



High Speed 3D Tomography on CPU, GPU, and FPGA

Nicolas Gac, Stéphane Mancini, Michel Desvignes, Dominique Houzet

► To cite this version:

Nicolas Gac, Stéphane Mancini, Michel Desvignes, Dominique Houzet. High Speed 3D Tomography on CPU, GPU, and FPGA. EURASIP Journal on Embedded Systems, 2008, 2008, <http://www.hindawi.com/GetArticle.aspx?doi=10.1155/2008/930250>. 10.1155/2008/930250 . hal-00367321

HAL Id: hal-00367321

<https://hal.science/hal-00367321>

Submitted on 13 Jan 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

High Speed 3D Tomography on CPU, GPU and FPGA

Nicolas GAC^{a,b}, Stéphane MANCINI^a, Michel DESVIGNES^a and Dominique HOUZET^a

Abstract—Back-Projection (BP) is a costly computational step in tomography image reconstruction such as Positron Emission Tomography (PET). To reduce the computation time, this paper presents a Pipelined, Pre-fetch and Parallelized Architecture for PET BP (3PA-PET). The key feature of this architecture is its original memory access strategy, masking the high latency of the external memory. Indeed, the pattern of the memory references to the data acquired hinder the processing unit.

The memory access bottleneck is overcome by an efficient use of the intrinsic temporal and spatial locality of the BP algorithm. A loop re-ordering allows an efficient use of general purpose processor's caches, for software implementation, as well as the 3D Predictive and Adaptive Cache (3D-AP Cache), when considering hardware implementations. Parallel hardware pipelines are also efficient thanks to a hierarchical 3D-AP Cache: each pipeline performs a memory reference in about one clock cycle to reach a computational throughput close to 100%.

The 3PA-PET architecture is prototyped on a System on Programmable Chip (SoPC) to validate the system and to measure its expected performances. Time performances are compared with a desktop PC, a workstation and a GPU (Graphic Processor Unit).

I. INTRODUCTION

Tomography consists of reconstructing an object from its projections via numerical methods [1]. This process is used in medical scanners, such as Computed Tomography (CT) or Positron Emission Tomography (PET) scanners. PET is a nuclear imaging modality; its goal is to measure the spatial and temporal distribution of a radio-tracer perfused in a patient's body. PET imaging is used in oncology, to detect, track and visualize tumors. After data acquisition, the 3D image of the radio-tracer

is reconstructed off line from the measures (called sinograms) to diagnose pathologies. Oncology and other clinical applications need a high quality reconstruction as fast as possible (few minutes at most) to reduce the device occupation and allow a patient repositioning¹. Also, dynamic PET is in need of even faster reconstruction.

Moreover, tomography is required in many other medical imaging techniques, such as 3D Magnetic Resonance Imaging and 3D Ultra-sound Imaging, or in other domains such as Synthetic Aperture Radar (SAR), contact-less control and industrial X-Ray applications. Therefore, the acceleration of the reconstruction algorithm is of great interest for various applications.

Due to the large amount of the acquired data and the complexity of the algorithms, reconstruction is a very time-consuming process. From a computing point of view, reconstruction methods can be classified into two main techniques: analytic (direct) reconstruction and iterative reconstruction. They both include a Back-Projection (BP) step that accounts for 50% to 70% of the processing time.

In 3D reconstruction, the computational complexity of the standard algorithm to reconstruct an N^3 data-set from N angles of projection is $O(N^4)$. In the previous decade, several algorithms have been proposed to reduce BP complexity. The lowest cost obtained is $O(N^3 \log N)$ but generally with a lower quality of reconstruction; also it doesn't take into account some required data management, which delay the process. Although CPUs have gained sufficient computing power for 2D reconstruction, with 3D reconstruction the increase of the amount of data for high quality images leads to higher computing times. Iterative reconstruction algorithms may reach several hours of processing [2].

a : GIPSA-lab, Grenoble Institute of Technology - INPG, BP 46, 38402 St Martin d'Hères, France

b : ETIS, CNRS, ENSEA, Univ Cergy-Pontoise, F-95000 Cergy-Pontoise, France

¹A patient can not experienced a radio-tracer twice in a short while and has to wait several months before a new examination, in case of bad camera positioning.

The algorithmic optimizations of reconstruction have reached some limits and it is becoming mandatory to reduce the computing time through architecture solutions. General purpose parallel computers benefit from recent competing technologies: the System on Programmable Chip (SoPC) and the GP-GPU (General Purpose Graphical Processing Unit).

This paper shows that a hardware implementation of the BP algorithm needs to overcome the memory bottleneck. This may be solved both by a loop reordering and the use of an efficient caching mechanism. Parallel hardware pipelines can be fed with a hierarchy of semi-general purpose cache such as the 3D-AP Cache [3]. The resulting architecture makes a better use of memory bandwidth than general purpose CPUs and GP-GPU.

The first parts of this paper present the use of the 3D BP algorithm and different solutions to accelerate it. Next, we present the memory bottleneck of a classical implementation of 3D BP to overcome. From this study, an efficient architecture is proposed: the Pipelined, Pre-fetch and Parallelized Architecture for 3D PET BP (3PA-PET). The quality, complexity and timing performance of the 3PA-PET architecture are also presented. Measures on its prototyping on a SoPC allows a comparison with the implementation of BP on CPU and GP-GPU.

II. 3D BP IN TOMOGRAPHY RECONSTRUCTION

In this section, we will first show that 3D PET BP and the 3D CT BP using respectively a parallel and a cone beam geometry, are close algorithms. Then, we will present some related works on acceleration of these two BPs on several architectures.

A. BP algorithms

1) *3D parallel beam BP for PET*: The detectors of a PET scanner are usually paving a cylinder and stacked in a set of rings of detectors [1]. The γ rays issued from the disintegration of a radio-tracer particle, are detected by a pair of sensors facing each others. The line which connects two sensors is called a Line Of Response (LOR) and the coincidence events counted on one LOR are stored in a *bin*. All the bins are stored in a sinogram as illustrated in figure 1. The reconstruction process attempts to estimate the image of the radio-tracer distribution f that has produced the sinogram.

The sinogram p_{PET} is a 4D space along $(\Delta, \psi, u_{\parallel}, v_{\parallel})$. The coordinates (Δ, v_{\parallel}) represent a

couple of rings : Δ is the axial distance between the two rings (segment number) and v_{\parallel} is the mean axial coordinate of the two rings (plane number). The coordinates (ψ, u_{\parallel}) represent one particular LOR between two rings: ψ is the azimuthal angle and u_{\parallel} is the tangential coordinate (bin number).

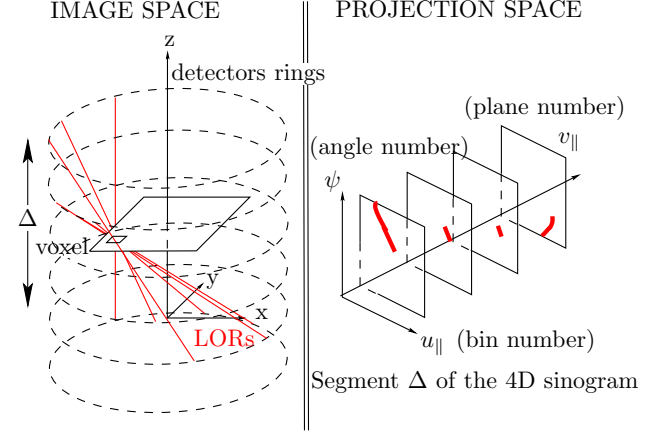


Figure 1. The acquired data are stored in a 4D sinogram, a sinogram bin corresponding to one particular LOR. To reconstruct one voxel, the data needed by the BP algorithm draw a 3D sinusoid in each segment Δ .

For each voxel (VOLUME piXEL) of coordinate $\vec{r} = (x, y, z)$, the BP algorithm sums up all the sinogram bins corresponding to that voxel projection:

$$f_{\text{PET}}^*(\vec{r}) = \iiint p_{\text{PET}}(\Delta, \psi, u_{\parallel}(\psi, \vec{r}), v_{\parallel}(\psi, \Delta, \vec{r})) J_{\Delta} d\psi d\Delta \quad (1)$$

J_{Δ} is a jacobian and the parallel beam coordinates $(u_{\parallel}, v_{\parallel})$ are computed as:

$$\begin{cases} u_{\parallel}(\psi, \vec{r}) &= x \cos \psi + y \sin \psi + \text{offset} \\ v_{\parallel}(\psi, \Delta, \vec{r}) &= \frac{\Delta}{2R_a} \cdot (x \sin \psi - y \cos \psi) + z + \text{offset} \end{cases} \quad (2)$$

Using $a_{ij}(\psi, \Delta)$ coefficients, we get:

$$\begin{cases} u_{\parallel}(\psi, \vec{r}) &= a_{00}x + a_{01}y + a_{03} \\ v_{\parallel}(\psi, \Delta, \vec{r}) &= a_{10}x + a_{11}y + a_{12}z + a_{13} \end{cases} \quad (3)$$

2) *Cone beam BP for CT*: The cone beam BP used in CT imaging modalities uses a similar algorithm [4]. In CT, the data is the X ray intensity reaching an X camera that rotates around

the observed volume. It measures the attenuation due to the density of tissues. The density $f_{\text{CT}}^*(\vec{r})$ is computed from the measured data p_{CT} following:

$$f_{\text{CT}}^*(\vec{r}) = \int p_{\text{CT}}(\alpha, u_{\text{V}}(\alpha, \vec{r}), v_{\text{V}}(\alpha, \vec{r})) \cdot w(\alpha, \vec{r})^2 \cdot d\alpha \quad (4)$$

Where α is the trajectory parameter of the camera. The cone beam coordinates $(u_{\text{V}}, v_{\text{V}})$ are computed as:

$$\begin{cases} u_{\text{V}}(\alpha, \vec{r}) &= (c_{00}x + c_{01}y + c_{02}z + c_{03}) \cdot w(\alpha, \vec{r}) \\ v_{\text{V}}(\alpha, \vec{r}) &= (c_{10}x + c_{11}y + c_{12}z + c_{13}) \cdot w(\alpha, \vec{r}) \end{cases} \quad (5)$$

where c_{ij} depends on α (i.e. $c_{ij} = c_{ij}(\alpha)$) and

$$w(\alpha, \vec{r}) = \frac{1}{c_{20} \cdot x + c_{21} \cdot y + c_{22} \cdot z + c_{23}} \quad (6)$$

3) *Comparison of CT and PET*: Although the CT BP is more complex due to the perspective transformation (eq 6), these algorithms are quite similar. Indeed, the summation over α (trajectory parameter) for CT BP, is equivalent to the summation over ψ and Δ for PET BP. Moreover, in these loops, both these BPs compute very similar projection coordinates $(u_{\parallel}, v_{\parallel})$ and $(u_{\text{V}}, v_{\text{V}})$. Nevertheless, the computation of the projection coordinates for CT BP needs a division by a distance weight $w(\alpha, \vec{r})$. Thus, the CT BP kernel has more arithmetic operations than has PET BP.

Supposing that one is able to design a pipeline that computes a sum update at each clock cycle, both for CT and PET BP, then the challenge is to fetch data along a complex path (a 3D sinusoid) in the acquired data (3D CT data or 4D PET sinogram). The method presented in this paper for solving the case of PET BP (parallel beam) could be transposed to solve the CT BP (cone beam).

B. Acceleration of reconstruction

Different computer architectures coupled with dedicated memory access strategies are used to accelerate the BP step of an analytic or iterative reconstruction, including: general purpose processors [5], [6], [7], graphical processors [8], [9], [10], [11], [12], [13], [14], the Cell processor [4], [15] or ASIC/FPGA architectures [16], [17], [18], [19], [20]. While most of these works have investigated cone

beam BP, only a few of them have investigated 3D parallel beam BP [2], [5], [8], [9].

The parallelisation of reconstruction algorithms on shared memory parallel general purpose computer [5] stays efficient only up to 4 processing units, because of conflictual accesses on the memory bus. Considering clusters of heterogeneous PCs [6], [7], efficiency of parallelisation drops down quickly because of the costly communication between PCs. After 10 PCs, parallelisation is not worthy. Yet on a distributed memory parallel computer, parallelisation works very well. for 3D PET iterative reconstruction, Jones *et al.* [2] succeeded to get an acceleration factor of about 30 with 32 processors. Ni *et al.* [21] achieved an excellent acceleration factor of 300, when they parallelized the Katsevitch algorithm, an exact cone beam BP with 300 CPUs.

Besides parallelisation on several nodes of general purpose processors, more efficient engines such as the GPU (Graphical Processing Unit) or the IBM Cell can be used. Current GPUs can be used either as a graphical pipeline, which is originally designed for [8], [9], [10], or as a multiprocessor chip thanks to the CUDA interface from Nvidia [10], [12], [13], [14], [15]. For both options, the acceleration factor of GPU is high, about an order of magnitude for cone beam BP. Xu *et al.* [10] have observed that an implementation of the cone beam BP using the graphic pipeline is 3 times faster than the one made with the CUDA interface. Kachelriess *et al.* [4] and Scherl *et al.* [15] present good result of acceleration of cone beam BP using the Cell processor. With its 1+8 cores, this architecture is an intermediate solution between general purpose parallel processors and GPU. The 8 vector engines have to be specifically programmed. Nevertheless, Scherl *et al.* [11] have measured that a GPU with CUDA is 3 times faster than the Cell for the BP alone.

FPGA technology is an alternative to processors, allowing designers to make a customized architecture. Most often, it is used to prototype ASIC implementations. In this context, FPGA implementations of 2D parallel beam BP [16] and 3D cone beam BP [17], [18], [19] have been investigated. These architectures are made of several pipelines working in parallel. Moreover, like the imageProX by Siemens [18], several FPGA chips can be used in a single board to raise the computational power.

Two memory access strategies have been applied

for all these architectures. In case the processor already has a memory cache, developers rely on it to optimize the external memory accesses. Otherwise, developers set up custom memory strategies in order to hide the memory access time. The most common approach is to use double buffering: the next required projection data is loaded from external memory, meanwhile the ongoing loaded projection data are back-projected. In this case, CPU and GPU memory strategies are based on an extensive use of the cache. For example Yang *et al.* [13] have observed that enabling a GPU cache is more competitive than software pre-fetching. On the other hand, the Cell and FPGA memory strategies have to be taken in charge by the software designers.

III. OVERCOMING THE MEMORY BOTTLENECK

In this section, we focus our study on finding out the best appropriate memory strategy to get the best fit between the 3D BP algorithm and a hardware architecture.

A. Memory access strategy

As the sinogram is kept in a SDRAM like external memory, we need an efficient memory management to overcome its latency and allow a high level of parallelism. The main difficulty is to deal with the high strides of addresses due to the sinusoidal pattern of references in the 4D space. A cache would help to hide the high latency of the external memory despite these strides. Standard caches therefore are inefficient as they exploit temporal and address locality of references. Hence they are used at their best when the references follow a 1D linear pattern as a cache line is loaded when a miss occurs.

Indeed, as shown earlier, the reconstruction of a single voxel $f(\vec{x})$ needs to follow a 3D sinusoid in the 4D sinogram. Such a pattern is of poor address locality but has a high index locality. Moreover, because of the $v_{||}$ dimension, the memory accesses for 3D BP have higher strides and are more distributed in the memory space than in the 2D BP case [22]. The challenge is to speed-up these memory accesses in a 4D data structure.

Therefore a new cache mechanism is needed. Estimating which bins would be referenced would help the cache to download the needed bins during the computing process.

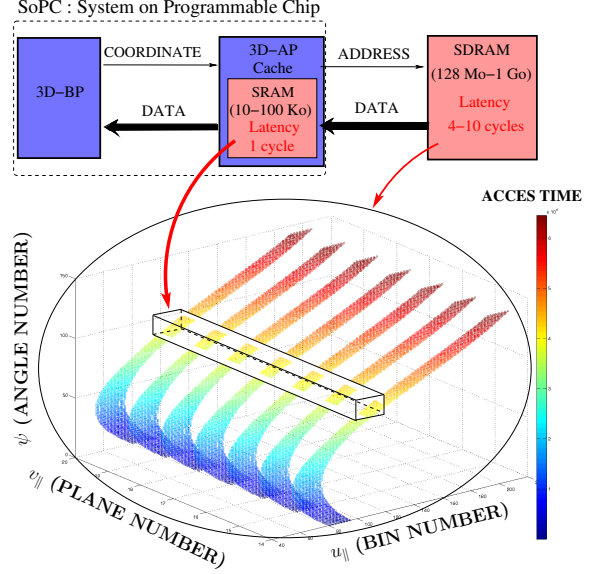


Figure 2. The memory access strategy is based on a fast and small cache memory inside the SoPC. The cache predicts the needs of the 3D BP unit and therefore succeeds to mask the high latency (4-10 cycles @200 Mhz) of the slower and bigger external SDRAM memory.

B. Improvement of spatial and temporal locality

A reconstruction with the Voxel-driven Bilinear Interpolation (VBI) standard BP is made of three loops, as described in algorithm III-B: the first loop is on the voxels \vec{r} , the second on the segments **Delta** and the third on the angles **psi**. Since voxels can be reconstructed independently, the loop on voxels can be split into two parts: one loop on blocks of voxels ($0 \dots n_{\max}$) and the other loop on the voxels of a block n ($\vec{r}_{\min}(n) \dots \vec{r}_{\max}(n)$).

Algorithm 1 The loop reordering of the 3D BP improves spatial and temporal locality

```

for  $n = 0$  to  $n_{\max}$  do
  for  $\Delta = 0$  to  $\Delta_{\max}$  do
    for  $\psi = 0$  to  $\psi_{\max}$  do
      for  $\vec{r} = \vec{r}_{\min}(n)$  to  $\vec{r}_{\max}(n)$  do
         $f(\vec{r}) += \text{bin}(\psi, \Delta, u_{||}(\psi, \vec{r}), v_{||}(\psi, \vec{x}))$ 

```

A loop reordering increases the temporal and spatial locality of memory accesses. Indeed, for given **psi** and **Delta** values, the data $\text{bin}(\psi, \vec{r})$ will be used several times for different voxels since the

projection of a 3D block of voxels is a 2D plane in the 4D space of the sinogram.

The figure 2 shows that the proposed loop re-ordering allows to cache a part of the sinogram. The BP of a block of voxels makes the references to follow a coherent 3D sinusoid in the sinogram along the time.

C. Mean Bin Reuse Rate (MBRR)

To give a theoretical estimation of the best achievable cache efficiency, the Mean Bin Reuse Rate (MBRR) is defined as the ratio between the number of bins accessed in cache memory by the processing units and the number of bins loaded in cache memory from the external memory. The ideal MBRR can be computed analytically. It depends on the shape of the block of reconstructed voxels. Figure 3, presents this optimal MBRR computed for each segment versus the size of the block.

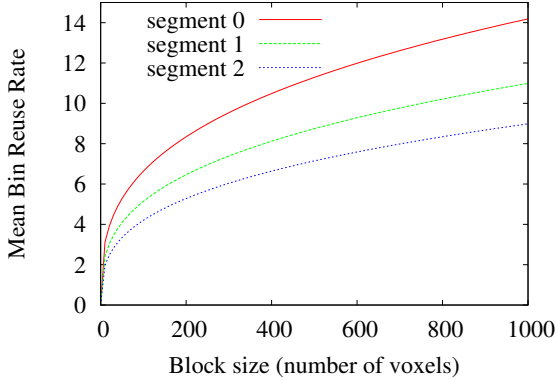


Figure 3. Mean Bin Reuse Rate (MBRR) estimated for a 3D BP without bilinear interpolation versus the size of reconstructed blocks of voxels for each segment of a Siemens HR+ sinogram (span 9 with 96 angles of projection)

IV. A 3P ARCHITECTURE FOR PET

In this section, we present the Pipelined, Pre-fetched and Parallelized Architecture for PET (3PA-PET). The 3PA-PET architecture is made of a high performance pipeline connected to a 3D Adaptive and Predictive Cache (3D-AP Cache). It allows to perform an update of a voxel value up to 1 operation per clock cycle (100% pipeline utilization), even for high latency memories.

A. Pipelined Architecture

The pipeline in figure 4, implements the different steps of the VBI standard BP: the computation of $u_{\parallel}(\mathbf{psi}, \vec{r})$ and $v_{\parallel}(\mathbf{psi}, \vec{r})$, the bilinear interpolation of the bin, and finally the accumulation of the voxel value. The forward flow control is done by packets passing through each stage of the pipeline.

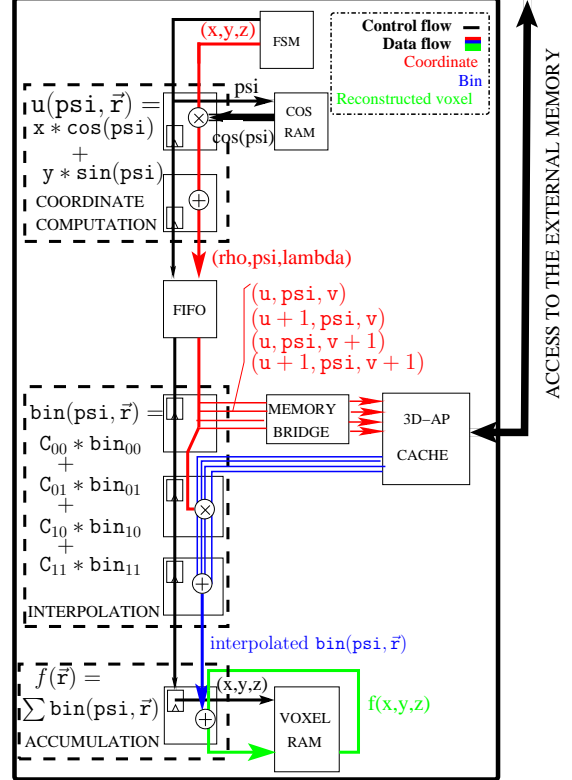


Figure 4. Pipeline of 3PA-PET

The 4 bins needed for the bilinear interpolation are fetched through the memory bridge. This bridge controls the 3D-AP Cache and can freeze (or not) the pipeline depending on the requested data availability. A backward flow control synchronizes the pipeline and the 3D-AP Cache.

B. Pre-fetch Architecture

The 3D-AP cache [3] masks the latency of the external memory so that the pipeline is no more systematically stalled. The memory bridge gets four bins from the cache at each clock cycle.

The 3D-AP Cache is a semi-generic cache memory mechanism that pre-fetches references following a continuous path into a 3D memory space. It was originally designed as a cache for a computer vision lip-tracking application [3] but it targets a large class of multi-dimensional processing algorithms. In the 3PA-PET architecture, the 3D-AP Cache is tuned to follow the references needed to reconstruct a block of voxels, as shown in figure 5. The pipeline issues spatial coordinates of the requested bin, here $(\text{psi}, u_{\parallel}, v_{\parallel})$, to the 3D-AP Cache. A part of the sinogram, namely the *cached zone*, is copied in an embedded memory. A tracking mechanism tries to maintain the center of the cached zone in the mean coordinate of the referenced data.

The 3D-AP Cache estimates dynamically which data is likely to be requested in the future. This is done by a statistical analysis on each axis of the previous references. Moreover, the 3D-AP cache masks the data transfer between the external memory and the internal cache memory. In the mean time, the cache grabs new data from the external memory, the data shared by the old and the new cached zone stay available for the processing unit.

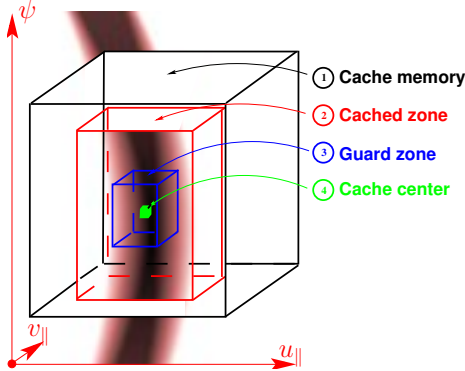
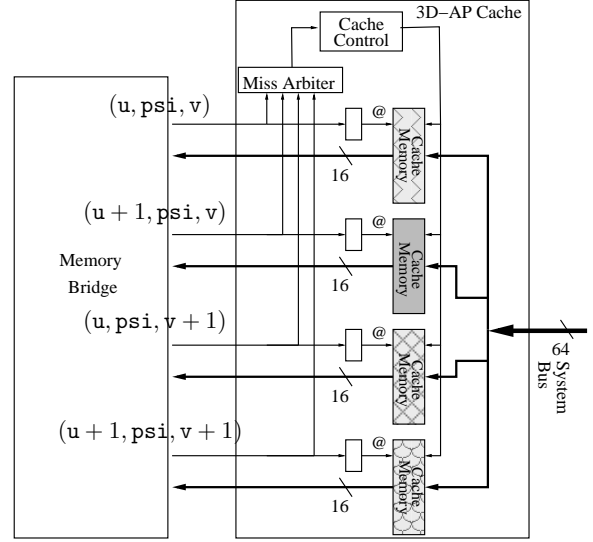


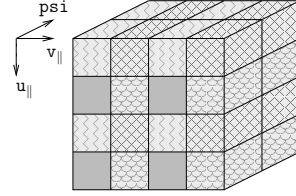
Figure 5. 3A-AP Cache zones

The 3D-AP Cache parameters have to be set beforehand by the user. In this study, we set for each dimension the value of five parameters:

- *cut-off and sampling frequencies*: the mean coordinate is computed by a first order low-pass IIR filter configured by these two frequencies.
- *cached zone size*: this zone is notified to the memory bridge to be available in cache. In this study, this size is a static parameter.



(a) Customized concurrent 3D-AP Cache architecture



(b) Mapping of bins to memory buffers

Figure 6. Memory architecture for bilinear interpolation

- *guard zone size* when the mean coordinate is out of this zone, the cache zone is updated.
- *cache speed*: it has to be set according to the speed of the data accesses performed by the application on each spatial dimension.

The cache is customized to allow four concurrent accesses to the bins needed to perform a bilinear interpolation. Figure 6 gives a simplified view of the 3D-AP cache to illustrate the involved memory architecture. The cache control unit grabs data from the external memory and splits the incoming data words to the different embedded memories. The cache control unit also manages the cache misses that could occur for some requested bins.

C. Parallelized Architecture

To increase the computing power, several pipelines are parallelized. A hierarchical cache reduces the memory bus occupation, when BP units work in parallel. In this hierarchical design,

one leaf cache is associated to one 3D BP unit while a root cache is feeding each of these leaf caches.

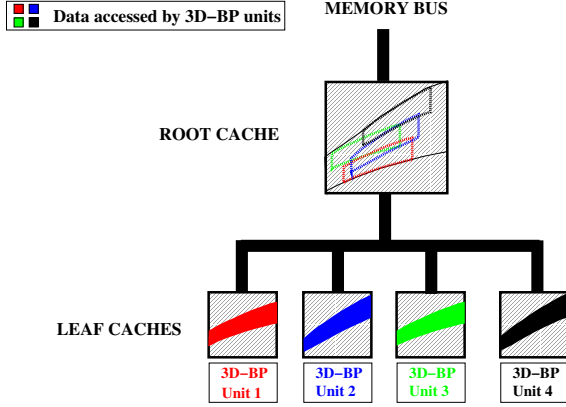


Figure 7. Each leaf cache is feeding by the root cache.

The spatial locality of the references of the pipelines is enabled by reconstructing a set of neighbor blocs. A pipeline reconstructs one bloc of voxels and all the pipeline share a loop over ψ_i . Some of the sinogram data are shared by the pipelines in the same way they are shared to reconstruct one bloc of data. The bins needed during the reconstruction of a set of bloc draw a 3D sinusoid. The cache concept presented previously with one unit, applies here in the same manner. Each leaf cache stores a 3D sinusoid needed to reconstruct a bloc. A higher level cache stores the union of these sinusoids as presented in figure 7.

V. 3PA-PET PERFORMANCES

A. Accuracy of reconstruction

The implemented VBI standard BP is a fixed point version of the original algorithm. Moreover the sinogram data is converted from float to short int (16 bits). The accuracy of reconstruction of 3PA-PET is measured between a reference reconstruction software and a software bit true model of 3PA-PET.

The reference data-set used is a sinogram of a 3D Shepp Logan volume of $128 \times 128 \times 63$ voxels. This phantom is a standard volume used in tomography to measure the accuracy of reconstruction. The sinogram is obtained from the STIR open source tool kit [23]. The volumes reconstructed by STIR and by 3PA-PET are shown on figures 8 and 9.

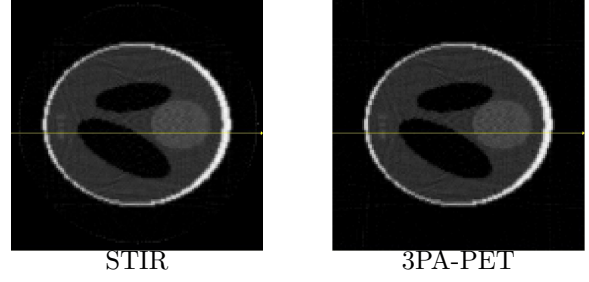


Figure 8. A slice of the 3D Shepp Logan phantom reconstructed by STIR and 3PA-PET BP.

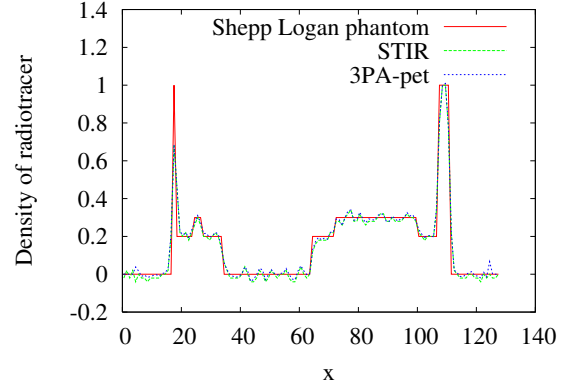


Figure 9. Profile of the 3D Shepp Logan phantom slices corresponding to the lines on figure 8

The accuracy of reconstruction of the 3PA-PET BP is measured with two metrics: the Mean Absolute Percentage Error (MAPE) and the Peak Signal-to-Noise Ratio (PSNR). Both compare a volume f_1 with a volume of reference f_{ref} . The PSNR corresponds to the ratio between the maximum of f_{ref} (dynamic range) and the mean squared error (MSE) of f_1 compared to f_{ref} :

$$PSNR = 20 \cdot \log_{10} \frac{\max(f_{\text{ref}})}{\sqrt{MSE(f_{\text{ref}}, f_1)}} \quad (7)$$

In table I, we compared the reference volume and the volumes reconstructed with STIR, with VBI floating-point arithmetic (VBI-flt) or with VBI fixed-point arithmetic (VBI-fix). All of the reconstruction methods have an intrinsic error around 3.9% with a PSNR of 10.5 dB when compared with the original volume. The floating-point and the fixed-point implementations have a MAPE of 0.13% and a PSNR of 23 dB. With different data type (*short int* versus *float*), the MAPE is about 1.1% and the PSNR of 19 dB. Thus we can conclude that

compared volumes		data	MAPE	PSNR
Accuracy of reconstruction				
STIR	/ original	float	3.89 %	10.5 dB
VBI-Flt	/ original	float	3.88 %	10.5 dB
VBI-Fix	/ original	float	3.88 %	10.5 dB
VBI-Flt	/ original	int16	3.97 %	10.5 dB
VBI-Fix	/ original	int16	3.97 %	10.5 dB
Compared reconstructions				
STIR	/ VBI-Flt	float	0.35 %	21.5 dB
VBI-Fix	/ VBI-Flt	float	0.13 %	26.2 dB
VBI-Fix	/ VBI-Flt	int16	0.13 %	23.0 dB
VBI-Fix	/ VBI-Flt	int16/flt	1.1 %	19.0 dB

Table I
ACCURACY OF RECONSTRUCTION AND COMPARED
RECONSTRUCTIONS FOR THE SHEPP LOGAN PHANTOM

the 3PA-PET implementation of the VBI BP is an accurate reconstruction system.

B. 3PA-PET complexity

The hardware resources used by the 3PA-PET architecture are presented on table II. The main BP FSM and the root cache control are shared between all of the units of the 3PA-PET architecture. Therefore the cost of an additional pipeline is only 800 slices. The sizes of a leaf and the root caches are respectively 2 KB and 18 KB. Hence, 9 BP units fit in a Xilinx Virtex 2 Pro VP30 chip and 16 units in a virtex 4 FX100.

	1 unit	4 units	9 units
<i>3D BP</i>			
CLB slices	573 (4.2%)	1817 (13.3%)	3924 (28.6%)
Multipliers	12 (9%)	48 (35%)	108 (79%)
<i>3D-AP Cache</i>			
CLB slices	672 (4.9%)	2830 (20.6%)	4804 (35.1%)
RAMs	2 kB (0.6%)	24 kB (7.8%)	36 kB (11.7%)
<i>3D BP + 3D-AP Cache</i>			
CLB slices	1245 (9.1%)	4637 (32.9%)	8728 (63.7%)

Table II
HARDWARE RESOURCES USED BY THE 3PA-PET ON A XILINX
VIRTEX 2 PRO VP30.

C. Efficiency of reconstruction

In order to assess the efficiency of the 3PA-PET architecture, we have measured the BP time on an Avnet development board connected to a PC

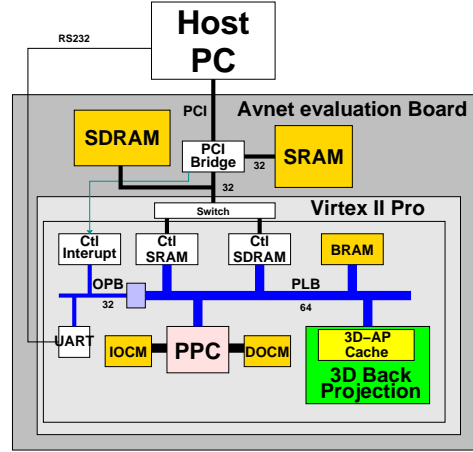


Figure 10. Evaluation system

host through a PCI interface (figure 10). The board contains an external SDRAM memory and a Xilinx SoPC (System on Programmable Chip). In order to investigate the 3PA-PET behavior with respect to the memory features, we have plugged it with a fake memory bus which could be set with different values of memory latency (l_{mem}) and memory bandwidth (BW_{mem}). The memory simulator estimates the time to access N_{line} lines of S_{line} Bytes following the relationship:

$$t_{\text{mem}} = N_{\text{line}} \cdot (l_{\text{mem}} + \frac{S_{\text{line}} - 1}{BW_{\text{mem}}}) \quad (8)$$

The times of reconstruction presented in this section are in clock cycles and scaled to one operation. An operation corresponds to one update of a voxel. The number of voxel's updates is equal to the number of voxels multiplied by the number of segments times the number of angles.

The results presented in figure 11 are achieved with one BP unit for the segment +2 which represents the worse case because the memory accesses draw the most incurvated 3D sinusoid. 3PA-PET is robust to high latencies and low bandwidth: the pipeline computes a voxel update in about 1 clock cycle, even for a memory latency of 30 cycles. This shows that the 3D-AP Cache succeeds to take advantage of the high spatial and temporal locality of the BP algorithm presented in section 3. The 3D-AP Cache follows the 3D memory path drawn during the BP process rather well. The cache miss-rate stays low (about 0.05% with $l_{\text{mem}} = 5$ cycles and

$BW_{\text{mem}} = 8 \text{ Bytes/cycle}$) which means that the 3D-AP Cache prediction is satisfactory and manages to hide the external memory latency.

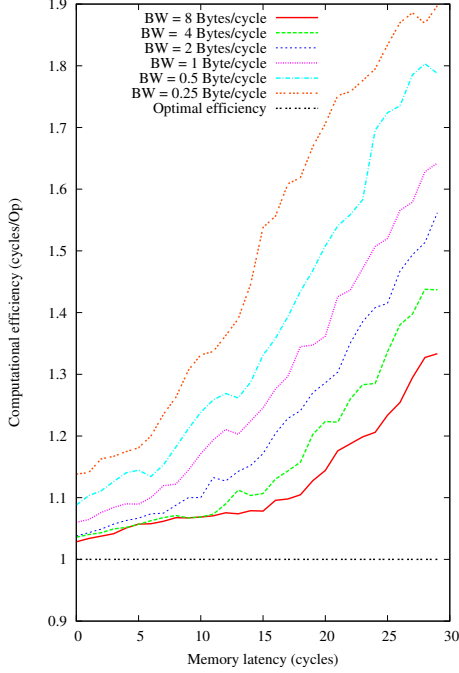


Figure 11. Cycles per operation for one unit of BP with respect to the latency and BandWidth (BW) of the external memory.

As illustrated in figure 12, the parallel 3PA-PET performances are not plenty satisfactory. Indeed, the efficiency of parallelisation decreases with the number of BP units. For instance, with a memory latency of 5 and for a complete BP, 4 units allow an acceleration of 3,2 (1.25 cycle/op per pipeline) and 8 units an acceleration of 4,7 (1.7 cycle/op per pipeline). Because the more units are working in parallel, the more busy the memory bus is. However, the hierarchical cache allows to make parallelisation a little bit more efficient thanks to the exploitation of the spatial and temporal locality existing between the data retrieved by each BP unit. Moreover, the measured MBRR for 8 units between the leaf loads and the leaf requests stay close to the MBRR measured for a single unit. This MBRR is about 8 for 8 units and 9 for 1 unit.

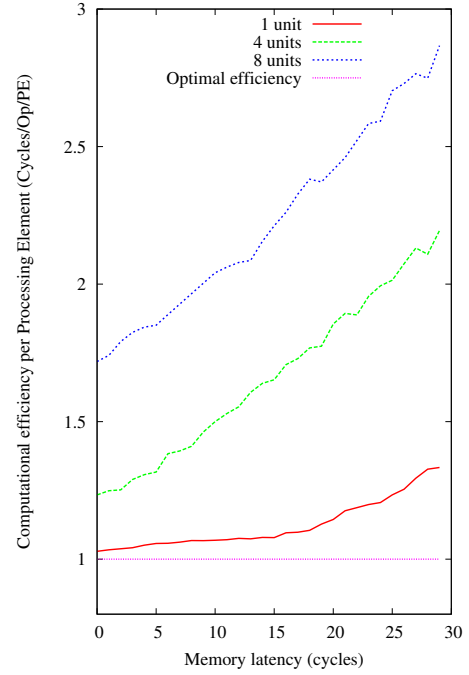


Figure 12. Cycles per operation per processing units for 1, 4 and 8 units of BP with respect to the latency and bandwidth of the external memory.

VI. COMPARISON WITH GENERAL PURPOSE AND GRAPHICS PROCESSORS

In table III, the 3PA-PET execution times are compared with STIR and the ones from software VBI BP on a desktop PC, a workstation and a GPU.

A. CPU implementation

Different software versions of BP, non optimized (v1) and optimized (v2 and v3) have been tested and compared to the STIR one on a Pentium 4 and on a bi-Xeon dual core. Two techniques of optimization have been applied with an extensive use of the cache memory and a reduction of the arithmetical operations.

First, an acceleration factor of 3 is obtained due to the reconstruction through blocks of voxels. This software loop reordering increases the use of the L1 cache (16 Ko). Indeed, the time of reconstruction with and without introduction of data locality, is respectively 54,7 s (v1) and 17,4 s (v2).

Secondly, the reduction of the arithmetic operations to compute the projection coordinates allows an acceleration by a factor 7 (software v3/software

v2 on table III). Indeed, the time performance is improved by a factor 2 due to an incremental computation of the coordinates as done by Kachelriess [4] *et al.* and again by a factor 3.5 when the inner loop is over z . The optimized code (VBI-flt(v3)) is presented on algorithm 2.

Algorithm 2 Reduction of operations to compute $u_{||}$ (same techniques used for $v_{||}$)

```

for  $n = 0$  to  $n_{\max}$  do
  for  $\delta = 0$  to  $\delta_{\max}$  do
    for  $\psi = 0$  to  $\psi_{\max}$  do
       $u_{||} = x_{n0} \cdot \cos \psi + y_{n0} \cdot \sin \psi$ 
      for  $x_n = x_{n0}$  to  $x_{n_{\max}}$  do
         $u_{||} = u_{||} + \cos \psi$ 
        for  $y_n = y_{n0}$  to  $y_{n_{\max}}$  do
           $u_{||} = u_{||} + \sin \psi$ 
          for  $z_n = z_{n0}$  to  $z_{n_{\max}}$  do
             $f(x_n, y_n, z_n) += \text{bin}(\delta, \psi, u_{||}, v_{||})$ 

```

Finally, this code has been parallelized using the *pthread* C-library to use the four cores of a bi-Xeon dual core workstation. One thread is associated to the reconstruction of one block.

B. GPU implementation

Current GPUs are cost effective solutions for the implementation of 3D tomography reconstruction because of their high level of parallelism. Moreover, the Nvidia GPUs are efficiently and easily programmed with the CUDA environment.

The Nvidia Geforce 8880 family has 2 to 16 vector processors (12 in our case), each one having 8 stream processors. It is programmable using standard C language with a few extensions without any knowledge about graphics pipeline. A non-incremental code is parallelized to run efficiently on these 12×8 multi-threaded stream processors. One thread is associated to one voxel reconstruction. Threads are grouped in blocks (16×16 in our case) which are scheduled at run-time one block per vector processor. Each couple of vector processors are associated with a 8 KB L1 2D texture read-only cache memory with 1D and 2D hard-wired interpolation. Moreover, the GPU offers a high memory bandwidth ($BW_{mem} = 64$ GB/s) and uses floating point computation. This makes it possible to efficiently parallelize the BP loops, as blocks of voxels correspond to 2D blocks of threads having

access to the read-only sinogram organized in 2D arrays through the 2D cache memory. The voxels are also organized in 2D arrays, each divided in $64 \times 16 \times 16$ blocks associated to a grid of $64 \times 16 \times 16$ blocks of threads. Thus each thread is responsible of 63 voxels considering a $63 \times 128 \times 128$ volume.

Two versions of thread code has been implemented. In the VBI-flt(v4) thread code, the loop over ψ is the inner loop, while in VBI-flt(v5) thread code, the loop over z is the inner loop as it is done for the VBI-flt(v3) CPU code. This allows a reduction of the number of projection coordinate computation. A speed-up factor of 2 is obtained with this code optimization (table III).

C. Discussion

The reconstruction times and efficiencies, global and per Processing Element (PE), are presented in table III for CPU, GPU and our 3PA-PET. To fairly compare our architecture with other technologies, the time measured on a Virtex 2 Pro has been scaled to a Virtex 4. Indeed, this technology is the same generation as the CPU and GPU used in this study. We have scaled the 35 MHz results to the GPU frequency (1.2 GHz) as well. This higher frequency could be reached through the design of a customized integrated circuit like an ASIC. Moreover, as Nvidia GTS 8800 GPU has five memory banks, we also present a prospective ASIC architecture which would have also five memory banks coupled with 5 processing blocks of 8 BP units each.

For the 200 MHz and the 1.2 GHz 3PA-PET, a memory latency (l_{mem}) of 25 ns has been used for the simulated memory bus. It corresponds respectively to a latency of 5 and 30 clock cycles.

On one hand, the GPU is the fastest hardware solution with a final reconstruction time of 50 ms. The ratio of computation over memory access is high enough for the GPU to allow the automatic overlapping of memory accesses with computations by the thread scheduling mechanism. Thus, due to its greatest computational power (96 PEs and hard-wired interpolation), Nvidia 8800GTS graphic processor is 10 times faster than 3PA-PET mapped on a Virtex 4, 10 times faster than a Xeon dual core and 50 times faster than a Pentium 4.

On the other hand, 3PA-PET is the most efficient architecture with a computational efficiency per Processing Element (PE) of about 2 cycles per

3D-BP Algorithm	PE (threads)	Time	Cycles/Op	
			/PE	total
Desktop PC : Pentium 4 (core freq.=3.2 Ghz,BW _{mem} =6.4 GB/s)				
STIR ¹	1	11.13 s	70.4	70.4
VBI-flt(v1)	1	54,7 s	355	355
VBI-flt(v2)	1	17,4 s	113	113
VBI-flt(v3)	1	2,5 s	16	16
Workstation : bi-Xeon dual core (core freq.=3 Ghz,BW _{mem} =10.6 GB/s)				
STIR ¹	1 (1)	5.74 s	34,5	34,5
VBI-flt(v3)	1 (1)	1.17 s	7,1	7,1
VBI-flt(v3)	2 (2)	583 ms	7,06	3,53
VBI-flt(v3)	4 (4)	294 ms	7,12	1,78
GPU : GTS8800 (shader freq.=1.2 Ghz,BW _{mem} =64 GB/s)				
VBI-flt(v4)	96 (192)	99 ms	25,9	0.27
VBI-flt(v5)	96 (192)	50 ms	13,0	0.14
FPGA ² : Virtex 4 (freq.=200 Mhz,BW _{mem} =0.8 GB/s,l _{mem} =25 ns)				
VBI-fix	1	2,5 s	1	1
VBI-fix	4	774 ms	1,25	0,31
VBI-fix	8	526 ms	1,7	0,21
ASIC ³ : one memory bank (freq.=1,2 Ghz,BW _{mem} =4,8 GB/s,l _{mem} =25 ns)				
VBI-fix	1	499 ms	1.21	1,21
VBI-fix	4	214 ms	2,07	0,517
VBI-fix	8	135 ms	2,62	0,328
ASIC ³ : five memory banks (freq.=1,2 Ghz,BW _{mem} =24 GB/s,l _{mem} =25 ns)				
VBI-fix	40	27 ms	2.62	0,065

¹ Time normalized to a 128 * 128 * 63 volume.

(STIR reconstructs $64^2\pi * 63$ cylindrical volumes)

² 35 Mhz results scaled to 200 Mhz ($l_{mem}=5$ cycles)

³ 35 Mhz results scaled to 1,2 Ghz ($l_{mem}=30$ cycles)

Table III

COMPARED TIME PERFORMANCE FOR THE 3D PET BP OF A 128×128×63 VOLUME FROM A SIEMENS HR+ SINOGRAM (5 SEGMENTS, SPAN 9, 96 ANGLES OF PROJECTION).

THROUGHPUT OF RECONSTRUCTION (CYCLES PER VOXEL UPDATE) IS PRESENTED FOR THE GLOBAL ARCHITECTURE AND PER PROCESSING ELEMENT (PE).

operation for a processing block made of 8 units coupled with one memory bank. Because of the fewer available computational and memory resources, the FPGA technology doesn't allow to have an efficient 3PA-PET system compared to GPU. Nevertheless, considering that a market would exist to justify it, an ASIC with five memory banks and five units of 3PA-PET (8 pipelines each) running at 1.2 Ghz, would be twice faster than the Nvidia GPU, 20 times than a Xeon dual core and 100 times than a Pentium 4. Furthermore, a better tuning of the 3D-AP Cache would allow to increase the available

parallelism and again increase the speed-up. Also, an ASIC implementation would be likely to have a lower consumption than a GPU (Nvidia 8800GTS needs 130 Watts).

All the studied architectures succeed to benefit from the spatial and temporal localities without any developer effort to set a double buffering memory strategy. This is only possible because of their own memory cache (1D cache for CPU, 2D texture cache for GPU and the semi-general purpose 3D-AP Cache for 3PA-PET). Nevertheless, 3PA-PET is the one that best exploits the memory throughput, as illustrated in figure 13. In this figure, all the 3D BP implementations are placed according to their computational efficiency (cycles/op) and to their available memory throughput (GB/s). The 'optimal architecture' in this figure, corresponds to a hardware architecture that would have an optimal balance between its computational and memory throughputs. Each PE of this optimal architecture computes one operation per cycle and its prefetching memory strategy only loads the necessary data in cache and delivers it in time to the processing units. Of course, the more PEs it has, the greater the memory throughput has to be. As one can observe, 3PA-PET is the architecture with a cache-based memory strategy that is the closest to the optimal one. This makes 3PA-PET the architecture with the best potential of acceleration.

VII. CONCLUSION

This paper presents several ways to speed-up the BP algorithm on different target architectures: general purpose CPU, GPU and FPGA/ASIC. These solutions exploit the temporal and 3D spatial locality that can be found in the BP algorithm. A suitable loop reordering shows to be efficient despite the high non-linearity of the algorithm. The 3PA-PET (pre-fetched and parallelized Architecture for PET) architecture is the one that makes the best use of this locality and allows a high level of parallelization with a high computational throughput.

Thanks to the 3D-AP Cache together with a loop reordering, 3PA-PET architecture proves to be an efficient parallel architecture that overcomes the memory bottleneck. Indeed, as it has been measured on a SoPC prototype, the pipelines are seldom stalled and the high latency and low bandwidth of memories can be overcome. Moreover, the com-

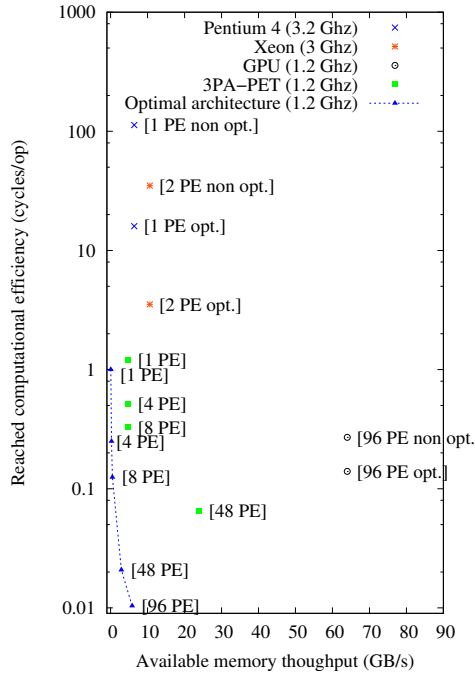


Figure 13. Memory throughput exploitation by CPU (optimized and non optimized code), GPU (optimized and non optimized code) and 3PA-PET implementations. 3PA-PET and optimal architecture results are obtained with $l_{mem} = 25ns$.

parison between the 3PA-PET architecture with a general purpose processor and a GPU highlights 3PA-PET efficiency. On one hand, the GPU has the best reconstruction time on a wall clock (followed by 3PA-PET and the CPU), on the other hand 3PA-PET makes the best use of the pipeline and clock cycles. An ASIC implementation with the same technological resources than of GPU would be of lower power consumption and faster: it would be twice faster than today's GPUs and 20 times faster than CPUs.

To conclude, the method of loop-reordering and use of an appropriate cache could be extended to other algorithms. The architecture principles presented in this article could be applied for the Cone beam BP needed in CT reconstruction.

REFERENCES

- [1] Kinahan P. E., et al. *Emission tomography : the fundamentals of PET and SPECT*, chapter Analytic image reconstruction methods. Elsevier Academic Press, 2004.
- [2] Jones J. P., et al. Data processing methods for a high throughput brain imaging pet research center. In *Nucl. Sci. Symp. 2006. IEEE*, volume 4, pages 2224–2228.
- [3] Mancini S. et al. An iir based 2d adaptive and predictive cache for image processing. In *Design of Circuits and Integrated Systems*, page 85, Bordeaux, France, 2004.
- [4] Kachelriess M., et al. Hyperfast parallel-beam and cone-beam backprojection using the cell general purpose hardware. *Medical Physics*, 34(4):1474–1486, April 2007.
- [5] Schellmann M., et al. Parallelization and runtime prediction of the listmode osem algorithm for 3d pet reconstruction. In *Nucl. Sci. Symp., 2006. IEEE*, volume 4, pages 2190–2195.
- [6] Shattuck D., et al. Internet2-based 3D PET image reconstruction using a PC cluster. *Phys. Med. Biol.*, 47(15):2785–2795, August 2002.
- [7] He T., et al. A heterogeneous windows cluster system for medical image reconstruction. In *IMSCCS '06.*, volume 1, pages 410–415, 2006.
- [8] Chidlow K. et al. Rapid emission tomography reconstruction. In *Proc. Int. Work. Volume Graphics (VG'03)*, Tokyo, Japan, July 2003.
- [9] Pratz G., et al. Fully 3-d list-mode osem accelerated by graphics processing units. In *Nucl. Sci. Symposium, 2006. IEEE*, volume 4, pages 2196–2202, Oct.
- [10] Xu F. et al. Real-time 3d computed tomographic reconstruction using commodity graphics hardware. *Physics in Medicine and Biology*, 52(12):3405–3419, 2007.
- [11] Scherl H., et al. Fast gpu-based ct reconstruction using the common unified device architecture (cuda). In *IEEE Nucl. Sci. Symp., NSS '07*, volume 6, pages 4464–4466.
- [12] Schiwietz T., et al. A fast and high-quality cone beam reconstruction pipeline using the gpu. In *Proc. SPIE Vol. 6510*, 2007.
- [13] Yang H., et al. Accelerating backprojections via cuda architecture. In *Proc. of Fully 3D*, pages 52–55, 2007.
- [14] Riabkov D., et al. Accelerated cone-beam backprojection using gpu-cpu hardware. In *Proc. of Fully 3D*, pages 68–71, 2007.
- [15] Scherl H., et al. On-the-fly-reconstruction in exact cone-beam ct using the cell broadband engine architecture. In *Proc. of Fully 3D*, pages 29–32, 2007.
- [16] Leiser M., et al. Parallel-beam backprojection: an fpga implementation optimized for medical imaging. *J. VLSI Signal Process. Syst.*, 39(3):295–311, 2005.
- [17] Li X. P2P-enhanced distributed computing in EM medical image reconstruction. In *PDPTA'04*, volume 2, pages 822–828.
- [18] Heigl B. et al. High-speed reconstruction for c-arm computed tomography. In *Proc. of Fully 3D*, pages 25–28, 2007.
- [19] Goddard I. et al. High-speed cone-beam reconstruction : An embedded systems approach. In *Proc. SPIE Medical Imaging Conf.*, pages 483–491, February 2002.
- [20] Terarecon . <http://www.terarecon.com/>.
- [21] Ni J., et al. Analysis of performance evaluation of parallel katsevich algorithm for 3-d ct image reconstruction. In *IMSCCS '06.*, volume 1, pages 258–265.
- [22] Gac N., et al. Hardware/software 2d-3d backprojection on a soc platform. In *Proc. of the 2006 ACM Symposium on Applied Computing (SAC)*, pages 222–228.
- [23] Thielemans K., et al. Stir: Software for tomographic image reconstruction release 2. In *Nucl. Sci. Symp. 2006. IEEE*, volume 4, pages 2174–2176.