

# A fast three-dimensional gamma evaluation using a GPU utilizing texture memory for on-the-fly interpolations

Lucas C. G. G. Persoon, Mark Podesta, Wouter J. C. van Elmpt, Sebastiaan M. J. J. G. Nijsten, and Frank Verhaegen<sup>a)</sup>

Department of Radiation Oncology (MAASTRO), GROW—School for Oncology and Developmental Biology, Maastricht University Medical Centre, Maastricht, The Netherlands

(Received 18 November 2010; revised May 2011; accepted for publication May 2011; published 20 June 2011)

**Purpose:** A widely accepted method to quantify differences in dose distributions is the gamma ( $\gamma$ ) evaluation. Currently, almost all  $\gamma$  implementations utilize the central processing unit (CPU). Recently, the graphics processing unit (GPU) has become a powerful platform for specific computing tasks. In this study, we describe the implementation of a 3D  $\gamma$  evaluation using a GPU to improve calculation time.

**Methods:** The  $\gamma$  evaluation algorithm was implemented on an NVIDIA Tesla C2050 GPU using the compute unified device architecture (CUDA). First, several cubic virtual phantoms were simulated. These phantoms were tested with varying dose cube sizes and set-ups, introducing artificial dose differences. Second, to show applicability in clinical practice, five patient cases have been evaluated using the 3D dose distribution from a treatment planning system as the reference and the delivered dose determined during treatment as the comparison. A calculation time comparison between the CPU and GPU was made with varying thread-block sizes including the option of using texture or global memory.

**Results:** A GPU over CPU speed-up of  $66 \pm 12$  was achieved for the virtual phantoms. For the patient cases, a speed-up of  $57 \pm 15$  using the GPU was obtained. A thread-block size of  $16 \times 16$  performed best in all cases. The use of texture memory improved the total calculation time, especially when interpolation was applied. Differences between the CPU and GPU  $\gamma$ s were negligible.

**Conclusions:** The GPU and its features, such as texture memory, decreased the calculation time for  $\gamma$  evaluations considerably without loss of accuracy. © 2011 American Association of Physicists in Medicine. [DOI: 10.1118/1.3595114]

Key words: radiotherapy, dose verification, gamma evaluation, 3D dosimetry, graphics processing unit (GPU), calculation speed, texture memory, interpolation

## I. INTRODUCTION

In radiotherapy, dose verification prior to and during treatment has become important due to complex dose delivery techniques. To compare measured to planned dose distributions, the widely accepted gamma ( $\gamma$ ) evaluation method is frequently used.<sup>1,2</sup> The  $\gamma$  evaluation is a method to quantify differences in dose distributions. A major disadvantage is that the method is slow.

Different  $\gamma$  implementations have been proposed in the literature<sup>3–6</sup> with the fastest calculation times reported of 0.6 s for cube sizes typically used in clinic. Predominantly, the computer's central processing unit (CPU) is used for the implementation. Recently, the graphics processing unit (GPU) has developed into a specialized device for specific computing tasks. The GPU is a powerful, massively parallel programmable architecture,<sup>7</sup> which can be used to decrease the calculation time for many applications. In radiotherapy, the GPU has already been used in a wide range of applications (e.g., ray-tracing, dose calculation, and CT reconstruction techniques)<sup>8–14</sup> reporting speed-up ratios from 6 to 908.

In this study, we implement a 3D  $\gamma$  evaluation on a GPU to reduce calculation time.

## II. METHODS AND MATERIALS

### II.A. Algorithm

The  $\gamma$  evaluation<sup>1,2</sup> combines a dose-difference and distance-to-agreement criterion for the comparison of dose distributions. A comparison dose ( $D_c(r_c)$ ) is compared with a reference dose ( $D_r(r_r)$ ). The dose difference is defined by  $\delta(r_c, r_r) = D_c(r_c) - D_r(r_r)$  and the spatial distance by  $r(r_c, r_r) = |r_c - r_r|$ . For all spatial grid point combinations, a function  $\Gamma(r_c, r_r)$  is calculated by

$$\Gamma(r_c, r_r) = \sqrt{((r^2(r_c, r_r)/\Delta d^2) + (\delta^2(r_c, r_r)/\Delta D^2))} \quad (1)$$

$\Delta D$  and  $\Delta d$  are the acceptance criteria.  $\Delta D$  being a percentage of the maximum reference dose and  $\Delta d$  being the distance to agreement. Then, for all points  $r$  within a spatial search volume  $\dot{U}$ , the  $\gamma$  is calculated by

$$\gamma(r) = \min\{\Gamma(r_c, r_r)\} \forall \{r_c \in \dot{U}\} \quad (2)$$

In our method, we apply a spherical volume  $\dot{U}$  of  $1 \times 1 \times 1$  cm<sup>3</sup> centered around  $r_r$  according to Wendling *et al.*,<sup>4</sup> and their early stopping criterion is applied, which calculates  $\Gamma(r_c, r_r)$  starting from the smallest  $r^2(r_c, r_r)$  to the largest

TABLE I. Patient characteristics.

Case	Reference dose cube size (voxels)	Zero padded dose cube size (voxels)	Treatment site	Treatment plan type	No. beams
A	166 × 124 × 88	192 × 128 × 88	Head and neck	3D CRT	4
B	139 × 131 × 88	160 × 160 × 88	Lung	SBRT	10
C	148 × 141 × 79	160 × 160 × 79	Lung	3D CRT	4
D	164 × 134 × 88	192 × 160 × 88	Head and neck	IMRT	7
E	141 × 133 × 78	160 × 160 × 78	Head and neck	IMRT	7

possible  $r^2(r_c, r_r)$  and stops when  $\sqrt{(r^2(r_c, r_r)/\Delta d^2)}$  exceeds the smallest found  $\Gamma(r_c, r_r)$ .

For this study, global dose differences with respect to the maximum  $D_r$  with criteria of  $\Delta D = 3\%$  and  $\Delta d = 3$  mm were used.

## II.B. GPU computing

An NVIDIA Tesla C2050 GPU (NVIDIA, Santa Clara, CA, USA) was used with double precision compute capability 2.0, 3 GB GDDR5 integrated memory, and NVIDIA's Compute Unified Device Architecture (CUDA) version 3.2 of the SDK. CUDA uses a parallel programming model and extends the C programming language by allowing the definition of kernels that are executed  $N$  times in parallel by  $N$  CUDA threads on the GPU.<sup>7</sup> To execute kernels, there are two important parameters to define: the thread-block size, a vector of three elements  $(x, y, z)$ , and the grid-block size, which is a vector of two elements  $(x, y)$ . The thread-block size cannot exceed a limit of 1024 threads for the TESLA C2050 and is hardware-specific. The thread-block size is the number of simultaneous launched threads. The total number of threads to be executed is the same as the product of the thread-block size and grid-block size. In our implementation, we executed for each voxel a thread keeping the product of the thread and grid block equal to the number of voxels. Memory access was conducted either using texture memory with the hardware built-in 3D texture functions to fetch voxels and immediate interpolate or by using global memory and an equivalent on-the-fly interpolation.<sup>7</sup> Texture memory is special hardware built on-top of the global memory to access data. CUDA supports a subset of the texturing hardware capabilities normally used for graphics and includes caching and interpolation (e.g., trilinear).

## II.C. CPU computing

The CPU simulations were executed on an Intel Xeon X5550 2.66 GHz quad-core CPU with hyperthreading. CPU and GPU calculations were implemented in C/C++ using Microsoft Visual Studio 2008 (Microsoft, Redmond, WA, USA). For benchmarking, single core performance was investigated, while multicore was also implemented using the parallel programming API OpenMP. Furthermore, the on-the-fly interpolation was implemented with an equivalent algorithm as implemented by NVIDIA.<sup>7</sup>

## II.D. Benchmarking and examples

Several cubic virtual phantoms were simulated.<sup>4</sup> The virtual phantom test is a setup with two dose cubes: a centered

one and a shifted one with magnitude changes. These phantoms were tested with varying cube sizes with a voxel size of  $0.1 \times 0.1 \times 0.1$  cm<sup>3</sup>. No interpolation of the comparison dose cubes ( $D_c$ ) was applied. The dose of the centered volume was set to a value of 100 cGy while outside the central volume it was set to 0 cGy. For benchmarking, three artificial errors were introduced in  $D_c$ .

Next, to show applicability in clinical practice, five patient cases (Table I) were evaluated using the 3D dose distribution from a treatment planning system as a reference  $D_r$  and the delivered dose measured during treatment. The measured dose cube  $D_c$  was on-the-fly interpolated to a grid size three times finer than the original. The  $D_r$  and  $D_c$  are, when necessary, zero padded to a multiple of 32 pixels in a slice.

To investigate the sensitivity of different thread-block sizes and the performance of texture memory, varying thread-block sizes were tested (i.e.,  $n \times n$  with  $n = 8, 16, 32$ ) and all calculations were executed using global memory and texture memory. The accuracy between the CPU and GPU implementations were evaluated by comparing several metrics.

## III. RESULTS

Table II shows the GPU over CPU speed-up for the virtual phantoms. The speed-up ranged from 49 to 82 times. Furthermore, it can be seen that the thread-block size of  $16 \times 16$  is optimal for all phantoms and shows the highest speed-up of  $66 \pm 12$ . The calculation time for the virtual phantoms using texture memory management decreased on average by 5%. No differences in accuracy between CPU and GPU results exceeding  $2 \cdot 10^{-7}\%$  were observed. For the patient cases, the speed-up ranged from 28 to 82 times. It was observed that the thread-block size of  $16 \times 16$  was also optimal with a speed-up of  $57 \pm 15$ . Using textures for the interpolation added an average performance gain of 30% in comparison with the GPU implementation using global memory and the on-the-fly interpolation. No differences exceeding  $9.1 \cdot 10^{-4}\%$  were observed.

## IV. DISCUSSION

In this work, we showed that the use of GPU architecture can speed-up the calculation of 3D  $\gamma$  evaluations considerably compared to CPU calculations. In our implementation, the interpolation step is executed on the GPU using texture memory. This is one of the major advantages of using GPUs to accelerate the calculation of 3D  $\gamma$ s. The benefit of using GPUs for the  $\gamma$  evaluation will increase further when new

TABLE II. Calculation time and speed-ups for the different virtual phantoms and patient cases. The introduced error is a shift and dose change in comparison with the reference dose. The different errors are given in parentheses for the  $128^3$  phantoms. The shift is introduced in each direction ( $x,y,z$ ). The calculation time is based on an average of five runs of the application. The standard deviation for the GPU kernel execution time was on average 2.5%, while for the CPU the average standard deviation was 1.5%. For the patient cases and  $384^3$  phantom the speed-up without using the texture memory is given in parentheses.

	CPU			GPU			Speed-up		
	Average time (s)	Thread-block size			Thread-block size				
		$8 \times 8$	$16 \times 16$	$32 \times 32$	$8 \times 8$	$16 \times 16$	$32 \times 32$		
Phantom $128^3$ , $5 \text{ cm}^3$ centered dose volume									
I (0.1 cm, 1%)	0.24	0.005	0.005	0.005	<b>49</b>	<b>49</b>	<b>49</b>		
II (0.2 cm, 3%)	0.80	0.140	0.010	0.011	<b>50</b>	<b>70</b>	<b>63</b>		
III (0.3 cm, 5%)	1.52	0.023	0.017	0.018	<b>65</b>	<b>89</b>	<b>84</b>		
Phantom $192^3$ , $7.5 \text{ cm}^3$ centered dose volume									
I	0.81	0.015	0.014	0.015	<b>54</b>	<b>58</b>	<b>54</b>		
II	1.95	0.042	0.037	0.039	<b>46</b>	<b>53</b>	<b>50</b>		
III	7.99	0.133	0.116	0.120	<b>60</b>	<b>69</b>	<b>67</b>		
Phantom $256^3$ , $10 \text{ cm}^3$ centered dose volume									
I	1.92	0.034	0.032	0.033	<b>56</b>	<b>59</b>	<b>58</b>		
II	5.74	0.085	0.084	0.090	<b>67</b>	<b>68</b>	<b>63</b>		
III	16.55	0.210	0.209	0.220	<b>79</b>	<b>79</b>	<b>75</b>		
Phantom $384^3$ , $15 \text{ cm}^3$ centered dose volume									
I	6.59	0.122	0.120	0.135	<b>54 (52)</b>	<b>55 (53)</b>	<b>48 (48)</b>		
II	17.07	0.257	0.241	0.277	<b>66 (61)</b>	<b>71 (67)</b>	<b>61 (58)</b>		
III	59.94	0.800	0.730	0.761	<b>74 (62)</b>	<b>82 (72)</b>	<b>78 (63)</b>		
Case									
A	0.36	0.009	0.009	0.011	<b>41 (28)</b>	<b>41 (30)</b>	<b>33 (28)</b>		
B	0.82	0.011	0.010	0.012	<b>74 (51)</b>	<b>82 (63)</b>	<b>68 (41)</b>		
C	0.59	0.012	0.010	0.015	<b>49 (46)</b>	<b>59 (50)</b>	<b>40 (35)</b>		
D	0.67	0.018	0.015	0.024	<b>37 (29)</b>	<b>44 (34)</b>	<b>28 (24)</b>		
E	0.70	0.015	0.013	0.021	<b>47 (33)</b>	<b>54 (39)</b>	<b>33 (25)</b>		

graphics hardware is introduced because GPUs evolve much faster than CPUs.<sup>7</sup>

There were a few considerations when choosing from the methods of speed increase available on the GPU architecture. The data volume sizes and search box sizes used in  $\gamma$  calculations exclude the option of effectively using the GPU shared memory to its potential. Coalesced global memory access is desirable and would offer the next best speed increase but much care must be taken to ensure a simple organized access pattern of the stored array. While an effective manipulate/copy/execute structure may be possible, it would be difficult to maintain as both the input volume resolution as well as the  $\gamma$  criterion vary from run to run. We decided to use texture fetching as it offers a good performance on memory reads that do not follow the access patterns required by global memory and has a higher bandwidth if there is locality to concurrent access patterns. This is suited to the style of a  $\gamma$  search box.<sup>7</sup> The memory access pattern is less linear and more local due to the specific ordered geographical search box. On the FERMI architecture, there is also the option of using the Level 1 (L1) or Level 2 (L2) cache which both use error checking correction (ECC) in contrast to texture memory, which is non-ECC. Using texture memory has the benefit of hardware accelerated trilinear interpolation, which is processed efficiently on a GPU, reducing the need to develop a fast interpolation or divided simplex function. This means a com-

paratively fast implementation of interpolation using texture memory is possible while keeping decent performance. We observed that using the texture memory is beneficial for the overall calculation time, on average 5% faster for the virtual phantoms and 30% faster for the patient cases. For the patient cases, the advantage of using texture memory versus global memory was higher because interpolation was executed on-the-fly.

A thread-block size of  $16 \times 16$  resulted in an optimal calculation time for all cases. NVIDIA provides several optimization tools; the *visual profiler* and the *occupancy calculator*. Occupancy can be seen as the number of concurrent threads in use and is a good measure of GPU efficiency. It was found that the thread-block size of  $8 \times 8$  has the lowest occupancy of 33%. This is because we use 30 registers per thread. The other two tested thread-block sizes have an occupancy of 66%. The best performance of  $16 \times 16$  can be explained because the number of thread-blocks per multiprocessor is six in comparison with only one thread-block for  $32 \times 32$  thread-block size. The  $16 \times 16$  thread-block size and the six thread-blocks per multiprocessor optimizes the number of warps (concurrently launched threads) processed on a multiprocessor and is, in this case, better able to fill the latency around expensive operations resulting in better overall performance.

The software and source-code using a fixed thread-block size of  $16 \times 16$  is freely available for non-commercial research

use. On the website<sup>15</sup>, both a separate C/C++ program as well as a MATLAB mex-file are downloadable. All compiled packages support only the windows platform.

## V. CONCLUSIONS

In this work, fast  $\gamma$  evaluation was implemented on GPU architecture. Comparison of large 3D dose matrices can be done in subsecond calculation times, with no loss of accuracy compared to CPU implementations. Further speed-ups may be realized by harnessing clusters of GPUs.

<sup>a)</sup>Author to whom correspondence should be addressed. Electronic mail: frank.verhaegen@maastro.nl; Telephone: +31-88-4455666; Fax: +31-88-4455667.

<sup>1</sup>D. A. Low, W. B. Harms, S. Mutic, and J. A. Purdy, "A technique for the quantitative evaluation of dose distributions," *Med. Phys.* **25**, 656–661 (1998).

<sup>2</sup>D. A. Low and J. F. Dempsey, "Evaluation of the gamma dose distribution comparison method," *Med. Phys.* **30**, 2455–2464 (2003).

<sup>3</sup>T. Ju, T. Simpson, J. O. Deasy, and D. A. Low, "Geometric interpretation of the gamma dose distribution comparison technique: interpolation-free calculation," *Med. Phys.* **35**, 879–887 (2008).

<sup>4</sup>M. Wendling, L. J. Zijp, L. N. McDermott, E. J. Smit, J. J. Sonke, B. J. Mijnheer, and M. van Herk, "A fast algorithm for gamma evaluation in 3D," *Med. Phys.* **34**, 1647–1654 (2007).

<sup>5</sup>J. Yuan and W. Chen, "A gamma dose distribution evaluation technique using the k-d tree for nearest neighbor searching," *Med. Phys.* **37**, 4868–4873 (2010).

<sup>6</sup>X. Gu, X. Jia, and S. B. Jiang, "GPU-based fast gamma index calculation," *Phys. Med. Biol.* **56**, 1431–1441 (2011).

<sup>7</sup>NVIDIA, "NVIDIA CUDA Programming Guide v3.2," NVIDIA, 2010.

<sup>8</sup>M. de Greef, J. Crezee, J. C. van Eijk, R. Pool, and A. Bel, "Accelerated ray tracing for radiotherapy dose calculations on a GPU," *Med. Phys.* **36**, 4095–4102 (2009).

<sup>9</sup>F. Xu, W. Xu, M. Jones, B. Keszthelyi, J. Sedat, D. Agard, and K. Mueller, "On the efficiency of iterative ordered subset reconstruction algorithms for acceleration on GPUs," *Comput. Meth. Progr. Biomed.* **98**, 261–270 (2009).

<sup>10</sup>R. Jacques, R. Taylor, J. Wong, and T. McNutt, "Towards real-time radiation therapy: GPU accelerated superposition/convolution," *Comput. Meth. Progr. Biomed.* **98**, 285–292 (2009).

<sup>11</sup>C. Men, X. Jia, and S. B. Jiang, "GPU-based ultra-fast direct aperture optimization for online adaptive radiation therapy," *Phys. Med. Biol.* **55**, 4309–4319 (2010).

<sup>12</sup>S. Hissoiny, B. Ozell, and P. Després, "A convolution-superposition dose calculation engine for GPUs," *Med. Phys.* **37**, 1029–1037 (2010).

<sup>13</sup>S. S. Samant, J. Xia, P. Muyan-Ozcelik, and J. D. Owens, "High performance computing for deformable image registration: towards a new paradigm in adaptive radiotherapy," *Med. Phys.* **35**, 3546–3553 (2008).

<sup>14</sup>P. Després, J. Rinkel, B. Hasegawa, and S. Prevrhal, "Stream processors: A new platform for Monte Carlo calculations," in *Third McGill Workshop on Monte Carlo Techniques in Radiotherapy Delivery and Verification*, edited by F. Verhaegen, *J. Phys.: Conf. Ser.* **102**, 012007 (2008).

<sup>15</sup><http://www.MISTIR.info>