



## Accelerating advanced MRI reconstructions on GPUs<sup>☆</sup>

S.S. Stone<sup>a,\*</sup>, J.P. Haldar<sup>b</sup>, S.C. Tsao<sup>a</sup>, W.-m.W. Hwu<sup>a</sup>, B.P. Sutton<sup>c</sup>, Z.-P. Liang<sup>b</sup>

<sup>a</sup> Center for Reliable and High-Performance Computing, University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA

<sup>b</sup> Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA

<sup>c</sup> Bioengineering Department and Biomedical Imaging Center, University of Illinois at Urbana-Champaign, Urbana, IL, 61801, USA

### ARTICLE INFO

#### Article history:

Received 11 March 2008

Received in revised form

10 May 2008

Accepted 10 May 2008

Available online 28 June 2008

#### Keywords:

GPU computing

MRI

Reconstruction

CUDA

### ABSTRACT

Computational acceleration on graphics processing units (GPUs) can make advanced magnetic resonance imaging (MRI) reconstruction algorithms attractive in clinical settings, thereby improving the quality of MR images across a broad spectrum of applications. This paper describes the acceleration of such an algorithm on NVIDIA's Quadro FX 5600. The reconstruction of a 3D image with 128<sup>3</sup> voxels achieves up to 180 GFLOPS and requires just over one minute on the Quadro, while reconstruction on a quad-core CPU is twenty-one times slower. Furthermore, for the data set studied in this article, the percent error exhibited by the advanced reconstruction is roughly three times lower than the percent error incurred by conventional reconstruction techniques.

© 2008 Elsevier Inc. All rights reserved.

## 1. Introduction

Mainstream microprocessors such as the Intel Pentium and AMD Opteron families have driven rapid performance increases and cost reductions in science and engineering applications for two decades. These commodity microprocessors have delivered GFLOPS to the desktop and hundreds of GFLOPS to cluster servers. This progress, however, slowed down in 2003 due to constraints on power consumption. Since that time, accelerators such as graphics processing units (GPUs) have led the advances in computational throughput for science and engineering applications. To wit, the peak throughput of programmable, floating-point, multiply-add operations on NVIDIA's GPUs has risen from 15 GFLOPS in early 2003 to nearly 400 GFLOPS in late 2007, a growth rate of roughly 1.9x per year. During the same period, the peak throughput of programmable FP MAD operations on Intel's CPUs has risen from 12 GFLOPS to roughly 100 GFLOPS, a growth rate of only 1.5x per year [27].

Recent advances in architecture have also increased the GPU's attractiveness as a platform for science and engineering applications. Prior to 2006, GPUs found very limited use in this

domain due to their limited support for both IEEE floating-point standards and arbitrary memory addressing. However, the recently released AMD R580 and NVIDIA G80 GPUs offer strong support for IEEE single-precision floating-point values (with double-precision soon to follow) and permit reads and writes to arbitrary addresses in memory [1,26]. Furthermore, modern GPUs use massive multithreading, fast context switching, and high memory bandwidth to tolerate ever-increasing latencies to main memory by overlapping long-latency loads in stalled threads with useful computation in other threads [22].

Increased programmability has also enhanced the GPU's suitability for science and engineering applications. For example, the G80 supports the single-program, multiple-data (SPMD) programming model, in which each thread is created from the same program and operates on a distinct data element, but all threads need not follow the same control flow path. As the SPMD programming model has been used on massively parallel supercomputers in the past, it is reasonable to expect that many high-performance applications will port easily to the G80 [22, 41]. Furthermore, general-purpose applications targeting the G80 are developed using ANSI C with simple extensions, rather than the cumbersome graphics application programming interfaces (APIs) [10,34] and high-level languages layered on graphics APIs [5, 7,40] that have been used in the past.

A wide variety of magnetic resonance imaging (MRI) applications, ranging from quantitative imaging of the brain to dynamic imaging of the beating heart, can benefit greatly from these increases in computational resources and advancements in

<sup>☆</sup> This work is based on an earlier work: Accelerating advanced MRI reconstructions on GPUs, in: Proceedings of the 2008 Conference on Computing Frontiers, (May 5–7, 2008) (c) ACM, 2008. <http://doi.acm.org/10.1145/1366230.13666274>.

\* Corresponding address: 224 Coordinated Science Lab, 1308 W. Main St., Urbana, IL 61801, USA.

E-mail address: [ssstone2@crhc.uiuc.edu](mailto:ssstone2@crhc.uiuc.edu) (S.S. Stone).

architecture and programmability. At present, many MRI experiments are specifically designed so that the image can be reconstructed quickly and efficiently on a standard CPU, often by acquiring the scan data on a uniform grid and applying a fast Fourier transform (FFT). However, in many applications the combination of tailored data acquisition and advanced image reconstruction significantly improves image quality. In particular, these techniques can increase signal-to-noise ratio, decrease scan time, and/or reduce imaging artifacts. However, advanced reconstruction algorithms often require several orders of magnitude more computation than conventional reconstruction algorithms. In this paper, we accelerate a reconstruction algorithm that can (1) generate MR images from arbitrary data sampling trajectories, and (2) incorporate prior anatomical knowledge into the reconstruction process, thereby increasing the signal-to-noise ratio while mitigating partial volume artifacts.

For these advanced reconstructions to be viable in clinical settings, dramatic and inexpensive computational acceleration is required. We find that advanced reconstructions from arbitrary scan trajectories are very well suited to acceleration on modern GPUs. In particular, an advanced reconstruction of an image comprising  $128^3$  voxels completes in just over one min on the G80, while the same reconstruction requires nearly 23 min on a quad-core CPU. Furthermore, relative to a conventional reconstruction of the data set examined in Section 5, the advanced reconstruction reduces the percent error in the reconstructed image by a factor of roughly 3X. The 21X acceleration achieved on the GPU makes the constrained reconstruction much more appealing in clinical settings.

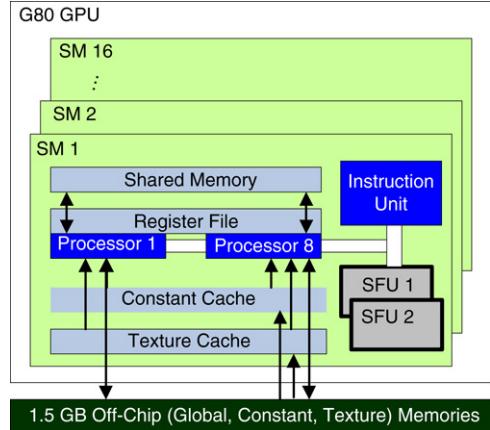
The remainder of this paper is organized as follows. Section 2 first describes the architecture of the Quadro FX 5600 and its G80 GPU, then discusses the advantages of advanced MRI reconstructions. Section 3 presents the GPU-based implementation of the advanced reconstruction algorithm. Section 4 describes experimental methodology. Section 5 presents results and discusses features of the Quadro that enable the advanced reconstruction to achieve up to 180 GFLOPS in performance. Section 6 discusses related work in GPU-based medical imaging. Section 7 concludes.

## 2. Background

### 2.1. The Quadro FX 5600 graphics card

The Quadro FX 5600 is a graphics card equipped with a G80 graphics processing unit (GPU). The Quadro has a large set of processor cores that can directly address a global memory. This architecture supports the single-program, multiple-data (SPMD) programming model, which is more general and flexible than the programming models supported by previous generations of GPUs, and which allows developers to easily implement data-parallel algorithms. In this section we discuss NVIDIA's Compute Unified Device Architecture (CUDA) and the architectural features of the G80 that are most relevant to accelerating MRI reconstructions. Similar descriptions are found in [29,30]. The interested reader may refer to [24,26] for additional details.

From the application developer's perspective, the CUDA programming model consists of ANSI C supported by several keywords and constructs. CUDA treats the GPU as a coprocessor that executes data-parallel kernel functions. The developer supplies a single source program encompassing both host (CPU) and kernel (GPU) code. NVIDIA's compiler, nvcc, separates the host and kernel codes, which are then compiled by the host compiler and nvcc,



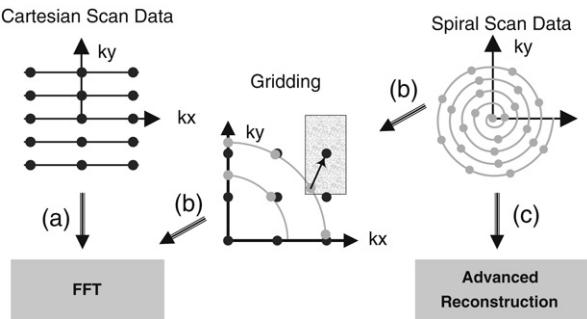
**Fig. 1.** Architecture of Quadro FX 5600.

respectively. The host code transfers data to and from the GPU's global memory via API calls, and initiates the kernel code by calling a function.

Fig. 1 depicts the Quadro's architecture. The G80 GPU consists of 16 *streaming multiprocessors* (SMs), each containing eight *streaming processors* (SPs), or processor cores, running at 1.35 GHz. Each SM has 8192 registers that are shared among all threads assigned to the SM. The threads on a given SM's cores execute in SIMD (single-instruction, multiple-data) fashion, with the instruction unit broadcasting the current instruction to the eight cores. Each core has a single arithmetic unit that performs single-precision floating point arithmetic and 32-bit integer operations. Additionally, each SM has two *special functional units* (SFUs), which perform more complex FP operations such as the trigonometric functions with low latency. Both the arithmetic units and the SFUs are fully pipelined. Thus, each SM can perform 18 FLOPS per clock cycle (one multiply-add operation per SP and one complex operation per SFU), yielding 388.8 GFLOPS (16 SM \* 18 FLOP/SM \* 1.35 GHz) of peak theoretical performance for the GPU.

The Quadro has 76.8 GB/s of bandwidth to its 1.5 GB, off-chip, global memory. Nevertheless, with computational resources supporting nearly 400 GFLOPS and each multiply-add instruction operating on up to 16 bytes of data, applications can easily saturate that bandwidth. Therefore, as depicted in Fig. 1, the G80 has several on-chip memories that can exploit data locality and data sharing to reduce an application's demands for off-chip memory bandwidth. For example, the Quadro has a 64 kB, off-chip *constant memory*, and each SM has an 8 kB constant memory cache. Because the cache is single-ported, simultaneous accesses of different addresses yield stalls. However, when multiple threads access the same address during the same cycle, the cache broadcasts that address's value to those threads with the same latency as a register access. This feature proves quite beneficial for the MRI reconstruction algorithm studied in this paper. In addition to the constant memory cache, each SM has a 16 kB *shared memory* for data that is either written and reused or shared among threads. Finally, for read-only data that is shared by many threads but not necessarily accessed simultaneously by all threads, the off-chip texture memory and the on-chip texture caches exploit 2D data locality to substantially reduce memory latency.

Threads executing on the G80 are organized into a three-level hierarchy. At the highest level, each kernel creates a single *grid*, which consists of many *thread blocks*. The maximum number of threads per block is 512. Each thread block is assigned to a single SM for the duration of its execution. Threads in the same block can share data through the shared memory and can perform barrier synchronization by invoking the `_syncthreads` primitive. Threads are otherwise independent, and synchronization across



**Fig. 2.** MRI reconstruction techniques. In (a) the scanner samples  $k$ -space on a uniform grid and reconstructs the image in one step via the FFT. In (b) the scanner samples  $k$ -space on a non-Cartesian (spiral) trajectory, then interpolates the samples onto a uniform grid and reconstructs the image in one step via the FFT. In (c) an advanced reconstruction algorithm is applied directly to the spiral scan data.

thread blocks is safely accomplished only by terminating the kernel. Finally, threads within a block are organized into warps of 32 threads. Each warp executes in SIMD fashion, with the SM's instruction unit broadcasting the same instruction to the eight cores on four consecutive clock cycles.

SMs can interleave warps on an instruction-by-instruction basis to hide the latency of global memory accesses and long-latency arithmetic operations. When one warp stalls, the SM can quickly switch to a ready warp in the same thread block or in some other thread block assigned to the SM. The SM stalls only if there are no warps with all operands available.

## 2.2. Advanced MRI reconstruction

Magnetic resonance imaging (MRI) is commonly used by the medical community to safely and non-invasively probe the structure and function of biological tissues from all regions of the body, and images generated using MRI have a profound impact in both clinical and research settings. MR imaging consists of two phases, acquisition (*scan*) and reconstruction. During the scan phase, the scanner samples data in the  $k$ -space domain (*i.e.* the spatial-frequency domain or Fourier transform domain) along a predefined trajectory. These samples are then transformed into the desired image during the reconstruction phase.

MRI is often limited by high noise levels, significant imaging artifacts, and/or long data acquisition times. In clinical settings, short scan times not only increase scanner throughput but also reduce patient discomfort, which tends to mitigate motion-related artifacts. High image resolution is equally important because it can enable earlier detection of pathology, leading to improved prognoses for patients. However, the goals of short scan time, high resolution, and high signal-to-noise ratio (SNR) often conflict; improvements in one metric tend to come at the expense of one or both of the others.

The sampling trajectory used by the MRI scanner can significantly affect the quality of the reconstruction. Fig. 2(a) and 2(c) depict a Cartesian scan trajectory and a non-Cartesian (spiral) scan trajectory, respectively. The Cartesian trajectory samples  $k$ -space on a uniform grid, which allows image reconstruction to be performed quickly and efficiently by applying a fast Fourier transform (FFT) directly to the acquired data. Although the reconstruction of Cartesian scan data is computationally efficient, non-Cartesian scan trajectories can be preferable because they are often faster and less sensitive to imaging artifacts caused by non-ideal experimental conditions. For these reasons, non-Cartesian trajectories with radial and spiral [3, Chapter 17] sampling patterns are becoming increasingly common in MRI.

Image reconstruction from non-Cartesian scan data presents both challenges and opportunities. In the most common approach, gridding, the samples are first interpolated onto a uniform

Cartesian grid and then reconstructed in one step via the FFT (see Fig. 2(b)) [18,33]. While gridding is computationally expedient, it satisfies no optimality criterion and cannot leverage prior information such as anatomical constraints. By contrast, statistically optimal image reconstructions can more accurately model imaging physics (*e.g.*, [12,28,39]) and can also incorporate additional prior information. For example, anatomically constrained reconstruction [14] incorporates anatomical information to reduce noise while preserving the resolution of known image features, enabling brief scans to yield high quality images. While such reconstructions have been impractical for large-scale 3D problems due to computational constraints, this paper shows that these reconstructions become viable in clinical settings when accelerated on GPUs. *Anatomically constrained reconstruction of non-Cartesian scan data enables brief scans to achieve high SNR, thereby decreasing imaging artifacts and increasing SNR simultaneously.* While such advanced reconstructions have been impractical for large-scale problems due to computational constraints, this paper shows that these reconstructions become viable in clinical settings when accelerated on GPUs.

We implemented the anatomically constrained reconstruction algorithm of [14]. This algorithm finds the solution to the following quasi-Bayesian estimation problem

$$\hat{\rho} = \arg \min_{\rho} \underbrace{\|\mathbf{F}\rho - \mathbf{d}\|_2^2}_{\text{data fidelity}} + \underbrace{\|\mathbf{W}\rho\|_2^2}_{\text{prior info}}, \quad (1)$$

where  $\hat{\rho}$  is a vector containing voxel values for the reconstructed image,  $\mathbf{F}$  is a matrix that models the imaging process,  $\mathbf{d}$  is a vector of data samples, and  $\mathbf{W}$  is a matrix that can incorporate prior information such as anatomical constraints. In clinical settings, these anatomical constraints are derived from one or more high-resolution, high-SNR scans of the patient, which reveal features such as the location of anatomical structures. The matrix  $\mathbf{W}$  is derived from these reference images. The first term in the above cost function imposes that data simulated from the reconstructed image should match somewhat closely with the real acquired data; the second term is used to impose prior information regarding the image statistics.

Because Eq. (1) defines a linear least squares problem, the solution is

$$\hat{\rho} = (\mathbf{F}^H \mathbf{F} + \mathbf{W}^H \mathbf{W})^{-1} \mathbf{F}^H \mathbf{d}. \quad (2)$$

However, the size of the matrix  $(\mathbf{F}^H \mathbf{F} + \mathbf{W}^H \mathbf{W})$  makes direct matrix inversion impractical for high-resolution reconstructions. For the  $128^3$ -voxel reconstructions examined in this paper, the inverted matrix contains well over four trillion complex-valued elements (the number of elements in the inverted matrix equals the square of the number of voxels in the reconstructed image). An iterative method for matrix inversion, such as the conjugate gradient (CG) algorithm [17], is therefore preferred.

The conjugate gradient algorithm reconstructs the image by iteratively solving for  $\hat{\rho}$ . During each iteration, the CG algorithm updates the current image estimate  $\rho$  to improve the value of the quasi-Bayesian cost function (Eq. (1)). The computational efficiency of the CG technique is largely determined by the efficiency of matrix–vector multiplication operations involving  $\mathbf{F}^H \mathbf{F}$  and  $\mathbf{W}^H \mathbf{W}$ , as these operations are required during each iteration of the CG algorithm. Fortunately, matrix  $\mathbf{W}$  often has a sparse structure that permits efficient multiplication by  $\mathbf{W}^H \mathbf{W}$ , and matrix  $\mathbf{F}^H \mathbf{F}$  has a convolutional structure [12,42] that enables efficient matrix multiplication via the FFT.

The advanced reconstruction algorithm described in this paper therefore consists of three primary computations. First, the algorithm computes each element of  $\mathbf{Q}$ , given by

$$Q(\mathbf{x}_n) = \sum_{m=1}^M |\phi(\mathbf{k}_m)|^2 e^{(i2\pi \mathbf{k}_m \cdot \mathbf{x}_n)}, \quad (3)$$

where  $\mathbf{Q}$  is the convolution kernel that facilitates multiplication operations involving  $\mathbf{F}^H \mathbf{F}$  and  $\phi(\cdot)$  is the Fourier transform of the voxel basis function. There are  $M$   $k$ -space sampling locations, with  $\mathbf{k}_m$  denoting the location of the  $m$ th sample. Likewise, there are  $N$  voxel coordinates, with  $\mathbf{x}_n$  denoting the coordinates of the  $n$ th voxel. Because  $\mathbf{Q}$  depends only on the scan trajectory (not the scan data) and the size of the image, it can be computed before the scan occurs and can be reused during any reconstruction that shares the same scan trajectory and image size.

Second, the algorithm computes the vector  $\mathbf{F}^H \mathbf{d}$ , defined as

$$[\mathbf{F}^H \mathbf{d}]_n = \sum_{m=1}^M \phi^*(\mathbf{k}_m) \mathbf{d}(\mathbf{k}_m) e^{(i2\pi \mathbf{k}_m \cdot \mathbf{x}_n)}. \quad (4)$$

Although Eq. (3) and Eq. (4) are quite similar, the former necessitates significantly more computation because the  $\mathbf{Q}$  algorithm oversamples the image space by a factor of two in each dimension. Therefore, during a 3D reconstruction, Eq. (3) is evaluated at  $8N$  values of  $\mathbf{x}_n$ , while the Eq. (4) is evaluated at only  $N$  values of  $\mathbf{x}_n$ . Finally, the CG solver performs iterative matrix inversion to solve Eq. (2).

The complexity of the advanced reconstruction far exceeds the complexity of a conventional, gridded reconstruction. Given a reconstruction problem of  $N$  voxels and  $M$  scan data points, the computations of  $\mathbf{Q}$  and  $\mathbf{F}^H \mathbf{d}$  have  $O(MN)$  complexity, compared to  $O(N \log N)$  complexity for reconstructions based on gridding and the FFT. For this reason, advanced reconstruction of high-resolution, three-dimensional images has been impractical in clinical settings, despite the technique's clear advantages over conventional reconstructions. Our work demonstrates that these advanced reconstructions can be performed quickly and efficiently on modern GPUs, increasing their viability in clinical settings.

### 2.3. Example application: 3D full-brain multi-echo acquisition

We illustrate one potential application of this work with real experimental data acquired using a novel acquisition scheme. The new multiparametric 3D structural imaging sequence provides several volumes with varying contrast in a multi-echo acquisition to assist in automatic brain segmentation, and we obtain volumes with  $T_1$ -weighting and  $T_2$ -weighting simultaneously and in complete registration.

Specifically, a 3D stack of spirals sequence was designed using the method of [13] with a  $256 \times 256 \times 176$  matrix size, 1 mm isotropic resolution, 17 spiral shots per slice, and with a TR of 350 ms. The sequence was acquired with multiple echoes to obtain a range of different contrasts during a single acquisition. In particular, we acquired a short echo time (2.2 ms) gradient echo spiral-out acquisition (GRE,  $T_1$ -weighted), followed by spiral-in/spiral-out acquisitions centered around two spin echo times at 46 ms (SE1) and 92 ms (SE2), respectively. All three images were acquired simultaneously with the same data acquisition trajectories and bandwidth, and therefore are coregistered. Subjects were scanned on a Siemens 3 T Allegra headscanner in accordance with the institutional review board using a single-channel head coil.

Advanced reconstructions and gridded reconstructions were performed with this data,<sup>1</sup> and the results are shown in Fig. 3. The advanced reconstruction's noise variance is more than 3 times better than that of the gridded reconstruction. The constraints used in the advanced reconstruction were obtained similarly to the technique described in [15], utilizing the whole image sequence to estimate a shared anatomical structure.

<sup>1</sup> Before processing, this data was filtered and resampled so that it could be reconstructed on a  $256 \times 256 \times 32$  voxel grid. The purpose of this preprocessing was to reduce the size of the CG solver's working set by a factor of 8 so that it could reside in the Quadro's 1.5 GB DRAM.

## 3. Advanced MRI reconstruction

The advanced MRI reconstruction algorithm described in Section 2.2 consists of three steps: computing the data structure  $\mathbf{Q}$  (which depends only on the scan trajectory), computing the vector  $\mathbf{F}^H \mathbf{d}$  (which depends on the scan trajectory and the scan data), and finding the image iteratively via a conjugate gradient linear solver. As Fig. 4 shows, the algorithms for  $\mathbf{F}^H \mathbf{d}$  and  $\mathbf{Q}$  are quite similar.<sup>2</sup> The most significant difference is that the  $\mathbf{Q}$  algorithm requires more computation because its outer loop executes  $8N$  iterations, compared to  $N$  iterations for  $\mathbf{F}^H \mathbf{d}$ . Otherwise,  $\mathbf{Q}$  suffers from the same bottlenecks and benefits from the same code transformations as  $\mathbf{F}^H \mathbf{d}$ .

Because  $\mathbf{Q}$  can be computed prior to acquiring an image's scan data, the critical path for a given reconstruction consists only of computing  $\mathbf{F}^H \mathbf{d}$  and executing the linear solver. Therefore, the remainder of this section describes the algorithms for  $\mathbf{F}^H \mathbf{d}$  and the linear solver, focusing on the implementation of the  $\mathbf{F}^H \mathbf{d}$  algorithm on the GPU. The interested reader may refer to [37] for more detailed discussion of  $\mathbf{Q}$ .

### 3.1. $\mathbf{F}^H \mathbf{d}$

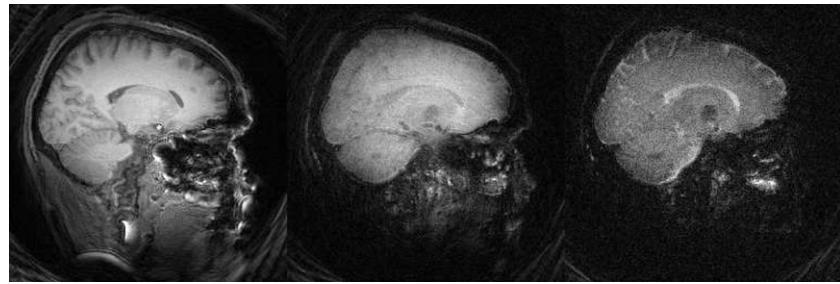
As Fig. 4(b) shows, the algorithm for  $\mathbf{F}^H \mathbf{d}$  is an excellent candidate for acceleration on the GPU because it contains substantial data-parallelism. The algorithm first computes the real and imaginary components of  $\mathbf{mu}$  at each sample point in the trajectory space ( $k$ -space), then computes the real and imaginary components of  $\mathbf{F}^H \mathbf{d}$  at each voxel in the image space. The value of  $\mathbf{F}^H \mathbf{d}$  at any voxel depends on the values of all sample points, but no elements of  $\mathbf{F}^H \mathbf{d}$  depend on any other elements of  $\mathbf{F}^H \mathbf{d}$ . Therefore, all elements of  $\mathbf{F}^H \mathbf{d}$  can be computed independently and in parallel.

Despite the algorithm's inherent parallelism, potential performance bottlenecks are evident. First, in the loop that computes the elements of  $\mathbf{F}^H \mathbf{d}$ , the ratio of floating-point operations to memory accesses is at best 3:1 and at worst 1:1. The best case assumes that the **sin** and **cos** operations are computed using five-element Taylor series that require 13 and 12 floating-point operations, respectively. The worst case assumes that each trigonometric operation is computed as a single operation in hardware. In either case, the GPU-based implementation of the algorithm must conserve memory bandwidth and tolerate memory latency. Second, the ratio of FP arithmetic to FP trigonometry is only 13:2. Thus, GPU-based implementation must tolerate or avoid stalls due to long-latency **sin** and **cos** operations.

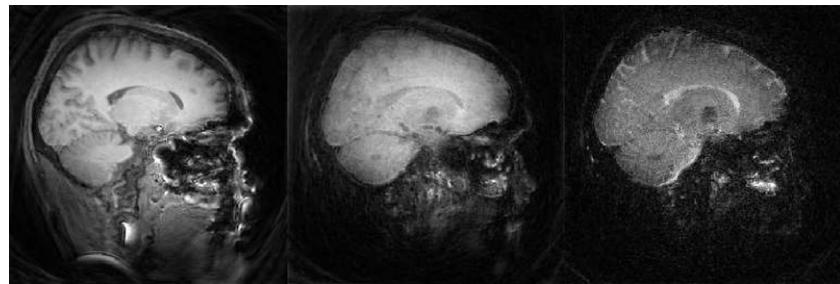
The GPU-based implementation of the  $\mathbf{F}^H \mathbf{d}$  algorithm (see Fig. 4(c)) uses the G80's constant memory caches to eliminate the potential bottleneck posed by memory bandwidth and latency. To overcome the memory bottleneck, the scan data is divided into many tiles, with each tile containing a distinct subset of sample points. For each tile, the host CPU loads the corresponding subset of sample points into constant memory before executing the **cmpPhD** function. Each thread then computes a partial sum for a single element of  $\mathbf{F}^H \mathbf{d}$  by iterating over all the sample points in the tile. This optimization significantly increases the ratio of FP operations to global memory accesses.

Likewise, the G80's special functional units (SFUs) enable the algorithm to avoid the potential bottleneck of long latency trigonometric operations. When the **use\_fast\_math** compiler option is invoked, the **sin** and **cos** operations are not linked to long-latency library calls, but rather are executed as individual, low-latency instructions on the SFUs. The speed of the SFU comes at the expense of some loss in accuracy when the argument to the **sin** or **cos** is very small, but, as we show in Section 5, this optimization does not necessarily decrease the overall accuracy of the algorithm.

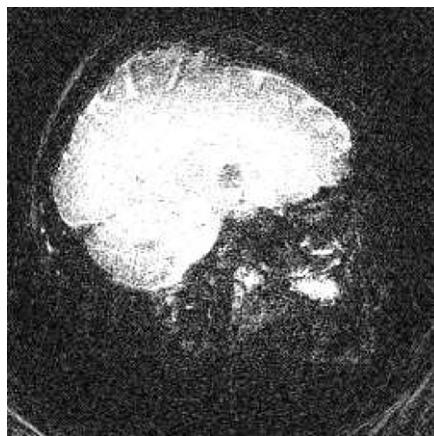
<sup>2</sup> In this work, we compute  $\mathbf{Q}$  and  $\mathbf{F}^H \mathbf{d}$  exactly, excluding numerical effects. These quantities have previously been calculated using fast approximations (e.g. [12,42]) due to the past impracticality of solving the exact problem.



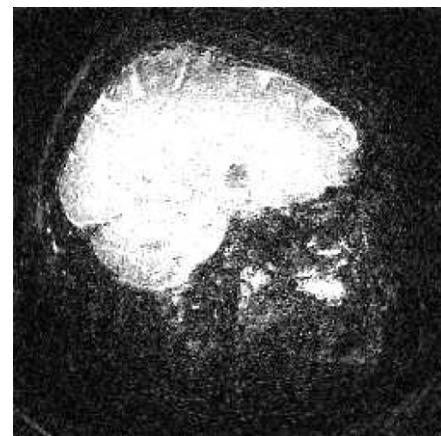
(a) Gridded Reconstructions. Left to right: GRE, SE1, and SE2.



(b) Constrained Reconstructions. Left to right: GRE, SE1, and SE2.



(c) Rescaled Gridded SE2.



(d) Rescaled Constrained SE2.

**Fig. 3.** Application of advanced MRI reconstruction to human brain data. (a) One slice from the gridding reconstructions from each of the three different images. (b) The corresponding slice from the constrained reconstructions. The SE2 image is shown with a modified colorscale in (c) and (d) to illustrate the significant SNR advantage of the advanced reconstruction. The constrained reconstruction's noise variance is more than 3 times lower than that of the gridded reconstruction.

### 3.2. Linear solver

The final phase of the image reconstruction consists of a linear solver that implements the preconditioned conjugate gradient (PCG) algorithm [9,21]. As described in Section 2, the solver iteratively solves Eq. (2) to find the desired image  $\hat{\rho}$ . When the iterations converge or the number of iterations exceeds a threshold, the solver terminates. During each iteration, the solver performs a large FFT and inverse FFT, several BLAS and sparse BLAS operations (including multiplication of sparse matrices and vectors, as well as addition, scaling, and scalar multiplication of vectors), and several other computations (such as summation reduction, shifting, and sampling).

We ported the PCG algorithm from MATLAB to C/CUDA, using NVIDIA's CUDA CUFFT Library [25] for the FFT and inverse FFT operations, implementing the BLAS and sparse BLAS operations in CUDA, and orchestrating the control flow and data marshaling in C. Complex-valued objects are represented using CUDA's **cufftComplex** data type, as required by the CUFFT Library. Sparse matrices are stored in compressed row format [11] to facilitate

efficient GPU-based execution of the expression  $\mathbf{A} * \mathbf{x}$ , where  $\mathbf{A}$  is a sparse matrix and  $\mathbf{x}$  is a vector. Although we have not rigorously analyzed the performance of the CUDA-based solver, it is roughly 25 times faster than a MATLAB-based version of the same algorithm. We use the CUDA-based solver for all experiments presented in Section 5 and view its performance as acceptable.

### 4. Methodology

To quantify the effects of the Quadro's architectural features on the performance and quality of the reconstruction, we implemented seven versions of the algorithm for  $\mathbf{F}^H \mathbf{d}$ , five of which are depicted in Fig. 5. The base version (GPU.Base, see Fig. 5(a)) simply executes in data-parallel fashion on the GPU, without using even the simplest optimizations to conserve memory bandwidth or tolerate long latency loads and trigonometric operations. The second version (GPU.RegAlloc, see Fig. 5(b)) register allocates the voxel data, thereby conserving some memory bandwidth and reducing the latency of all voxel accesses. GPU.Layout (Fig. 5(c)) register allocates the voxel data and changes the layout of the

<pre> // calc  φ(k<sub>m</sub>) ^2 // at each sample point m for (m = 0; m &lt; M; m++) {     phiMag[m] = rPhi[m]*rPhi[m] +         iPhi[m]*iPhi[m]; }  // calc Q at each voxel n for (n = 0; n &lt; 8*N; n++) {     for (m = 0; m &lt; M; m++) {         // e<sup>2πk<sub>m</sub>x<sub>n</sub></sup>         exp = 2*PI*(kx[m] * x[n] +             ky[m] * y[n] +             kz[m] * z[n]);         // ae<sup>ic</sup> = a*cos(c)+ia*sin(c)         rQ[n] += phiMag[m]*cos(exp);         iQ[n] += phiMag[m]*sin(exp);     } } </pre>	<pre> // calc mu = φ*(k<sub>m</sub>)d(k<sub>m</sub>) // at each sample point m for (m = 0; m &lt; M; m++) {     rMu[m] = rPhi[m]*rD[m] +         iPhi[m]*iD[m];     iMu[m] = rPhi[m]*iD[m] -         iPhi[m]*rD[m]; }  // calc FHd at each voxel n for (n = 0; n &lt; N; n++) {     for (m = 0; m &lt; M; m++) {         // e<sup>2πk<sub>m</sub>x<sub>n</sub></sup>         exp = 2*PI*(kx[m] * x[n] +             ky[m] * y[n] +             kz[m] * z[n]);         cArg = cos(exp);         sArg = sin(exp);         // (a+bi)e<sup>ic</sup> = a*cos(c)-b*sin(c)         // + i(b*cos(c)+a*sin(c))         rFhD[n] += rMu[m]*cArg -             iMu[m]*sArg;         iFhD[n] += iMu[m]*cArg +             rMu[m]*sArg;     } } </pre>	<pre> // calc mu at each sample point m __global__ void cmpMu(float* rPhi, iPhi, rD, iD, rMu, iMu, int M) {     int m = blockIdx.x * MU_THREADS_PER_BLOCK + threadIdx.x;     if (m &lt; M) {         rMu[m] = rPhi[m]*rD[m] + iPhi[m]*iD[m];         iMu[m] = rPhi[m]*iD[m] - iPhi[m]*rD[m];     } }  // calc FHd at one voxel n __global__ void cmpFhD(float* gx, gy, gz, grFhD, giFhD) {     // find the index of the voxel assigned to this thread     int n = blockIdx.x * FH_THREADS_PER_BLOCK + threadIdx.x;      // register allocate voxel inputs and outputs     x = gx[n]; y = gy[n]; z = gz[n];     rFhD = grFhD[n]; iFhD = giFhD[n];      // loop over all the sample points in the current tile     for (int m = 0; m &lt; SAMPLE PTS PER TILE; m++) {         // s (sample data) is held in constant memory         float exp = 2 * PI * (s[m].kx * x +             s[m].ky * y +             s[m].kz * z);         cArg = cos(exp);         sArg = sin(exp);         rFhD += s[m].rMu*cArg - s[m].iMu*sArg;         iFhD += s[m].iMu*cArg + s[m].rMu*sArg;     }      grFhD[n] = rFhD;     giFhD[n] = iFhD; } </pre>
---	--	---

(a) Q algorithm

(b) FH<sup>H</sup>d algorithm(c) FH<sup>H</sup>d algorithm in CUDA

**Fig. 4.** Data-parallel phases of advanced MRI reconstruction. Panels (a) and (b) show simplified C code for the algorithms that compute **Q** and **F<sup>H</sup>d**, respectively. Panel (c) depicts the **F<sup>H</sup>d** algorithm in CUDA.

scan data in the Quadro's global memory so that accesses to the scan data make more efficient use of the memory bandwidth. GPU.ConstMem (Fig. 5(d)) register allocates the voxel data and places the scan data in the Quadro's constant memory so that accesses to the scan data are cached. The fifth version (GPU.FastTrig, see Fig. 5(e)) additionally uses the G80's special functional units to compute fast, approximate versions of the trigonometric operations. The sixth version, GPU.Tune, also uses experimentally-tuned settings for three code transformations: loop unrolling, data tiling (scan points per thread), and number of threads per block. The tuned settings balance allocation of GPU resources to improve hardware utilization and thread efficiency. Finally, GPU.Multi executes the tuned version on multiple Quadros.

To obtain a reasonable baseline, we implemented two versions of **F<sup>H</sup>d** on the CPU. Version CPU.DP uses double-precision for all floating-point values and operations, while version CPU.SP uses single-precision. Both CPU versions are compiled with Intel's icpc (version 10.1) using flags -O3 -msse3 -axT -vec-report3 -fp-model fast = 2, which (1) vectorizes the algorithm's dominant loops using instructions tuned for the Core 2 architecture, and (2) links the trigonometric operations to fast, approximate functions in the math library. Based on experimental tuning with a smaller data set, the inner loops are unrolled by a factor of four and the scan data is tiled to improve locality in the L1 cache.

Each GPU version of **F<sup>H</sup>d** is compiled using nvcc -O3 (CUDA version 1.1) and executed on a 1.35 GHz Quadro FX 5600. The Quadro card is housed in a system with a 2.4 GHz dual-socket, dual-core Opteron 2216 CPU. Each core has a 1 MB L2 cache. The CPU versions use pthreads to execute on all four cores of 2.66 GHz Core 2 Extreme quad-core CPU, which has peak theoretical capacity of 21.2 GFLOPS per core and a 4 MB L2 cache. The CPU versions perform substantially better on the Core 2 Extreme quad-core than on the dual-socket, dual-core Opteron.

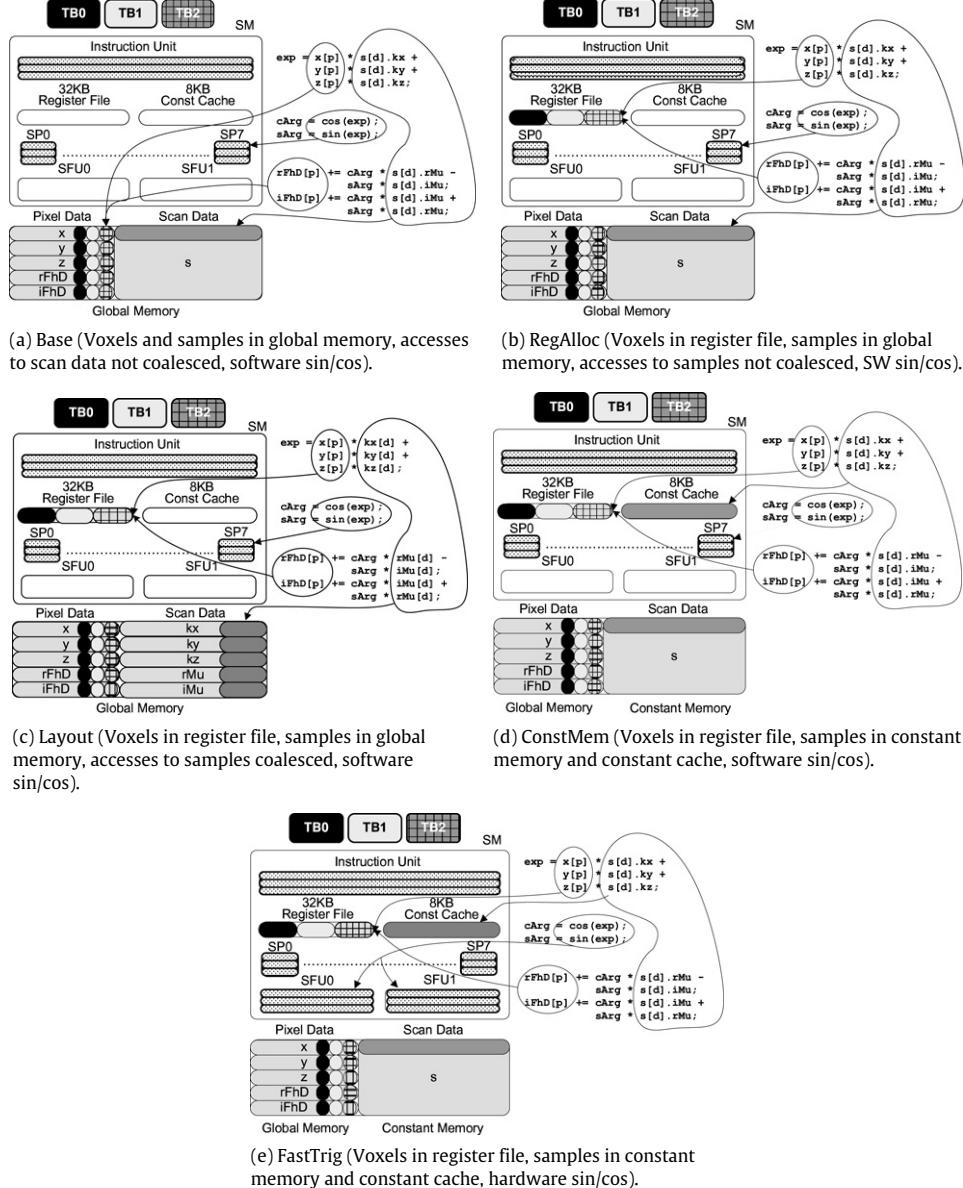
All reconstructions use the GPU version of the linear solver, which executes 60 iterations on the Quadro FX 5600. Two versions of **Q** were computed on the Core 2 Extreme, one using double-precision and the other using single-precision. The single-precision **Q** was used for all GPU-based reconstructions and for

the reconstruction involving CPU.SP, while the double-precision **Q** was used only for the reconstruction involving CPU.DP. As the computation of **Q** is not on the reconstruction's critical path, we give **Q** no further consideration.

To facilitate comparison of the advanced reconstruction with a conventional reconstruction, we also evaluated a reconstruction based on gridding and the FFT [18]. Our version of the gridded reconstruction is not optimized for performance, but it is fair to assume that an optimized implementation would execute in several seconds [37].

All reconstructions are performed on sample data obtained from a simulated, three-dimensional, non-Cartesian scan of a phantom image [20]. There are 284,592 sample points in the scan data set, and the image is reconstructed at 128<sup>3</sup> resolution, for a total of 2<sup>21</sup> voxels. In the first set of experiments, the simulated data contains no noise. In the second set of experiments, we added complex white Gaussian noise to the simulated data. When determining the quality of the reconstructed images, the *percent error* and *peak signal-to-noise ratio* metrics are used. The percent error is the root-mean-square (RMS) of the voxel error divided by the RMS voxel value in the true image (after the true image has been sampled at 128<sup>3</sup> resolution). To permit fair comparison of the gridded and advanced reconstructions, we adjusted the scale of each gridded image to match the scale of the true image before computing the gridded image's percent error and PSNR.

The data (runtime, GFLOPS, and images) presented in Section 5 were obtained by reconstructing each image once with each of the 11 implementations of the **F<sup>H</sup>d** algorithm described above. There are two exceptions to this policy. For GPU.Tune and GPU.Multi, the time required to compute **F<sup>H</sup>d** is so small that run-time variations in performance may become non-negligible. Therefore, for these configurations we computed **F<sup>H</sup>d** three times and reported the average performance. Also, when performing multiple reconstructions of the same data set back-to-back on the same computer, we do not clear the caches between successive calculations of **F<sup>H</sup>d** or successive executions of the linear solver. In the case of the **F<sup>H</sup>d** algorithm, which has a relatively small working set, the runtime increases by roughly 10% when the caches are cold.

Fig. 5. Versions of the  $F^H d$  algorithm on the GPU.

By contrast, the linear solver, which has a relatively large working set, exhibits a 30% increase in runtime when the caches are cold. However, given that (1) some data in the working sets depend only on the scan trajectory, and (2) clinicians are likely to use the same computer to perform several successive reconstructions on the same patient or with the same scan trajectory, it is difficult to determine the extent to which cold caches are an accurate reflection of clinical conditions. Finally, the reconstruction times reported in Section 5 exclude the time required to create an image from the reconstructed data (roughly one second).

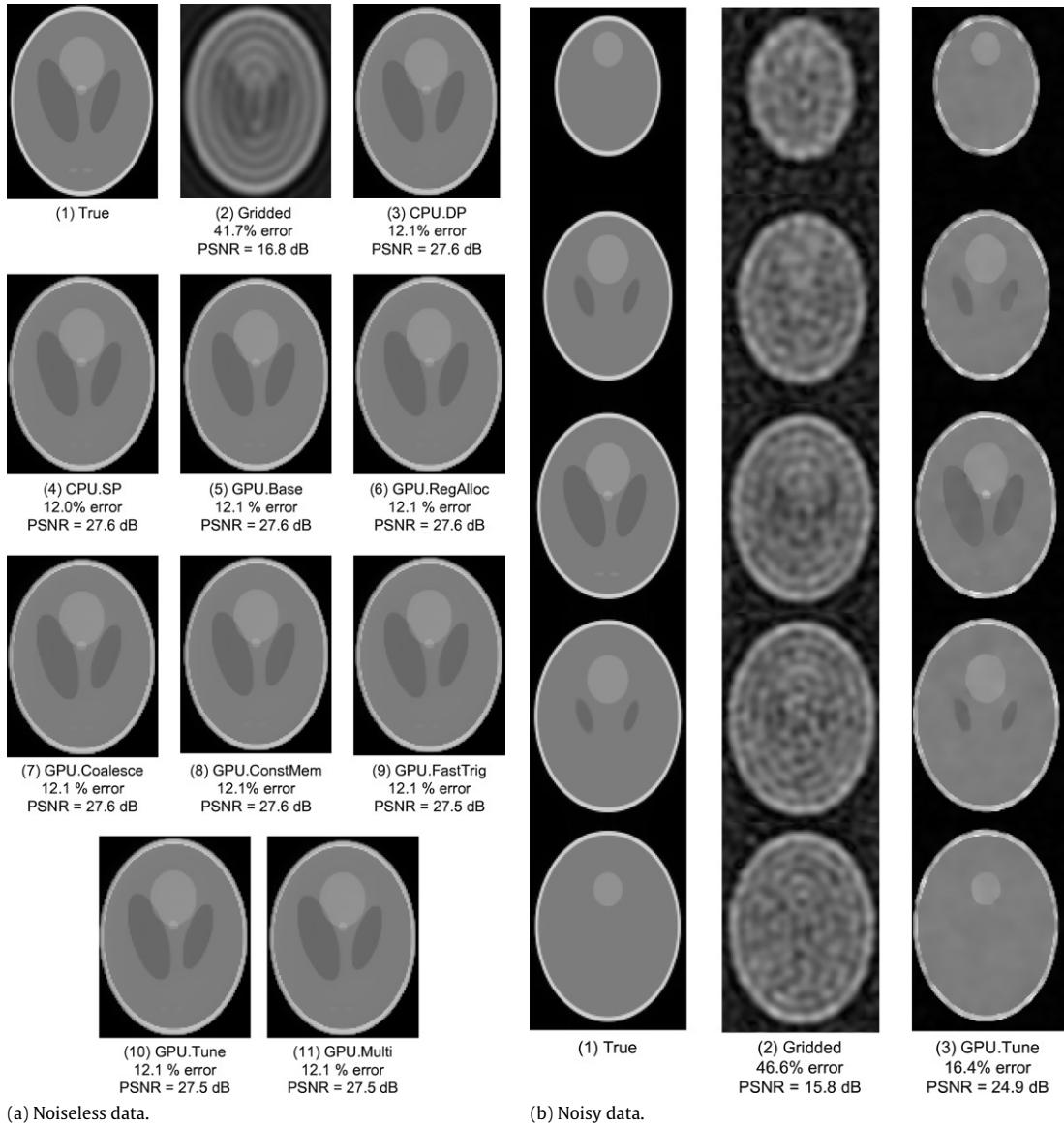
Finally, the advanced reconstruction leverages two optimizations that are not evident elsewhere in our discussion. First, the scan trajectory is symmetric, and the advanced reconstruction uses prior knowledge of that symmetry to mitigate the effects of numerical imprecision on the reconstruction's accuracy. Second, we manually balanced the resolution and the noise in the advanced reconstruction by performing the reconstruction multiple times while adjusting a regularization parameter. Adjustment of the regularization parameter can be performed prior to acquiring the

sample data, given the sampling trajectory, noise levels, and other readily available prior information [14].

## 5. Evaluation

To be useful in clinical settings, the advanced reconstruction must satisfy two criteria. First, the quality of an image obtained via the advanced reconstruction should significantly exceed the quality of an image obtained via a gridded reconstruction. Second, the reconstruction must complete quickly. After image acquisition, the patient typically remains in the scanner during image reconstruction. The scanner operator then decides whether the image is acceptable or whether it should be acquired again. Any delays therefore increase patient discomfort and decrease scanner throughput. Also, when the administration of a medical treatment depends on the MR images, any delay is at best frustrating and at worst harmful to the patient's health.

Our experiments indicate that the advanced reconstruction definitely satisfies the first criterion. As Fig. 6(a) shows, advanced reconstruction of the noiseless data yields significantly better



**Fig. 6.** Phantom images. (a) Noiseless data: One 2D slice of the 3D image. The percent error and PSNR values in each sub-figure caption are calculated over the entire 3D image. (b) Noisy data: Three 2D slices of the 3D image. The percent error and PSNR values in each sub-figure caption are calculated over the entire 3D image.

images than gridded reconstruction. Relative to the true image (Fig. 6(a)(1)), the advanced reconstructions (Fig. 6(a)(3–11)) exhibit 12% to 13% error and 27 dB to 28 dB PSNR, compared to 42% error and 17 dB PSNR for the gridded reconstruction (Fig. 6(a)(2)). There are no significant differences among the images obtained from the advanced reconstruction, despite the use of single-precision floating point in Fig. 6(a)(4–11) and approximate trigonometric operations in Fig. 6(a)(3, 4, and 9–11).

The images reconstructed from the noisy data (Fig. 6(b)) further demonstrate the superiority of the advanced reconstruction. Relative to the true image, the advanced reconstruction exhibits 16% error and 25 dB PSNR, while the error and PSNR for the gridded reconstruction are 47% and 16 dB, respectively. Again, there are no significant differences among the images obtained from the various versions of the advanced reconstruction.

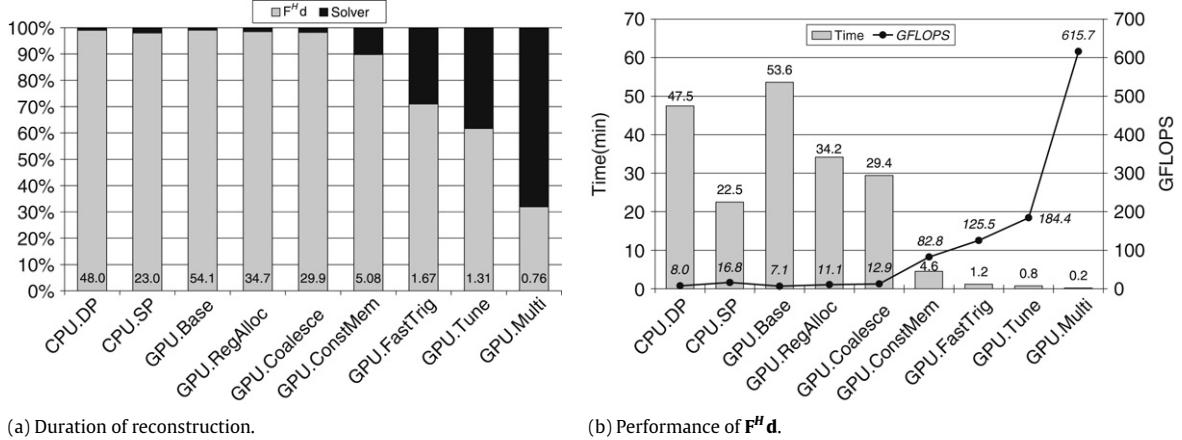
In addition to producing significantly better images than the gridded reconstruction, the GPU-accelerated advanced reconstruction arguably satisfies the second criterion for clinical use: speed. As Fig. 7(a) shows, the fastest single-GPU version of the advanced reconstruction completes in 66 s. This reconstruction time is clearly much more appealing for clinical applications than

the fastest CPU-based reconstruction, which completes in nearly 23 min.

The fastest single-GPU version of the advanced reconstruction computes  $\mathbf{F}^H \mathbf{d}$  in 49 s, compared to 22.5 min for the fastest CPU-based reconstruction. The remainder of this section describes how the advanced reconstruction leverages the GPU's resources to achieve such impressive acceleration when computing  $\mathbf{F}^H \mathbf{d}$ . We find that the constant memory caches are quite effective in reducing the number of accesses to global memory, while the special functional units provide substantial acceleration for the trigonometric computations in the algorithm's inner loops. We also find that experimentally-tuned code transformations have a significant impact on the algorithm's performance. Specifically, the algorithm's performance increases by 47% when the tiling factor, the number of threads per block, and the loop unrolling factor are experimentally tuned.

### 5.1. GPU.Base

As Fig. 7(b) shows, GPU.Base is significantly slower than CPU.SP, the optimized, single-precision, quad-core implementation of  $\mathbf{F}^H \mathbf{d}$ .



(a) Duration of reconstruction.

(b) Performance of  $F^H \mathbf{d}$ .

**Fig. 7.** Performance of advanced MRI reconstruction. (a) The reconstruction time includes the time to compute  $F^H \mathbf{d}$  and the time to run 60 iterations of the linear solver. The number at the bottom of each bar is the reconstruction time in minutes. (b) Performance of  $F^H \mathbf{d}$  computation. The first six configurations (CPU.DP – GPU.ConstMem) compute the trigonometric functions in software, using approximately 13 and 12 FLOPS for the **sin** and **cos** operations, respectively. The remaining configurations compute the trigonometric operations in hardware; therefore, each **sin** or **cos** accounts for a single FLOP.

In GPU.Base (see Fig. 5(a)), the inner loops are not unrolled. There are 256 threads per block and 256 scan points per tile. Because GPU.Base leverages neither the constant memory nor the shared memory, memory bandwidth and latency are significant performance bottlenecks. With one 4-byte global memory accesses for every three FP operations, and with memory bandwidth of 76.8 GB/s, the upper limit on the kernel's performance is only 57.6 GFLOPS. Due to other performance bottlenecks, the kernel actually achieves only 7.0 GFLOPS, less than half of the 16.8 GFLOPS achieved by CPU.SP.

### 5.2. GPU.RegAlloc

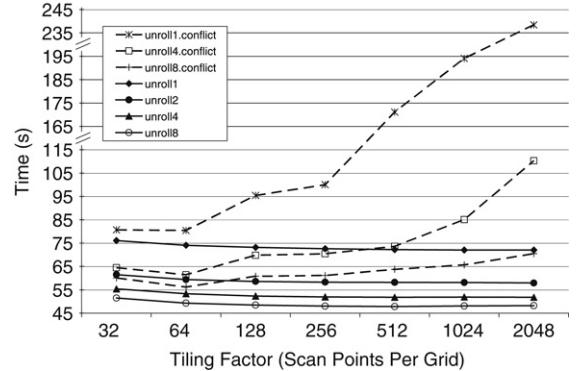
Relative to GPU.Base, GPU.RegAlloc (see Fig. 5(b)) decreases the time required to compute  $F^H \mathbf{d}$  from 53.6 min to 34.2 min. In short, register allocating the voxel data increases the computation intensity (the ratio of FP operations to off-chip memory accesses) from 3:1 to 5:1. This substantial reduction in required off-chip memory bandwidth translates into increased performance. Eliminating the two stores to global memory during every loop iteration is particularly beneficial.

### 5.3. GPU.Layout

By changing the layout of the scan data in global memory (see Fig. 5(c)), GPU.Layout achieves an additional speedup of 16% over GPU.RegAlloc. The underlying causes of GPU.Layout's improved performance relative to GPU.RegAlloc are difficult to identify. GPU.Layout does require less overhead to marshal the data into struct-of-arrays format than GPU.RegAlloc requires to marshal the data into array-of-structs format, which may account for a small fraction of the improvement in performance. Furthermore, we speculate that the Quadro's memory controller may provide some form of buffering that the struct-of-arrays layout leverages more successfully than does the array-of-structs layout. Section 5.4 offers additional insight into the relative merits of these two data layouts.

### 5.4. GPU.ConstMem

GPU.ConstMem (Fig. 5(d)) achieves speedup of 6.4X over GPU.Layout by placing each tile's scan data in constant memory rather than global memory. GPU.ConstMem therefore benefits from each SM's 8 kB constant memory cache. At 4.6 min and 82.8



**Fig. 8.** Constant cache conflicts. The three dashed lines represent hybrid versions of GPU.Layout and GPU.ConstMem, such that the scan data is placed in constant memory using the struct-of-arrays layout depicted in Fig. 5(c). Each dashed line corresponds to a different loop unrolling factor. Likewise, the four solid lines correspond to versions of GPU.ConstMem with different loop unrolling factors. Configurations with loop unrolling factor of 2 are excluded from the hybrid versions because some of those configurations hang for unknown reasons.

GFLOPS, this version of  $F^H \mathbf{d}$  is 4.9X faster than the optimized CPU version.

The GPU.ConstMem and GPU.Layout configurations underscore an important trade-off between optimizing latency to off-chip memory and optimizing latency to the constant caches. As discussed above, the array-of-structs data layout in Fig. 5(c) improves performance relative to the struct-of-arrays data layout in Fig. 5(b). However, as Fig. 8 shows, the struct-of-arrays layout yields very bad performance if the scan data is placed in constant memory, presumably because this layout leads to excessive conflicts in the constant cache. In particular, as the tiling factor increases, the time required for these hybrid versions (with scan data in constant memory but with struct-of-arrays layout) to compute  $F^H \mathbf{d}$  steadily increases [30]. This phenomenon most likely occurs because elements of nearby scan points ( $\mathbf{k}_x$ ,  $\mathbf{k}_y$ , etc.) map to the same cache line, so that the warps continually contend for the same cache lines. By contrast, the various versions of GPU.ConstMem actually require less time to compute  $F^H \mathbf{d}$  as the tiling factor increases; cache conflicts are clearly much less common.

We now analyze the off-chip memory accesses on a single SM during the execution of three thread blocks. With 7 global memory accesses per thread, 256 threads per thread block, and 3 thread blocks per SM, there are 5376 accesses to global memory.

**Table 1**  
CUDA 1.0 vs CUDA 1.1

Config	$\mathbf{F}^H \mathbf{d}$				Reconstruction (All times in min:s)					
	CUDA 1.0		CUDA 1.1		CUDA 1.0			CUDA 1.1		
	Registers	GFLOPS	Registers	GFLOPS	FhD	Solver	Total	FhD	Solver	Total
GPU.Base	14	7.02	18	7.05	53:51	0:25	54:17	53:36	0:17	53:53
GPU.RegAlloc	16	11.07	20	11.07	34:08	0:25	34:33	34:09	0:17	34:27
GPU.Layout	16	13.10	20	12.85	28:51	0:25	29:16	29:25	0:17	29:42
GPU.ConstMem	15	85.82	19	82.79	4:24	0:25	4:50	4:34	0:18	4:52
GPU.FastTrig	13	127.54	13	125.49	1:10	0:25	1:36	1:11	0:17	1:29
GPU.Tune	19	151.45	13	184.40	0:59	0:25	1:25	0:49	0:17	1:06
GPU.Multi	19	502.22	13	615.75	0:18	0:25	0:43	0:15	0:18	0:32

The Registers column refers to registers per thread.

Assuming no constant memory cache evictions due to conflicts, there are also 1280 accesses to constant memory (256 data points per tile, with 5 floating-point values per data element), yielding a total of 6656 off-chip memory accesses. The number of floating-point computations performed by the 3 thread blocks is  $3^*256^*256^*38 = 7471,104$ . Thus, the ratio of FP operations to off-chip memory accesses has increased by over two orders of magnitude, from 3:1 to 1100:1. However, GPU.ConstMem still achieves only 82.8 GFLOPS (roughly 20% to 25% of the Quadro's peak theoretical throughput), which implies the existence of another bottleneck.

### 5.5. GPU.FastTrig

GPU.FastTrig (Fig. 5(e)) achieves acceleration of nearly 4X over GPU.ConstMem by using the special functional units (SFUs) to compute each trigonometric operation as a single operation in hardware. When compiled without the **use\_fast\_math** compiler option, the algorithm uses implementations of **sin** and **cos** provided by an NVIDIA math library. Assuming that the library computes **sin** and **cos** using a five-element Taylor series, the trigonometric operations require 13 and 12 floating-point operations, respectively. By contrast, when compiled with the **use\_fast\_math** option, each **sin** or **cos** computation executes as a single floating-point operation on an SFU. The SFU achieves low latency at the expense of some accuracy. In our experiments (not shown), the images reconstructed by GPU.FastTrig always had lower or only slightly higher percent error than images reconstructed by GPU.ConstMem. Thus, the SFU's approximate implementations of **sin** and **cos** often have negligible impact on the reconstruction's accuracy. However, further experimentation is necessary to determine whether there are experimental conditions under which these instructions might decrease the quality of a reconstruction.

### 5.6. GPU.Tune

While GPU.FastTrig overcomes the potential bottlenecks related to off-chip memory accesses and trigonometric computations, the algorithm still performs at only 125.5 GFLOPS, which is roughly one-third of the Quadro's peak theoretical performance. To determine the impact of experiment-driven code transformations, we conducted an exhaustive search that varied the number of threads per block from 32 to 512 (by increments of 32), the tiling factor from 32 to 2048 (by powers of 2), and the loop unrolling factor from 1 to 8 (by powers of 2).<sup>3</sup> Recent work has demonstrated that this type of experimental tuning can be performed quickly and accurately using static analysis techniques, as long as the code is parameterized correctly [30]. For reference, all previous configurations (GPU.Base – GPU.FastTrig) performed no loop unrolling and

set both the number of threads per block and the tiling factor to 256. The exhaustive, experiment-driven search selects 128 threads per block, a tiling factor of 512, and a loop unrolling factor of 8. This configuration increases the algorithm's performance by 47%, with the runtime decreasing to 49 s and the throughput increasing to 184 GFLOPS.

Although the code is now well-optimized and tuned, the achieved throughput is just under 50% of the Quadro's peak throughput. To understand the remaining constraints on performance, we examined GPU.Tune's ptx code (the assembly-like code generated by nvcc and consumed by the CUDA runtime). In short, the unrolled loop contains 3 integer instructions at the top and 7 integer instructions at the bottom, with the original loop body replicated 8 times in between. The original loop body consists of five loads (to constant memory) and 14 FP instructions (11 simple arithmetic, 2 trigonometric, and 1 MAD). Thus, the unrolled loop computes 120 FLOPS using 162 instructions (74% efficiency), while the peak theoretical throughput requires 2 FLOPS per instruction (200% efficiency). Clearly, GPU.Tune performs somewhat better than this analysis suggests it would, because 100% efficiency is required to reach 50% of the Quadro's peak throughput. We assume that the CUDA runtime is responsible for this performance boost. For example, the ptx code uses 3 MUL and 2 ADD instructions to compute the quantity **exp** in Fig. 4(c). As **exp** is expressed in standard sum-of-products form, 1 MUL and 2 MAD instructions are clearly preferred. This transformation alone would boost GPU.Tune's efficiency to 82%.

### 5.7. GPU.Multi

In this final experiment, the voxels are divided into four distinct subsets, with one of four Quadros computing  $\mathbf{F}^H \mathbf{d}$  for each subset. This optimization decreases the time required to compute  $\mathbf{F}^H \mathbf{d}$  to 14.5 s and increases the throughput to over 600 GFLOPS. The acceleration is slightly sub-linear because the overheads (I/O, data marshaling, etc.) represent a significant fraction the time required to compute  $\mathbf{F}^H \mathbf{d}$ . With  $\mathbf{F}^H \mathbf{d}$ 's runtime reduced to just 14.5 s, Amdahl's law is beginning to assert itself.

### 5.8. Improvements in CUDA 1.1

With CUDA 1.1, the performance of the advanced MRI reconstruction is roughly 20% better than with CUDA 1.0. Enhancements to nvcc's register allocation policy are partially responsible for the improved performance. As Table 1 shows, GPU.Tune uses 13 registers per thread in CUDA 1.1 (when the loop is unrolled 8 times), compared to 19 registers per thread in CUDA 1.0 (when the loop is unrolled only 5 times). Additional loop unrolling in CUDA 1.0 is counter-productive, as the number of registers per thread steadily increases from 19 to 29 as the loop unrolling factor increases from 5 to 8. Increasing the per-thread register usage causes a corresponding decrease in utilization, because the number of threads that can execute simultaneously is inversely proportional to the number of registers per thread.

<sup>3</sup> Configurations with non-power-of-2 loop unrolling factors routinely hang in CUDA 1.1 for unknown reasons.

Likewise, enhancements to nvcc's code generation also contribute to improved performance. In CUDA 1.0, nvcc generates four additional integer instructions each time the inner loop of  $\mathbf{F}^H \mathbf{d}$  is unrolled. These instructions compute the base address of the next data sample in constant memory. Each load instruction then uses an integer offset to access the desired element of the data sample (e.g.,  $\mathbf{k}_x$  or  $\mathbf{k}_y$ ). By contrast, in CUDA 1.1, the base address of the next data sample is computed once at the top of the unrolled loop, and the integer offsets are adjusted so that each load accesses the correct data sample. Given that the original loop body contains only 19 ptx instructions, the unnecessary overhead of 4 additional instructions per loop iteration is significant.

## 6. Related work

General-purpose computing on graphics processing units (often termed *GPGPU* or *GPU computing*) supports a broad range of scientific and engineering applications, including physical simulation, signal and image processing, database management, and data mining [27]. Medical imaging was one of the first GPU computing applications. In 1994 Cabral et al. observed that volume rendering essentially performs a generalized Radon transform, while the filtered backprojection algorithm for computed tomography (CT) reconstruction performs an inverse Radon transform. The CT reconstruction based on filtered backprojection achieved a speedup of two orders of magnitude on the SGI RealityEngine [6].

A wide variety of CT reconstruction algorithms have since been accelerated on graphics processors [8,22,23,43], and the Cell Broadband Engine [4,31]. Filtered backprojection algorithms receive further attention in [22,43], while [8] studies the performance of two iterative algorithms for CT reconstruction (the Maximum Likelihood Expectation Maximization and Ordered Sub-set Expectation Maximization algorithms) on the GPU. In [23] the GPU is used to accelerate Simultaneous Algebraic Reconstruction Technique (SART), an algorithm that increases the quality of image reconstruction relative to the conventional filtered backprojection algorithm under certain conditions. SART, which requires significantly more computation than backprojection, becomes a viable clinical option when executed on the GPU. Finally, [4,31] accelerate CT reconstructions based on cone-beam backprojection on the Cell/BE.

By contrast, MRI reconstruction on the GPU has not been studied extensively. Research in this area has focused on accelerating the fast Fourier transform (FFT), which is a key component of many MRI reconstruction algorithms. Speedups on the order of 2x-9x have been reported [19,32,38]. In [35], Sørensen et al. use a GPU to accelerate a gridding algorithm for MRI reconstruction, achieving a substantial speedup over the baseline implementation. GPU-based parallel imaging shows promising results in [16]. Finally, the acceleration of the advanced reconstruction algorithm described in this paper builds on our earlier work with the same algorithm [36,37]. Baskaran et al. have independently observed that the  $\mathbf{F}^H \mathbf{d}$  algorithm can be efficiently mapped to the GPU using a parallelizing compiler [2].

## 7. Conclusions and future work

In many applications, magnetic resonance imaging is limited by high noise levels, imaging artifacts, and long scan times. Advanced image reconstructions, which can operate on arbitrary scan trajectories and incorporate anatomical constraints, can mitigate these limitations at the expense of substantial computation. The computational resources, architectural features, and programmability of the Quadro FX 5600 reduce the time for an advanced reconstruction of non-uniform MRI scan data from nearly 23 min on a quad-core CPU to just over one min on the Quadro, making the reconstruction practical for many clinical applications.

The single-precision floating-point arithmetic and approximate trigonometric operations that help accelerate the advanced reconstruction may, under certain conditions, degrade the quality of the reconstructed image. While we did not observe this phenomenon during our reconstructions of the 3D phantom image, we view further investigation of the advanced reconstruction algorithm's sensitivity to numerical approximations as important future work.

## Acknowledgments

The authors wish to thank Keith Thulborn and Ian Atkinson of the Center for MR Research at the University of Illinois at Chicago for assisting with an earlier version of this paper and for providing the scan trajectory used in some of our experiments. We also thank the Bioengineering Department of the University of Illinois at Urbana-Champaign (UIUC) for providing the *in vivo* data used in our experiments, and thank the National Center for Supercomputing Applications at UIUC for donating time on its Quadro Plex cluster. We acknowledge the support of the Gigascale Systems Research Center, one of five research centers funded under the Focus Center Research Program, a Semiconductor Research Corporation program. Experiments were made possible by generous donations of hardware from NVIDIA and Intel and by NSF CNS grant 05-51665. This work was supported in part by research grants NIH-P41-EB03631-16 and NIH-R01-CA098717. This material is based on work supported under two National Science Foundation Graduate Research Fellowships (Sam Stone, Justin Haldar). Any opinions, findings, conclusions, or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the NSF.

## Appendix. Supplementary data

Supplementary data associated with this article can be found, in the online version, at doi:10.1016/j.jpdc.2008.05.013.

## References

- [1] AMD stream processor. <http://ati.amd.com/products/streamprocessor/index.html>.
- [2] M.M. Baskaran, U. Bondhugula, J. Ramanujam, A. Rountev, P. Sadayappan, A compiler framework for optimization of affine loop nests for general purpose computations on gpus, Technical Report OSU-CISRC-12/07-TR78, Ohio State University, Dec. 2007.
- [3] M.A. Bernstein, K.F. King, X.J. Zhou, Handbook of MRI Pulse Sequences, Elsevier Academic Press, Burlington, MA, 2004.
- [4] O. Bockenbach, M. Knaup, M. Kachelrieß, Implementation of a cone-beam backprojection algorithm on the Cell Broadband Engine processor, in: SPIE Medical Imaging 2007: Physics of Medical Imaging, 2007.
- [5] I. Buck, Brook specification v0.2, October 2003.
- [6] B. Cabral, N. Cam, J. Foran, Accelerated volume rendering and tomographic reconstruction using texture mapping hardware, in: 1994 Symposium on Volume Visualization, 1994.
- [7] Cg. [http://developer.nvidia.com/page/cg\\_main.html](http://developer.nvidia.com/page/cg_main.html).
- [8] K. Chidlow, T. Möller, Rapid emission tomography reconstruction, in: Int'l Workshop on Volume Graphics, 2003.
- [9] P. Concus, G. Golub, D. O'Leary, A Generalized Conjugate Gradient Method for the Numerical Solution of Elliptic Partial Differential Equations, Academic Press, New York, 1976.
- [10] DirectX developer center. <http://www.microsoft.com/directx/>.
- [11] J. Dongarra, Compressed Row Storage CRS. <http://netlib.org/utk/papers/templates/node91.html>.
- [12] J.A. Fessler, S. Lee, V.T. Olafsson, H.R. Shi, D.C. Noll, Toeplitz-based iterative image reconstruction for MRI with correction for magnetic field inhomogeneity, IEEE Transactions on Signal Processing 53 (9) (2005) 3393–3402.
- [13] G. Glover, Simple analytic spiral  $k$ -space algorithm, Magnetic Resonance in Medicine 42 (2) (1999) 412–415.
- [14] J. Haldar, D. Hernando, S.-K. Song, Z. Liang, Anatomically constrained reconstruction from noisy data, Magnetic Resonance in Medicine 59 (2008) 810–818.
- [15] J.P. Haldar, Z. Liang, Joint reconstruction of noisy high-resolution MR image sequences, in: IEEE International Symposium on Biomedical Imaging: From Nano to Macro, 2008, pp. 752–755.

- [16] M. Hansen, D. Atkinson, T. Sørensen, Cartesian SENSE and  $k\text{-}t$  SENSE reconstruction using commodity graphics hardware, *Magnetic Resonance in Medicine* 59 (3) (2008) 463–468.
- [17] M. Hestenes, E. Stiefel, Methods of conjugate gradients for solving linear systems, *Journal of Research of the National Bureau of Standards* 49 (6) (1952) 409–436.
- [18] J.I. Jackson, C.H. Meyer, D.G. Nishimura, A. Macovski, Selection of a convolution function for Fourier inversion using gridding, *IEEE Transactions on Medical Imaging* 10 (3) (1991) 473–478.
- [19] T. Jansen, B. von Rymon-Lipinski, N. Hanssen, E. Keeve, Fourier volume rendering on the GPU using a split-stream FFT, in: 9th International Fall Workshop on Vision, Modeling, and Visualization, 2004.
- [20] C. Koay, J. Sarlls, E. Ozarslan, Three dimensional analytical magnetic resonance imaging phantom in the Fourier domain, *Magnetic Resonance in Medicine* 58 (2007) 430–436.
- [21] J. Meijerink, H. van der Vorst, An iterative solution method for linear systems of which the coefficient matrix is a symmetric M-matrix, *Mathematics of Computation* 31 (137) (1977) 148–162.
- [22] K. Mueller, F. Xu, N. Neophytou, Why do commodity graphics hardware boards (GPUs) work so well for acceleration of computed tomography?, in: SPIE Electronic Imaging 2007, Computational Imaging V Keynote, 2007.
- [23] K. Mueller, R. Yagel, Rapid 3-D cone-beam reconstruction with the simultaneous algebraic reconstruction technique (SART) using 2-D texture mapping hardware, *IEEE Transactions on Medical Imaging* 19 (12) (2000) 1227–1237.
- [24] J. Nickolls, I. Buck, NVIDIA CUDA software and GPU parallel computing architecture, Microprocessor Forum, May 2007.
- [25] NVIDIA corporation, CUDA CUFFT Library, version 1.1, 2007.
- [26] NVIDIA corporation, NVIDIA CUDA Programming Guide, version 1.1, 2007.
- [27] J. Owens, D. Luebke, N. Govindaraju, M. Harris, J. Kruger, A. Lefohn, T. Purcell, A survey of general-purpose computation on graphics hardware, *Computer Graphics Forum* 26 (1) (2007) 80–113.
- [28] K.P. Pruessmann, M. Weiger, P. Börnert, P. Boesiger, Advances in sensitivity encoding with arbitrary  $k$ -space trajectories, *Magn. Res. Med.* 46 (4) (2001) 638–651.
- [29] S. Ryoo, C. Rodrigues, S. Baghsorkhi, S. Stone, D. Kirk, W. Hwu, Optimization principles and application performance evaluation of a multithreaded GPU using CUDA, in: Symposium on Principles and Practice of Parallel Programming, PPoPP, 2008.
- [30] S. Ryoo, C. Rodrigues, S. Stone, S. Baghsorkhi, S. Ueng, J. Stratton, W. Hwu, Optimization space pruning for a multithreaded GPU, in: International Symposium on Code Generation and Optimization, CGO, 2008.
- [31] M. Sakamoto, M. Murase, Parallel implementation for 3-D CT image reconstruction on Cell Broadband Engine, in: International Conference on Multimedia and Expo, 2007.
- [32] T. Schiwietz, T. Chang, P. Speier, R. Westermann, MR image reconstruction using the GPU, in: SPIE Medical Imaging 2006, 2006.
- [33] H. Schomberg, J. Timmer, The gridding method for image reconstruction by Fourier transformation, *IEEE Transactions on Medical Imaging* 14 (3) (1995) 596–607.
- [34] M. Segal, K. Akeley, The OpenGL Graphics System: A Specification (Version 2.0), Silicon Graphics, Inc., October 2004.
- [35] T. Sørensen, T. Schaeffter, K. Noe, M. Hansen, Accelerating the non-equispaced fast Fourier transform on commodity graphics hardware, *IEEE Transactions on Medical Imaging* 27 (4) (2008) 538–547.
- [36] S. Stone, J. Haldar, S. Tsao, W. Hwu, B. Sutton, Z. Liang, Accelerating advanced MRI reconstructions on GPUs, in: Proceedings of International Conference on Computing Frontiers CF, 2008.
- [37] S. Stone, H. Yi, J. Haldar, W. Hwu, B. Sutton, Z. Liang, How GPUs can improve the quality of magnetic resonance imaging, in: First Workshop on General Purpose Processing on Graphics Processing Units, GPGPU, 2007.
- [38] T. Sumanaweera, D. Liu, Medical image reconstruction with the FFT, in: M. Pharr (Ed.), *GPU Gems 2: Programming Techniques for High-Performance Graphics and General-Purpose Computation*, Addison-Wesley, 2005, pp. 765–784.
- [39] B.P. Sutton, D.C. Noll, J.A. Fessler, Fast, iterative image reconstruction for MRI in the presence of field inhomogeneities, *IEEE Transactions on Medical Imaging* 22 (2) (2003) 178–188.
- [40] D. Tarditi, S. Puri, J. Oglesby, Accelerator: Using data parallelism to program GPUs for general-purpose uses, in: Int'l Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS-XII, 2006.
- [41] P. Trancoso, M. Charalambous, Exploring graphics processor performance for general purpose applications, in: Euromicro Symposium on Digital System Design, Architectures, Methods, and Tools, DSD 2005, 2005.
- [42] F.T.A.W. Wajer, Non-Cartesian MRI scan time reduction through sparse sampling, PhD thesis, Technische Universiteit Delft, Delft, Netherlands, 2001.
- [43] X. Xue, A. Cheryauka, D. Tubbs, Acceleration of fluoro-CT reconstruction for a mobile C-Arm on GPU and FPGA hardware: A simulation study, in: SPIE Medical Imaging 2006, 2006.

- [2] J.A. Fessler, B.P. Sutton, Nonuniform fast Fourier transforms using min-max interpolation, *IEEE Transactions on Signal Processing* 51 (2) (2003) 560–574.
- [3] P.C. Lauterbur, Image formation by induced local interactions: Examples employing nuclear magnetic resonance, *Nature* 242 (1973) 190–191.



**S.S. Stone** received his M.S. degree in Electrical and Computer Engineering in 2007 from the University of Illinois at Urbana-Champaign, where he was supported by a National Science Foundation Graduate Research Fellowship. He received the B.S. degree in Computer Engineering in 2003 from Virginia Tech, where he was supported by the Harry Lynde Bradley Scholarship. While at the University of Illinois, his research interests included computer architecture and magnetic resonance image reconstruction. Sam enrolled at Harvard Law School in the fall of 2008.



**J.P. Haldar** is a Ph.D. candidate in electrical and computer engineering at the University of Illinois at Urbana-Champaign. He received the B.S. degree in 2004 and the M.S. degree in 2005, both in electrical engineering from the University of Illinois at Urbana-Champaign. His research interests include statistical signal processing and biomedical inverse problems, with focus on data acquisition, image reconstruction, and parameter estimation schemes for magnetic resonance imaging and spectroscopy.



**S.C. Tsao** is an M.S. candidate in Electrical and Computer Engineering and a graduate research assistant in the IMPACT research group at the University of Illinois at Urbana-Champaign. She received her B.S. in Computer Science and Information Engineering from National Chung Cheng University in Taiwan in 2006. Stephanie is interested in studying and investigating the parallelization of sequential code on GPUs to gain insight into parallelization techniques for many-core systems.



**W.-m.W. Hwu** is a Professor and holds the Walter J. ("Jerry") Sanders III-Advanced Micro Devices Endowed Chair in Electrical and Computer Engineering of the University of Illinois at Urbana-Champaign. His research interests are in the area of architecture, implementation, and compilation for parallel computer systems. He is the director of the IMPACT research group ([www.crcr.uiuc.edu/Impact](http://www.crcr.uiuc.edu/Impact)). For his contributions in research and teaching, he received the ACM SigArch Maurice Wilkes Award, the ACM Grace Murray Hopper Award, the Tau Beta Pi Daniel C. Drucker Eminent Faculty Award, and the ISCA Most Influential Paper Award. He is a fellow of IEEE and ACM. Hwu leads the GSRC ([www.gigascale.org](http://www.gigascale.org)) Concurrent Systems Theme. Hwu received his Ph.D. degree in Computer Science from the University of California, Berkeley in 1987.



**B.P. Sutton** joined the Bioengineering Department at the University of Illinois at Urbana-Champaign in January, 2006. Dr. Sutton received a B.S. in General Engineering from the University of Illinois at Urbana-Champaign in 1998. He earned M.S.'s in Biomedical and Electrical Engineering (2001) and a Ph.D. in Biomedical Engineering from the University of Michigan in 2003. Brad's research focuses on development of novel acquisition techniques in structural and functional brain imaging.



**Z.-P. Liang** received his Ph.D. degree in Biomedical Engineering from Case Western Reserve University in 1989. He is currently Professor of Electrical and Computer Engineering at the University of Illinois at Urbana-Champaign. He is also affiliated with the Beckman Institute for Advanced Science and Technology, the Computational Biophysics Program, and the Department of Bioengineering. Dr. Liang's research interests include magnetic resonance imaging, superresolution image reconstruction using a priori constraints, statistical and learning-based methods for biomedical image analysis, and their application to functional brain mapping, cancer imaging, and cardiac imaging.

## Further reading

- [1] C.B. Ahn, J.H. Kim, Z.H. Cho, High-speed spiral-scan echo planar NMR imaging, *IEEE Transactions on Medical Imaging* 5 (1) (1986) 2–7.