

# How To Write Fast Numerical Code: A Small Introduction

Srinivas Chellappa, Franz Franchetti, and Markus Püschel

Electrical and Computer Engineering  
Carnegie Mellon University  
{schellap, franzf, pueschel}@ece.cmu.edu

**Abstract.** The complexity of modern computing platforms has made it increasingly difficult to write numerical code that achieves the best possible performance. Straightforward implementations based on algorithms that minimize the operations count often fall short in performance by an order of magnitude. This tutorial introduces the reader to a set of general techniques to improve the performance of numerical code, focusing on optimizations for the computer’s memory hierarchy. Two running examples are used to demonstrate these techniques: matrix-matrix multiplication and the discrete Fourier transform.

## 1 Introduction

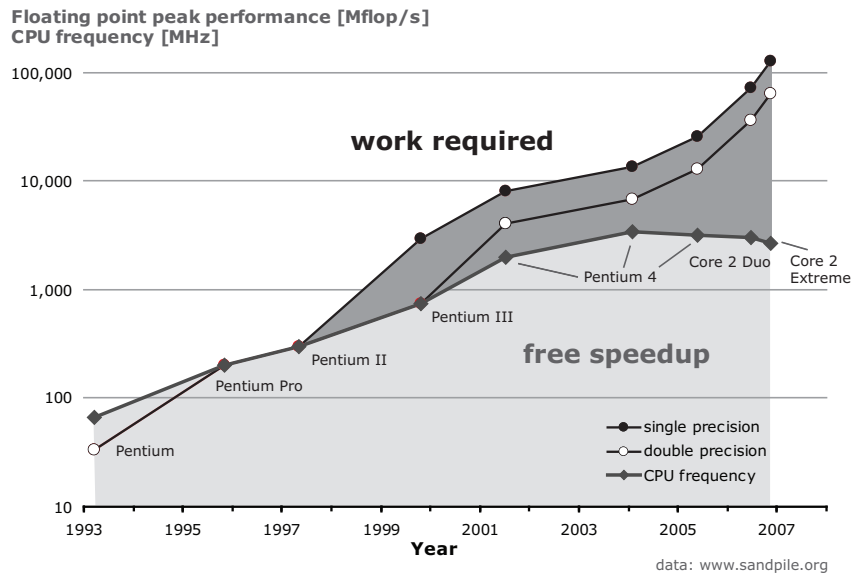
The growth in the performance of computing platforms in the past few decades has followed a reliable pattern usually referred to as Moore’s Law. Moore observed in 1965 [53] that the number of transistors per chip roughly doubles every 18 months and predicted—correctly—that this trend would continue. In parallel, due to the shrinking size of transistors, CPU frequencies could be increased at roughly the same exponential rate. This trend has been the big supporter for many performance demanding applications in scientific computing (such as climate modeling and other physics simulations), consumer computing (such as audio, image, and video processing), and embedded computing (such as control, communication, and signal processing). In fact, these domains have a practically unlimited need for performance (for example, the ever growing need for higher resolution videos), and it seems that the evolution of computers is well on track to support these needs.

However, everything comes at a price, and in this case it is the increasing difficulty of writing the fastest possible software. In this tutorial, we focus on *numerical* software. By that we mean code that mainly consists of floating point computations.

**The problem.** To understand the problem we investigate Fig. 1, which considers various Intel architectures from the first Pentium to the (at the time of this writing) latest Core 2 Extreme. The  $x$ -axis shows the year of release. The  $y$ -axis, in log-scale, shows both the CPU frequency (in MHz) and the single/double precision theoretical peak performance (in Mflop/s = Mega Floating point Operations per Second) of the respective machines. First we note, as explained above, the exponential increase in CPU frequency. This results in a “free” speedup for numerical software. In other words, legacy code written for

an obsolete predecessor will run faster without any extra programming effort. However, the theoretical performance of computers has evolved at a faster pace due to increases in the processors' parallelism. This parallelism comes in several forms, including pipelining, superscalar processing, vector processing and multi-threading. Single-instruction multiple-data (SIMD) vector instructions enable the execution of 2, 4, or more operations in parallel. The latest generations are also "multicore," which means 2, 4, or more processing cores<sup>1</sup> exist on a single chip. Exploiting parallelism in numerical software is not trivial, it requires implementation effort. Legacy code typically neither includes vector instructions, nor is it multi-threaded to take advantage of multiple processor cores or multiple processors. Ideally, compilers would take care of this problem by automatically vectorizing and parallelizing existing source code. However, while much outstanding compiler research has attacked these problems (e.g., [50, 64, 30]), they are in general still unsolved. Experience shows that this is particularly true for numerical problems. The reason is, for numerical problems, the platform's available parallelism often requires an algorithm structured differently than the one that would be used in the corresponding sequential code. Compilers are inherently not capable of changing or restructuring algorithms since doing so requires knowledge of the algorithm domain.

### Evolution of Intel Platforms



**Fig. 1.** The evolution of computing platform's peak performance versus their CPU frequency explains why high performance software development becomes increasingly harder.

<sup>1</sup> At the time of this writing 8 cores per chip is the maximum commonly available CPU configuration.

Similar problems are caused by the computer’s memory hierarchy, independently of the available parallelism. The fast processor speeds have made it increasingly difficult to “feed all floating point execution units” at the necessary rate to keep them busy. Moving data from and to memory has become the bottleneck. The memory hierarchy, consisting of registers, and multiple levels of cache, aims to address this problem, but can only work if data is accessed in the proper order. One cache miss may incur a penalty of 20–100s CPU cycles, a time in which 100 or more floating point operations could have been performed. Again, compilers are inherently limited in optimizing for the memory hierarchy since optimization may require an algorithm restructuring or an entirely different choice of algorithm to begin with.

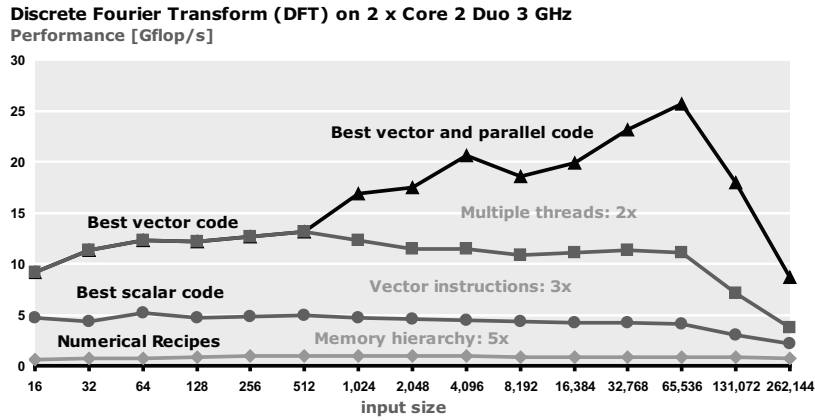
Adding to these problems is the fact that CPU frequency scaling is approaching its end due to limits to the chip’s possible power density (see Fig. 1): since 2004 it has hovered around 3GHz. This implies *the end of automatic speed-up*; future performance gains will be made exclusively by increasing parallelism.

To quantify the problem we look at two representative examples, which are among the most important numerical kernels used: the discrete Fourier transform (DFT) and the matrix-matrix multiplication (MMM).

It is well-known that the complexity of the DFT, for input size  $n$  is  $O(n \log(n))$  due to the availability of fast Fourier transform algorithms (FFTs) [68]. Fig. 2 shows the performance of four FFT implementations on an Intel Core platform with 4 cores. The  $x$ -axis is the input size  $n = 2^4, \dots, 2^{18}$ . The  $y$ -axis is the performance in Gflop/s. For all implementations, the operations count is estimated as  $5n \log_2(n)$ , so the numbers are proportional to inverse runtime. The bottom line shows the performance of the implementation by Numerical Recipes [58] compiled with the best available compiler (the Intel vendor compiler `icc` in this case) and all optimizations enabled. The next line (best scalar) shows the performance of the fastest standard C implementation for the DFT and is roughly 5 times faster due to optimizations for the memory hierarchy. The next line (best vector) shows the performance, when in addition vector instructions are used for a further gain of a factor of 3. Finally, for large sizes, another factor of 2 can be gained by writing multithreaded code to use all processor cores. Note that all four implementations *have roughly the same operations count* for a given size but the performance difference is a factor of 12 for small sizes, and a factor of up to 30 for large sizes. The uppermost three lines were obtained using Spiral [60, 63]; similar performance can be achieved using FFTW [32, 33, 24].

Fig. 3 shows a similar plot for MMM (assuming square matrices), where the bottom line corresponds to a standard, triple loop implementation. Here the performance difference to the best code can be as much as 160 times, with 5-20 times due to optimizations for the memory hierarchy. All the implementations have exactly the same floating point operations count of  $2n^3$ . The top two lines are for GotoBLAS [36]; the best scalar code is produced using ATLAS [70].

To summarize the above discussion, the task of achieving the highest performance with an implementation lies to a great extent with the programmer. For a given problem, he or she has to carefully consider different algorithms and possibly restructure them to



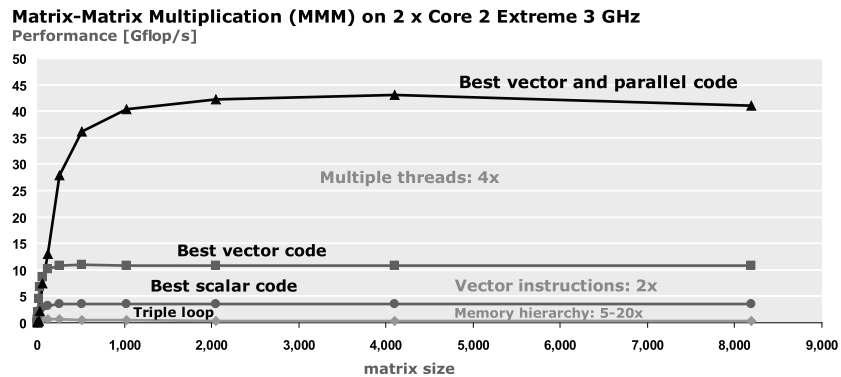
**Fig. 2.** Performance of four single precision implementations of the discrete Fourier transform. The operations count is roughly the same.

adapt to the given platform’s memory hierarchy and available parallelism. This is very difficult, time-consuming, and requires interdisciplinary knowledge about algorithms, software optimizations, and the hardware architecture. Further, the tuning process is platform-dependent: an implementation optimized for one computer will not necessarily be the fastest one on another. Consequently, to achieve highest performance, tuning has to be repeated with the release of each new platform. Since the times of a free speedup (due to frequency scaling) are over, this retuning has become mandatory if any performance gains are desired. Needless to say, the problem is not merely an academic one, but one that affects the software industry as a whole.

**Automatic performance tuning.** A number of research efforts have started to address this problem in a new area called “automatic performance tuning” [54]. The general idea is to at least partially automate the implementation and optimization procedure. Two basic approaches have emerged so far in this area: adaptive libraries and source code generators.

Examples of adaptive libraries include FFTW [33] for the discrete Fourier transform and adaptive sorting libraries [7, 49]. In both cases, the libraries are highly optimized, and beyond that, have degrees of freedom with regard to the divide-and-conquer strategy to choose (both DFT and sorting are done recursively in these libraries). This strategy is determined at runtime, on the given platform, using a search mechanism. This way, the library can dynamically adapt to the computer’s memory hierarchy. Sparsity and OSKI from the BeBOP group [42, 20, 6, 69] is another example of such a library, used to compute sparse matrix-vector products.

On the other hand, source code generators produce algorithm implementations from scratch. They are used to generate either crucial components, or libraries in their entirety. For instance, ATLAS [71, 20] generates the kernel code for matrix-matrix multi-



**Fig. 3.** Performance of four double precision implementations of matrix-matrix multiplication. The operations count is exactly the same.

plication. It does so by generating many different variants arising from different choices of blocking, loop unrolling, and instruction ordering. These are all measured and the fastest one is selected using search methods.

FFTW also uses a generator to produce small size DFT kernels [31]. Here, no search is used, but many optimizations are performed before the actual code is output. Spiral [60, 27] is a library generator for arbitrary sized linear transform including the DFT, filters, and others. Besides enumerating alternatives, Spiral uses an internal mathematical language to optimize algorithms at a high level of abstraction before source code is generated. This includes algorithm restructuring for the memory hierarchy, vector instructions, and multi-threaded code [27–29].

Other automatic performance tuning efforts include [42] for sparse linear algebra and [5] for tensor computations.

This new research is promising but much more work is needed to automate the implementation and optimization of a large set of library functionality.

**Summary.** We summarize the main points of this section:

- *End of free-speedup for legacy code.* CPU frequencies have hit the power wall and stalled. Future performance gains in computers will be obtained through increasing parallelism. This means code has to be rewritten to take advantage of the available parallelism and performance.
- *Minimizing operations count does not mean maximizing performance.* Floating-point operations are much cheaper than cache misses. Fastest performance requires code that is adapted to the memory hierarchy, uses vector instructions and multiple cores (if available). As a consequence, we have the following problem.

- *The performance different between a straightforward implementation and the best possible can be a factor of 10, 20, or more.* This is true even if the former is based on an algorithm that is optimal in its operations count.
- *It is very difficult to write the fastest possible code.* The reason is that performance-optimal code has to be carefully optimized for the platform's memory hierarchy and available parallelism. For numerical problems, compilers cannot perform these optimizations, or can only perform them to a very limited extent.
- *Performance is in general non-portable.* The fastest code for one computer may perform poorly on another.

**About this tutorial.** The goal of this tutorial is to provide the reader with a small introduction to the performance optimization of numerical problems. The computers considered in this tutorial are COTS (commercial off-the-shelf) desktop computers with the latest microarchitectures such as Core 2 Duo or the Pentium from Intel, the Opteron from AMD, or the PowerPC from Apple/Motorola. We assume the reader to have the level of knowledge of a junior (third year) student in computer science or engineering. This includes basic knowledge of computer architecture, algorithms, and a solid familiarity with C programming.

Section 2 provides some basic background information on algorithm analysis, the MMM and the DFT, features of modern computer systems relevant to this tutorial, and compilers and their correct usage. It also identifies data access patterns that are necessary for obtaining high performance on modern computer systems. Section 3 first introduces the basics of benchmarking numerical code and then provides a general high-level procedure for attacking the problem of performance optimization given an existing program that has to be tuned for performance. This procedure reduces the problem to the optimization of performance-critical kernels, which is then studied in Section 4 using MMM and DFT as examples.

Along with the explanations, we provide programming exercises to provide the reader with hands-on experience.

## 2 Background

In this section we provide the necessary background for this tutorial. We briefly review algorithm analysis, introduce MMM and the DFT, discuss the memory hierarchy of off-the-shelf microarchitectures, and explain the use of compilers. The following standard books provide more information on algorithms [16], MMM and linear algebra [19], the DFT [68, 66], and computer architecture and systems [40, 12].

### 2.1 Cost Analysis Of Algorithms

The starting point for any implementation of a numerical problem is the choice of algorithm. Before an actual implementation, algorithm analysis, based on the number of

operations performed, can give a rough estimate of the performance to be expected. We discuss the floating point operations count and the degree of reuse.

**Cost: asymptotic, exact, and measured.** It is common in algorithm analysis to represent the asymptotic runtime of an algorithm in  $O$ -notation as  $O(f(n))$ , where  $n$  is the input size and  $f$ , a function [16]. For numerical algorithms,  $f(n)$  is typically determined from the number of floating point operations performed. The  $O$ -notation neglects constants and lower order terms; for example,  $O(n^3 + 100n^2) = O(5n^3)$ . Hence it is only suited to describe the performance *trend* but not the *actual* performance itself. Further, it makes a statement only about the asymptotic behavior, i.e., behavior as  $n$  goes to infinity. Thus it is in principle possible that an  $O(n^3)$  algorithm performs better than an  $O(n^2)$  algorithm for all practically relevant input sizes  $n$ .

A better form of analysis for numerical algorithms is to compute the *exact* number of floating point operations, or at least the exact highest order term. However, this may be difficult in practice. In this case, profiling tools can be used on an actual implementation to determine the number of operations actually performed. The latter can also be used to determine the computational bottleneck in a given implementation.

However, even if the exact number of operations of an algorithm and its implementation is known, it is very difficult to determine the actual runtime. As an example consider Fig. 3: all four implementations require exactly  $2n^3$  operations, but the runtime differs by up to two orders of magnitude.

**Reuse: CPU bound vs. memory bound.** Another useful measure of an algorithm is the degree of reuse. The asymptotic reuse for an  $O(f(n))$  algorithm is given by  $O(f(n)/n)$ . Intuitively, the degree of reuse measures how often a given input value is used in a computation during the algorithm. A high degree of reuse implies that an algorithm may perform better (in terms of operations per second) on a computer with memory hierarchy, since the number of computations dominates the number of data transfers from memory to CPU. In this case we say that the algorithm is *CPU bound*. A low degree of reuse implies that the number of data transfers from memory to CPU is high compared to the number of operations and the performance (in operations per second) may deteriorate: in this case we say that the algorithm is *memory bound*.

A CPU bound algorithm will run faster on a machines with a faster CPU. A memory bound algorithm will run faster on a machine with a faster memory bus.

## 2.2 Matrix-Matrix Multiplication

Matrix-matrix multiplication (MMM) is arguably the most important numerical kernel functionality. It is used in many linear algebra algorithms such as solving systems of linear equations, matrix inversion, eigenvalue computations, and many others. We will use MMM, and the DFT (Section 2.3) as examples to demonstrate optimizations for performance.

**Definition.** Given a  $k \times m$  matrix  $A = [a_{i,j}]$  and a  $m \times n$  matrix  $B = [b_{i,j}]$ , the product  $C = AB$  is a  $k \times n$  matrix with entries

$$c_{i,j} = \sum_{k=1}^m a_{i,k} b_{k,j}.$$

For actual applications, usually  $C = C + AB$  is implemented instead of  $C = AB$ .

**Complexity and analysis.** Given two  $n \times n$  matrices  $A, B$ , MMM computed as  $C = C + AB$  by definition requires  $n^3$  multiplications and  $n^3$  additions for a total of  $2n^3 = O(n^3)$  floating point operations. Since the input data (the matrices) have size  $O(n^2)$ , the reuse is given by  $O(n^3/n^2) = O(n)$ .

Asymptotically better MMM algorithms do exist. Strassen's algorithm [65] requires only  $O(n^{\log_2 7}) \approx O(n^{2.808})$  operations. The actual crossover point (i.e., when it requires less operations than computation by definition) is  $n = 655$ . However, the more complicated structure of Strassen's algorithm and its weaker numerical stability reduce its applicability. The best-known algorithm for MMM is due to Coppersmith-Winograd and requires  $O(n^{2.376})$  [15]. The large hidden constant and a complicated structure have so far made this algorithm impractical for real applications.

**Direct implementation.** A direct implementation of MMM is the triple loop shown below.

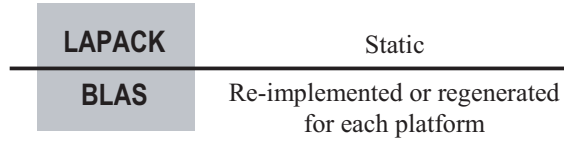
```
// MMM - direct implementation
for(i=0; i<m; i++)
  for(j=0; j<n; j++)
    for(k=0; k<n; k++)
      c[i][j] += a[i][k] * b[k][j];
```

**BLAS and LAPACK.** MMM is an example of a Basic Linear Algebra Subroutine (BLAS) [10]. These routines are used as kernels in fundamental linear algebra algorithms such as linear equation solving, eigenvalue computations, singular value decompositions, LU/Cholesky/QR decompositions, and others. These functions are implemented in the library LAPACK (Linear Algebra PACKage) [3] using MMM and other BLAS routines as kernels (see Fig. 4). The idea behind this design is to reimplement or reoptimize BLAS for new hardware architectures while leaving LAPACK unmodified. The performance improvements for BLAS will then translate into performance gains for LAPACK. This design has proven very successful until the release of multicore systems which may require a redesign of LAPACK.

#### Further Reading.

- *Linear Algebra.* General information about numerical linear algebra can be found in [19, 12].
- *BLAS.* Further information on ATLAS is available in [70, 71, 4]. Details on Goto-BLAS can be found at [36, 37].





**Fig. 4.** LAPACK is implemented on top of BLAS

- *Linear Algebra Libraries.* LAPACK is described in [48, 3]. The distributed memory extension ScaLAPCK is described in [62, 9]. An alternative approach is pursued by PLAPACK [57, 13] and FLAME [39, 8, 25].

### 2.3 Discrete Fourier Transform

In this section we provide background for the DFT and fast algorithms to compute the DFT (Fast Fourier Transforms, or FFTs). We explain these algorithms using the Kronecker product formalism.

**Definition.** The discrete Fourier transform (DFT) of an input vector  $x$  of length  $n$  is defined as the matrix-vector product

$$y = \text{DFT}_n x, \quad \text{DFT}_n = [\omega_n^{k\ell}]_{0 \leq k, \ell < n}, \quad \omega_n = e^{-2\pi i/n}.$$

In this tutorial we assume that  $n$  is a two-power - i.e.,  $n = 2^k$ .

**Kronecker product formalism.** We describe fast algorithms for the DFT using the Kronecker product formalism [68]. There are several reasons for using the Kronecker product formalism to represent these algorithms. First, it is easy to translate algorithms expressed this way directly into code, as we shall see later. Second, representing algorithms at this level of abstraction allows us to *match* the algorithms to a specific hardware architecture. For instance, the algorithms can be mapped to vector architectures by identifying and manipulating structures at the formula level [29, 28]. Third, this representation allows for easy manipulation of the algorithm to derive variations, which helps with exploring the algorithm search space.

We define  $I_n$  as the  $n \times n$  identity matrix. The tensor (or Kronecker) product of matrices is defined as

$$A \otimes B = [a_{i,j} B]_{i,j} \quad \text{with} \quad A = [a_{i,j}]_{i,j}.$$

In particular,

$$I_n \otimes A = \begin{bmatrix} A & & & \\ & A & & \\ & & \ddots & \\ & & & A \end{bmatrix}.$$

We also introduce the iterative direct sum

$$\bigoplus_{i=0}^{n-1} A_i = \begin{bmatrix} A_0 & & & \\ & A_1 & & \\ & & \ddots & \\ & & & A_{n-1} \end{bmatrix},$$

which generalizes  $I_n \otimes A$ . The stride permutation matrix  $L_m^{mn}$  permutes an input vector  $x$  of length  $mn$  as

$$in + j \mapsto jm + i, \quad 0 \leq i < m, \quad 0 \leq j < n.$$

If  $x$  is viewed as an  $n \times m$  matrix, stored in row-major order, then  $L_m^{mn}$  performs a transposition of this matrix.

**Recursive FFT.** Using the above formalism, the well-known Cooley-Tukey FFT in its recursive form can be written as a factorization of the  $\text{DFT}_n$  matrix into a product of sparse matrices:

$$\text{DFT}_{mn} = (\text{DFT}_m \otimes I_n) D_{m,n} (I_m \otimes \text{DFT}_n) L_m^{mn}. \quad (1)$$

Here  $D_{m,n}$  is the diagonal “twiddle” matrix defined as

$$D_{m,n} = \bigoplus_{i=0}^n \text{diag}(\omega_{mn}^0, \omega_{mn}^1, \dots, \omega_{mn}^{n-1})^i.$$

Equation (1) computes a DFT of size  $mn$  in four steps. First, the input vector is permuted by  $L_m^{mn}$ . Second,  $m$  DFTs of size  $n$  are computed recursively on segments of the vector. Third, the vector is scaled element wise by  $D_{m,n}$ . Lastly,  $n$  DFTs of size  $m$  are computed recursively at stride  $m$ .

The recursively called smaller DFTs are computed similarly until the base case  $n = 2$  is reached and computed by definition using an addition and a subtraction:

$$\text{DFT}_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}. \quad (2)$$

In summary, (1) and (2) are sufficient to compute DFTs of arbitrary two-power sizes. To compute DFTs of other sizes, other FFT algorithms are required [68].

**Algorithms and formulas.** There is a degree of freedom in applying (1) to recursively compute a DFT, namely in factoring the given DFT input size  $n$ . For instance one can factor  $8 \rightarrow 2 \times 4 \rightarrow 2 \times (2 \times 2)$  using two recursive applications of (1). The complete FFT algorithm for this factorization could then be written as the following *formula*:

$$\text{DFT}_8 = (\text{DFT}_2 \otimes I_4) D_{8,4} (I_2 \otimes (\text{DFT}_2 \otimes I_2) D_{4,2} (I_2 \otimes \text{DFT}_2) L_2^4) L_2^8. \quad (3)$$

**Direct implementation.** A straightforward implementation of (1) can be easily obtained since the occurring matrix formulas have a direct interpretation in terms of code

formula	code
$y = (A_n B_n)x$	<code>t[0:1:n-1] = B(x[0:1:n-1]); y[0:1:n-1] = A(t[0:1:n-1]);</code>
$y = (I_m \otimes A_n)x$	<code>for(i=0;i&lt;m;i++)   y[i*n:1:i*n+n-1] =     A(x[i*n:1:i*n+n-1]);</code>
$y = (A_m \otimes I_n)x$	<code>for(i=0;i&lt;m;i++)   y[i:n:i+m-1] =     A(x[i:n:i+m-1]);</code>
$y = (\bigoplus_{i=0}^{m-1} A_n^i)x$	<code>for(i=0;i&lt;m;i++)   y[i*n:1:i*n+n-1] =     A(i, x[i*n:1:i*n+n-1]);</code>
$y = D_{m,n}x$	<code>for(i=0;i&lt;m*n;i++)   y[i] = Dmn[i]*x[i];</code>
$y = L_m^{mn}x$	<code>for(i=0;i&lt;m;i++)   for(j=0;j&lt;n;j++)     y[i+m*j]=x[n*i+j];</code>

**Table 1.** Translating formulas to code.  $x$  denotes the input and  $y$  the output vector. The subscript of  $A$  and  $B$  specifies the size of the (square) matrix. We use Matlab-like notation:  $x[b:s:e]$  denotes the subvector of  $x$  starting at  $b$ , ending at  $e$  and extracted at stride  $s$ .

as shown in Table 1. The implementation of (1) would hence have four steps corresponding to the four factors in (1).

Observe in Table 1 that the multiplication of a vector by a tensor product containing an identity matrix can be computed using loops. The working set for each of the  $m$  iterations of  $y = (I_m \otimes A_n)x$  is a contiguous block of size  $n$  and the base address is increased by  $n$  between iterations. In contrast, the working sets of size  $m$  of the  $n$  iterations of  $y = (A_m \otimes I_n)x$  are interleaved, leading to stride  $n$  within one iteration and a unit stride base update across iterations.

**Cost analysis.** Computing the DFT using (1) requires, independent of the recursion strategy,  $O(n \log n)$  floating-point operations on  $O(n)$  data elements. Hence, the degree of reuse is  $O(\log n)$ .

The exact number of real operations depends on the chosen factorizations of  $n$  and is between 4 and  $5n \log_2 n$  operations.

Currently, the best FFT algorithm (in terms of the arithmetic cost) for 2-power sizes requires  $\frac{34}{9}n \log_2 n + O(n)$  many operations [47].

**Iterative FFTs.** The original FFT by Cooley and Tukey [14] was not the recursive algorithm (1), but its iterative equivalent. It can be obtained by expanding the DFT recursively by always using the factorization  $n = 2 \cdot n/2$ , and then rearranging the

parentheses and fusing adjacent permutations. The result is the iterative FFT

$$\text{DFT}_n = \left( \prod_{i=1}^k (I_{2^{i-1}} \otimes \text{DFT}_2 \otimes I_{2^{n-i}}) D'_{n,i} \right) R_n, \quad (4)$$

where the  $D_{n,i}$  are diagonal matrices and  $R_{p,n}$  is the bit-reversal permutation [68].

This algorithm underlies the implementation in Numerical Recipes [58] whose performance was shown in Fig. 2. The corresponding code is shown below.

```
#include <math.h>

#define SWAP(a,b) tempr=a;a=b;b=tempr
void four1(float *data, int *nn, int *isign)
{ /* altered for consistency with original FORTRAN.
  /* Press, Flannery, Teukolsky, Vetterling "Numerical
  * Recipes in C" tuned up ; Code works only when *nn is
  * a power of 2 */
  int n, mmax, m, j, i;
  double wtemp, wr, wpr, wpi, wi, theta, wpin;
  double tempr, tempi, datar, datai,
    datalr, datali;
  n = *nn * 2;
  j = 0;
  for(i = 0; i < n; i += 2)
  { if (j > i) { /* could use j>i+1 to help
                * compiler analysis */
      SWAP(data[j], data[i]);
      SWAP(data[j + 1], data[i + 1]);
    }
    m = *nn;
    while (m >= 2 && j >= m) {
      j -= m;
      m >>= 1;
    }
    j += m;
  }
  theta = 3.141592653589795 * .5;
  if (*isign < 0)
    theta = -theta;
  wpin = 0; /* sin(+PI) */
  for(mmax = 2; n > mmax; mmax *= 2)
  { wpi = wpin;
    wpin = sin(theta);
    wpr = 1 - wpin * wpin - wpin * wpin;
    /* cos(theta*2) */
    theta *= .5;
    wr = 1;
    wi = 0;
    for(m = 0; m < mmax; m += 2)
```

```

{ j = m + mmax;
  tempr = (double) wr *(datalr = data[j]);
  tempi = (double) wi *(datali = data[j + 1]);
  for(i = m; i < n - mmax * 2; i += mmax * 2)
  { /* mixed precision not significantly more
     * accurate here; if removing double casts,
     * tempr and tempi should be double */
    tempr -= tempi;
    tempi = (double) wr *datali + (double) wi *datalr;
    /* don't expect compiler to analyze j > i+1 */
    datalr = data[j + mmax * 2];
    datali = data[j + mmax * 2 + 1];
    data[i] = (datar = data[i]) + tempr;
    data[i + 1] = (datai = data[i + 1]) + tempi;
    data[j] = datar - tempr;
    data[j + 1] = datai - tempi;
    tempr = (double) wr *datalr;
    tempi = (double) wi *datali;
    j += mmax * 2;
  }
  tempr -= tempi;
  tempi = (double) wr *datali + (double) wi *datalr;
  data[i] = (datar = data[i]) + tempr;
  data[i + 1] = (datai = data[i + 1]) + tempi;
  data[j] = datar - tempr;
  data[j + 1] = datai - tempi;
  wr = (wtemp = wr) * wpr - wi * wpi;
  wi = wtemp * wpi + wi * wpr;
}
}
}

```

### Further Reading.

- *FFT algorithms*. [56, 66] give an overview of FFT algorithms. [68] uses the Kronecker product formalism to describe many different FFT algorithms, including parallel and vector algorithms. [46] uses the Kronecker formalism to parallelize and vectorize FFT algorithms.
- *FFTW*. FFTW can be downloaded at [24]. The latest version, FFTW3, is described in [33]. The previous version FFTW2 is described in [32] and the codelet generator *genfft*, in [31].
- *SPIRAL*. Spiral is a program generation system for transforms. The core system is described in [60] and on the web at [63]. Using Kronecker product manipulations, SIMD vectorization is described in [26, 27], shared memory (SMP and multicore) parallelization in [28], and message passing (MPI) in [11].
- *Open source FFT libraries*. FFTPACK [23] is a mixed-radix Fortran FFT library. The GNU Scientific library (GSL) [38] contains a C port of FFTPACK. UHFFT

[51, 67] is an adaptive FFT library. Numerical Recipes [58] contains the radix-2 FFT implementation shown above. FFTW [22] provides a parallel FFT library for distributed memory machines.

- *Proprietary FFT libraries.* The AMD Core Math Library (ACML) [1] is the vendor library for AMD processors. Intel provides fast FFT implementations as a part of their Math Kernel Library (MKL) [52] and Integrated Performance Primitives (IPP) [45]. IBM’s IBM Engineering and Scientific Software Library (ESSL) [21] and the parallel version (PESSL) contain FFT functions optimized for IBM machines. The vDSP library contains FFT functions optimized for AltiVec. The libraries of the Numerical Algorithms Group (NAG) [55] and the International Mathematical and Statistical Library (IMSL) [43] also contain FFT functionality.

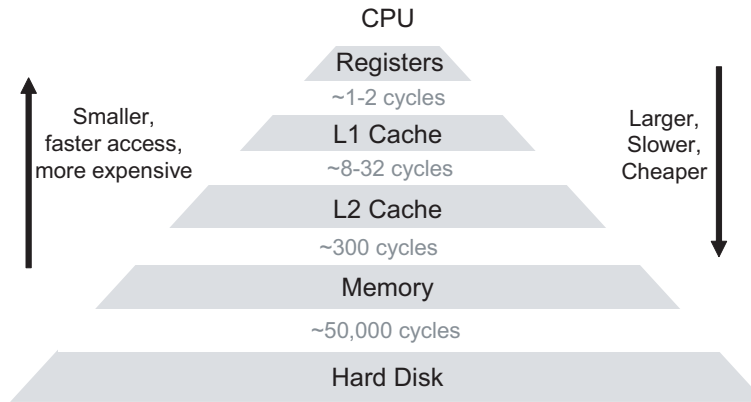
## 2.4 State-Of-The-Art Desktop and Laptop Computer Systems

Modern computers include several performance enhancing microarchitectural features like cache systems, a memory hierarchy, virtual memory, and CPU features like vector and parallel processing. While these features usually increase the achievable performance, they also make the optimization process more complex. This section introduces several microarchitectural features relevant to this tutorial. For further reading, refer to [40, 12].

**Memory hierarchy.** Most computer systems use a *memory hierarchy* to bridge the speed gap between the processor(s) and its connection to main memory. As shown in Fig. 5, the highest levels of the memory hierarchy contain the fastest and the smallest memory systems, and vice versa.

A hierarchical memory enables the processor to take advantage of the memory locality of computer programs. Optimizing numerical programs for the memory hierarchy is one of the most fundamental approaches to producing fast code, and the subject of this tutorial. Programs typically exhibit temporal and spatial memory locality. Temporal locality is the concept that a memory location that is referenced by a program will likely be referenced again in the near future. Spatial locality is the concept that the likelihood of referencing a memory location by a program is higher if a nearby location was recently referenced. High performance computer software must be designed so that the hardware can easily take advantage of locality. Thus, this tutorial focuses on writing fast code by designing programs that exhibit maximal temporal and spatial localities.

**Registers.** Registers inside the processor are the highest level of the memory hierarchy. Any value (address or data) that is involved in computation has to eventually be placed in a register. Registers may be designed to hold only a specific type of value (special purpose registers), or only floating point values (eg., double FP registers), vector values (vector registers), or any value (general purpose registers). The number of registers in a processor varies by architecture. A few examples are provided in Table 2. When an active computation requires more values to be held than the register space will allow, a *register spill* occurs, and the register contents are written to lower levels of memory from which they will be reloaded again. Register spills are expensive. To avoid them



**Fig. 5.** Memory hierarchy. Typical latencies for data transfers from the CPU to each of the levels are shown. The numbers shown here are only an indication, and the actual numbers will depend on the exact architecture under consideration and the access sequence of the program.

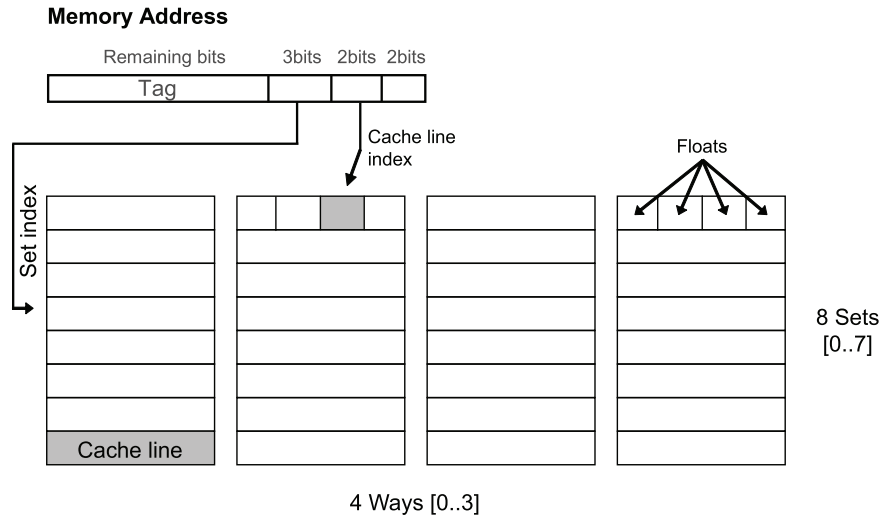
and speed up computation, a processor might make use of internal registers that are not visible to the programmer. Many optimizations that work on other levels of the memory hierarchy can typically also be extended to the register level.

Processor	Integer Registers	Double FP Registers
Core 2 Extreme	16	16
Itanium 2	128	128
UltraSPARC T2	32	32
POWER6	32	32

**Table 2.** Sample scalar register space (per core) in various architectures. In addition to integer and FP registers the Core 2 Extreme also has 16 multimedia registers.

**Cache memory.** Cache memory is a small, fast memory that resides between the main memory and the processor. It reduces average memory access times by taking advantage of spatial and temporal locality. When the processor initially requests data from a memory location (called a cache miss), the cache fetches and stores the requested data and data spatially close. Subsequent accesses, called *hits*, can be serviced by the cache without needing to access main memory. A well designed cache system has a low miss to hit ratio (also known as just the miss ratio or miss rate).

Caches, as shown in Fig. 6 are divided into cache lines (also known as blocks) and sets. Data is moved in and out of cache memory in chunks equal to the line size. Cache lines exist to take advantage of spatial locality. Multiple levels of caches and separate data



**Fig. 6.** 4-way set associative cache with cache line size of 4 single precision floats (4 bytes per float) per line, and cache size of 128 floats (total cache size is 512 bytes). The figure also illustrates the parts of a memory address used to index into the cache. Since each data element under consideration is 4 bytes long, the two least significant bits are irrelevant in this case. The number of bits used to address into the cache line would be different for double precision floats.

Level/Type	Size	Associativity
L1 Data (per core)	32 kB	8-way set associative
L1 Instruction (per core)	32 kB	8-way set associative
L2 Unified (common)	4 MB	8-way set associative

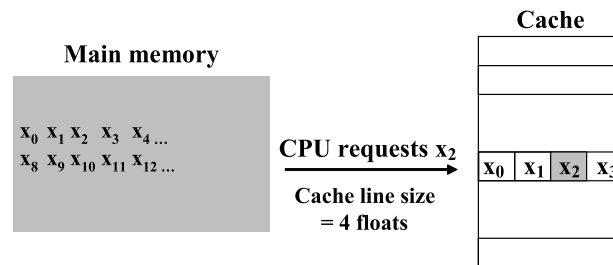
**Table 3.** Cache system example: Intel Core2 Duo, Merom Notebook processor.

and instruction caches may exist, as shown in Table 3. Caches may be direct mapped (every main memory location is mapped to a specific cache location) or  $k$ -way set associative (every main memory location is mapped to  $k$  possible cache locations). In addition to misses caused due to data being brought in for the first time and those due to cache capacity constraints, caches that are not fully associative might incur conflict misses [41].

Since cache misses are typically expensive, writing fast code involves designing programs to have a low miss rates. This is accomplished using two important guiding principles, illustrated in Fig. 7 and described below:

- **Reuse.** Once data is brought into the cache, the program should reuse it as much as possible before it gets evicted. In other words, programs must try to avoid scattering computations made on a particular data location throughout the execution of the program. Otherwise, the same data (or data location) has to go through several





**Fig. 7.** Neighbor use and reuse: When the CPU requests  $x_2$ ,  $x_0$ ,  $x_1$ , and  $x_3$  are also brought into the cache since the cache line size holds 4 floats.

cycles of being brought into the cache and subsequently evicted, which increases runtime.

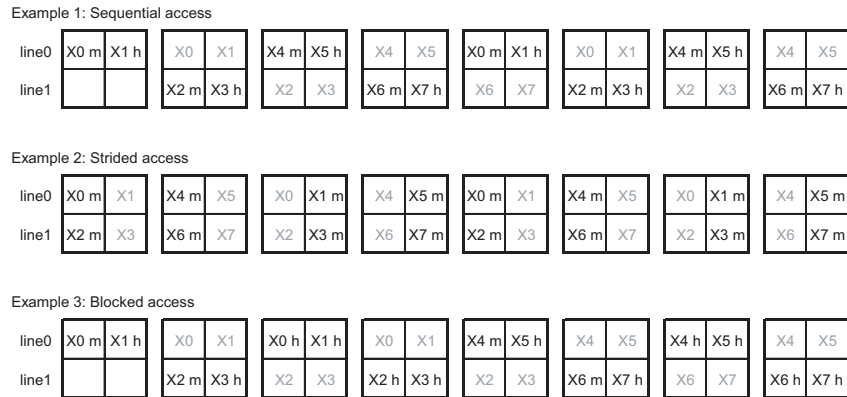
- **Neighbor use (Using all data brought in).** Data is always brought into the cache in chunks the size of a cache line. This is seen in Fig. 7, where one data element  $x_2$  was requested, and three others are also brought in since they belong to the same cache line. To take advantage of this, programs must be designed to perform computations on neighboring data (physically close in memory) before the line is evicted. This might involve reordering loops, for instance, to work on data in small chunks.

These two principles work at multiple levels. For instance, code can be designed to use and reuse all data within a single cache block, as well as within an entire cache level. In fact, these principles hold throughout the memory hierarchy, and thus can be used at various cache and memory levels. Depending on the computation being done, techniques that use these principles may not be trivial to design or implement.

In scientific or numerical computing, data typically consists of floating point numbers. Therefore, it helps to view the cache organization, lines, and sets in terms of the number of floating point numbers that can be held. For instance, the cache shown in Fig. 6 is a 512 byte, 4-way set associative cache with a line size of 16 bytes. There are a total of 32 lines (512 bytes / 16 bytes per line), and 4 sets (32 lines / 8 lines per set). If we note that each cache line can hold 4 floats (16 bytes / 4 bytes per float), we can immediately see that the cache can hold a total of 128 floats. This means that datasets larger than 128 floats will not fit in the cache. Also, if we make an initial access to 128 consecutive floats, there will be a total of 32 cache misses and 96 cache hits (since 4 floats in a line are loaded on each cache miss). This gives us a rough estimate of the runtime of such a set of accesses, which is useful both in designing programs and in performing sanity checks.

**Cache analysis.** We now consider three examples of accessing an array in various sequences, and analyze their effects on the cache.

Consider a simple direct mapped 16 byte data cache with two cache lines, each of size 8 bytes (two floats per line). Consider the following code sequence, in which the array



**Fig. 8.** Cache access analysis: The state of the complete cache for each example is shown after every two accesses, along with whether the two accesses resulted in hits or misses (shown by ‘h’ or ‘m’). The two requests just made are shown in black, while the remaining parts of the cache are shown in gray. To save space, square brackets are not shown: ‘X0’ refers to X[0].

X is cache-aligned (that is, X[0] is always loaded into the beginning of the first cache line) and accessed twice in consecutive order:

```
float X[8];
for(int j=0; j<2; j++)
  for(int i=0; i<8; i++)
    access(X[i]);
```

The top row on Fig. 8 shows the states of the cache after every two (out of the total of sixteen) accesses for this example. To analyze the cache footprint and pattern of this code sequence, we first observe that the size of the array (8 floats) exceeds the size of the cache (4 floats). We then observe that a total of 16 accesses are made to the array. To calculate how many result in hits, and how many in misses, we observe the cache access pattern of the code. The pattern is “0123456701234567” (only the indices of X accessed are shown). We note that an access to any even index of X results in that element and the subsequent element being loaded since they are in the same cache line. Thus, accessing X[0] loads X[0] and X[1] into the cache. We can then compute the hit/miss pattern to be: “MHMHMHMHMHMHMHMH”. So in all, there are 8 hits and 8 misses.

We now look at another code sequence that again accesses the same array twice (similar to the last example), albeit with a stride of 2:

```
float X[8];
for(int j=0; j<2; j++)
{
  for(int i=0; i<7; i+=2)
    access(X[i]);
  for(int i=1; i<8; i+=2)
    access(X[i]);
}
```

The middle row on Fig. 8 shows the corresponding cache states for this example. The access pattern here is "0246135702461357". A similar analysis shows us that the miss ratio is even worse: every single access in this pattern results in a miss (with a total of 16 misses and 0 hits). This example illustrates an important point: strided accesses generally result in poor cache efficiency, since they effectively "make the cache smaller."

Finally, let us consider a third code sequence that again accesses the same array twice:

```
float X[8];
for(i=0; i<2; i++)
  for(k=0; k<2; k++)
    for(j=0; j<4; j++)
      access(X[j+(i*4)]);
```

The bottom row on Fig. 8 shows the corresponding cache states for this example. The access pattern here is "0123012345674567". Counting the hits and misses, ("MHMH-HHHMHMHMHMH"), we observe that there are 12 hits and 4 misses. We also note that if this rearrangement is legal, it is a cache optimized version of the original code sequence. In fact, this rearrangement is an example of both of the previously mentioned principles behind optimizing for the memory hierarchy: reuse and neighbor use. Unlike the first example, the "0123" block is reused here before being evicted. Unlike the second example, every time an even-indexed element is accessed, the succeeding odd-indexed element which is a part of the same cache line is also immediately accessed. Thus, analyzing the cache can help us estimate and improve the cache performance of a program.

**CPU features.** Modern microprocessors also contain other performance enhancing features. Most processors contain pipelined superscalar out-of-order cores with multiple execution units. Pipelining is a form of parallelism where different parts of the processor work simultaneously on different components of different instructions. Superscalar cores can retire more than one instruction per processor clock cycle. Out-of-order processing cores can detect instruction dependencies and reschedule the instruction sequence for performance. The programmer has to be cognizant of these features in order to be able to optimize for a particular architecture.

Most such aggressive cores also contain multiple execution units (for instance, floating point units) for increased performance. This means that a processor might be able to, for instance, simultaneously retire one floating point add instruction every cycle, and one floating point multiplication instruction every other cycle. It is up to the programmer and the compiler to keep the processor's execution units adequately busy (primarily via instruction scheduling and memory locality) in order to achieve maximum performance.

The theoretical rate at which a processor can perform floating point operations is known as the processor's *theoretical peak performance*. This is measured in flop/s (Floating point Operations per Second). For instance, a processor running at 1GHz that can retire one addition every cycle, and one multiplication every other cycle has a theoretical peak of 1.5 Gflop/s. The theoretical peak of a Core 2 Extreme processor operating under various modes is shown in Table 4.

	1 core	2 cores	4 cores
x87 double	2	4	8
SSE2 double	4	8	16
x87 float	2	4	8
SSE float	8	16	32

**Table 4.** Core 2 Extreme: Floating-point operations per cycle for different operation modes.

In practice, cache misses, pipeline stalls due to dependencies, branches, branch mispredictions, and the fact that meaningful programs contain instructions other than floating point instructions, do not allow a processor to perform at its theoretical peak performance. Further, the achievable performance also depends on the inherent limitations of the algorithm, such as reuse. For example, MMM, with a reuse degree of  $O(n)$  can achieve close to the peak performance of 48 Gflop/s (as seen in Fig. 3), whereas the DFT with a reuse degree of  $O(\log(n))$  reaches only about 50% (as seen in Fig. 2).

In summary, knowing a processor’s theoretical peak and the an algorithm’s degree of reuse gives us a rough estimate of the extent to which a program could potentially be improved.

Modern processors also contain two major explicit forms of parallelism: vector processing and multicore processing, which are important for writing fast code, but beyond the scope of this tutorial.

#### Further reading.

- *General computer architecture.* [40, 12]
- *CPU/architecture specific.* [44, 2]

## 2.5 Using Compilers

To produce fast code it is not sufficient to write and optimize source code—the programmer must also ensure that the code that is written gets compiled into an efficient binary executable. This involves careful selection and use of compiler flags, use of language extensions, and monitoring and analyzing the compiler’s output. Furthermore, in some situations, it is best to let the compiler know of all the degrees of freedom it has, so it can optimize well. In other situations, it is best to direct the compiler to do exactly what is required. This section goes over the compile process, what to keep in mind before, while, and after compiling, and some of the common pitfalls related to the compiling process.

**Variable declaration: memory allocation.** Understanding how C handles the allocation of space for variables is beneficial. C assigns variables to different *storage class specifiers* by default, based on where in the source code they appear. The default stor-

age class for a variable can be overridden by preceding a variable declaration with the desired storage class specifier.

Variables that are shared among source files use the `extern` storage class. Global variables belong to the `static` storage class, and typically exist in static memory (as do extern variables), which means that they exist as long as the program executes. Local variables belong to the `auto` (automatic) storage class, which means that they are allocated on the stack automatically upon entering the local block within which they are defined, and destroyed upon exit. The `register` storage class requests that the compiler allocates space for the variable directly in the CPU registers. These are useful to eliminate load/store latencies on heavily used variables. Keep in mind that depending on the compiler being used, care should be taken to initialize variables before usage.

**Variable declaration: qualifiers.** Most compilers provide further means to specify variable attributes through *qualifiers*. A `const` qualifier specifies that a variable's value will never change. A `volatile` qualifier is used to refer to variables whose values might be influenced by sources external to the compiler's knowledge. Operations involving volatile variables will not be optimized by the compiler, in order to preserve correctness. A `restrict` qualifier is especially useful to writing fast code, since it tells the compiler that a certain memory address will be restricted to access via the specified pointer. This allows for effective compiler optimization.

Finally, memory alignment can also be specified by qualifiers. Such qualifiers are specific to the compiler being used. For instance, `__attribute__((aligned(128)))` requests a variable to be aligned at the specified 128-byte memory boundary. Such requests allow variables to be aligned to cache line boundaries or virtual memory pages as desired. Similar qualifiers can be used to tell the compiler that the address pointed to by a pointer is memory aligned.

**Dynamic memory allocation.** Dynamic memory allocation, using `malloc` for example, allocates memory in the *heap*, and a pointer to the allocated memory is returned. If alignment is of importance, many libraries provide a `memalign` function (the Intel equivalent is `_mm_malloc`) to allocate memory aligned to a specified boundary. The alternative is to allocate more memory than required, and to then check and shift the returned pointer adequately to achieve the required alignment.

**Inline assembly and intrinsics.** Sometimes, it is best to write assembly code to access powerful features of the machine which may not be available via C. Assembly can be included as a part of any program in C using inline assembly. However, inline assembly use must be minimized as it might interfere with compiler optimizations. Architecture vendors typically provide C language extensions to allow programmers to access special machine instructions. These extensions, called *intrinsics*, are similar to function calls that allow the programmer to avoid writing inline assembly. Importantly, intrinsics allow the compiler to understand what data and/or control the programmer is manipulating, thus allowing for better optimization. As an example, Intel's MMX and SSE extensions to the x86 ISA can be accessed via C intrinsics provided by Intel.

**Compiler flags.** Most compilers are highly configurable via a plethora of command line options and flags. In fact, finding the right set of compiler options that yield optimal performance is non-trivial. However, there are some basic ideas to keep in mind while using a compiler, as listed below. Note that these ideas apply to most compilers.

- *C standards.* A compiler can be set to follow a certain C standard such as C99. Certain qualifiers and libraries might need specific C standards to work. By switching to a newer standard, the programmer can typically communicate more to the compiler, thus enabling it to work better.
- *Architecture specifications.* Most compilers will compile and optimize by default for a basic architecture ISA standard to maximize compatibility. Machine and architecture specific optimizations may not be performed as a result. For instance, a compiler on an AMD Athlon processor may compile to the x86 standard by default, and not perform Athlon-specific optimizations. Instructing the compiler to compile for the correct target architecture may result in considerable performance gains. Additional flags may be required for these optimizations. For example, gcc requires the `-sse` flag to include vector instructions.
- *Optimization levels.* Most compilers usually define several optimization levels that can be selected. Determining the optimization level that yields maximum performance is a black art usually done by trial and error. A more aggressive optimization level doesn't necessarily yield better performance. Optimization levels are usually a shortcut to turn on or off a large set of compiler flags (discussed next).
- *Specialized compiler options.* Compilers typically perform numerous optimizations, many which can be selectively turned on or off and configured through command line flags. Loop unrolling, function inlining, instruction scheduling, and other loop optimizations are only some of the available configurable optimizations. Usually, finding the right optimization level is sufficient, but sometimes, inspection of assembly code provides insights that can be used to fine-tune compiler optimizations.

**Compiler output.** The output of the compiler is usually an executable binary. As mentioned earlier, the compiler can also be used to produce various intermediate stages, including the preprocessed source, assembly code, and the object code. Sometimes, it is important and useful to visually inspect the assembly code to better understand both the performance of an executable and the behavior of the compiler.

Compilers also output warnings, which can be controlled through compiler flags. Sometimes, a seemingly innocuous warning might provide excellent insights into the source of a bug, which makes these warnings a significant debugging tool.

Optimization reports are an important part of the compiler output that must be inspected. For instance, a vectorizing compiler will inform the programmer of whether it was able to successfully vectorize or not. A failure to vectorize a program that was expected to be vectorized is a reason for examining the program carefully, and modifying or annotating the code as necessary.

In conclusion, it is important for programmers to be knowledgeable about the compiler that they use in order to be able to use the compiler efficiently, and to ensure that poor compiler usage does not diminish the results of code designed for high performance.

**Further reading.**

- *Gnu Compiler Collection (gcc)*. [34]
- *Intel compiler*. [61]

**2.6 Exercises**

1. **Operations count.** The Walsh-Hadamard transform (WHT) is related to the DFT but has a simpler structure and simpler algorithms. The WHT is defined only for 2-power input sizes  $n = 2^k$ , as given by the matrix

$$\text{WHT}_{2^k} = \underbrace{\text{DFT}_2 \otimes \text{DFT}_2 \otimes \dots \otimes \text{DFT}_2}_{k \text{ factors}},$$

where  $\text{DFT}_2$  is as defined in (2).

- (a) How many entries of the WHT are zeros and why? Determine the number of additions and the number of multiplications required when computing the WHT by definition.
- (b) The WHT of an input vector can be computed iteratively or recursively using the following formulas:

$$\text{WHT}_{2^k} = \prod_{i=0}^{k-1} (I_{2^{k-i-1}} \otimes \text{DFT}_2 \otimes I_{2^i}) \quad (\text{iterative}) \quad (5)$$

$$\text{WHT}_{2^k} = (\text{DFT}_2 \otimes I_{2^{k-1}})(I_2 \otimes \text{WHT}_{2^{k-1}}) \quad (\text{recursive}) \quad (6)$$

- (c) Determine the exact operations counts (again, additions and multiplications separately) of both algorithms. Also determine the degree of reuse as defined in Section 2.1.
2. **Direct implementations.** Implement and execute generic versions of the following:
    - the direct implementation of MMM (code snippet given in Section 2.2),
    - the Numerical Recipes code for the DFT as given in [58],
    - a recursive implementation of the WHT based on (6)

This code will also be used in the exercises of later sections.

3. **Determining hardware information.** In this exercise, you will determine the relevant hardware configuration of your computer. You will use this information in later exercises.

Determine the following information about your computer:

- CPU type and clock speed
- For each cache: size, associativity, and cache line size
- Size of main memory
- System bus speed

Here are a few tips on how to determine this information:

- Look in the computer’s manual.
- Look in the CPU manufacturer’s manual.
- For CPU information in Linux, execute `cat /proc/cpuinfo`.
- For cache information in Linux, search for lines with the word ‘cache’ in the kernel ring buffer. You can do so by typing:  
`dmesg | grep '^CPU.*cache'` on most systems.

### 3 Performance Optimization: The Basics

In this section we will review the basic steps required to assess the performance of a given implementation, also known as “benchmarking.” We focus on runtime benchmarking as the most important case. (Other examples of benchmarking includes assessing the usage of memory or other resources.)

For a given program, the basic procedure consists of three steps:

1. finding the hotspots (hotspots are the most frequently executed code regions),
2. timing the hotspots, and
3. analyzing the measured runtimes.

It is essential to find the parts of the program that perform the bulk of the computation and restrict further investigation to these *hotspots*. Optimizing other parts of the program will have little to no effect on the overall runtime. In order to obtain a meaningful runtime measurement, one has to build a test environment for each hotspot that exercises and measures it in the correct way. Finally, one has to assess the measured data and relate it to the cost analysis of the respective hotspot. This way one can make efficiency statements and target the correct (inefficient) hotspot for further optimization.

#### 3.1 Finding The Hotspots

The first step in benchmarking is to find the parts of the program where most time is spent. Most development platforms contain a *profiling* tool. For instance, the development environment available on the GNU/Linux platform contains the GNU `gprof`



profiler. On Windows platforms, the Intel VTune tool [17] that plugs into Microsoft's Visual Studio [18] can be used to profile applications.

If no profiling tool is available, obtaining first-order profiling information can be obtained by inserting statements throughout the program that print out the current system time. In this case, less is more, as inserting too many time points may have side effects on the measured program.

**Example: GNU tool chain.** We provide a small example of using the GNU tool chain to profile a sample program.

Consider the following program:

```
#include <stdio.h>

float function1()
{ int i; float retval=0;
  for(i=1; i<1000000; i++)
    retval += (1/i);
  return(retval);
}

float function2()
{ int i; float retval=0;
  for(i=1; i<10000000; i++)
    retval += (1/(i+1));
  return(retval);
}

void function3() { return; }

int main()
{ int i;
  printf("Result: %.2f\n", function1());
  printf("Result: %.2f\n", function2());
  if (1==2) function3();
  return(0);
}
```

Our final objective is to optimize this program. In order to do so, we first need to find where the program spends most of its execution time, using `gprof`.

As specified in the `gprof` manual [35], three steps are involved in profiling using `gprof`:

1. Compile and link with profiling enabled:

```
gcc -O0 -lm -g -pg -o ourProgram ourProgram.c
```

The resulting executable is instrumented. This means that in addition to executing your program, it will also write out profiling data when executed. (Note: We use the `-O0` flag to prevent the compiler from inlining our functions and performing

other optimizing transforms that might make it difficult for us to make sense of the profile output. For profiling to provide us with meaningful information, we would need to compile at the level of optimization that we intend to finally use, with the understanding that mapping the profiler's output back to the source code in this case might involve some effort.)

2. Execute the program to generate the profile data file

```
./ourProgram
```

The program executes and writes the profile data to `gmon.out`.

3. Run `gprof` on the profile data file to analyze the profile data

```
gprof ourProgram gmon.out > profile.txt
```

The analysis is now contained in `profile.txt`. This file shows you how many times each function was executed, and how much time was spent in each function, and plenty of other detail. For our example program, we obtain:

% time	cumulative seconds	self seconds	calls	self ms/call	total ms/call	name
92.68	0.38	0.38	1	380.00	380.00	function2
7.32	0.41	0.03	1	30.00	30.00	function1

We can see that most of the program runtime was spent in executing `function2`, with relatively little spent on `function1`. This tells us that it is most important to optimize the runtime of `function2`.

Further down in `profile.txt`, we see that `gprof` also tells us if the time taken by a function was spent inside the function or inside other function calls made by the function. Note that `gprof` can take several other arguments to produce different kinds of profiling analyses for the executable, including the number of times a certain line in the source code was executed.

### 3.2 Timing a Hotspot

Once the hotspots have been found, we need to measure their runtime for further analysis. Each hotspot must be timed separately with an appropriate timing routine. The general idea is the following:

1. Read the current time (start time) from the appropriate time source.
2. Execute the kernel/hotspot. Iterate an adequate number of times to obtain a meaningful value off the time source.
3. Read the current time (end time) from the appropriate time source.

4. Execution time of the kernel/hotspot =  $\frac{\text{End time} - \text{Start time}}{\text{Number of iterations}}$ .

We first discuss time sources, and reading the time from them; then we explain how to write a timing routine to get meaningful results.

**Time functions.** Depending on the system one is using, a variety of time sources to “get the current time” may be available:

- Most Unix systems define `gettimeofday()` to portably query the current time (as defined in IEEE Std 1003.1).
- ANSI C defines `ctime()` and `clock()` as portable ways of obtaining the current time.
- On Intel processors, the `rdtsc` instruction reads the time stamp counter which allows near-cycle accurate timing. On PowerPC processors, the `mf spr` instruction reads the time-base register.

Generally, portable time functions have much less precision than cycle-counter-based methods. The pros and cons of various timing methods are listed below:

Timer type	Advantages	Disadvantages
Wall clock; unix: <code>gettimeofday()</code>	Simple to use, highly portable	Low resolution, does not account for background tasks
System timer; unix: <code>time</code> command	Gives wall clock, user-cpu, and system-cpu times	Relatively low resolution
Hardware times-tamp counter (discussed below)	High resolution, most precise and accurate	Does not account for background system load (effectively, wall clock time), best for kernels with short runtimes; non-portable

We give a simplified example of a timing macro based on `rdtsc` (a hardware times-tamp counter) for a 32-bit Intel processor to be used with Microsoft VisualStudio:

```
typedef union
{
  __int64 int64;
  struct {__int32 lo, hi;} int32;
} tsc_counter;

#define RDTSC(cpu_c) \
{ __asm rdtsc \
  __asm mov (cpu_c).int32.lo,eax \
  __asm mov (cpu_c).int32.hi,edx \
}
```

The corresponding code sequence in GNU C looks slightly different:

```
typedef union
{ unsigned long long int64;
```

```

    struct {unsigned int lo, hi;} int32;
} tsc_counter;

#define RDTSC(cpu_c)          \
    __asm__ __volatile__ ("rdtsc" :      \
    "=a" ((cpu_c).int32.lo),          \
    "=d" ((cpu_c).int32.hi))

```

**Timing routine.** A timing routine calls the function that is to be timed without executing the original program. The objective is to isolate the kernel and measure the runtime of the kernel with the least disturbance and highest accuracy possible. A timing routine consists of the following steps:

- Initialize kernel-related data structures.
- Initialize kernel input data.
- Call kernel a few times to warm up microarchitectural components.
- Read current time.
- Call kernel multiple times to obtain an adequately precise value from the timing source used.
- Read current time.
- Divide the time difference by the number of kernel calls.

To obtain more stable timing results, one often has to run multiple timings and take the average or minimum value.

We give an example timing routine for a matrix multiplication function, using the RDTSC macro defined above.

```

double time_MMM(int N, double *src, double *dst)
{ // init
  for(i=0; i<N; i++)
    src[i] = 0.0;
  init_MMM(A,B,C,N,K,M); // if needed

  // warm up
  MMM(A,B,C,N,K,M);

  // time
  RDTSC(t0);
  for(int i=0; i<TIMING_REPETITIONS; i++)
    MMM(A,B,C,N,K,M);
  RDTSC(t1);

  // compute runtime in cycles
  return (double)((t1.int64-t0.int64)/TIMING_REPETITIONS);
}

```

**Known problems.** The following problems may occur when timing numerical kernels:

- Too few iterations of the function to be timed are executed between the two time stamp readings, and the resulting timing is inaccurate due to poor timer resolution.
- Too many iterations are executed between the two time stamp readings, and the resulting timing is affected by system events.
- The machine is under load and the load has side effects on the measured program.
- Multiple timing jobs are executed concurrently, and they interfere with one another.
- Data alignment of input and output triggers cache problems.
- Virtual-to-physical memory translation makes timing irreproducible.
- The time stamp counter overflows and either triggers an interrupt or produces a meaningless value.
- Reading the timestamp counters requires hundred(s) of cycles, which itself affects the timing.
- The linking order of object files changes locality of static constants and this produces cache interference.
- The machine was not rebooted in a long time and the operating system state causes problems.
- The control flow in the numerical kernel being timed is data-dependent and the test data is not representative.
- The kernel is in-place (i.e., all computations are performed directly on the original data locations without creating copies), and the norm of the output is larger than the norm of the input. Repetitive application of the kernel leads to exponential growth of the norm and finally triggers floating-point exceptions which interfere with the timing.
- The transform is timed with a zero vector, and the operating system is “smart,” and responds to a request for a large zero-vector dynamic memory allocation by returning a special zero-valued copy-on-write virtual memory page. Read accesses to this “page” would be much faster than accesses to a page that is actually allocated, since this page is a special one maintained by the operating system for efficiency.

One needs to be very careful when timing numerical kernels to rule out these problems. Getting highly accurate, reproducible, stable timing results for the full range of problem sizes is often nontrivial. Note that small problem sizes may suffer from timer resolution issues, while large problem sizes with longer runtimes may suffer from intervening processes.

### 3.3 Analyzing the Measured Runtime

We now know how to calculate the theoretical peak performance and the memory bandwidth for our target platform, and how to obtain the operations count and the runtime for our numerical kernel. The next step is to use these to conduct a performance analysis that answers two questions:

- What is the limiting resource, i.e., is the kernel CPU-bound or memory-bound? This provides an idea of the various optimization methods that can be used.
- How efficient is the implementation with respect to the limiting resource? This shows the potential performance increase we can expect through optimization.

**Normalization.** To assess the runtime behavior of a kernel as function of problem size, the runtime (or inverse runtime) has to be normalized with the asymptotic or exact operations count. For instance, FFT performance is usually reported in pseudo Mflop/s. This value is computed as  $5n \log_2 n / \text{runtime}$  for  $\text{DFT}_n$ , using the operations count of the radix-2 FFT algorithm as normalization factor. For MMM, the situation is easier, since all currently relevant implementations have the exact operations count  $2n^3$ .

Let us now take a look at at Fig. 2. The Numerical Recipes FFT program achieves almost the same pseudo Mflop/s value, independently of the problem size. This means that all problem sizes run approximately at the same level of (in)efficiency. In contrast, the best code shows a wide variation of performance, generally at a much higher pseudo Mflop/s level. In particular, the performance ramps up to 25 Gflop/s and then drops dramatically. It seems that the kernel becomes more and more efficient with larger problems, but only up to a certain size. Analysis shows that the breakdown occurs once the whole working set of the computation does not fit into the L2 cache any more and the problem switches from being CPU-bound to memory-bound, since the FFT's reuse degree is only  $O(\log(n))$ .

In contrast, Fig. 3 shows that MMM maintains the performance even for out-of-cache sizes. This is possible since MMM has a reuse of  $O(n)$ , higher than the FFT.

Fig. 2 shows that performance plots for high-performance implementations can feature unanticipated characteristics. That is especially true if the kernel changes behavior, for instance, if it slowly changes from being CPU-bound to memory-bound as the kernel size is varied.

**Relative performance.** Absolute performance only tells a part of the story. Comparing the measured performance to the theoretical peak performance shows how efficient the implementation is. A low efficiency for an algorithm with sufficiently high reuse means there is room for optimization.

We continue examining our examples from Fig. 2 and Fig. 3, with the target machine being a Core 2 Extreme at 3 GHz.

In Fig. 2, Numerical Recipes is a single-core single-precision x87 implementation and thus the corresponding peak performance is 6 Gflop/s (see Table 4). As Numerical Recipes reaches around 1 pseudo Gflop/s it runs at about 16% of the peak. Note that if

SSE (4-way vector) instructions and all four cores are used, the peak performance goes up by a factor of 16. (see Table 4). The best scalar code achieves around 4 Gflop/s or about 60% of the x87 peak. The fastest overall code uses SSE and 4 cores and reaches up to 25 Gflop/s or 25% of the quad-core SSE peak.

In Fig. 2, the overall fastest code reaches and sustains about 42 Gflop/s or about 85% of the quad-core SSE2 peak. This is much higher than the FFT due to the higher degree of reuse in the MMM as compared with the FFT.

### 3.4 Exercises

1. In this exercise, you will measure and analyze the performance of the naive implementations of matrix multiplication, the WHT, and the DFT. The steps you will need to follow to complete this exercise are given below. For this exercise, use the hardware configuration of your computer as you determined in Exercise 3 in Section 2.
  - (a) **Determine your computer's theoretical peak performance.** The theoretical peak performance is the number of floating point operations that can be done in a second. This is found by determining the CPU clock speed, and examining the microarchitecture to look at the throughput of floating point operations. For instance, a CPU running at 900 MHz that can retire 2 floating point operations per cycle, has a theoretical peak performance of 1800 Mflop/s. If the type of instructions that the CPU can retire at the same rate includes FMA (fused multiply add) instructions, the theoretical peak would be 3600 Mflop/s (2 multiplies, 2 adds per cycle = 4 operations per cycle, at 900 million cycles per second =  $900 \times 4$  Mflop/s). For this exercise, do not consider vector operations.
  - (b) **Measure runtimes.** Use your implementations of the MMM, WHT, and DFT as completed in Exercise 2. Use the techniques described in Section 3.2 to measure the runtimes of your implementations.
  - (c) **Determine performance and interpret results.**
    - Performance: The performance of your implementation is its number of floating point operations per unit time, measured in flop/s. For the DFT, for instance, the number of operations is taken to be  $5n \log(n)$ .
    - Percentage peak performance: This is simply the percentage of theoretical peak performance. For instance, if your measured code runs at 1.2 Gflop/s on a machine with a peak performance of 3.6 Gflop/s, this implies that your implementation achieves 33.3% of peak performance.
2. **Micro-benchmarks: Mathematical functions.** We assume a Pentium compatible machine. Determine the runtime (in cycles) of the following computations ( $x, y$  are doubles) as accurately as possible:

- $y = x$
- $y = 7.12x$
- $y = x + 7.12$
- $y = \sin(x), x \in \{0.0, 0.2, 4.1, 170.32\}$
- $y = \log(x), x \in \{0.001, 1.00001, 10.65, 2762.32\}$
- $y = \exp(x), x \in \{-1.234e - 17, 0.101, 3.72, 1.234e25\}$

There are a total of 15 runtimes. Explain the results. The benchmark setup should be as follows:

- (a) Allocate two vector doubles  $x[N]$  and  $y[N]$  and initialize all  $x[i]$  to be one of the values from above.
- (b) Use

```
for(i=0; i<N; i++)
    y[i] = f(x[i]);
```

to compute  $y[i] = f(x[i])$ , with  $f()$  being one of the functions above and time this `for`-loop.

- (c) Choose  $N$  such that all data easily fits into L1 cache but there are enough iterations to obtain a reasonable amount of work.
- (d) Use the x86 time stamp counter via the interface provided by `rdtsc.h`, as listed in Section 3.2.

To accurately measure these very short computations, use the following guidelines:

- Only time the actual work, leave everything else (initializations, timing related computations, etc.) outside the timing loop.
- Use the C preprocessor to produce a parameterized implementation to easily check different parameters.
- You may have to run your `for(N)` loop multiple times to obtain reasonable timing accuracy.
- You may have to take the minimum across multiple such measurements to obtain stable results. Thus, you might end up with three nested loops.
- You must perform a warm up for the experiment: variables where you store the timing results, the timed routine and the data vectors should all be loaded into the L1 cache, since cache misses might result in inaccurate timing results.
- Alignment of your data vectors on cache line sizes or page sizes can influence the runtime significantly.



- The use of `CPUID` to serialize the CPU before reading the RDTSC as explained in the Intel manual produces a considerable amount of overhead and may be omitted for this exercise.

## 4 Optimization for the Memory Hierarchy

In this section we describe methods for optimizations targeted at the memory hierarchy of a state-of-the-art computer system. We divide the discussion into four sections:

- Performance-conscious programming.
- Optimizations for cache.
- Optimizations for the registers and CPU.
- Parameter-based performance tuning.

We first overview the general concepts, and then apply them to the MMM and the DFT.

### 4.1 Performance-Conscious Programming

Before we discuss specific optimizations, we need to ensure that our code does not yield poor performance because it violates certain procedures fundamental to writing fast code. Such procedures are discussed in this section. It is important to note that programming for high performance may go to some extent against standard software engineering principles. This is justified if performance is critical.

**Language: C.** For high performance implementations, C is a good choice, as long as one is careful with the language features used (see below). The next typical choice for high-performance numerical code is Fortran, which tends to be more cumbersome to use than C when dynamic memory and dynamic data structures are used.

Object-oriented programming (C++) must be avoided for performance-critical parts since using object oriented features such as operator overloading and late binding incurs significant performance overhead. Languages that are not compiled to native machine code (like Java) should also be avoided.

**Arrays.** Whenever possible, one-dimensional arrays of scalar variables should be used. Higher dimensional objects should be linearized: an  $m \times n$  matrix should be represented by a vector of length  $mn$ , with the matrix element  $(i, j)$  mapped to the vector element  $in + j$ .

**Records.** Complicated `struct` and `union` data types should be avoided. Multiple arrays should be favored over one array with `struct` entries. To represent vectors of complex numbers, vectors of real numbers of twice the size should be used, with the real and imaginary parts appearing as pairs along the vector.

**Dynamic data structures.** Dynamically generated data structures like linked lists and trees must be avoided if the algorithm using them can be implemented on array structures instead. Heap storage must be allocated in large chunks, as opposed to separate allocations for each object.

**Control flow.** `while` loops and loops with complicated termination conditions must be avoided. `for` loops with loop counters and loop bounds known at compile-time must be used whenever possible. `switch`, `?:`, and `if` statements must be avoided in hot spots and inner loops. For small, repetitive tasks, macros are a better choice than functions. Macros are expanded before compilation while the compiler must perform analysis on inline functions.

## 4.2 Cache Optimization

For lower levels in the memory hierarchy (L1, L2, L3 data cache, TLB) the overarching optimization goal is to reuse data as much as possible once brought in. The architecture of a set-associative cache (Fig. 6) suggests three major optimization methods that target different hardware restrictions.

- Blocking: working on data in chunks that fit into the respective cache level, to overcome restrictions due to cache capacity,
- Loop merging: merging consecutive loops that sweep through data into one loop to reuse data in the cache and hence make the best use of the restricted memory bandwidth, and,
- Buffering: copying data into contiguous temporary buffers to overcome conflict cache misses due to cache associativity.

The actual optimization process applies one or more of these ideas to some of the levels of the memory hierarchy. It is not always a good idea to apply all methods to all levels, as code complexity may increase dramatically.

Finally, the correct parameters for blocking and/or buffering on the targeted computer system have to be found. A good approach is to write the program parameterized, i.e., collect all parameters as named constants. Then it is easy to try different parameter settings by hand or using a script to find the variant that achieves the highest performance.

**Blocking.** The basic idea of blocking is to perform the computation in “blocks” that operate on a subset of the input data to achieve memory locality. This can be achieved in different ways. For example, loops in loop nests, like the triple-loop MMM in Section 2.2 may be split and swapped (a transformation called tiling) so that the working set of the inner loops fits into the targeted memory hierarchy level, whereas the outer loop jumps from block to block. Another way to achieve blocking is to choose a recursive algorithm to start with. Recursive algorithms naturally divide a large problem into smaller problems that typically operate on subsets of the data. If designed and parameterized well, at some level all sub-problems fit into the targeted memory level and blocking is

achieved implicitly. An example of such an algorithm is the recursive Cooley-Tukey FFT in (1).

**Loop merging.** Numerical algorithms often have multiple stages. Each stage accesses the whole data set before the next stage can start, which produces multiple sweeps through the working set. If the working set does not fit into the cache this can dramatically reduce performance.

In some algorithms the dependencies do not require that *all* operations of a previous stage are completed before *any* operation in a later stage can be started. If this is the case, loops can be merged and the number of passes through the working set can be reduced. This optimization is essential for implementing high-performance DFT functions.

**Buffering.** When working on multi-dimensional data like matrices, logically close elements can be far from each other in linearized memory. For instance, matrix elements in one column are stored at a distance equal to the number of columns of that matrix. Cache associativity and cache line size get into conflict if one wants to hold, for instance, a small rectangular section of such a matrix in cache, leading to cache thrashing.

One simple solution is to copy the desired block into a contiguous temporary buffer. That incurs a one-time cost but alleviates cache thrashing. This optimization is often called buffering.

### 4.3 CPU and Register Level Optimization

Optimization for the highest level in the memory hierarchy, the registers, is to some extent similar to optimizations for the cache. However it also needs to take into account microarchitectural properties of the target CPU. Current high-end CPUs are superscalar, out-of-order, deeply pipelined, feature complicated branch prediction units, and many other performance enhancing technologies. From a high-level point of view, one can summarize the optimization goals for a modern CPU as follows. An efficient C program should:

- have inner loops with adequately large loop bodies,
- have many independent operations inside an inner loop body,
- use automatic variables whenever possible,
- reuse loaded data elements to the extent possible,
- avoid math library function calls inside an inner loop if possible.

Some of these goals might conflict with others, or are constrained by machine parameters. The following methods help us achieve the stated goals:

- Blocking
- Unrolling and scheduling

- Scalar replacement
- Precomputation of constants

We now discuss these methods in detail.

**Blocking.** Register-level blocking partitions the data into chunks on which the computation can be performed within the register set. Only initial loads and final stores but no register spilling is required. Sometimes a small amount of spilling can be tolerated. We show the blocking of a single loop as example. Consider the example code below.

```
for(i=0; i<8; i++)
{ y[2*i]   = x[2*i] + x[2*i+1];
  y[2*i+1] = x[2*i] - x[2*i+1];
}
```

We block the *i* loop, obtaining the following code.

```
for(i1=0; i1<4; i1++)
  for(i2=0; i2<2; i2++)
  { y[4*i1+2*i2] = x[4*i1+2*i2] + x[4*i1+2*i2+1];
    y[4*i1+2*i2+1] = x[4*i1+2*i2] - x[4*i1+2*i2+1];
  }
```

On many machines registers are only addressable by name but not indirectly via other registers (holding loop counters). In this case, once the data fits into registers, either loop unrolling or software pipelining with register rotation (as supported by Itanium) is required to actually take advantage of register-blocked computation.

**Unrolling and scheduling.** Unrolling produces larger basic blocks. That allows the compiler to apply strength reduction to simplify expressions. It decreases the number of conditional branches thus decreasing potential branch mispredictions and condition evaluations. Further it increases the number of operations in the basic block and allows the compiler to better utilize the register file. However, too much unrolling may increase the code size too much and overflow the instruction cache. The following code is the code above with unrolled inner loop *i2*.

```
for(i1=0; i1<4; i1++)
{ y[4*i1]   = x[4*i1] + x[4*i1+1];
  y[4*i1+1] = x[4*i1] - x[4*i1+1];
  y[4*i1+2] = x[4*i1+2] + x[4*i1+3];
  y[4*i1+3] = x[4*i1+2] - x[4*i1+3];
}
```

Unrolling exposes an opportunity to perform instruction scheduling. With unrolled code, it becomes easy to determine data dependencies between instructions. Issuing an instruction right after a preceding instruction that it is dependent upon will lead to the CPU pipeline being stalled until the former instruction completes. Instruction scheduling is the process of rearranging code to include independent instructions in between two dependent instructions to minimize pipeline stalls.

Scheduling large basic blocks with complicated dependencies may be too challenging for the compiler. In this case source scheduling may help. Source scheduling is the (legal) reordering of statements in the unrolled basic block. Different scheduling algorithms apply different rules, aiming at, e.g., minimizing distance between producer and consumer (which may potentially not be too short), and/or minimizing the number of live variables for each statement in the basic block. It is sometimes better to source schedule basic blocks and turn off aggressive scheduling by the compiler.

The number of registers, quality of the C compiler, and size of the instruction cache limit the amount of unrolling, that increases performance. Experiments show that on current machines, roughly 1,000 operations are the limit. Note, that unrolling always increases the size of the loop body, but not necessarily the instruction-level parallelism. Depending on the algorithm, more complicated loop transformations may be required. One example is the MMM, discussed later.

**Scalar replacement.** In C compilers, pointer analysis is complicated, and using even the simplest pointer constructs can prevent “obvious” optimizations. This observation extends to arrays with known size. It is very important to replace arrays that are fully inside the scope of an innermost loop by one automatic, scalar variable per array element. This can be done as the array access pattern does not depend on any loop variable and will help compiler optimization tremendously. As an example, consider the following code:

```
double t[2];
for(i=0; i<8; i++)
{ t[0] = x[2*i] + x[2*i+1];
  t[1] = x[2*i] - x[2*i+1];
  y[2*i]   = t[0] * D[2*i];
  y[2*i+1] = t[0] * D[2*i];
}
```

The variable `t` is now scalarized, resulting in code the compiler will be able to better optimize:

```
double t0, t1;
for(i=0; i<8; i++)
{ t0 = x[2*i] + x[2*i+1];
  t1 = x[2*i] - x[2*i+1];
  y[2*i]   = t0 * D[2*i];
  y[2*i+1] = t1 * D[2*i];
}
```

The difference is that `t0` and `t1` are automatic variables and can be held in registers whereas the array `t` will most likely be allocated in memory, and loaded and stored from memory for each operation.

**Precomputation of constants.** In a CPU-bound kernel, all constants that are known ahead of time should be precompute at compile time or initialization time and stored in a data array. At execution time, the kernel simply loads the precomputed data instead of needing to invoke math library functions. Consider the following example.

```

for(i=0; i<8; i++)
    y[i] = x[i] * sin(M_PI * i / 8);

```

The program contains an function call to the math library in the inner loop. Calling `sin()` can cost multiple thousands of cycles on modern CPUs. However, all the constants are known before entering the kernel and thus can be precomputed.

```

static double D[8];
void init()
{ for(int i=0; i<8; i++)
    D[i] = sin(M_PI * i / 8);
}

...
// in the kernel
for(i=0; i<8; i++)
    y[i] = x[i] * D[i];

```

The initialization function needs to be called only once. If the kernel is used over and over again, precomputation results in enormous savings. If the kernel is used only once, chances are that performance does not matter.

#### 4.4 Parameter-Based Performance Tuning

Many of the optimizations for the memory hierarchy discussed above have inherent degrees of freedom such as the block size for blocking or the degree of unrolling the code. While it may be possible to derive a reasonable estimation of these parameters, the complexity of modern microarchitecture makes an exact prediction impossible. In fact, often the best value may be a surprise to the programmer. As a consequence, it makes sense to perform an empirical search to find those parameters. These means creating the variants, ideally through a script or parameterized coding (for instance, defining all parameters as C preprocessor constants in a separate header file), and measuring their performance to find the best. The result may depend on the target platform, so the search should be repeated for every new platform.

This parameter-based performance optimization is one of the techniques used in recent research on automatic performance tuning. An overview of some prominent efforts can be found in [54].

## 5 MMM

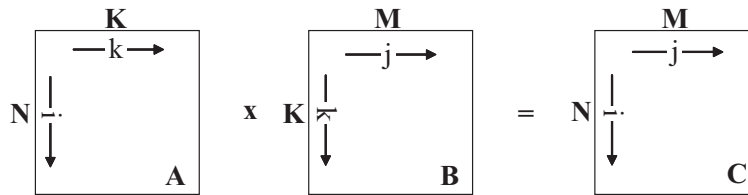
In this section, we optimize matrix-matrix multiplication (MMM) for the memory hierarchy. We explain the optimizations implemented by the ATLAS [70], and organize the steps as in Section 4. ATLAS is a program generator for MMM and other BLAS routines and also performs other optimizations not discussed here. Our presentation closely follows the one in Yotov et al. [73], which presents a model-based version of ATLAS.

For the rest of this section, we will assume the dimensions of the input matrices  $A$  and  $B$  to be  $N \times K$  and  $K \times M$  respectively, which implies an  $N \times M$  output matrix  $C$ . For simplicity, we will further assume that various optimization parameters are perfectly divisible by these dimensions whenever such a division is necessary. The computation considered is  $C = C + AB$ .

**Naive Implementation.** Matrix-matrix multiplication (MMM), as defined in Section 2.2, is naively implemented using the triple loop shown below. We use 2D array notation (for instance,  $C[i][j]$ ) to keep the code more readable. However, in an implementation where the matrix sizes are not known at compile time, one should resort to a linearized representation of  $C$ ,  $A$ , and  $B$  (see Section 4.1).

```
// K, M, N are compile-time constants
double C[N][M], A[N][K], B[K][M];
// Assume C is initialized to zero
for(i=0; i<N; i++)
  for(j=0; j<M; j++)
    { for(k=0; k<K; k++)
      C[i][j] += A[i][k] * B[k][j];
    }
}
```

The C language stores two-dimensional arrays in row-major order. Therefore, a cache miss to a (memory aligned) matrix element causes that element and adjacent elements in the same row being loaded into one cache line of the cache (see Fig. 7). Thus, accessing a large matrix by rows is cache efficient, while accessing it by columns is not.

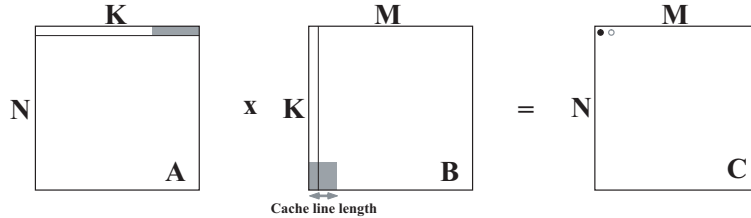


**Fig. 9.** Data access pattern for the naive MMM

Fig. 9 illustrates the data access pattern of the naive implementation. From this figure, we see the output locality of the computation: all accesses to each element in  $C$  are consecutive, and  $C$  is completed element by element, row by row. However, unless all input and output arrays fit into the cache, the naive implementation has poor locality with respect to  $A$  and  $B$ .

We analyze the naive implementation by counting the number of cache misses. We assume a cache line size of 64 bytes, or 8 (double precision) floating point values, and that  $N$  is large with respect to the cache size. To compute the first entry in  $C$ , we need to access the entire first row of  $A$  and the entire first column of  $B$ . Accessing a row of  $A$  results in  $N/8$  misses (one for each group of 8) due to the row-major storage order,

while accessing a column of  $B$  results in a full  $N$  misses, yielding a total of  $(9/8)N$  misses for the first entry in  $C$ .



**Fig. 10.** The state of the cache at the end of computation of the first element of  $C$  (marked by a solid circle) is shown. Areas of the input matrices marked in gray are cache resident at this point. The next element of  $C$  to be computed is shown as a hollow circle.

To analyze the computation of the second entry of  $C$ , we first observe that the parts of  $A$  and  $B$  that will be accessed first are not in the cache. Namely, since  $N$  is much larger than the cache, the first few elements of the first row of  $A$  were in cache but meanwhile have been overwritten. Similarly, the first elements of the second column of  $B$  were already in cache (each element shared a cache line with its neighbor in the first column) but also have been overwritten. This is illustrated in Fig. 10, which shows in gray the parts of  $A$  and  $B$  that are in cache after the first entry of  $C$  is computed. Consequently, the number of misses involved in computing the second entry (and every subsequent entry of  $C$ ), produces also  $(9/8)N$  misses. Therefore, the total number of misses generated by this algorithm (for the  $N^2$  entries in  $C$ ) is  $(9/8)N^3$ . In summary, there is no reuse and no neighbor use, a problem resolved to the extent possible by the optimizations in the next sections.

## 5.1 Cache Optimization

**Blocking.** One of the most important MMM optimizations is blocking, as introduced in Section 4.2. Blocking involves performing the addition and multiplication operations on *blocks* of the original matrix, instead of individual elements. The idea is to increase locality by restricting the computation at any point to work on small chunks that fit entirely into the cache. We will also see that blocking essentially increases reuse and neighbor use, the concepts previously presented in Section 2.4.

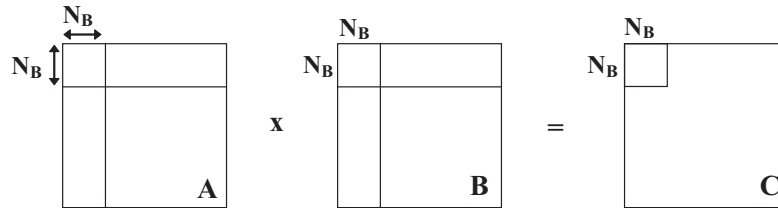
The compiler loop transformation that implements blocking is known as *tiling* [72, 70, 73]. Blocking or tiling the MMM for each level of the memory hierarchy involves adding three more nested loops to the basic triple-loop implementation. The code for the MMM blocked for one memory level with block size  $N_B$  follows.



```

// MMM loop nest (j, i, k)
for(i=0; i<N; i+=NB)
  for(j=0; j<M; j+=NB)
    for(k=0; k<K; k+=NB)
      // mini-MMM loop nest (i0, j0, k0)
      for(i0=i; i0<(i + NB); i0++)
        for(j0=j; j0<(j + NB); j0++)
          for(k0=k; k0<(k + NB); k0++)
            C[i0][j0] += A[i0][k0] + B[k0][j0];

```



**Fig. 11.** Blocking for the cache: mini-MMMs

Fig. 11 shows the data access pattern of blocking for the cache. The three additional innermost loops cause each matrix to be divided into blocks of size  $N_B \times N_B$ . Notice the similarity in the access pattern to the naive implementation, except at the block level instead of at the element level.

We now analyze this version of the MMM to determine the impact on the number of cache misses. We assume that the block size is larger than the cache line size, and for now that several blocks can fit into the cache. This implies that accessing a block results only in  $N_B^2/8$  misses, regardless of the access sequence.

Computing the first *block* of  $C$  results in accessing all the topmost blocks of  $A$ , and all the leftmost blocks of  $B$ . This results in  $(N_B^2/8 + N_B^2/8)(N/N_B)$  cache misses. Similar to the reasoning used in the analysis of the naive version, computing each block of  $C$  results in the same amount of misses, and therefore, the total number of misses generated by this algorithm (for the  $(N/N_B)^2$  blocks in  $C$ ) is  $N^3/4N_B$ , which is significantly less than the  $(9/8)N^3$  misses in the naive version.

The smaller blocks are called the mini-MMMs, following [74].  $N_B$  is an optimization parameter that must be chosen such that the working set of the mini-MMM fits entirely into the cache. A simple translation of our assumption that blocks from the two input and output matrices (our *working set*) fit into a fully associative cache is expressed by the following equation:  $3N_B^2 \leq C_s$ , where  $C_s$  is the cache size. ATLAS determines  $N_B$  by searching and trying different arbitrary values and picking the one that results in the best performance.

Yotov et al. [73] use a model based approach, and choose  $N_B$  based directly on the cache size. A careful examination of the data access pattern of the blocked MMM reveals that

the working set at a finer granularity consists only of a single element in  $C$  (since each element in  $C$  is reused completely by the innermost  $k_0$  loop before it moves on to the next element), a single row in the  $A$  matrix (since a row is fully reused before the program moves on to the next row), and the entire block in  $B$ . Therefore, Yotov et al. note that the complete mini blocks from all three matrices need not fit into the cache, but only that the following relationship needs to hold true:  $N_B^2 + N_B + 1 \leq C_s$ . The authors thus pick the largest value of  $N_B$  that will satisfy this inequality.

Blocking for MMM works because it increases cache reuse and neighbor use, our guiding principles discussed in Section 2. Cache reuse is increased because once a block is brought into the cache, it is used several times before being overwritten. Neighbor use is increased for the input matrix  $B$ , since in the naive implementation, it is accessed in the column major order, which uses only one element in each cache block once before evicting the block. In the blocked implementation, all elements in the cache block are used before eviction.

Blocking can be done for any or all of the cache levels. Since the latency between L2 and L3 or main memory is typically much higher than the latency between L1 and L2, it is necessary to block at least for the cache closest to memory. It pays to additionally block for L1 cache only if it compensates for the extra loop overhead.

An additional optimization that can be done for the cache is to exchange the  $i$  and the  $j$  loops, depending upon the relative sizes of the  $A$  and  $B$  matrices.

**Loop merging.** Loop merging is not applicable to the MMM.

**Buffering.** Buffering (also known as copying) for MMM is applicable for large sizes. The basic idea behind buffering is to copy tiles of the input and output matrices into sequential order in memory to minimize cache conflict misses (and TLB misses if the matrices span multiple pages), inside each mini-MMM. The following listing illustrates buffering. The matrix  $B$  is fully buffered at the beginning since it is accessed in full during each iteration of the outermost  $i$  loop. Vertical panels of  $A$  are used during each iteration of  $j$ , and are buffered just before the  $j$  loop begins. Finally, in some cases, it might be beneficial to copy a single tile of  $C$  before the  $k$  loop, since a single tile is reused by each iteration of the  $k$  loop. Note that the benefits of buffering have to outweigh the costs, which might not hold true for very small or very large matrices.

```

for(i=0; i<M; i+=NB)
    // Buffer a panel of A here
    for(j=0; j<N; j+=NB)
        // Copy a block (tile) of C here
        for(k=0; k<K; k+=NB)
            // mini-MMM loop nest as before (i0, j0, k0)
            ...
            ...

```

## 5.2 CPU and Register Level Optimization

We now look at optimizing the MMM for the CPU. We continue with our MMM example from the previous section.

**Blocking.** Blocking for the registers looks similar to blocking for the cache. Another set of nested triple loops is added. The resulting code is shown below:

```
// MMM loop nest (j, i, k)
for(i=0; i<N; i+=NB)
  for(j=0; j<M; j+=NB)
    for(k=0; k<K; k+=NB)
      // mini-MMM loop nest (i0, j0, k0)
      for(i0=i; i0<(i + NB); i0+=MU)
        for(j0=j; j0<(j + NB); j0+=NU)
          for(k0=k; k0<(k + NB); k0+=KU)
            // micro-MMM loop nest (j00, i00)
            for(k00=k0; k00<=(k0 + KU); k00++)
              for(j00=j0; j00<=(j0 + NU); j00++)
                for(i00=i0; i00<=(i0 + MU); i00++)
                  C[i00][j00]=C[i00][j00]+A[i00][k00]*B[k00][j00];
```

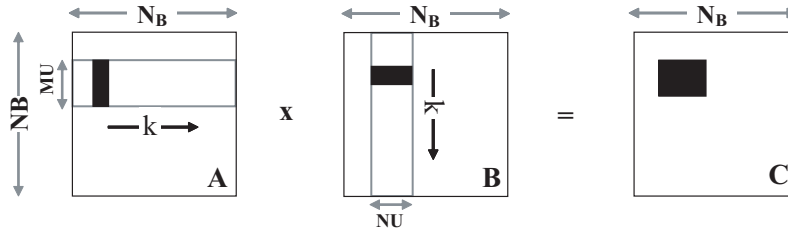


Fig. 12. mini-MMMs and micro-MMMs (from [73])

As Fig. 12 shows, each mini-MMM is now computed by blocking it into a sequence of micro-MMMs. Each micro-MMM multiplies an  $M_U \times 1$  block of  $A$  by a  $1 \times N_U$  block of  $B$ , with the output being a  $M_U \times N_U$  block of  $C$ . At this level of blocking, we have a degree of freedom in choosing  $M_U$  and  $N_U$ . These parameters must be chosen so that a micro-MMM fits into register space (thus avoiding register spills).

ATLAS searches over arbitrary values for these parameters to choose the ones that result in the fastest code. Yotov et al., with a reasoning that is similar to the one used in choosing  $N_B$  in the previous section, select these parameters based on the inequality  $M_U + N_U + (M_U \times N_U) \leq N_R$ , where  $N_R$  is the number of data (integer or floating point) registers. They further refine the reasoning and this inequality in [73].

Locality is not the only objective of blocking for register space. Note that in the code above, the micro-MMM have a loop order of  $kij$ . While this reduces output locality, it

also provides better instruction level parallelism (all the  $M_U N_U$  addition/multiplication pairs are independent) when combined with loop unrolling discussed next.

**Unrolling and scheduling.** Loop unrolling and scheduling, as discussed in Section 4.3 on page 35, can be used to further optimize MMM. We unroll the two innermost loops to get  $M_U \times N_U$  additions and multiplications. Note that these instructions are of the form  $C+ = AB$ . As mentioned in [71], such an instruction will not execute well on machines without a fused multiply-add unit, since the addition is dependent on the multiplication, and will cause a pipeline stall until the multiplication is completed. Thus, it may be beneficial to separate the addition and the multiplication operations here, and schedule them with unrelated intervening instructions to minimize pipeline stalls.

**Scalar replacement.** When the innermost loops are unrolled, each array element appears multiple times in the unrolled code. For the reasons discussed earlier in Section 4.3, replacing array references by scalar variables in unrolled code enables compiler optimizations to work better.

**Precomputation of constants.** Since the MMM does not have constants that can be precomputed, this optimization does not apply.

### 5.3 Parameter-Based Performance Tuning

The above discussion identifies several parameters that can be used for tuning. ATLAS performs this tuning automatically by generating the variants and selecting the fastest using a search procedure.

**Blocking for cache.**  $N_B$  is the main optimization parameter used to control the block size of the mini-MMMs. If several levels of blocking are desired, additional blocking parameters arise.

**Blocking for registers.** When blocking for the registers,  $M_U$  and  $N_U$  are the main tunable parameters, and must be chosen such that the micro-MMM does produce register spills.

### 5.4 Exercises

1. **Mini-MMM.** The goal of this exercise is to implement a fast mini-MMM to multiply two square  $N_B \times N_B$  matrices ( $N_B$  is a parameter), which is used within an MMM.
  - (a) (Based on definition - *program0*): Use your naive implementation of the MMM as completed in Exercise 2 on page 23.
  - (b) (Blocking - *program1*): Block into micro MMMs with  $M_U = N_U = 2, K_U = 1$ . The inner triple loop should have the  $kij$  order. You can use the code given in Section 5.2 for your implementation. Unroll (manually) the innermost  $i$  and  $j$ -loop such that multiplications and additions alternate. Perform scalar replacement (manually).

- (c) (Unrolling - *program2*, *program3*) Unroll the innermost  $k$ -loop by a factor of 2 (*program2*), and by a factor of 4 (*program3*), ( $K_U = 2$  and 4, this doubles and quadruples the loop body) again performing scalar replacement. Assume that 4 divides  $N_B$ .
- (d) (Performance plot, used to search for best block size  $N_B$ ) Determine the L2 cache size  $C$  (in doubles, i.e., 8 byte units) of your computer. Measure the performance (in Mflop/s) of your four codes for all  $N_B$  with  $16 \leq N_B \leq \min(80, \sqrt{C})$  for values of  $N_B$  divisible by 4. Create a plot with  $N_B$  on the  $x$ -axis, and the performance on the  $y$ -axis. The plot will contain 4 lines - one line for each of the programs. Discuss the plot, including answers to the following questions: which  $N_B$  and which code yields the maximum performance? What is the percentage of peak performance in this case?
- (e) Does it improve your best program so far you switch the outermost loop order from  $ijk$  to  $jik$ ?

## 2. MMM.

- (a) Implement an MMM for multiplying two square  $N \times N$  matrices assuming  $N_B$  divides  $N$ , blocked into  $N_B \times N_B$  blocks. Use your best mini-MMM code from Exercise 1.
- (b) Create a performance plot comparing this implementation and *program0* (based on definition) above for an interesting range of  $N$  (up to sizes where the matrices do not fit into the L2 cache). Plot the size ( $N$ ) on the  $x$ -axis, against the performance (in Mflop/s) on the  $y$ -axis.
- (c) Analyze and discuss the plot.

## 6 DFT

In this section we describe the design and implementation of a high-performance function to compute the FFT. The approach we have to take is different from the one taken to optimize the MMM in Section 5: we do not start with a naive implementation that is transformed into an optimized form, but design the code from scratch. This is due to the more complex structure of the available FFT algorithms. Note that, in contrast to MMM, an implementation based on the definition of the DFT in Section 2.3 is not competitive.

The first main problem is the choice of a suitable FFT algorithm, since many different variants are available that differ vastly in structure. It makes no sense to start with the wrong FFT algorithm and optimize the implementation step by step. In particular, when targeting a machine with a memory hierarchy, choosing an implementation based on the iterative radix-2 FFT used in Numerical Recipes (Section 2.3) is suboptimal.

In our discussion below we design a recursive radix-4 FFT implementation. Generalization to a mixed-radix recursive implementation conceptually relatively straight-forward

in concept, but technically complex. The optimization steps taken follow to a large extent, the design of FFTW 2.x [32]. FFTW uses a program generator in addition, to automatically implement optimized unrolled basic blocks [31].

In all our DFT code examples the (complex) data is assumed to be stored in interleaved complex double-precision arrays. We pass around pointers of type `double`, and two neighboring `double` elements are one complex number. All strides are relative to complex numbers.

## 6.1 Cache Optimization

In this section we derive the recursive skeleton and the kernel specification for our DFT implementation.

**Blocking.** Blocking a DFT algorithm is done by choosing the recursive Cooley-Tukey FFT algorithm (1) as starting point instead of the iterative FFT used by the Numerical Recipes code in Section 2.3. The block size is the chosen *radix*  $m$  in (1), which is a degree of freedom. We assume a radix-4 implementation with  $mn = 4^k$ , i.e., we pick  $m = 4$  and  $n = 4^{k-1}$ . The corresponding recursion is

$$\text{DFT}_{4^k} = (\text{DFT}_4 \otimes I_{4^{k-1}}) D_{4,4^{k-1}} (I_4 \otimes \text{DFT}_{4^{k-1}}) L_4^{4^k}. \quad (7)$$

For  $k > 1$  our implementation will recursively apply (7) to the terms  $\text{DFT}_{4^{k-1}}$  in the right side of (7). The terms  $\text{DFT}_4$  are recursion leaves and not implemented using (7). We will discuss their implementation in Section 6.2.

This recursion is *right-expanded*—the first stage gets recursively expanded while the second stage uses radix-4 kernels. Right-expanded recursive implementations have superior data locality as only a small amount of temporary storage is needed and the second stage can be implemented in-place.

**Loop merging.** A naive implementation of (7) leads to a recursive function with 4 stages (corresponding to the four matrix factors) and thus 4 sweeps through the data. However, the stride permutation  $L_4^{4^k}$  is just a data reordering and thus is a candidate for loop merging. Similarly, the twiddle factor matrix  $D_{4,4^{k-1}}$  is a diagonal matrix and can be merged with the subsequent stage.

We now sketch the derivation of a recursive implementation of (7). We partition (7) into two expressions as

$$\text{DFT}_{4^k} = ((\text{DFT}_4 \otimes I_{4^{k-1}}) D_{4,4^{k-1}}) \cdot ((I_4 \otimes \text{DFT}_{4^{k-1}}) L_4^{4^k}). \quad (8)$$

which become two stages in the recursive function

```
void DFT_rec(int N, int n, double *Y, double *X, int s);
```

that implements (7). `DFT_rec` computes  $y = \text{DFT}_N x$  for  $N = 4^n$  out-of-place with input data starting at `*X` and output data at `*Y`. The output is accessed at unit stride

while the input is accessed at stride  $s$ . It may sound counterintuitive to have different strides for input and output, but we will soon see that this choice allows us to obtain a high performance implementation. At the top level, the function needs to be called with  $s=1$ , as facilitated by the following C macro (`log4( )`, not shown here, is a function to compute  $k$ ).

```
#define DFT(N, Y, X) DFT_rec(N, log4(N), Y, X, 1)
```

For  $N = 4$ , we reach the leaf of the recursion and call a base case kernel function.

```
void DFT4_base(double *Y, double *X, int s);
```

`DFT4_base` computes  $y = \text{DFT}_4 x$  out-of-place with input data starting at `*X` and output data at `*Y`. The input is accessed at stride  $s$ . The output is contiguous and thus accessed at unit stride.

The first stage of (8),  $y = (I_4 \otimes \text{DFT}_{4^{k-1}})L_4^{4^k} x$ , is handled as follows. According to Table 1 the tensor product is translated into a loop with 4 iterations, with the iterations implementing  $y = \text{DFT}_{4^{k-1}} x$ . The tensor product  $I_4 \otimes \text{DFT}_{4^{k-1}}$  partitions the input and output into contiguous chunks. Thus,  $\text{DFT}_{4^{k-1}}$  operates in each iteration on one of these contiguous arrays of size  $4^{k-1}$ . The input to  $\text{DFT}_{4^{k-1}}$  is readdressed by  $y = L_4^{4^k} x$ . This means that  $\text{DFT}_{4^{k-1}}$  inside the tensor product actually reads from  $x_j, x_{j+4}, x_{j+8}, \dots, x_{j+4^{k-1}}$ , where  $j$  is the loop iteration number. Thus, the function implementing  $\text{DFT}_{4^{k-1}}$  needs to be called with input stride 4 but unit output stride. Further, the first stage needs to be implemented out-of-place, i.e.,  $x$  and  $y$  need to be different memory regions.

The second stage,  $y = (\text{DFT}_4 \otimes I_{4^{k-1}})D_{4,4^{k-1}}x$ , first scales the input by a diagonal matrix and then sweeps with a radix-4 DFT kernel over it. According to Table 1 the tensor product is translated into a loop with  $4^{k-1}$  iterations, where each iteration computes  $y = \text{DFT}_4 x$ . The tensor product  $I_{4^{k-1}} \otimes \text{DFT}_4$  partitions the input and output into  $4^{k-1}$  interleaved chunks of size 4. Thus,  $\text{DFT}_4$  operates on  $x_j, x_{j+4^{k-1}}, x_{j+2 \cdot 4^{k-1}}$ , and  $x_{j+3 \cdot 4^{k-1}}$ , where  $j$  is the loop iteration number. The diagonal  $D_{4,4^{k-1}}$  is fused by simply premultiplying each input to  $\text{DFT}_4$  with the correct complex number right after loading from memory but before using the code of  $\text{DFT}_4$ . This in effect implements  $y = (\text{DFT}_4 D_j)x$  for a suitable, loop-dependent  $D_j$ . If no output of this stage is written before all inputs are used, this stage can be implemented in-place, i.e., both  $x$  and  $y$  can be in the same memory region.

Hence we need to implement the following kernel function for the second stage

```
void DFT4_twiddle(double *Y, int s, int n, int j);
```

as base case. `DFT4_twiddle` computes  $y = (\text{DFT}_4 D_j)x$  in-place with input and output data starting at `*Y`, accessed at stride  $s$ .  $D_j$  is a diagonal matrix derived from the twiddle diagonal  $D_{4^k,4}$ ,  $s$ ,  $n$ , and  $j$ . Details on  $D_j$  are beyond the scope of this tutorial.

The final recursive function is given below. There are some address multiplications by 2, required to implement arrays of complex numbers as arrays (of twice the size) of real numbers.

```
// recursive radix-4 DFT implementation

// top-level call to DFT function
int log4(int);
#define DFT(N, Y, X) DFT_rec(N, log4(N), Y, X, 1)

// DFT kernels
void DFT4_base(double *Y, double *X, int s);
void DFT4_twiddle(double *Y, int s, int N, int j);

// recursive radix-4 DFT function
// N: problem size
// Y: output vector
// X: input vector
// s: stride to access X
void DFT_rec(int N, int n, double *Y, double *X, int s)
{ int k;

  if (N==4)
    // Y = DFT_4 X
    DFT4_base(Y, X, s);
  else {
    // Y = (I_4 x DFT_{N/4})(L^{N_4}) X
    for(j=0; j<4; j++)
      DFT_rec(N/4, n-1, Y+8*j, X+2*j*s, s*4);
  }
  // Y = (DFT_4 x I_{N/4})(D_{N,4}) Y
  for(j=0; j<N/4; j++)
    DFT4_twiddle(Y+2*j*s, N/4, n, j);
}
```

**Buffering.** The kernel `DFT4_twiddle` accesses both input and output in a stride. If (for larger DFT sizes) these strides become large 2-powers and the kernel size  $n$  for  $\text{DFT}_n$  is larger than the cache associativity, cache thrashing occurs. In this situation, each iteration of the  $k$ -loop has to load  $n$  cache lines and all these cache lines get evicted before the next iteration of the  $k$ -loop could use the already loaded remainder of the cache lines.

Typical machines have a cache associativity of 4-way or 8-way. This means one should use small DFT kernels for large DFT sizes which introduce large strides and do not fit into cache. Small DFT kernels means more passes through the data set, which in turn hurts performance by requiring higher memory bandwidth. The combination of these issues is one reason why the performance of the DFT code in Fig. 2 drops dramatically when  $n$  exceeds  $2^{16}$ .



Buffering alleviates these problems to a certain degree. An initial and final copy operation costs overhead, but all intermediate steps are done on contiguous data, preventing cache thrashing.

As an example, buffering is performed on the second loop of the preceding code, leading to the following code. We assume a cache line size of 2 complex numbers (= 4 doubles). To implement buffering, we first split the  $j$ -loop into  $N/8 \times 2$  iterations. We add copying to the body of the *outer* tiled ( $j_1$ ) loop. Our copy operation handles cache lines and thus data for multiple DFTs. In particular, we copy 4 sets of 4 consecutive elements (4 cache lines) into a local buffer. The inner tiled ( $j_2$ ) loop performs 2 DFTs on the local contiguous buffer. The large, performance degrading stride (4) in the original  $j$ -loop gets replaced by a small stride (2) in the  $j_2$ -loop at the cost of two copy operations that copy whole cache lines. The threshold parameter  $th$  controls for which sizes the second loop gets buffered.

```
// recursive radix-4 DFT function with buffering
// N: problem size
// Y: output vector
// X: input vector
// s: stride to access X
// th: threshold size to stop buffering
void DFT_buf(int N, int n, double *Y, double *X, int s, int th)
{ int i, j1, j2, k;
  // local buffer
  double buf[16];

  if (N==4)
    // Y = DFT_4 X
    DFT4_base(Y, X, s);
  else {
    // Y = (I_4 x DFT_{N/4})(L^{N_4}) X
    if (N > th)
      for(k=0; k<4; k++)
        DFT_buf(N/4, Y+8*j, X+2*j*s, s*4);
    else
      for(j=0; j<4; j++)
        DFT_rec(N/4, Y+8*k, X+2*j*s, s*4);
  }
  // Y = (DFT_4 x I_{N/4})(D_{N,4}) Y, buffered
  // tiled j loop
  for(j1=0; j1<N/8; j1++)
  { // copy 4 chunks of 4 to local buffer
    for(i=0; i<4; i++)
      for(k=0; k<4; j++)
        buf[4*i+k] = Y[4*j1+8*i+k];

    // perform 2 DFT_4 on contiguous data
    for(j2=0; j2<2; j2++)
      DFT4_twiddle(buf, N/4, n-1, j);
  }
}
```

```

// copy 4 chunks of 4 to output
for(i=0; i<4; i++)
    for(k=0; j<4; j++)
        Y[4*j1+8*i+k] = buf[4*i+k];
    }
}

```

## 6.2 CPU and Register Level Optimization

This section describes the design and implementation of optimized DFT base cases (kernels). We again restrict the discussion to the recursive radix-4 FFT algorithm. Extensions to mixed-radix implementations requires different kernel sizes, all implemented following the ideas presented in this section. High-performance implementations may use kernels of up to size 64 [31, 59].

**Blocking.** We apply (1) to the  $DFT_4$ :

$$DFT_4 = (DFT_2 \otimes I_2) D_{4,2} (I_2 \otimes DFT_2) L_2^4. \quad (9)$$

As (9) is a recursive formula, an implementation based on (9) is automatically blocked.

**Unrolling and scheduling.** We implement (9) according to the rules summarized in Table 1. We aim at implementing recursion leafs. Thus the code needs to be unrolled. Due to the recursive nature of (1), kernels derived from (1) are automatically scheduled.

For DFT kernels, larger unrolled kernels lead to slightly less operations, as more twiddle factors are known at optimization time and one can take better advantage of trivial complex multiplications. However, larger kernels do not increase the available instruction level parallelism as much as in MMM, since the DFT data flow is more complicated and imposes stronger constraints on the operation ordering.

**Scalar replacement.** We next apply scalar replacement as described in Section 4.3. This leads to the following code for `DFT4_base`. From the discussion in Section 6.1 we know that this function loads at complex stride `s` from `*x` and writes at unit stride to `*Y`. We obtain the following code.

```

// DFT4 implementation
void DFT4(double *Y, double *X, int s)
{
    double t0, t1, t2, t3, t4, t5, t6, t7;
    t0 = (X[0] + X[4*s]);
    t1 = (X[2*s] + X[6*s]);
    t2 = (X[1] + X[4*s+1]);
    t3 = (X[2*s+1] + X[6*s+1]);
    t4 = (X[0] - X[4*s]);
    t5 = (X[2*s+1] - X[6*s+1]);
    t6 = (X[1] - X[4*s+1]);
    t7 = (X[2*s] - X[6*s]);
    Y[0] = (t0 + t1);
    Y[1] = (t2 + t3);
}

```

```

Y[4] = (t0 - t1);
Y[5] = (t2 - t3);
Y[2] = (t4 - t5);
Y[3] = (t6 + t7);
Y[6] = (t4 + t5);
Y[7] = (t6 - t7);
}

```

**Precomputation of constants.** The kernel `DFT4_twiddle` computes  $y = (\text{DFT}_4 D_j)x$ , which contains multiplication with the complex diagonal  $D_j$ . The entries of  $D_j$  are complex roots of unity (twiddle factors) that depend on the recursion level and the loop counter  $j$ . Computing the actual entries of  $D_j$  needs evaluations of  $\sin \frac{k\pi}{N}$  and  $\cos \frac{k\pi}{N}$  for suitable values of  $k$  and  $N$ , which requires expensive calls to the math library. Hence these numbers should be precomputed.

We introduce an initialization function `init_DFT` that precomputes all diagonals required for size  $N$  and stores pointers to the tables (one table for each recursion level) in the global variable `double **DN`. We do not show the function as the actual computation of the constants is beyond the scope of the tutorial. We show the C code for the function `DFT4_twiddle` below.

```

// twiddle table, initialized by init_DFT(N)
void init_DFT(int N);
double **DN;

// C macro for complex multiplication
#define CMPLX_MULT(cr, ci, a, b, idx) \
{ double ar, ai, br, bi; \
  ar = a[2*idx]; ai = a[2*idx+1]; \
  br = b[2*idx]; bi = b[2*idx+1]; \
  cr = ar*br - ai*bi; \
  ci = ar*bi + ai*br; \
}

// DFT4*D_j implementation
void DFT4_twiddle(double *Y, int s, int n, int j)
{ double t0, t1, t2, t3, t4, t5, t6, t7;
  X0, X1, X2, X3, X4, X5, X6, X7;
  double *Dj;

  // complex multiplications from D_N
  Dj = DN[n]+8*j;
  CMPLX_MULT(X0, X1, Y, Dj, 0);
  CMPLX_MULT(X2, X3, Y, Dj, 1);
  CMPLX_MULT(X4, X5, Y, Dj, 2);
  CMPLX_MULT(X6, X7, Y, Dj, 3);

  // operations from DFT4
  t0 = (X0 + X4);
  t1 = (X2 + X6);

```

```

    t2 = (X1 + X5);
    t3 = (X3 + X7);
    t4 = (X0 - X4);
    t5 = (X3 - X7);
    t6 = (X1 - X5);
    t7 = (X2 - X6);
    Y[0] = (t0 + t1);
    Y[1] = (t2 + t3);
    Y[4*s] = (t0 - t1);
    Y[4*s+1] = (t2 - t3);
    Y[2*s] = (t4 - t5);
    Y[2*s+1] = (t6 + t7);
    Y[6*s] = (t4 + t5);
    Y[6*s+1] = (t6 - t7);
}

```

### 6.3 Parameter-Based Performance Tuning

We now discuss the parameters in our DFT implementation that can be tuned to the memory hierarchy.

**Base case sizes.** The most important parameter tuning is the selection of base cases. To allow for multiple base cases `DFT_base` and `DFT_twiddle`, the program structure must become more complex, as a data structure describing the recursion and containing function pointers to the appropriate kernels replaces the two parameters `N` and `n` in `DFT_rec`. The resulting program would be very similar to FFTW 2.x.

After this infrastructural change the system can apply any function `DFT_twiddle` in the second stage of the recursion and any function `DFT_base` as recursion leaf. The tuning process needs to find for each recursion step the right kernel size. FFTW uses both a cost estimation and runtime experiments based on dynamic programming to find good parameter choices [33]. Showing the full implementation is beyond the scope of this tutorial.

**Threshold for buffering.** The second parameter decides the sizes for which buffering should be applied. This depends on the cache size of the target machine, as buffering only becomes beneficial for problem sizes that are not resident in the L2 cache.

**Buffer size.** Finally, we need to set the buffer size based on the cache line size of the target machine to prevent cache thrashing. The cache line size can be either looked up or found experimentally.

### 6.4 Exercises

1. In Exercise 2 in Section 2, you completed a recursive implementation of the WHT. In this exercise, implement the triple loop, iterative version of the WHT using (5).

Create a performance plot (size versus Mflop/s) for sizes  $2^1$ – $2^{20}$  for both versions. Discuss the plot.

2. Create unrolled WHTs of sizes 4 and 8 based on the recursive WHT algorithm. (The number of operations should match the cost computed in Exercise 1c in Section 2).
3. Now implement recursive radix-4 and radix-8 implementations of the WHT based on the formulas

$$\text{WHT}_{2^k} = (\text{WHT}_4 \otimes I_{2^{k-2}})(I_4 \otimes \text{WHT}_{2^{k-2}}) \quad (\text{radix-4})$$

$$\text{WHT}_{2^k} = (\text{WHT}_8 \otimes I_{2^{k-3}})(I_8 \otimes \text{WHT}_{2^{k-3}}) \quad (\text{radix-8})$$

In these implementations, the left hand side WHT (of size 4 or 8) should be your unrolled kernel (which then has to handle input data at a stride) called in a loop; the right hand side is a recursive call (also called in a loop). Further, in both implementations, you may need one step with a different radix to handle all input sizes.

Measure the performance of both implementations, again for sizes  $2^1$ – $2^{20}$  and add it to the previous plot (4 lines total).

4. Try to further improve the code or perform other interesting experiments. For example, what happens if one considers more general algorithms based on

$$\text{WHT}_{2^k} = (\text{WHT}_{2^i} \otimes I_{2^{k-i}})(I_{2^i} \otimes \text{WHT}_{2^{k-i}})$$

The unrolled code could be the WHT on the left hand side of the above equation. Alternatively, one could run a search to find the best radix in each step independently.

## 7 Acknowledgement

This work was supported by DARPA through the Department of Interior grant NBCH1050009 and by NSF through awards 0234293, 0325687, and 0702386.

## References

1. ACML web site. <http://developer.amd.com/acml.jsp>.
2. Advanced Micro Devices (AMD) Inc. *Software Optimization Guide for AMD Athlon 64 and AMD Optero Processors*, 2005. <http://developer.amd.com/devguides.jsp>.
3. E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 3rd edition, 1999.
4. ATLAS web site. [math-atlas.sourceforge.net](http://math-atlas.sourceforge.net).

5. G. Baumgartner, D. Bernholdt, D. Cociovora, R. Harrison, M. Nooijen, J. Ramanujan, and P. Sadayappan. A performance optimization framework for compilation of tensor contraction expressions into parallel programs. In *Proc. Int'l Workshop on High-Level Parallel Programming Models and Supportive Environments (held in conjunction with IEEE Int'l Parallel and Distributed Processing Symposium (IPDPS))*, 2002.
6. BeBOP web site. <http://bebop.cs.berkeley.edu/>.
7. Eran Bida and Sivan Toledo. An automatically-tuned sorting library. *Software: Practice and Experience*, 2007.
8. Paolo Bientinesi, John A. Gunnels, Margaret E. Myers, Enrique Quintana-Orti, and Robert van de Geijn. The science of deriving dense linear algebra algorithms. *TOMS*, 31(1):1–26, March 2005.
9. L. S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users' Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, 1997.
10. L. S. Blackford, J. Demmel, J. Dongarra, I. Duff, S. Hammarling, G. Henry, M. Heroux, L. Kaufman, A. Lumsdaine, A. Petitet, R. Pozo, K. Remington, and R. C. Whaley. An updated set of Basic Linear Algebra Subprograms (BLAS). *ACM Transactions on Mathematical Software*, 2001.
11. Andreas Bonelli, Franz Franchetti, Jürgen Lorenz, Markus Püschel, and Christoph W. Ueberhuber. Automatic performance optimization of the discrete Fourier transform on distributed memory computers. In *Proc. International Symposium on Parallel and Distributed Processing and Applications (ISPA)*, 2006.
12. Randal E. Bryant and David R. O'Hallaron. *Computer Systems: A Programmer's Perspective*. Prentice Hall, 2003.
13. Almadena Chchelkanova, John Gunnels, Greg Morrow, James Overfelt, and Robert van de Geijn. Parallel implementation of blas: General techniques for level 3 blas. *Concurrency: Practice and Experience*, 9(9):837–857, 1997.
14. J. W. Cooley and J. W. Tukey. An algorithm for the machine calculation of complex Fourier series. *Math. of Computation*, 19:297–301, 1965.
15. Don Coppersmith and Shmuel Winograd. Matrix multiplication via arithmetic progressions. *Journal of Symbolic Computation*, 9:251–280, 1990.
16. Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Cliff Stein. *Introduction to algorithms*. MIT Press, Cambridge, MA, USA, 2001.
17. Intel Corp. Intel vtune.
18. Microsoft Corp. Microsoft visual studio.
19. James W. Demmel. *Applied numerical linear algebra*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1997.
20. Jim Demmel, Jack Dongarra, Victor Eijkhout, Erika Fuentes, Antoine Petitet, Rich Vuduc, Clint Whaley, and Katherine Yelick. Self adapting linear algebra algorithms and software. *Proceedings of the IEEE*, 93(2):293–312, 2005. Special issue on "Program Generation, Optimization, and Adaptation".
21. ESSL and PESSL web site. <http://www-03.ibm.com/systems/p/software/essl.html>.
22. FFTE web site. <http://www.ffte.jp>.
23. FFTPACK web site. <http://www.netlib.org/fftpack/>.
24. FFTW web site. <http://www.fftw.org>.
25. FLAME web site. <http://www.cs.utexas.edu/users/flame/>.
26. F. Franchetti and M Püschel. Short vector code generation for the discrete Fourier transform. In *Proc. IEEE Int'l Parallel and Distributed Processing Symposium (IPDPS)*, pages 58–67, 2003.

27. F. Franchetti, Y. Voronenko, and M. Püschel. Loop merging for signal transforms. In *Proc. Programming Language Design and Implementation (PLDI)*, pages 315–326, 2005.
28. F. Franchetti, Y. Voronenko, and M. Püschel. FFT program generation for shared memory: SMP and multicore. In *Proc. Supercomputing*, 2006.
29. 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.
30. Bjorn Franke and Michael F. P. O’Boyle. A complete compiler approach to auto-parallelizing c programs for multi-dsp systems. *IEEE Trans. Parallel Distrib. Syst.*, 16(3):234–245, 2005.
31. M. Frigo. A fast Fourier transform compiler. In *Proc. ACM SIGPLAN conference on Programming Language Design and Implementation (PLDI)*, pages 169–180, 1999.
32. M. Frigo and S. G. Johnson. FFTW: An adaptive software architecture for the FFT. In *Proc. IEEE Int’l Conf. Acoustics, Speech, and Signal Processing (ICASSP)*, volume 3, pages 1381–1384, 1998.
33. Matteo Frigo and Steven G. Johnson. The design and implementation of FFTW3. *Proceedings of the IEEE*, 93(2):216–231, 2005. Special issue on “Program Generation, Optimization, and Adaptation”.
34. Gcc:options that control optimization. <http://gcc.gnu.org/onlinedocs/gcc/Optimize-Options.html>.
35. GNU. Gnu gprof manual. [http://www.gnu.org/software/binutils/manual/gprof-2.9.1/html\\_mono/gprof.html](http://www.gnu.org/software/binutils/manual/gprof-2.9.1/html_mono/gprof.html).
36. K. Goto and R. van de Geijn. On reducing TLB misses in matrix multiplication, FLAME working note 9. Technical Report TR-2002-55, The University of Texas at Austin, Department of Computer Sciences, Nov. 2002.
37. GotoBLAS web site. <http://www.tacc.utexas.edu/general/staff/goto/>.
38. GSL web site. <http://www.gnu.org/software/gsl/>.
39. John A. Gunnels, Fred G. Gustavson, Greg M. Henry, and Robert A. van de Geijn. FLAME: Formal linear algebra methods environment. *TOMS*, 27(4):422–455, December 2001.
40. John L. Hennessy and David A. Patterson. *Computer Architecture: A Quantitative Approach*. Morgan Kaufmann, May 2002.
41. M. D. Hill and A. J. Smith. Evaluating associativity in cpu caches. *IEEE Trans. Comput.*, 38(12):1612–1630, 1989.
42. E.-J. Im, K. Yelick, and R. Vuduc. Sparsity: Optimization framework for sparse matrix kernels. *Int’l J. High Performance Computing Applications*, 18(1), 2004.
43. IMSL web site. <http://www.vni.com/products/ims/>.
44. Intel Corporation. *Intel 64 and IA-32 Architectures Optimization Reference Manual*, 2007. <http://www.intel.com/products/processor/manuals/index.htm>.
45. IPP web site. <http://www.intel.com/cd/software/products/asmo-na/eng/perflib/ipp/302910.htm>.
46. J. R. Johnson, R. W. Johnson, D. Rodriguez, and R. Tolimieri. A methodology for designing, modifying, and implementing FFT algorithms on various architectures. *Circuits Systems Signal Processing*, 9:449–500, 1990.
47. Steven G. Johnson and Matteo Frigo. A modified split-radix FFT with fewer arithmetic operations. *IEEE Trans. Signal Processing*, 55(1):111–119, 2007.
48. LAPACK web site. <http://www.netlib.org/lapack/>.
49. Xiaoming Li, Mara Jess Garzar, and David Padua. A dynamically tuned sorting library. In *Proc. International Symposium on Code Generation and Optimization (CGO)*, pages 111–124, 2004.
50. L. Meadows, S. Nakamoto, and V. Schuster. A vectorizing, software pipelining compiler for LIW and superscalar architecture. In *Proceedings of Risc ’92*, San Jose, CA, February 1992.

51. D. Mirković and S. L. Johnsson. Automatic performance tuning in the UHFFT library. In *Proc. Int'l Conf. Computational Science (ICCS)*, volume 2073 of *LNCS*, pages 71–80. Springer, 2001.
52. MKL web site. <http://www.intel.com/cd/software/products/asm-na/eng/307757.htm>.
53. Gordon E. Moore. Cramming more components onto integrated circuits. *Readings in computer architecture*, pages 56–59, 2000.
54. José M. F. Moura, Markus Püschel, David Padua, and Jack Dongarra. Scanning the issue: Special issue on program generation, optimization, and platform adaptation. *Proceedings of the IEEE, special issue on "Program Generation, Optimization, and Adaptation"*, 93(2):211–215, 2005.
55. NAG web site. <http://www.nag.com/>.
56. H. J. Nussbaumer. *Fast Fourier Transformation and Convolution Algorithms*. Springer, 2nd edition, 1982.
57. PLAPACK web site. <http://www.cs.utexas.edu/users/plapack/>.
58. W. H. Press, B. P. Flannery, Teukolsky S. A., and Vetterling W. T. *Numerical Recipes in C: The Art of Scientific Computing*. Cambridge University Press, 2nd edition, 1992.
59. M. Püschel, B. Singer, J. Xiong, J. M. F. Moura, J. Johnson, D. Padua, M. Veloso, and R. W. Johnson. SPIRAL: A generator for platform-adapted libraries of signal processing algorithms. *Int'l Journal of High Performance Computing Applications*, 18(1):21–45, 2004.
60. Markus Püschel, José M. F. Moura, Jeremy Johnson, David Padua, Manuela Veloso, Bryan W. Singer, Jianxin Xiong, Franz Franchetti, Aca Gačić, Yevgen Voronenko, Kang Chen, Robert W. Johnson, and Nick Rizzolo. SPIRAL: Code generation for DSP transforms. *Proc. of the IEEE*, 93(2):232–275, 2005. Special issue on *Program Generation, Optimization, and Adaptation*.
61. Quick-reference guide to optimization with intel compilers version 10.x. [http://cache-www.intel.com/cd/00/00/22/23/222300\\_222300.pdf](http://cache-www.intel.com/cd/00/00/22/23/222300_222300.pdf).
62. ScaLAPACK web site. <http://www.netlib.org/scalapack/>.
63. Spiral web site, 1998. [www.spiral.net](http://www.spiral.net).
64. Stanford SUIF Compiler Group. SUIF: A parallelizing & optimizing research compiler. Technical Report CSL-TR-94-620, Computer Systems Laboratory, Stanford University, May 1994.
65. V. Strassen. Gaussian elimination is not optimal. *Numerische Mathematik*, 14(3):354–356, 1969.
66. R. Tolimieri, M. An, and C. Lu. *Algorithms for discrete Fourier transforms and convolution*. Springer, 2nd edition, 1997.
67. UHFFT web site. <http://www2.cs.uh.edu/~mirkovic/fft/parfft.htm>.
68. C. Van Loan. *Computational Framework of the Fast Fourier Transform*. SIAM, 1992.
69. Richard Vuduc, James W. Demmel, and Katherine A. Yelick. OSKI: A library of automatically tuned sparse matrix kernels. In *Proceedings of SciDAC 2005*, Journal of Physics: Conference Series, San Francisco, CA, USA, June 2005. Institute of Physics Publishing. (*to appear*).
70. R. C. Whaley and J. Dongarra. Automatically Tuned Linear Algebra Software (ATLAS). In *Proc. Supercomputing*, 1998.
71. R. C. Whaley, A. Petitet, and J. J. Dongarra. Automated empirical optimization of software and the ATLAS project. *Parallel Computing*, 27(1-2):3–35, 2001. Also available as University of Tennessee LAPACK Working note #147, UT-CS-00-448, 2000 ([www.netlib.org/lapack/lawns/lawn147.ps](http://www.netlib.org/lapack/lawns/lawn147.ps)).
72. M. Wolfe. Iteration space tiling for memory hierarchies. In *Third SIAM Conference on Parallel Processing for Scientific Computing*, December 1987.



73. K. Yotov, X. Li, G. Ren, M. Garzaran, D. Padua, K. Pingali, and P. Stodghill. Is search really necessary to generate high-performance blas. In *Proceedings of the IEEE, special issue on Program Generation, Optimization, and Adaptation.*, 2005.
74. Kamen Yotov, Xiaoming Li, Gang Ren, Maria Garzaran, David Padua, Keshav Pingali, and Paul Stodghill. A comparison of empirical and model-driven optimization. *Proceedings of the IEEE*, 93(2), 2005. special issue on "Program Generation, Optimization, and Adaptation".