

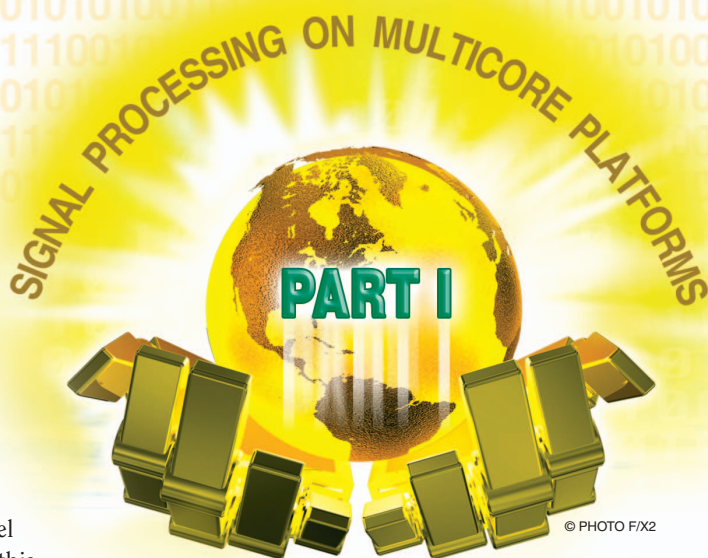
Discrete Fourier Transform on Multicore

A review of optimizations necessary
for good multicore performance

This article gives an overview on the techniques needed to implement the discrete Fourier transform (DFT) efficiently on current multicore systems. The focus is on Intel-compatible multicores, but we also discuss the IBM Cell and, briefly, graphics processing units (GPUs). The performance optimization is broken down into three key challenges: parallelization, vectorization, and memory hierarchy optimization. In each case, we use the Kronecker product formalism to formally derive the necessary algorithmic transformations based on a few hardware parameters. Further code-level optimizations are discussed. The rigorous nature of this framework enables the complete automation of the implementation task as shown by the program generator Spiral. Finally, we show and analyze DFT benchmarks of the fastest libraries available for the considered platforms.

INTRODUCTION

The evolution of computing platforms is at an historic inflection point: after years of exponential growth, central processing unit (CPU) frequency has stalled due to physical limitations. However, the theoretical floating point peak performance, a critical measure for the processing abilities of a platform, continues to increase at the pace predicted by Moore's law. The reason is an increase in parallelism in the form of multiple processor cores and vector processing abilities [1].



© PHOTO F/X2

The consequences for signal processing software are dramatic: it means the end of free speedup for legacy software and a dramatically increased difficulty of writing high-performance code since the programmer now has to use multiple threads and vector instruction sets and tune the code to the memory hierarchy. Unfortunately, this process is platform specific: performance does not port easily. Failure to apply these optimizations by hand can result in dramatic performance losses, as shown in Figure 1. The figure shows the performance (in gigafloating point operations per second (Gflop/s); higher is better) for four implementations of the DFT, arguably the most important signal processing function and focus of this article. All implementations use fast Fourier

transform (FFT) algorithms with roughly the same number of operations. Yet, the performance difference between the best and worst is a factor of 12–35. The bottom line shows the implementation from Numerical Recipes [2] based on a standard radix-2 iterative FFT. The next line is one of the best scalar (standard C code) implementations and is five times faster, since it is tuned to the memory hierarchy. The next one gains up to another factor of three by using streaming single instruction, multiple data (SIMD) extensions (SSE) vector instructions. Finally, the fastest code uses up to four processor cores and speeds up the code by another factor of three for larger sizes. The performance drop for very large sizes is due to inevitable cache misses once the working set cannot be held in the largest cache.

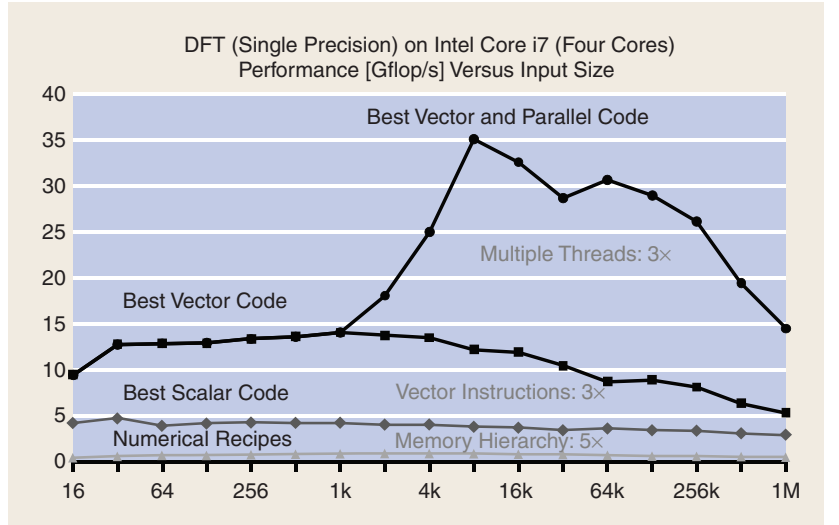
The main reasons for the achieved speedup are algorithmic: all top three lines are based on nonstandard FFT variants whose structure matches the architectural constraints imposed by this multicore platform. Equally notable, the source code for the fast implementations were not written by a human but automatically produced by a tool called Spiral [3]–[5] that we developed and whose framework underlies this article.

We present a tutorial-style overview of the optimizations needed to achieve good performance on Intel multicore platforms (top line in Figure 1), the IBM Cell, and briefly touch on GPUs. Starting from a standard recursive FFT, we use the Kronecker product formalism to derive the necessary structural or algorithmic optimizations. Necessary code-level optimizations are also discussed. In this way, we address all three key problems identified in Figure 1: multiple cores, vector instruction sets, and the memory hierarchy.

Finally, we show and analyze performance benchmarks of the fastest libraries for Intel Core, IBM Cell, and GPUs. Specifically, we show benchmarks for the FFT in the west (FFTW) [6]–[8], Intel’s Integrated Performance Primitives (IPP), Spiral-generated libraries, FFT for the IBM Cell (FFTC) Broadband Engine [9], and [10] for the GPU. Other libraries and related work will be introduced throughout the article.

FRAMEWORK

There are two fundamentally different ways of representing linear transforms like the DFT: as matrix-vector products or using summations. Correspondingly, fast algorithms are represented either with a matrix formalism as in [11] and [12] or as nested summations as in most signal processing books. In this section, we introduce the matrix formalism and use it to express various classical FFTs. Later in the article, we use this formalism as a tool for structural manipulation of FFTs to derive variants that can be efficiently mapped to current multicore systems. The formal nature of the approach makes computer generation of FFT libraries possible, as we have demonstrated with the tool Spiral [13] that we developed and briefly discuss later.



[FIG1] Naive DFT implementations based on only minimizing the operations count underperform considerably on modern multicore CPUs.

We focus on the DFT but note that the framework extends to a large class of linear transforms [3].

DFT

The DFT of n input samples x_0, \dots, x_{n-1} is defined in summation form as

$$y_k = \sum_{0 \leq \ell < n} \omega_n^{k\ell} x_\ell, \quad 0 \leq k < n, \quad (1)$$

with $\omega_n = \exp(-2\pi j/n)$. Stacking the x_ℓ and y_k into vectors x and y yields the equivalent form of a matrix-vector product

$$y = \text{DFT}_n x, \quad \text{DFT}_n = [\omega_n^{k\ell}]_{0 \leq k, \ell < n}. \quad (2)$$

We drop x and y and simply think of the matrix DFT_n as the transform, implicitly assuming that it is multiplied to x . Fast algorithms are now expressed as factorizations of DFT_n using the formalism described next.

MATRIX FORMALISM

We denote with I_n the $n \times n$ identity matrix, and the butterfly matrix with

$$\text{DFT}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}. \quad (3)$$

The Kronecker product of matrices A and B is defined as

$$A \otimes B = [a_{k,\ell} B], \quad \text{for } A = [a_{k,\ell}].$$

It replaces every entry $a_{k,\ell}$ of A by the matrix $a_{k,\ell} B$. Most important are the cases where A or B is the identity. As examples, we show in Figure 2 the structure and dataflow of the “dual” constructs $I_4 \otimes \text{DFT}_2$ and $\text{DFT}_2 \otimes I_4$. The former is obviously parallel, the latter has vector structure: it can be viewed as a single DFT_2 operating on vectors of length four instead of scalars. This will be crucial later.

The stride permutation matrix L_m^{mn} permutes the elements of the input vector as $in + j \rightarrow jm + i$, $0 \leq i < m$, $0 \leq j < n$. If the vector x is viewed as an $n \times m$ matrix, stored in row-major order, then L_m^{mn} performs a transposition of this matrix. Further, if P is a permutation (matrix), then $A^P = P^{-1}AP$ is the conjugation of A with P . Note that $y = A^P x$ implies $Py = APx$, i.e., A is performed on vectors permuted with P .

There are various identities connecting these constructs [14] as shown in Table 1. For example, (9) shows how to translate the duals in Figure 2 into each other. The Kronecker product naturally arises in two-dimensional (2-D) and three-dimensional (3-D) DFTs that, respectively, can be written as

$$\text{DFT}_{m \times n} = \text{DFT}_m \otimes \text{DFT}_n, \quad (13)$$

$$\text{DFT}_{k \times m \times n} = \text{DFT}_k \otimes \text{DFT}_m \otimes \text{DFT}_n. \quad (14)$$

RECURSIVE DIT FFT

The matrix formalism can be used to express FFTs as factorizations of the matrix DFT_n in (2). As an example, the recursive general-radix decimation in time (DIT) Cooley-Tukey FFT for $n = km$ is

$$\text{DFT}_{km} = (\text{DFT}_k \otimes I_m) T_m^n (I_k \otimes \text{DFT}_m) L_k^n. \quad (15)$$

Here, T_m^n is a diagonal matrix containing the twiddle factors. Figure 3 shows the special case $16 = 4 \times 4$. Using this algorithm on an input vector x means (reading Figure 3 from right to left):

- L_4^{16} : Read x at stride four.
- $I_4 \otimes \text{DFT}_4$: Apply four DFT_4 s on consecutive chunks.
- T_4^{16} : Scale by the diagonal elements of T_4^{16} .
- $\text{DFT}_4 \otimes I_4$: Apply four DFT_4 s at stride four.

RECURSIVE FFT VARIANTS

Looking at Table 1, it becomes clear that many variants of (15) can be derived. For example, simple transposition using that

$\text{DFT}_n^\top = \text{DFT}_n$ yields the decimation in frequency (DIF) Cooley-Tukey FFT

$$\text{DFT}_{km} = L_m^n (I_k \otimes \text{DFT}_m) T_m^n (\text{DFT}_k \otimes I_m). \quad (16)$$

The four step algorithm [11], [15], [16]

$$\text{DFT}_{km} = (\text{DFT}_k \otimes I_m) T_m^n L_k^n (\text{DFT}_m \otimes I_k) \quad (17)$$

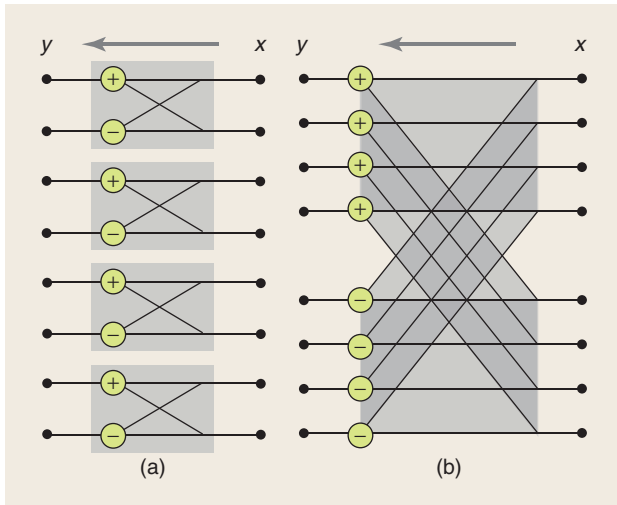
was originally developed for vector computers. It produces the longest possible unit stride vector operations at the cost of a transposition. The six step algorithm [11], [17]

$$\text{DFT}_{km} = L_k^n (I_m \otimes \text{DFT}_k) L_m^n T_m^n (I_k \otimes \text{DFT}_m) L_k^n. \quad (18)$$

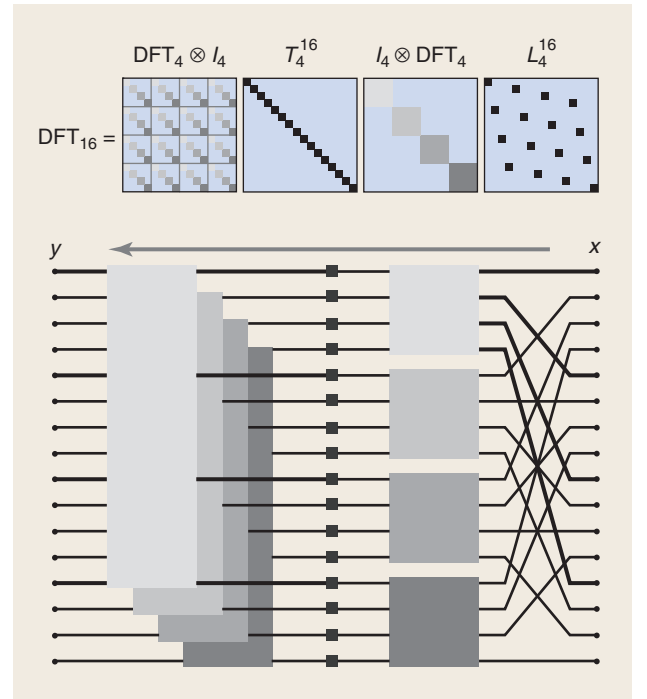
was originally developed for distributed memory machines. It produces fully local computation at the cost of three global

[TABLE 1] FORMULA IDENTITIES TO MANIPULATE FFT ALGORITHMS. A IS $n \times n$, AND B AND C ARE $m \times m$. A^\top IS THE TRANSPOSE OF A .

$$\begin{aligned} (BC)^\top &= C^\top B^\top & (4) \\ (A \otimes B)^\top &= A^\top \otimes B^\top & (5) \\ I_{mn} &= I_m \otimes I_n & (6) \\ A \otimes B &= (A \otimes I_m)(I_n \otimes B) & (7) \\ A \otimes (BC) &= (A \otimes B)(A \otimes C) & (8) \\ A \otimes B &= L_n^{mn}(B \otimes A)L_m^{mn} & (9) \\ (L_m^{mn})^{-1} &= L_n^{mn} & (10) \\ L_n^{kmn} &= (L_n^{kn} \otimes I_m)(I_k \otimes L_n^{mn}) & (11) \\ L_{km}^{kn} &= (I_k \otimes L_m^{mn})(L_k^{kn} \otimes I_m) & (12) \end{aligned}$$



[FIG2] Dataflow (right to left) of a parallel and its "dual" vector construct. (a) $y = (I_4 \otimes \text{DFT}_4)x$. (b) $y = (\text{DFT}_4 \otimes I_4)x$.



[FIG3] Cooley-Tukey FFT (15) for $16 = 4 \times 4$ as matrix formula and as (complex) data-flow graph (from right to left). Some lines are bold to emphasize the strided access.

transpositions (all-to-all data exchanges). Both (17) and (18) are derived from (15) using (9) and (10). The required matrix transposition can be blocked using (11) and (12) for more efficient data movement.

For a 2-D DFT, applying (7) to (13) yields the row-column algorithm

$$\text{DFT}_{m \times n} = (\text{DFT}_m \otimes I_n)(I_m \otimes \text{DFT}_n). \quad (19)$$

The 2-D vector radix algorithm [18] is derived from (6)–(10)

$$\text{DFT}_{mn \times rs} = (\text{DFT}_{m \times r} \otimes I_{ns})^{I_m \otimes L_r^m \otimes I_s} (T_n^{mn} \otimes T_s^{rs}) (I_{mr} \otimes \text{DFT}_{n \times s})^{I_m \otimes L_r^m \otimes I_s} (L_m^{mn} \otimes L_r^{rs}).$$

Higher-dimensional versions are derived similarly, and the associativity of \otimes gives rise to more variants.

ITERATIVE FFT ALGORITHMS

A second class (historically earlier) of FFT algorithms based on (15) are the iterative FFT algorithms, obtained by recursively expanding (15) and again using Table 1. The simplest are radix- r forms (usually $r = 2, 4, 8$), which require an FFT size of $n = r^k$ (more complicated mixed-radix variants always exist).

The DIT triple-loop FFT [11], [19],

$$\text{DFT}_{r^k} = R_r^{r^k} \prod_{i=0}^{k-1} D_i^{r^k} (I_{r^{k-i-1}} \otimes \text{DFT}_r \otimes I_{r'}) \quad (20)$$

is the simplest iterative algorithm. $R_r^{r^k}$ is the radix- r digit reversal permutation and $D_i^{r^k}$ contains the twiddle factors in the i th stage. A radix-2 version is implemented by Numerical Recipes [2]. A variant of (20) is the Pease FFT [20], [11], which has constant geometry, i.e., the control flow is independent of the stage; however, it also requires the digit reversal permutation. It was originally developed for parallel computers, and its regular structure makes it a good choice for field-programmable gate arrays (FPGAs) or application-specific integrated circuits (ASICs).

The Stockham FFT [11], [21]

$$\text{DFT}_{r^k} = \prod_{i=0}^{k-1} (\text{DFT}_r \otimes I_{r^{k-i-1}}) D_i^{r^k} (L_r^{r^{k-i}} \otimes I_{r'}) \quad (21)$$

is self-sorting, i.e., it does not have a digit reversal permutation. It was originally developed for vector computers.

DFT VARIANTS AND OTHER FFTs

In practice, several variants of the DFT in (2) are needed including forward/inverse, interleaved/split complex format, for complex/real input data, in place/out-of-place ($y = x$ or not), and others [22]. Fortunately, most of these variants are close to the standard DFT in (2), so fast code for the latter can be adapted. An exception is the DFT for real input data, which has its own class of FFTs (see [4] for an overview using the above formalism). This article focuses on the standard one-dimensional (1-D) interleaved (alternating real and imaginary parts) complex DFT in (2).

DFT algorithms fundamentally different from (15) include prime-factor (n is a product of coprime factors), Rader (n is prime), and Bluestein or Winograd (any n) FFTs and can also be

expressed in the above formalism [11]. In practice, these are mostly used for small sizes < 32 , which then serve as building blocks for large composite sizes via (15).

FROM MATRIX FORMULAS TO IMPLEMENTATIONS

Table 2 shows how to translate matrix formulas into basic sequential loop code. However, strictly applying the upper part of the table will lead to low performance. The last three entries show optimized translations. The last two show how a diagonal scaling is always fused with subsequent loops and how permutations are (almost) always done as readdressing in the subsequent loop. The shown translation of $y = (L_m^{mn} \otimes I_k)x$ replaces every scalar in the shown translation of $y = L_m^{mn}x$ by a vector of length k , similar as in Figure 2(b).

Further, different types of implementations are possible, briefly discussed next.

SINGLE-SIZE KERNELS

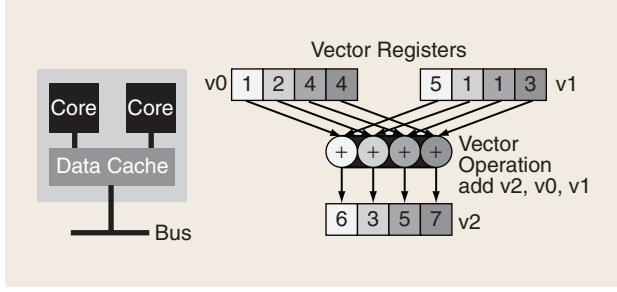
For small input sizes ≤ 32 or 64, matrix formulas are often implemented as fully unrolled code blocks. In this case array scalarization and, to a lesser extent, algebraic optimizations and scheduling are used to achieve best performance and can be completely automated [3], [7], [23], [24]. For example, a DFT_8 kernel is implemented in a few tens of lines of code.

SINGLE-SIZE LOOP CODE

If the input size is known in advance, a fixed formula or data flow can be chosen and implemented using nested loops arising from

[TABLE 2] FROM MATRIX FORMULAS TO CODE. THE SUBSCRIPT OF A , B SPECIFIES THE (SQUARE) MATRIX SIZE. $x[b:s:e]$ DENOTES THE SUBVECTOR OF x STARTING AT b , ENDING AT e , EXTRACTED AT STRIDE s . THE DIAGONAL ELEMENTS OF D ARE STORED IN AN ARRAY WITH THE SAME NAME.

MATRIX FORMULA	MATLAB PSEUDOCODE
$y = (A_n B_n)x$	<pre>t[0:1:n-1] = B(x[0:1:n-1]); y[0:1:n-1] = A(t[0:1:n-1]);</pre>
$y = \left(\prod_{i=0}^{k-1} A_i\right)x$	<pre>y = x; for (i=0; i<k; i++) x = y; y = A(i, x);</pre>
$y = (I_m \otimes A_n)x$	<pre>for (i=0; i<m; i++) y[i*n:1:i*n+n-1] = A(x[i*n:1:i*n+n-1]);</pre>
$y = (A_m \otimes I_n)x$	<pre>for (i=0; i<n; i++) y[i:n:i+m*n-n] = A(x[i:n:i+m*n-n]);</pre>
$y = D_n x$	<pre>for (i=0; i<n; i++) y[i] = Dn[i]*x[i];</pre>
$y = L_m^{mn} x$	<pre>for (i=0; i<m; i++) for (j=0; j<n; j++) y[i+m*j:1:i+m*j] = x[n*i+j:1:n*i+j];</pre>
$y = L_m^{mn} \otimes I_k x$	<pre>for (i=0; i<m; i++) for (j=0; j<n; j++) y[k*(i+m*j):1:k*(i+m*j)+k-1] = x[k*(n*i+j):1:k*(n*i+j)+k-1];</pre>
$y = (A_m \otimes I_n) D_{mn} x$	<pre>for (i=0; i<n; i++) t= Dmn[i:n:i+m*n-n] * x[i:n:i+m*n-n]; y[i:n:i+m*n-n] = A(t);</pre>
$y = (I_m \otimes A_n) L_m^{mn} x$	<pre>for (i=0; i<m; i++) y[i*n:1:i*n+n-1] = A(x[i:m:i+m*n-m]);</pre>



[FIG4] Shared cache in a multicore CPU and SIMD vector extensions.

tensor products and iterative products. As discussed before, scaling and readdressing is merged into kernels for high performance. These implementations can be generated automatically using Spiral [3], [23], [25]. One such single-size implementation can require up to a few hundred lines of code.

GENERAL-SIZE LOOP CODE

Only iterative algorithms lend themselves to general-size loop code implementation. Again, readdressing is folded into the computational kernel. An example is the Numerical Recipes code [2], which implements (20) for $r = 2$ as a triple loop in about one page of C code.

GENERAL-SIZE RECURSIVE CODE

Translating formulas into recursive code is complicated but is at the heart of several high-performance portable general-size FFT libraries [8], [26]–[28]. A tutorial for a recursive radix-4 FFT is given in [29], leading to about two pages of C code. Extension to vector and multicore platforms considerably increases the code size: for example, FFTW contains more than 200,000 lines of code. The implementation of such libraries was also automated using Spiral [4], [28].

MAPPING FFTs TO MULTICORE CPUs

Historically, the Kronecker product formalism was used to develop FFTs for parallel target platforms such as small-scale and massive multiprocessors and vector computers [11], [16], [30]. We now discuss how to extend this approach to state-of-the-art multicore CPUs. The new hardware characteristics that need to be captured are: 1) multiple cores communicating through shared caches or explicit messages, 2) SIMD short vector instructions, and 3) the memory hierarchy and its transfer restrictions, such as caches and DMA-based streaming memory. We will address each of these features in three steps. First, we identify relevant hardware parameters. Second, we identify a set of matrix formulas that can be mapped efficiently for these parameters. Third, we derive a variant of the recursive Cooley-Tukey FFT (15) that is a member of this set. In each case, we also briefly discuss the mapping to actual code including further relevant code-level optimizations. In the section “Benchmarks on Multicore CPUs,” we then instantiate the concepts and algorithms to an Intel Core and the Cell BE and briefly discuss GPUs and FPGAs.

Choosing recursive algorithms is not a requirement. It is possible to start from iterative algorithms or combine one or two steps of recursion with an iterative algorithm and achieve reasonable performance (as demonstrated by several vendor libraries). However, many current high-performance libraries for cache-based machines implement the recursive FFT algorithms [3], [4], [6], [26], [27], [31] discussed here.

PARALLELISM: MULTIPLE CORES

The multicore CPUs we target may have shared caches (Figure 4), private caches, or scratchpads (local stores) with data being transferred in packets. Cache coherent architectures transfer data implicitly between private and shared caches as required. Data transfer between scratchpads has to be managed explicitly by the programmer. In each case, to obtain best performance it is crucial to ensure that the whole data content in a transfer (e.g., cache line or DMA packet) is used by the receiver (spatial locality) and that the number of transfers is minimized (temporal locality).

MACHINE MODEL

We assume that the packet size is a multiple of an atomic packet size of μ complex numbers. For instance, on a cache-based memory hierarchy, a cache coherency event, a cache miss, or an eviction always transmits a whole cache line (e.g., 64 B translates into $\mu = 8$ for complex single-precision). In scratchpad-based systems like the Cell, DMA packets need to be of sufficient size: to yield reasonable performance, a Cell DMA transfer between SPEs should be at least 128 B ($\mu = 16$ for complex single-precision numbers) and a multiple of 16.

We consider CPUs with p cores. Well-designed parallel code is load balanced (all cores have the same amount of work), with minimal data transmission between cores, performed in packets of size μ . On shared memory multicores this implies that the code is free of false sharing (two cores accessing different elements in the same cache line).

Matrix formulas solely built from

$$I_p \otimes A, D_n, P \otimes I_\mu \text{ (} P \text{ permutation, } D_n \text{ diagonal)}$$

with A a $m \times n$ matrix and $u | m, n$ can be implemented efficiently as parallel code; we call them parallel constructs. Namely, $I_p \otimes A$ is load balanced and embarrassingly parallel [see Figure 2(a)], i.e., it does not require any communication. The same holds for scaling by D_n . Finally, the communication pattern $P_n \otimes I_\mu$ transmits entire packets of size μ between cores. Note that products of parallel constructs are again parallel constructs.

MULTICORE COOLEY-TUKEY FFT

We now state a multicore FFT built exclusively from parallel constructs, derived using Table 1 [5]

$$\text{DFT}_{mm} = (I_p \otimes (\text{DFT}_m \otimes I_{n/p}))^{((L_p^m \otimes I_{n/p\mu}) \otimes I_\mu)} T_m^{mm} (I_p \otimes (I_{m/p} \otimes \text{DFT}_n) L_{m/p}^{m/p}) ((L_p^m \otimes I_{m/p\mu}) \otimes I_\mu). \quad (22)$$

Implementation of (22) on a cache-based system relies on the cache coherency protocol to transmit cache lines of length μ between cores and requires a global barrier. Implementation on a scratchpad-based system requires explicit sending and receiving of the data packets, and depending on the communication interface additional synchronization may be required.

Equation (22) can be used as outermost recursion to enable multicore parallelization. The smaller DFTs are then expanded using the short vector Cooley-Tukey FFT (23) or the vector recursion (29) shown later in this section.

Historically, the Pease and the six step FFT (18) were starting points for parallel iterative or recursive implementations, but due to changed tradeoffs these algorithms are no longer a good choice in many cases.

MAPPING TO C CODE

OpenMP [32] is a good choice for parallel code if it is well supported by the target platform's compiler. OpenMP allows the programmer to declare certain loops to be parallel and to specify variables as shared or private. It enables the programming of sophisticated parallel software without needing to deal with lower-level threading details. As example, the formula $I_4 \otimes \text{DFT}_2$ in Figure 2(a) is translated into the following OpenMP program snippet:

```
double x[8], y[8];
#pragma omp parallel for
for (int i=0; i<4; i++)
{
  y[2*i] = x[2*i] + x[2*i+1];
  y[2*i+1] = x[2*i] - x[2*i+1];
}
```

Note that only a C `#pragma` is inserted to instruct the compiler to parallelize the `for` loop. If OpenMP is turned off, these pragmas are ignored and the program becomes sequential. Because of the regular structure of FFTs only a few more issues have to be addressed for efficient OpenMP parallelization. For correctness, the sharing and privatization of variables with the OpenMP `shared` and `private` clauses has to be done properly to avoid race conditions and other problems. For performance, scheduling hints can be provided [e.g., `schedule(static)`]. In addition, if only a subset of the available cores is to be used, affinity can be set to choose the subset. For example, two threads operating on the same data should be physically close, i.e., share a high level of the memory hierarchy.

On some target platforms no OpenMP compiler may be available, in which case one must use threading libraries like the portable POSIX threads (pthreads) [33] library or use operating system threading interfaces to build the required parallel loops and barriers. While this approach can yield a slight performance advantage, it requires understanding of the target architecture, its memory consistency model, and cache coherency protocols.

On the Cell processor the programmer needs to manage and synchronize threads on the PPE and the SPEs and perform data movements via DMA transfers. Due to the Cell's

unconventional architecture, libraries for it must be adapted to take advantage of its features. While there have been some programming paradigms ported to the Cell (including some function off-loading interfaces), for the class of programs discussed in this article, performance is best obtained by avoiding the overhead of such interfaces.

SIMD VECTORIZATION

Most multicore CPUs include vector instruction sets. SIMD vector extensions add vector registers (two-way double or four-way float on the Core i7 and the Cell), and much longer vectors will be available in the near future (e.g., 16-way single precision on Intel's upcoming Larrabee GPU, and four-way double and eight-way single precision in AVX on the next generation of Intel multicore CPUs). Vector instructions then operate on these registers in parallel, providing high potential speed-up. A four-way vector addition is shown in Figure 4.

MACHINE MODEL

To obtain best performance on vector extensions, data should be loaded and stored with vector memory operations that transfer complete, naturally aligned vectors. Unaligned and subvector accesses are expensive. All operations on the vector registers should be vector operations (vector addition, subtraction, and multiplication). Data reorganization within registers (shuffles) are needed for FFTs but should be minimized.

For this article, we restrict ourselves to what we call complex vectorization. We denote the machine vector length with ν , meaning ν complex numbers are packed into a vector register of length 2ν ; e.g., for four-way float SSE, $\nu = 2$.

All formulas built solely from

$$A \otimes I_\nu, D_n \text{ (complex diagonal), and } L_\nu^2$$

can be implemented efficiently with vector instructions; we call them vector constructs. Moreover, if A and B are vector constructs, then AB and $I_n \otimes A$ are vector constructs.

First, $A \otimes I_\nu$ is naturally vectorized [e.g., Figure 2(b)]: vector code for $y = (A \otimes I_\nu)x$ can be obtained from scalar code implementing $y = Ax$ by simply replacing all scalar operations by the corresponding vector operations, and all scalar variables by vector variables. Second, we assume that $y = D_n x$ can be implemented efficiently for $\nu|n$. This is a reasonable assumption: for instance, the SSE 4.2 instruction set implemented by the Core i7 contains instructions for the efficient mapping of complex multiplications. Third, $y = L_\nu^2 x$ can always be implemented with a small number of vector instructions [34].

SHORT VECTOR COOLEY-TUKEY FFT

We show a short vector FFT algorithm that is built from vectorizable constructs, derived from (15) using Table 1. It requires only a small number of in-register shuffles [35]

$$\text{DFT}_{mn} = ((\text{DFT}_m \otimes I_{n/\nu}) \otimes I_\nu) T_n^{mn} (I_{m/\nu} \otimes (I_{n/\nu} \otimes L_\nu^2) (L_{n/\nu}^n \otimes I_\nu) (\text{DFT}_n \otimes I_\nu) (L_{m/\nu}^{mn/\nu} \otimes I_\nu). \quad (23)$$

This FFT is composable with memory hierarchy optimized FFTs. Namely, inserting (23) into the vector recursion (29) shown later yields again a vector construct.

A somewhat more complicated real (using a real representation of matrices) short vector FFT is derived in [36] and [37].

Traditional vector algorithms like the Four Step algorithm (17) or the Stockham algorithm (21) were designed for traditional vector computers with much longer vectors. Due to the expensive permutations, they are not a good choice for short-vector SIMD architectures.

MAPPING TO C CODE

The most convenient way to efficiently use SIMD extensions is through intrinsic function interfaces provided by most high-performance compilers. For instance, the Intel C++ compiler, Microsoft's VisualStudio C compiler, IBM's XL C compiler, and the GNU C compiler provide such an interface for the supported SIMD extensions.

The programmer uses a data type and function abstraction of the SIMD extensions to implement C code. The compiler understands the data types and special functions and maps the C program to the respective instructions. In this scenario, the programmer must select the appropriate instructions and make sure machine restrictions like data alignment are met. However, the programmer does not have to directly use assembly and thus is spared from register allocation and instruction scheduling. For example, the formula $\text{DFT}_2 \otimes I_4$ in Figure 2(b) is implemented by the following C program snippet using intrinsics for the Intel C++ compiler. `__m128` is a built-in data type to abstract XMM vector registers, and `_mm_add_ps()` abstracts the SSE instruction `addps` through a function call

```
__m128 x[2], y[2];
y[0] = _mm_add_ps(x[0], x[1]);
y[1] = _mm_sub_ps(x[0], x[1]);
```

MEMORY HIERARCHY

Our target multicore CPUs have a memory hierarchy with multiple cores sharing the off-chip bandwidth. Machines with memory hierarchies present algorithm designers with the following two challenges:

- *Temporal locality*: Faster memory levels are smaller, and the working set must be blocked to fit into that level to minimize data transfers.
- *Spatial locality*: Data transfer between memory hierarchy levels happens in packets. This implies that transferred packets should be fully used to avoid wasting memory bandwidth.

MACHINE MODEL

A memory hierarchy can have multiple levels. For a given level, we call the capacity N if it can hold the working set for the computation of $y = Ax$ for an $N \times N$ matrix A . This implies that the input vector x , the output vector y , and all necessary temporary arrays and constants fit into that cache level. For instance, if we

consider double-precision, one (complex) value is 16 B. If A is a DFT, N is the cache size divided by 64 (assuming a factor of four space overhead). As before, we assume that data is transferred between the current level and the next lower level in the memory hierarchy in packets of μ complex numbers. Moreover, if it is a set-associative cache, it can hold α lines of μ elements in the same set, and hence there are $\sigma = N/(\alpha\mu)$ sets.

We call formula A a memory construct if, during the computation of $y = Ax$, elements of x are loaded once and never stored, and elements of y are never loaded and stored once. Obvious memory constructs are

$$A_n (n \leq N), P \otimes I_\mu (P \text{ a permutation}), D_n. \quad (24)$$

The first has a sufficiently small working set. The second loads or stores complete packets. The last, diagonal scaling, poses no problems.

One problematic construct in (15) has the form $A_m \otimes I_n$. As Table 2 shows, the loop body accesses data at stride n yielding poor spatial locality unless $m \leq \alpha$, which is very restrictive. However, at the expense of some overhead, this condition can be relaxed to $m \leq N/\mu = \sigma\alpha$ through buffering. It is done by first tiling the loop by μ and then copying the working set for the innermost loop into contiguous memory. The tiled loop corresponds to the formula

$$A_m \otimes I_n = (I_{n/\mu} \otimes (A_m \otimes I_\mu))^{(I_{n/\mu} \otimes I_\mu)}. \quad (25)$$

Buffering means that in the above formula the conjugation is implemented using actual copy operations (in contrast to translating them into reindexing) based on the third-last entry in Table 2. The resulting pseudocode snippet is shown next. On the Cell the copy operations are translated into DMA instructions instead

```
double x[m*n], y[m*n];
for (j=0; j<n/mu; j++)
{ // allocate buffers
  double u[m*mu], v[m*mu];
  // copy into buffer
  for (k=0; k<m; k++)
    u[k*mu:1:k*(mu+1)-1] =
      x[j*mu+k*n:1:j*(mu+1)-1+k*n];
  // compute A on buffered contiguous data
  for (i=0; i<mu; i++)
    v[i:mu:i+m*mu-mu] =
      A(u[i:mu:i+m*mu-mu]);
  // copy data back
  for (k=0; k<m; k++)
    y[j*mu+k*n:1:j*(mu+1)-1+k*n] =
      v[k*mu:1:k*(mu+1)-1];
}
```

The other problematic construct in (15) has the form $(I_m \otimes A_n)L_m^{mm}$ and can be handled similarly: the loop is again

tilled but only the load side has a strided access and needs to be buffered. Formally,

$$(I_m \otimes A_n) L_m^{mn} = (I_{m/\mu} \otimes L_\mu^{n\mu} (A_n \otimes I_\mu)) (L_{m/\mu}^{mn/\mu} \otimes I_\mu) \quad (26)$$

expresses the tiling and again the permutation $L_{m/\mu}^{mn/\mu} \otimes I_\mu$ is implemented using explicit copy operations.

MEMORY HIERARCHY OPTIMIZATIONS

To obtain FFT algorithms suitable for the memory hierarchy (i.e., the algorithm is a memory construct), we start with (15)

$$\text{DFT}_{mn} = (\text{DFT}_m \otimes I_n) T_n^{mn} (I_m \otimes \text{DFT}_n) L_m^{mn}. \quad (27)$$

$\text{DFT}_m \otimes I_n$ is a memory construct for $m \leq \alpha$ or $m \leq \sigma\alpha$ if buffering is applied. This explains why relatively small values of m (the radix) work well in practice. $I_m \otimes \text{DFT}_n$ is a memory construct for $n \leq N$. If $n > N$, recursive application of (27) will eventually yield $n \leq N$, producing another memory construct $\text{DFT}_m \otimes I_n$ at each recursion step along the way. This suggests that the largest possible radix m is a good choice in each step. We show a two-level recursion for further discussion

$$\begin{aligned} \text{DFT}_{kmn} &= (\text{DFT}_k \otimes I_{mn}) T_{mn}^{kmn} \\ &\quad (I_k \otimes (\text{DFT}_m \otimes I_n) T_n^{mn} (I_m \otimes \text{DFT}_n) L_m^{mn}) L_k^{kmn}. \end{aligned} \quad (28)$$

In (27), the only nonmemory construct is one L_m^{mn} at every recursion level with $mn > N$. Using buffering at each step incurs too much overhead. Buffering the right-most L s in (28) jointly is not possible. One solution is to give up on spatial locality: all right-most L s are fused and merged into the first loop as explained before. A better solution is to translate it into a memory construct, which is indeed possible and done, e.g., in FFTW [6]. Namely, the entire second line in (28) is translated into the vector recursion

$$\begin{aligned} & (I_k \otimes (\text{DFT}_m \otimes I_n) T_n^{mn}) \\ & (I_k^{km} \otimes I_n) (I_m \otimes (I_k \otimes \text{DFT}_n) L_k^{kn}) (L_m^{mn} \otimes I_k) \end{aligned} \quad (29)$$

with $\mu \mid n, k$. This is repeated until the problematic $(I_k \otimes \text{DFT}_n) L_k^{kn}$ is small enough to be a memory construct. The formula manipulation leading to (29) manifests itself as loop splitting and loop exchange in the equivalent code [6]. Historically, the iterative triple loop algorithm (20) was used to compute FFTs on a single CPU. However, once the data set does not fit in cache, cache thrashing occurs and the performance drops drastically.

MAPPING TO C CODE

The structural optimization ensures that the algorithm has good cache locality. When mapping to code, the following additional issues have to be addressed: 1) how to create efficient basic blocks, 2) how to exploit degrees of freedom, and 3) how to handle constants (twiddle factors).

On modern, deeply pipelined superscalar processors, the recursive FFT has to be terminated with a basic block that is sufficiently

large but does not cause instruction cache misses. Experiments show that a DFT of a size between 4ν and 32ν (ν is the SIMD vector length) is a good choice. The basic block is obtained by unrolling an FFT with minimal operations count and performing scalar replacement to enable efficient register allocation and scheduling. Additional small gains may be achievable by C code scheduling, reduction of the needed constants, and a few other techniques. This process was automated in [7] and [23]. The downside is considerably increased code size. For example, FFTW [6] requires several megabytes of C code to implement 1-D FFTs based on (15).

The FFTs (22), (23), (25), (26), (29) contain degrees of freedom (mainly the respective radix, if and where to buffer, and when to terminate with a basic block) that can be searched over for further platform adaptation. Dynamic programming has been proven efficient in most cases [3], [6].

Finally, the twiddle constants are usually precomputed except for possibly very large sizes for which the FFT becomes memory bound. In this case, computation on the fly can yield considerable improvements. The decision is again handled by search in FFTW and Spiral-generated libraries [4].

MAPPING FFTs TO GPUS AND FPGAs

On early machines, large programs with complex loop structures were expensive, and memory access was relatively cheap. Multiple passes through the entire data set were acceptable while recursive functions were hard to implement and expensive. Thus, iterative FFT algorithms like (20) and (21) were the best choice and were developed first. Current machines with streaming memory (GPUs) or small memories (DSP processors or embedded processors) operate in a similar tradeoff spot. On FPGAs, the simple loop structure of iterative algorithms makes them the preferred choice. For example, the regular Pease FFT and variants of the iterative FFT are good choices for latency and throughput optimized FPGA implementations, respectively [11], [20], [38]. Next, we very briefly discuss GPUs and provide references for more details.

Since the advent of programmable pixel shaders, general-purpose programming on GPUs became an increasingly viable option. Earlier GPUs like Nvidia's G79 series were a first step towards that goal. However, with the Nvidia G80 series GPUs have become truly programmable. While they are still optimized for graphics-like workloads, mapping nongraphics applications with similar structure can result in astonishing performance including for the DFT. The caveat is that the performance is often only achievable for data resident in GPU memory; data transfer between GPU and host CPU is still very expensive and may nullify any speedup obtained through the GPU's high raw performance (we discuss this issue to greater detail in the section "GPU"). Intel's upcoming Larrabee platform may improve this situation.

MACHINE MODEL

State-of-the-art GPUs like the Nvidia G80 series applies ideas from symmetrically multithreaded (e.g., Tera MTA) and vector computers (e.g., Cray T90) to achieve high streaming performance [39]. In addition, minimal control flow, small computational kernels, and spatial locality are a requirement to achieve high performance.

GPU FFT ALGORITHMS

The above analysis suggests the Stockham FFT algorithm (21), originally developed for vector computers. Indeed, most FFT implementations developed for GPUs [10], [40]–[42] are based on (21), and the radix is chosen to match the GPU's hardware parameters.

MAPPING TO CODE

In early GPU computing, the computation had to be mapped to pixel shaders using graphics languages like OpenGL and DirectX. The language Cg was a first step towards more general-purpose shaders that could be programmed in a C-like language.

With the G80 series, Nvidia introduced Compute Unified Device Architecture (CUDA) [43], which makes it possible to run more general compute-intensive algorithms on GPUs. The program has to be broken into host code (implemented in C with CUDA library calls) and CUDA kernels; these kernels are programmed in a special C dialect. The kernels are run in a data-parallel SPMD fashion on a grid; a sufficiently large grid and homogeneous kernel code allows for high throughput performance. It is paramount to structure the data access pattern such that kernels that are grid neighbors operate on contiguous data (spatial locality), to enable coalesced memory access. CUDA kernels are compiled to a platform-independent byte code that the CUDA driver translates into actual GPU code on its first invocation, introducing a significant overhead.

OpenCL [44] is an emerging open standard for parallel programming of heterogeneous systems; one of its targets is GPU computing and partitioning of computation across CPUs and GPUs.

COMPUTER GENERATION OF LIBRARIES

The framework presented in this article concisely describes FFTs and enables structural optimization to efficiently match algorithms to multicore platforms based on a few crucial parameters. The formal nature of the approach has another major advantage: it serves as a blueprint for the computer generation of transform libraries. We have demonstrated this with Spiral, a system that generates high-performance libraries for linear transforms including the DFT. Given only textbook algorithms (as in the section “Framework”), Spiral generates multithreaded, vectorized source code. In Spiral, the matrix formalism is used as domain-specific declarative language called SPL, on which structural optimizations are performed by rewriting systems. In addition, to perform all necessary optimizations, Spiral uses an extension called Σ -SPL [25], [28] not described here due to space restrictions. The source code-level optimizations are performed by Spiral's backend compiler (an extension of [23]). Spiral can generate code for transforms of fixed input size [3], [5], [37], [45] or general input size transform libraries [4], [28] that are similar to FFTW.

Major advantages of library generation include the efficient handling of the implementation complexity and the ability to quickly port to new platforms. The complexity is due to the combination of many nontrivial transformations including those in this article, the need for further code-level optimizations (such as the unrolling of small kernels as was briefly discussed), the degrees of freedom in

these optimizations (such as the choice of radix or the use or not use of buffering), the need for specialization (e.g., for small code sizes, FFTW provides many variants), the number of transforms needed (the DFT variants discussed in the section “Framework” and other transforms), and the set of available algorithms.

Problems with porting include different programming models (vector instructions, explicit DMA, OpenMP) and the difficulty of maintaining performance, both exacerbated by the fast evolution of platforms.

A library generator greatly alleviates these problems, and, as shown in the next section, often without losses and sometimes even gains in performance. For example, an increasing number of transform routines in Intel's IPP (starting with 6.0) are generated by Spiral, the main reason being superior performance.

BENCHMARKS ON MULTICORE CPUS

In this section, we show FFT benchmarks of the fastest libraries on multicores that are state of the art (at the time of this writing). The focus is on an Intel Core quadcore system and the IBM Cell BE with nine cores, but we also include results for the Nvidia GPU GTX280 with 240 cores. For the Core and the Cell, we consider Spiral-generated libraries, which implement the exact techniques discussed and the similar FFTW 3.2 library. FFTW implements the recursive Cooley-Tukey FFT, buffering, vector recursion, and SIMD vectorization, using algorithms similar but not equal to (22) and (23). For the Core, we also show Intel's IPP 6.0 [22] and for the Cell FFTC [9]. For the GPU, we extracted the results from [10].

The performance for input size n is computed as $5n \log_2(n)/t$, where t is the run time in seconds. This is a slight overestimate since the true flop count is closer to $4n \log_2(n)$ and depends on the exact algorithm and recursion strategy chosen.

INTEL MULTICORE

PLATFORM

We consider a 2.66 GHz Intel Core i7 quadcore processor (Nehalem microarchitecture, 45 nm process) with SSE 4.2 instruction set. It has three levels of cache and 25.6 GB/s memory bandwidth (using all three on-chip memory controllers). Each core supports hyperthreading but for programs with high arithmetic density (like FFTs), SMT does not provide any benefit, so in Spiral we use no more than four threads using OpenMP.

The Core i7 implements SSE 4.2, providing two-way double precision and four-way single precision floating point vector support. Moreover, it supports complex arithmetic operations packing one complex double-precision number or two complex single-precision numbers into vector registers. While unaligned memory access is supported, the best performance is achieved with 16-B aligned vector loads and stores.

The Core i7 implements the shared memory paradigm. Each core has a private 64 kB L1 cache and 256 kB L2. The 8 MB L3 cache is shared among cores. All caches have 64 B cache lines (four complex double-precision numbers or eight complex single-precision numbers).

The theoretical peak performance is 85.12 Gflop/s for single and 42.56 Gflop/s for double precision.

RESULTS

Figure 5(a) and (b) shows results for out-of-place double-precision and single-precision two-power FFTs. FFTW and the Spiral-generated library [4] are compiled with the Intel C++ compiler 11.0 and flags “-O3 -xS;” IPP is provided as binary. The measurements are with “warm” cache.

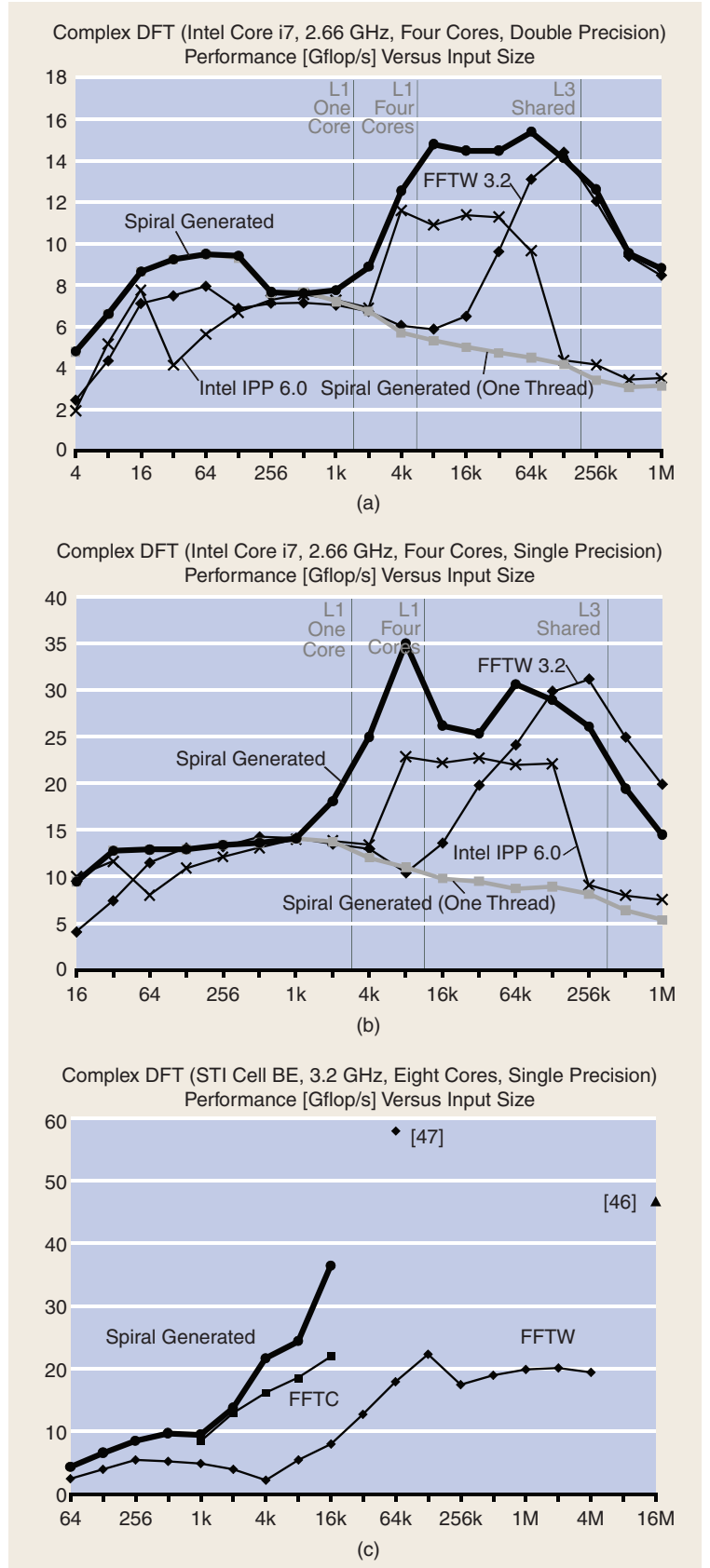
For the Spiral-generated libraries, the working set for input size n is $6n$ real numbers or $4 \cdot 1/16n$ real numbers if the twiddle factors are computed on the fly. Figure 5(a) and (b) indicates maximal cache resident sizes. For example, in Figure 5(a), an FFT is L3 cache resident up to $n = 2^{17}$.

Overall, in Figure 5(a) and (b), Spiral is about equal and often faster than the hand-written libraries. For small sizes, the performance reaches up to 10 Gflop/s (double precision) or 14 Gflop/s (single precision); in double precision, a slight drop occurs for the first size (256) that is not completely unrolled. This is due to the occurring index expressions that can be fully precomputed and inlined only if the code is unrolled. The first speedup through threading occurs already for a working set that fits into the L1 cache of one core. Subsequently, the performance ramps up as four cores are used on L1, L2, or L3 cache resident data yielding about 3x speedup over one thread. The performance peaks at 15 Gflop/s (double precision) and 35 Gflop/s (single precision) for vectorized code running on all four cores. The drop for single precision and input size 16 K may be an artifact of imperfect search. For sizes outside the shared L3 cache, the performance drops as the computation becomes memory bound. At this point buffering, vector recursion, and on-the-fly twiddle computation become crucial.

Since the FFTs used by the Spiral-generated library are already adapted as explained in this article, even a random choice of recursion will yield reasonable performance (say, within 2x). The detailed shape of the best recursion for a given DFT size on a given multi-core CPU is difficult to predict. However, we made the following observations.

The multicore Cooley-Tukey FFT (22) is used as top-level algorithm sizes that are large enough to benefit from parallelization. Typically, $m = \mu$ in (22) leads to a yields good performance. Further, the short vector Cooley-Tukey FFT (23) provides perfect SIMD vectorization. On the Core i7 $n = \nu$ is a good choice.

For cache-resident sizes, the standard Cooley-Tukey FFT (15) is a good choice with k small enough such that DFT_k can be implemented fully unrolled,



[FIG5] DFT performance on the (a), (b) 2.66 GHz Intel Core i7, and the (c) 3.2 GHz Cell BE. Higher is better.

and the machine has enough registers to support the computation. In practice, $8\nu \leq k \leq 32\nu$ (CPU-dependent) is a good choice. Once the working set no longer fits into the last cache level (or into the local store on the Cell) the involved tradeoffs become tricky. In addition to (15), the vector recursion (29) and buffering (25), (26) become fastest. On the Core i7 for large enough DFT_{kmm} , a typical out-of-cache decomposition applies the vector recursion (29) with $8 \leq k, m \leq 32$ until the working set fits into the last cache level. In addition, $\text{DFT}_k \otimes I_{mm}$ and $\text{DFT}_m \otimes I_n$ are buffered and the twiddle factors computed on the fly. This coincides with our abstract analysis in the section “Memory Hierarchy.”

CELL BE

PLATFORM

We consider a 3.2 GHz Cell BE with nine cores, including one traditional PowerPC core and eight SIMD vector cores [called synergistic processing elements (SPEs)]. Each SPE includes its own fast on-chip 256 KB local memory (local store) that is designed to be explicitly managed by the programmer. This means intercore and main memory-local store transfers must be performed via DMA; the achieved DMA bandwidth increases with DMA packet size. The Cell includes a set of four-way single precision SIMD instructions for the SPEs, accessible via C intrinsics. The vectorization, however, is very similar to Intel’s SSE. The peak performance of the Cell is 204.8 Gflop/s single precision (SPEs only) and 14.4 Gflop/s double precision.

RESULTS

Figure 5(c) shows the (latency) performance of Spiral-generated code (separate functions for every size [45] in this case), FFTC [9],

and FFTW 3.2, each using the interleaved-complex data format. Spiral-generated code is compiled with spu-gcc (flags: “-O2”), the other data is extracted from the respective papers. In addition we include the 2^{24} -sized DFT from [46] and the 2^{16} -sized DFT from [47] (both use split-complex data format). The latter has been estimated to achieve a throughput of 116 Gflop/s.

Spiral-generated code to date is limited to sizes for which the working set fits into the union of all local stores; the same seems to hold for FFTC. Both perform better than FFTW for these sizes. The excellent performance of both [47] and [46] is due to a highly optimized 2^8 -sized kernel.

All loads and stores from main memory, and all intercore communication (permutations in (22)) are performed explicitly using DMA instructions. We use the Cell’s intercore messaging mechanisms for synchronization barriers.

The Cell allows initiated DMA instructions to proceed in the background along with active computation. Although not currently used in our code in Figure 5(c), large out-of-chip DFT sizes can use a multibuffering technique based on (25) to (partially) hide memory costs. Data can be stored and loaded for the previous and next iterations in separate buffers while computation progresses for the current iteration. The explicit move operations in (25) would become DMA instructions.

For parallel code, in contrast to the Core i7, the best m, n in (22) found are both close to \sqrt{mn} , since this maximizes the packet size μ . The remaining choices found are similar to the Core i7; inside (22) and (23), a (15) with $m \approx 128$ is chosen.

Other results include [48], who achieve about 22 Gflop/s on a single SPE for DFTs of input sizes 2^{10} and 2^{13} resident in the SPE’s local store. [49] implement 2-D and 3-D parallel SPE-resident FFT kernels achieving up to 30 Gflop/s.

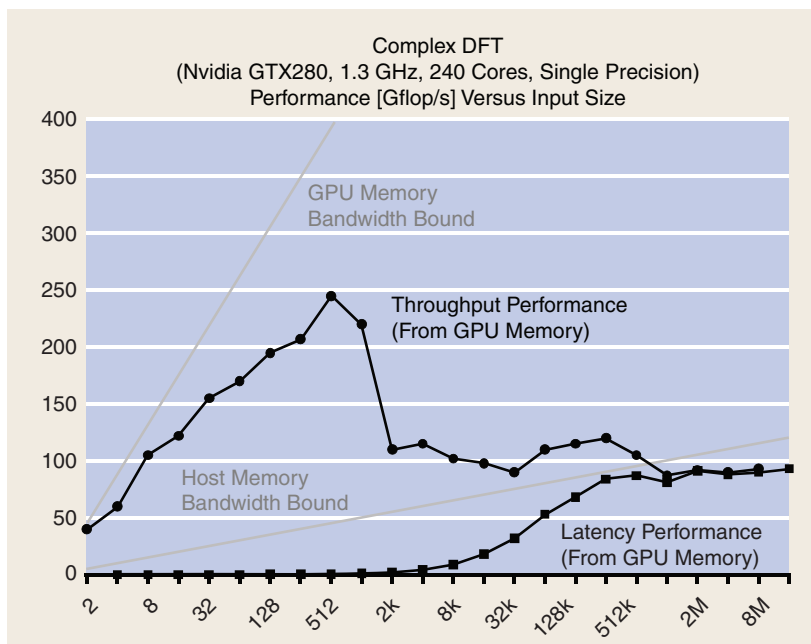
GPU

PLATFORM

We consider the Nvidia 280GTX with 240 cores grouped into 30 multiprocessors, 1 GB of on-GPU main memory and a GPU memory bandwidth of 140 GB/s. The shader clock is 1.3 GHz, and each core can perform 1 fused multiply-add and a multiply operation per cycle, leading to 936 Gflop/s peak performance. The connection between CPU memory and GPU uses PCIe 2.0, which for 16 lanes has a bandwidth of 16 GB/s (8 in each direction). The theoretical (single precision) peak performance is 936 Gflop/s.

RESULTS

The run time results in Figure 6 are taken from [10] and appear to be the fastest at the time of this writing. The memory configuration limits the achievable floating-point performance for FFTs to $43.75 \log_2 n$ Gflop/s out of GPU memory (obtained from 140 GB/s GPU



[FIG6] DFT performance on a Nvidia GTX 280 GPU. Higher is better.

memory bandwidth), and to $5 \log_2 n$ Gflop/s out of host memory (obtained from 16 GB/s PCIe bandwidth). These bounds are included as gray lines in Figure 6. The plot shows latency (one DFT is performed) and throughput ($2^{23}/n$ DFTs of size n are performed in batch mode), both with GPU memory resident data.

Throughput performance plays well to the strength of the GPU and ramps up almost parallel to the memory bandwidth bound until $n = 512$; then cache limitations cause the performance to drop to about 100 Gflop/s. We note that it is possible (not specified in [10]) that the batch mode computes $\text{DFT}_n \otimes I$ (interleaved DFTs) rather than $I \otimes \text{DFT}_n$. Not surprisingly, latency performance in contrast can amortize memory latency only for large sizes.

It seems attractive to transparently utilize the GPU as accelerator for batched single precision DFTs in CPU computations. The problem is the PCIe bandwidth between host memory and GPU memory, which yields the lower gray line as performance bound, and realistically probably half of that. Consequently, only very large sizes would benefit. However, if the entire application can be implemented on the GPU, the full performance can be harnessed, and can yield for the DFT up to an eightfold performance improvement compared to a CPU and up to fivefold compared to the Cell.

FPGA

While not directly in the scope of this article, we mention (for comparison) that the FFT in [13] and [38] achieves up to 40 Gflop/s throughput performance on a Virtex-4 (V4-FX140) for a DFT_{256} in single-precision floating point using all logic available (more data is not readily available). On Virtex-6, twice the performance is possible using twice the resources. This performance requires the data to be on chip; the off-chip bandwidth is typically about 10 GB/s in each direction. Note that usually fixed point is used on FPGAs and that most applications leave only a small part of the available logic for FFTs. In summary, the main appeal of FPGAs for DFTs is in power savings rather than in use as mere accelerators.

For commercial state-of-the-art FFTs on FPGAs see, e.g., [50] and [51].

CONCLUSIONS

The end of CPU frequency scaling and the advent of multicore systems has two major consequences for compute-intensive signal processing applications. First, it marks the end of free speedup for legacy software. Second, the software development skill required to achieve optimal performance is dramatically increased. As we have shown for the DFT, minimizing operations count alone does not yield optimal or even close-to-optimal performance. Instead, the structure of algorithms becomes crucial and has to be matched to the target architecture. Specifically, the checklist for high performance is efficient parallelization, vectorization, and memory hierarchy optimization. The necessary transformations are likely to stay out of reach for compilers since they require domain knowledge and the ability to assess the many available choices. To handle the implementation complexity, we believe it is important to develop rigorous

approaches that formalize algorithmic optimizations by connecting the algorithm structure with architecture parameters. We have presented such a framework for the DFT and used it to give an overview on FFTs and optimizations for current multicores. Further, as we demonstrated with Spiral, the rigorous nature of the framework enables automation: the computer generation of DFT libraries that often achieve excellent performance compared to their handwritten counterparts.

ACKNOWLEDGMENT

This work was supported by NSF through awards 0325687, 0702386, by DARPA (DOI grant NBCH1050009), the ARO grant W911NF0710416, and by Intel Corp. and Mercury Computer Systems, Inc.

AUTHORS

Franz Franchetti (franzf@ece.cmu.edu) is an assistant research professor with the Electrical and Computer Engineering (ECE) Department at Carnegie Mellon University. He received the M.Sc. degree in technical mathematics and the Ph.D. degree in computational mathematics from the Vienna University of Technology in 2000 and 2003, respectively. He was a postdoctoral research associate with the Institute for Analysis and Scientific Computing in 2003. From 2004 to 2005, he was a postdoctoral research associate with the ECE Department at Carnegie Mellon University and a recipient of the Austrian Science Fund Schrödinger Fellowship. In 2006, he was member of the team winning the Gordon Bell Prize (Peak Performance Award). From 2005 to 2008, he was systems scientist (special faculty) with Carnegie Mellon's ECE Department.

Markus Püschel (pueschel@ece.cmu.edu) is a professor in the Electrical and Computer Engineering Department at Carnegie Mellon University. He received his M.Sc. degree in mathematics and his Ph.D. degree in computer science in 1995 and 1998, respectively, both from the University of Karlsruhe, Germany. He was an associate editor for *IEEE Transactions on Signal Processing*, and *IEEE Signal Processing Letters* and was a guest editor of the *Proceedings of the IEEE* and the *Journal of Symbolic Computation*. He received the College of Engineering at Carnegie Mellon's Outstanding Research Award and the Eta Kappa Nu Award for Outstanding Teaching. He also holds the title of Privatdozent at the University of Technology, Vienna, Austria. He is a Senior Member of the IEEE.

Yevgen Voronenko (yvonen@ece.cmu.edu) is a project scientist in the Electrical and Computer Engineering Department at Carnegie Mellon University (CMU). He received a B.S. degree in computer science from Drexel University in 2003 and his Ph.D. degree from CMU in 2008, where he was awarded the A.G. Milnes Outstanding Ph.D. Dissertation Award. His research interests include scientific computing, software engineering, programming languages, and compiler design.

Srinivas Chellappa (schellap@ece.cmu.edu) is a Ph.D. candidate in the Electrical and Computer Engineering Department at Carnegie Mellon University. He is currently working on automatic program generation of DSP kernels for distributed

memory platforms. His research interests include parallel and distributed systems, automatic code generation, operating systems, virtualization, and computer architecture.

José M.F. Moura (moura@ece.cmu.edu) is a professor at Carnegie Mellon University (CMU), where he manages a large program between CMU and Portugal (www.icti.cmu.edu). He is also codirector of CMU's Center for Sensed Critical Infrastructures Research. His interests include statistical and algebraic signal and image processing. He was editor-in-chief of *IEEE Transactions in Signal Processing* and interim editor-in-chief of *IEEE Signal Processing Letters*. He is the president of the IEEE Signal Processing Society. He is a Fellow of the IEEE, a Fellow of AAAS, and corresponding member of the Academy of Sciences of Portugal.

REFERENCES

- [1] G. Blake, R. G. Dreslinski, and T. Mudge, "A survey of multicore architectures," *IEEE Signal Processing Mag.*, vol. 26, no. 6, pp. 26–37, 2009.
- [2] W. H. Press, B. P. Flannery, S. A. Teukolsky, and W. T. Vetterling, *Numerical Recipes in C: The Art of Scientific Computing*, 2nd ed. Cambridge, U.K.: Cambridge Univ. Press, 1992.
- [3] M. Püschel, J. M. F. Moura, J. Johnson, D. Padua, M. Veloso, B. W. Singer, J. Xiong, F. Franchetti, A. Gačić, Y. Voronenko, K. Chen, R. W. Johnson, and N. Rizzolo, "SPIRAL: Code generation for DSP transforms," *Proc. IEEE* (Special Issue on Program Generation, Optimization, and Adaptation), vol. 93, no. 2, pp. 232–275, 2005.
- [4] Y. Voronenko, F. de Mesmay, and M. Püschel, "Computer generation of general size linear transform libraries," in *Proc. Code Generation and Optimization (CGO)*, 2009, pp. 102–113.
- [5] F. Franchetti, Y. Voronenko, and M. Püschel, "FFT program generation for shared memory: SMP and multicore," in *Proc. Supercomputing (SC)*, 2006, p. 51.
- [6] M. Frigo and S. G. Johnson, "The design and implementation of FFTW3," *Proc. IEEE* (Special Issue on Program Generation, Optimization, and Adaptation), vol. 93, no. 2, pp. 216–231, 2005.
- [7] M. Frigo, "A fast Fourier transform compiler," in *Proc. Programming Language Design and Implementation (PLDI)*, 1999, pp. 169–180.
- [8] M. Frigo and S. G. Johnson, FFTW 3.2 [Online]. Available: www.fftw.org
- [9] D. A. Bader and V. Agarwal, "FFTC: Fastest Fourier transform for the IBM cell broadband engine," in *Proc. Int. Conf. High Performance Computing (HiPC)*, 2007, pp. 172–184.
- [10] N. K. Govindaraju, B. Lloyd, Y. Dotsenko, B. Smith, and J. Manferdelli, "High performance discrete Fourier transforms on graphics processors," in *Proc. Supercomputing (SC)*, 2008, pp. 1–12.
- [11] C. Van Loan, *Computational Framework of the Fast Fourier Transform*. Philadelphia, PA: SIAM, 1992.
- [12] R. Tolimieri, M. An, and C. Lu, *Algorithms for Discrete Fourier Transforms and Convolution*, 2nd ed. New York: Springer-Verlag, 1997.
- [13] M. Püschel, Spiral Web site [Online]. Available: www.spiral.net
- [14] J. Johnson, R. W. Johnson, D. Rodriguez, and R. Tolimieri, "A methodology for designing, modifying, and implementing Fourier transform algorithms on various architectures," *IEEE Trans. Circuits Syst.*, vol. 9, no. 4, pp. 449–500, 1990.
- [15] A. Norton and A. J. Silberberger, "Parallelization and performance analysis of the Cooley-Tukey FFT algorithm for shared-memory architectures," *IEEE Trans. Comput.*, vol. 36, no. 5, pp. 581–591, 1987.
- [16] M. Hegland, "Block algorithms for FFTs on vector and parallel computer," in *Parallel Computing: Trends and Applications*. Amsterdam, The Netherlands: Elsevier, 1994, pp. 129–136.
- [17] D. H. Bailey, "FFTs in external or hierarchical memory," *J. Supercomput.*, vol. 4, no. 1, pp. 23–35, Mar. 1990.
- [18] D. B. Harris, J. H. McClellan, D. S. K. Chan, and H. W. Schuessler, "Vector radix fast Fourier transform," in *Proc. Int. Conf. Acoustics, Speech, and Signal Processing (ICASSP)*, 1977, pp. 548–551.
- [19] J. W. Cooley and J. W. Tukey, "An algorithm for the machine calculation of complex Fourier series," *Math. Comput.*, vol. 19, pp. 297–301, Apr. 1965.
- [20] M. C. Pease, "An adaptation of the fast Fourier transform for parallel processing," *J. ACM*, vol. 15, no. 2, pp. 252–264, Apr. 1968.
- [21] P. N. Schwarztrauber, "Multiprocessor FFTs," *Parallel Comput.*, vol. 5, pp. 197–210, July 1987.
- [22] Intel Corp., *Intel integrated performance primitives (ipp) 6.0* [Online]. Available: <http://software.intel.com/en-us/intel-ipp>
- [23] J. Xiong, J. Johnson, R. Johnson, and D. Padua, "SPL: A language and compiler for DSP algorithms," in *Proc. Programming Language Design and Implementation (PLDI)*, 2001, pp. 298–308.
- [24] I. W. Selesnick and C. S. Burrus, "Automatic generation of prime length FFT programs," *IEEE Trans. Signal Processing*, vol. 44, no. 1, pp. 14–24, 1996.
- [25] F. Franchetti, Y. Voronenko, and M. Püschel, "Loop merging for signal transforms," in *Proc. Programming Language Design and Implementation (PLDI)*, 2005, pp. 315–326.
- [26] D. Mirković and S. L. Johnson, "Automatic performance tuning in the UHFFT library," in *Proc. Int. Conf. Computational Science (ICCS) (Lecture Notes in Computer Science)*, vol. 2073. New York: Springer-Verlag, 2001, pp. 71–80.
- [27] D. Takahashi, "An implementation of parallel 1-D FFT using SSE3 instructions on dual-core processors," in *Proc. Int. Workshop State-of-the-Art in Scientific and Parallel Computing (PARA)*, 2006, pp. 1178–1187.
- [28] Y. Voronenko, "Library generation for linear transforms," Ph.D. dissertation, Elect. Comput. Eng., Carnegie Mellon Univ., Pittsburgh, PA, 2008.
- [29] S. Chellappa, F. Franchetti, and M. Püschel, "How to write fast numerical code: A small introduction," in *Proc. Summer School on Generative and Transformational Techniques in Software Engineering (GTTSE)* (Lecture Notes in Computer Science, vol. 5235). Berlin: Springer-Verlag, 2008, pp. 196–259.
- [30] J. R. Johnson, R. W. Johnson, D. Rodriguez, and R. Tolimieri, "A methodology for designing, modifying, and implementing Fourier transform algorithms on various architectures," *IEEE Trans. Circuits, Syst., Signal Processing*, vol. 9, no. 4, pp. 449–500, 1990.
- [31] A. Ali, L. Johnson, and J. Subhlok, "Scheduling FFT computation on SMP and multicore systems," in *Proc. Int. Conf. Supercomputing (ICS)*, 2007, pp. 293–301.
- [32] OpenMP. (1998). OpenMP C and C++ application program interface, version 1.0 [Online]. Available: www.openmp.org
- [33] B. Gallmeister, *POSIX.4*. Sebastopol, CA: O'Reilly, 1994.
- [34] F. Franchetti and M. Püschel, "Generating SIMD vectorized permutations," in *Proc. Int. Conf. Compiler Construction (CC)* (Lecture Notes in Computer Science, vol. 4959). Berlin: Springer-Verlag, 2008, pp. 116–131.
- [35] F. Franchetti and M. Püschel, "Short vector code generation for the discrete Fourier transform," in *Proc. IEEE Int. Parallel and Distributed Processing Symp. (IPDPS)*, 2003, pp. 58–67.
- [36] F. Franchetti and M. Püschel, "A SIMD vectorizing compiler for digital signal processing algorithms," in *Proc. IEEE Int. Parallel and Distributed Processing Symp. (IPDPS)*, 2002, pp. 20–26.
- [37] F. Franchetti, Y. Voronenko, and M. Püschel, "A rewriting system for the vectorization of signal transforms," in *Proc. High Performance Computing for Computational Science (VECPAR)*, 2006, pp. 363–377.
- [38] P. A. Milder, F. Franchetti, J. C. Hoe, and M. Püschel, "Formal datapath representation and manipulation for implementing DSP transforms," in *Proc. Design Automation Conf. (DAC)*, 2008, pp. 385–390.
- [39] A. Snavely, L. Carter, J. Boisseau, A. Majumdar, K. S. Gatlin, N. Mitchell, J. Feo, and B. Koblenz, "Multi-processor performance on the Tera MTA," in *Proc. Supercomputing (SC)*, 1998, pp. 1–8.
- [40] K. Moreland and E. Angel, "The FFT on a GPU," in *Proc. ACM SIGGRAPH/EUROGRAPHICS Conf. Graphics Hardware*, 2003, pp. 112–119.
- [41] N. K. Govindaraju and D. Manocha, "Cache-efficient numerical algorithms using graphics hardware," *Parallel Comput.*, vol. 33, no. 10–11, pp. 663–684, 2007.
- [42] A. Nukada, Y. Ogata, T. Endo, and S. Matsuoaka, "Bandwidth intensive 3-d FFT kernel for GPUs using CUDA," in *Proc. Supercomputing (SC)*, 2008, pp. 1–11.
- [43] Nvidia Corp., Nvidia CUDA [Online]. Available: www.nvidia.com/cuda
- [44] Khronos Group, OpenCL [Online]. Available: www.khronos.org/opencl/
- [45] S. Chellappa, F. Franchetti, and M. Püschel, "Computer generation of fast Fourier transforms for the cell broadband engine," in *Proc. Int. Conf. Supercomputing (ICS)*, 2009, pp. 26–35.
- [46] A. C. Chow, G. C. Fossum, and D. A. Brokenshire, "A programming example: Large FFT on the cell broadband engine," IBM, Tech. Rep., May 2005 [Online]. Available: [https://www-01.ibm.com/chips/techlib/techlib.nsf/techdocs/0AA2394A505E0FB872570AB005BF0F1/\\$file/GSPx_FFT_paper_legal_0115.pdf](https://www-01.ibm.com/chips/techlib/techlib.nsf/techdocs/0AA2394A505E0FB872570AB005BF0F1/$file/GSPx_FFT_paper_legal_0115.pdf)
- [47] J. Greene and R. Cooper, "A parallel 64K complex FFT algorithm for the IBM/Sony/Toshiba cell broadband engine processor," in *Proc. Global Signal Processing Expo (GSPx)*, 2005.
- [48] L. Cico, R. Cooper, and J. Greene, "Performance and programmability of the IBM/Sony/Toshiba cell broadband engine processor," in *Proc. (EDGE) Workshop*, 2006.
- [49] P. Bientinesi, N. Pitsianis, and X. Sun, "Multi-dimensional array operations for signal processing algorithms," in *Proc. Int. Workshop State-of-the-Art in Scientific and Parallel Computing (PARA)*, 2008.
- [50] 4DSP Inc., 4DSP [Online]. Available: www.4dsp.com/fft.htm
- [51] Dillon Engineering, Dillon FFT [Online]. Available: www.dilloneng.com/fft_ip