimensional cross-section analysis, *Ann. Biomed. Eng.* 43 (5) (2015) 1223–1234.
- [13] E. Smistad, Seeded region growing, 2015, (<https://github.com/smistad/FAST/tree/master/source/FAST/Algorithms/>).
- [14] E. Smistad, A.C. Elster, F. Lindseth, GPU accelerated segmentation and centerline extraction of tubular structures from medical images, *Int. J. Comput. Assisted Radiol. Surg.* 9 (4) (2014) 561–575.
- [15] P. Harish, P.J. Narayanan, Accelerating large graph algorithms on the GPU using CUDA, in: S. Aluru, M. Parashar, R. Badrinath, V.K. Prasanna (Eds.), *High Performance Computing – HiPC 2007*, Springer Berlin Heidelberg, Berlin, Heidelberg, 2007, pp. 197–208.
- [16] X. Zhang, X. Li, Y. Feng, A medical image segmentation algorithm based on bi-directional region growing, *Optik* 126 (20) (2015) 2398–2404.
- [17] H. Jiang, B. He, D. Fang, Z. Ma, B. Yang, L. Zhang, A region growing vessel segmentation algorithm based on spectrum information, *Comput. Math. Methods Med.* 2013 (2013).
- [18] T.C. Pessoa, J. Gmys, N. Melab, F.H. de Carvalho Junior, D. Tuytens, A GPU-based backtracking algorithm for permutation combinatorial problems, in: J. Carretero, J. Garcia-Blas, R.K. Ko, P. Mueller, K. Nakano (Eds.), *Algorithms and Architectures for Parallel Processing*, Springer International Publishing, Cham, 2016, pp. 310–324.
- [19] B.A. Hechtman, A.D. Hilton, D.J. Sorin, TREES: a CPU/GPU task-parallel runtime with explicit epoch synchronization, arXiv:1608.00571. (2016).
- [20] M. Greiner, Stack implementation on programmable graphics hardware, in: *Vision, Modeling, and Visualization 2004: Proceedings*, 2004, p. 255.
- [21] V. Vineet, P.J. Narayanan, CUDA cuts: Fast graph cuts on the GPU, in: *2008 IEEE Computer Society Conference on Computer Vision and Pattern Recognition Workshops*, 2008, pp. 1–8, doi:10.1109/CVPRW.2008.4563095.
- [22] S. Xiao, W.C. Feng, Inter-block GPU communication via fast barrier synchronization, in: *2010 IEEE International Symposium on Parallel Distributed Processing (IPDPS)*, 2010, pp. 1–12, doi:10.1109/IPDPS.2010.5470477.
- [23] M. Harris, Cuda pro tip: write flexible kernels with grid-stride loops, 2015.
- [24] M. Sourouri, S.B. Baden, X. Cai, Panda: a compiler framework for concurrent CPU+GPU execution of 3d stencil computations on GPU-accelerated supercomputers, *Int. J. Parallel Program.* 45 (3) (2017) 711–729.
- [25] S. Park, J. Lee, H. Lee, J. Shin, J. Seo, K.H. Lee, Y.-G. Shin, B. Kim, Parallelized seeded region growing using CUDA, *Comput. Math. Methods Med.* 2014 (2014).
- [26] X. Tang, A. Pattnaik, H. Jiang, O. Kayiran, A. Jog, S. Pai, M. Ibrahim, M.T. Kandemir, C.R. Das, Controlled kernel launch for dynamic parallelism in GPUs, in: *2017 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2017, pp. 649–660, doi:10.1109/HPCA.2017.14.
- [27] Y. Komura, Y. Okabe, GPU-based single-cluster algorithm for the simulation of the ising model, *J. Comput. Phys.* 231 (4) (2012) 1209–1215.
- [28] T. Sorensen, H. Evrard, A.F. Donaldson, Cooperative kernels: GPU multitasking for blocking algorithms, in: *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*, ACM, 2017, pp. 431–441.
- [29] G. Rai, T. Nair, Gradient based seeded region grow method for ct angiographic image segmentation, arXiv:1001.3735. (2010).
- [30] J.E. Stone, D. Gohara, G. Shi, OpenCL: a parallel programming standard for heterogeneous computing systems, *IEEE Des. Test* 12 (3) (2010) 66–73, doi:10.1109/MCSE.2010.69.

- [31] Å.A. Fretland, V.J. Dagenborg, G.M.W. Bjørnelv, A.M. Kazaryan, R. Kristiansen, M.W. Fagerland, J. Hausken, T.I. Tønnessen, A. Abildgaard, L. Barkhatov, et al., Laparoscopic versus open resection for colorectal liver metastases, *Ann. Surg.* 267 (2) (2018) 199–207.
- [32] R. Naseem, F.A. Cheikh, A. Beghdadi, O.J. Elle, F. Lindseth, Cross modality guided liver image enhancement of CT using MRI, in: *2019 8th European Workshop on Visual Information Processing (EUVIP)*, IEEE, 2019, pp. 46–51.
- [33] Y. Zhao, H. Li, S. Wan, A. Sekuboyina, X. Hu, G. Tetteh, M. Piraud, B. Menze, Knowledge-aided convolutional neural network for small organ segmentation, *IEEE J. Biomed. Health Inf.* 23 (4) (2019) 1363–1373, doi:10.1109/JBHI.2019.2891526.
- [34] D.N.H. Thanh, D. Sergey, V.B.S. Prasath, N.H. Hai, Blood vessels segmentation method for retinal fundus images based on adaptive principal curvature and image derivative operators, *Int. Arch. Photogramm. Remote Sens. Spatial Inf. Sci. XLII-2/W12* (2019) 211–218, doi:10.5194/isprs-archives-xlii-2-w12-211-2019.
- [35] O. Zachariadis, A. Teatini, N. Satpute, J. Gómez-Luna, O. Mutlu, O. Jakob Elle, J. Olivares, Accelerating B-spline Interpolation on GPUs: Application to Medical Image Registration, *Comput. Methods Programs Biomed.* (2020), doi:10.1016/j.cmpb.2020.105431.