

Compile-Time Optimizations for Tensor Contraction Expressions*

Gerald Baumgartner, Daniel Cociorva, Chi-Chung Lam, P. Sadayappan
Dept. of Computer and Information Science
Ohio State University
{gb,cociorva,clam,saday}@cis.ohio-state.edu

J. Ramanujam
Dept. of Electrical and Computer Engineering
Louisiana State University
jxr@ee.lsu.edu

Abstract

This paper provides an overview of compile-time optimizations being developed in the context of a program synthesis system for a class of quantum chemistry computations. These computations are expressible as a set of tensor contractions and arise in electronic structure modeling. The synthesis system will take as input a high-level specification of the computation and generate high-performance parallel code for a number of target architectures. Several components of the synthesis system are described, focusing on compile-time optimization issues that they address.

1. Introduction

The development of high-performance parallel programs for scientific applications is often very complicated. The effect of the algorithm choice on memory access costs and communication overhead are often very complex. Currently, available tools for software development and performance optimization do not provide adequate support to the developers of high-performance scientific applications. Often, the time to develop an efficient parallel program for a computational model is the limiting factor in the rate of progress of the science. Therefore, approaches to automated synthesis of high-performance parallel programs from high-level specifications are very attractive. In general, automatic program synthesis is an intractable problem. However, for specific domains, this is feasible, as is being demonstrated by the SPIRAL project [23] for the domain of signal processing.

In this paper, we provide an overview of a project that is developing a program synthesis system to facilitate the rapid development of high-performance parallel programs for a class of scientific computations encountered in chemistry and physics - electronic structure calculations, where many computationally intensive components are expressible as a set of tensor contractions. Currently, manual development of accurate quantum chemistry models in this domain is very tedious and takes an expert several months to years to develop and debug. The synthesis tool aims to reduce the development time to hours/days, by having the chemist specify the computation in a high-level form, from which an efficient parallel program is automatically synthesized. This should enable the rapid synthesis of high-performance implementations of sophisticated ab-initio quantum chemistry models, including models that are too tedious for manual development by quantum chemists.

The computational domain that we consider is extremely compute-intensive and consumes significant computer resources at national supercomputer centers. Many of these codes are limited in the size of the problem that they can currently solve because of memory and performance limitations. The computational structures that we address are present in some computational physics codes modeling electronic properties of semiconductors and metals, and in computational chemistry codes such as ACES II, GAMESS, Gaussian, NWChem[9], PSI, and MOLPRO. In particular, they comprise the bulk of the computation with the coupled cluster approach to the accurate description of the electronic structure of atoms and molecules [18, 19]. Computational approaches to modeling the structure and interactions of molecules, the electronic and optical properties of molecules, the heats and rates of chemical reactions, etc., are crucial to the understanding of chemical processes in real-world systems. Examples of applications include combustion and atmospheric chemistry, chemical vapor depo-

*Supported in part by the National Science Foundation through awards CHE-0121676 and EIA-9986052, and the U. S. Department of Energy through award DE-AC05-00OR22725.

sition, protein structure and enzymatic chemistry, and industrial chemical processing. Computational chemistry and materials science account for significant fractions of super-computer usage at national centers (for example, approximately 85% of total usage at Pacific Northwest National Laboratories, 30% at NERSC, and about 50% of one of SDSC’s systems).

2. The computational context

In the class of computations considered, the final result to be computed can be expressed in terms of tensor contractions, essentially a collection of multi-dimensional summations of the product of several input arrays. Due to commutativity, associativity, and distributivity, there are many different ways to compute the final result, and they could differ widely in the number of floating point operations required. Consider the following expression:

$$S_{abij} = \sum_{cdefkl} A_{acik} \times B_{befl} \times C_{dfjk} \times D_{cdel}$$

If this expression is directly translated to code (with ten nested loops, for indices $a - l$), the total number of arithmetic operations required will be $4 \times N^{10}$ if the range of each index $a - l$ is N . Instead, the same expression can be rewritten by use of associative and distributive laws:

$$S_{abij} = \sum_{ck} \left(\sum_{df} \left(\sum_{el} B_{befl} \times D_{cdel} \right) \times C_{dfjk} \right) \times A_{acik}$$

This corresponds to the formula sequence shown in Fig. 1(a) and can be directly translated into code as shown in Fig. 1(b). This form only requires $6 \times N^6$ operations. However, additional space is required to store temporary arrays $T1$ and $T2$. Often, the space requirements for the temporary arrays poses a serious problem. For this example, abstracted from a quantum chemistry model, the array extents along indices $a - d$ are the largest, while the extents along indices $i - l$ are the smallest. Therefore, the size of temporary array $T1$ would dominate the total memory requirement.

The operation minimization problem encountered here is a generalization of the well known matrix-chain multiplication problem, where a linear chain of matrices to be multiplied is given, e.g., ABCD, and the optimal order of pair-wise multiplications is sought, i.e., ((AB)C)D versus (AB)(CD), etc. In contrast, for computations expressed as sets of matrix contractions, there is additional freedom in choosing the pair-wise products. For the above example, instead of forcing a single chain order, e.g., ABCD, other orders are possible, such as the BDCA order shown for the operation-reduced form above.

The problem of determining the operator tree with minimal operation count is NP-complete, and an efficient pruning search procedure has been developed [16, 17]. For the

above example, although the latter form is far more economical in terms of the number of operations, its implementation will require the use of temporary intermediate arrays to hold the partial results of the parenthesized array subexpressions. Sometimes, the sizes of intermediate arrays needed for the “operation-minimal” form are too large to even fit on disk.

A systematic way to explore ways of reducing the memory requirement for the computation is to view it in terms of potential loop fusions. Loop fusion merges loop nests with common outer loops into larger imperfectly nested loops. When one loop nest produces an intermediate array which is consumed by another loop nest, fusing the two loop nests allows the dimension corresponding to the fused loop to be eliminated in the array. This results in a smaller intermediate array and thus reduces the memory requirements. For the example considered, the application of fusion is illustrated in Fig. 1(c). This way, $T1$ can be reduced to a scalar and $T2$ to a 2-dimensional array, without changing the number of operations.

For a computation comprising of a number of nested loops, there are often many fusion choices, that are not all mutually compatible. This is because different fusion choices could require different loops to be made the outermost. In prior work, we addressed the problem of finding the choice of fusions for a given operator tree that minimized the total space required for all arrays [13, 14, 15].

3. An example from CCSD(T)

One of the most computationally intensive components of many quantum chemistry packages is the CCSD(T) scheme. It is a coupled cluster approximation that includes single and double excitations from the Hartree-Fock wavefunction plus a perturbative estimate for the *connected* triple excitations. For molecules well described by a Hartree-Fock wave function, this method predicts bond energies, ionization potentials, and electron affinities to an accuracy of ± 0.5 kcal/mol, bond lengths accurate to ± 0.0005 Å, and vibrational frequencies accurate to ± 5 cm^{-1} . This level of accuracy is adequate to answer many of the questions that arise in studies of chemical systems.

As a concrete example from quantum chemistry, we discuss a component of the CCSD(T) calculation. The following representative equation arises in the Laplace factorized expression for linear triples perturbation correction:

$$\begin{aligned} A3A &= X_{ce,af} Y_{ae,cf} + X_{a\bar{e},c\bar{f}} Y_{c\bar{e},a\bar{f}} + X_{a\bar{e},\bar{c}f} Y_{\bar{c}\bar{e},af} \\ &+ X_{\bar{a}e,c\bar{f}} Y_{ce,\bar{a}\bar{f}} + X_{\bar{a}e,\bar{c}f} Y_{\bar{c}\bar{e},\bar{a}\bar{f}} + X_{\bar{a}\bar{e},\bar{c}f} Y_{\bar{c}\bar{e},\bar{a}\bar{f}}, \end{aligned}$$

where X and Y are of the form $X_{ae,cf} = t_{ij}^{ae} t_{ij}^{cf}$ and $Y_{ce,af} = \langle cb || ek \rangle \langle ab || fk \rangle$, respectively.

Integrals with two vertical bars have been antisymmetrized and may be expressed as: $\langle pq || rs \rangle = \langle pq |$

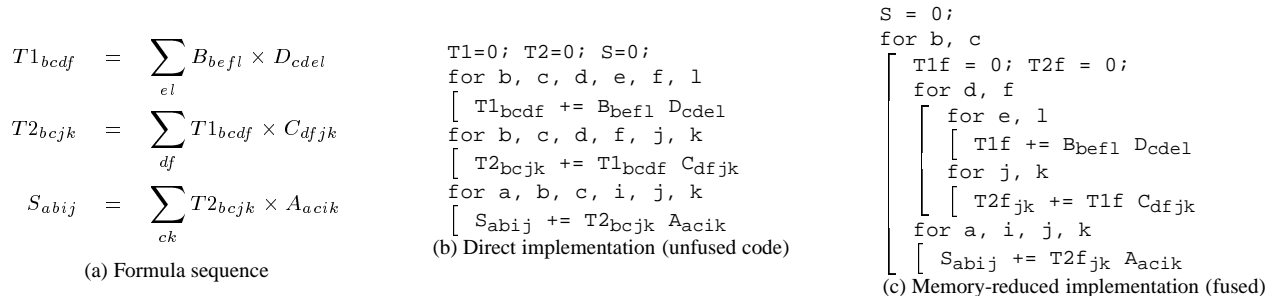


Figure 1. Example illustrating use of loop fusion for memory reduction.

$rs\rangle - \langle pq | sr\rangle$), where integrals with one vertical bar are of the form $\langle \mu\nu | \omega\lambda \rangle = \int \int dr^3 ds^3 \phi_\mu(\underline{\mathbf{r}})\phi_\nu(\underline{\mathbf{s}})|\underline{\mathbf{r}} - \underline{\mathbf{s}}|^{-1}\phi_\omega(\underline{\mathbf{r}})\phi_\lambda(\underline{\mathbf{s}})$ and are quite expensive to compute (requiring on the order of 1000 arithmetic operations). Electrons may have either up or down (or alpha/beta) spin. Down spin is denoted here with an over-bar. The indices i, j, k, l, m, n refer to occupied orbitals, of number O between 30 and 100. The indices a, b, c, d, e, f refer to unoccupied orbitals of number V between 1000 and 3000. The integrals are written in the molecular orbital basis, but must be computed in the underlying atom-centered Gaussian basis, and transformed to the molecular orbital basis. We omit these details in our discussion here.

A3A is one of many contributions to the energy, and among the most expensive, scaling as $O(OV^5)$. Here, we assume that we have already computed the amplitudes t_{ij}^{ae} , and they must be read as necessary, and contracted to form a block of X . The integrals $\langle cb || ek \rangle$ must be recomputed as necessary, contracted to form a block of Y corresponding to X , and the two contracted to form the scalar contribution to the energy.

Fig. 2 shows pseudo-code for the computation of one of the energy components E for A3A. Temporary arrays $T1$ and $T2$ are used to store the integrals of form $\langle ab || ek \rangle$, where the functions f_1 and f_2 represent the integral calculations (in reality, f_1 and f_2 represent the same array/function; but it is more convenient to treat them as distinct initially, to simplify our explanation about the space-time trade-off problem addressed by the synthesis system). The intermediate quantities X_{aecf} are computed by contracting over (i.e., summing over products of) input array T , while the intermediate quantities Y_{ceaf} are obtained by contracting over $T1$ and $T2$. The final result is a single scalar quantity E , that is obtained by adding together the $O(OV^3)$ pair-wise products $X_{aecf}Y_{ceaf}$.

The cost of computing each integral f_1, f_2 is represented by C_f , and in practice is of the order of hundreds or a few thousand arithmetic operations. The pseudo-code form shown in Fig. 2 is computationally very efficient in minimizing the number of expensive integral function evaluations f_1 and f_2 , and maximizing the reuse of

the stored integrals in $T1$ and $T2$ (each element of $T1$ and $T2$ is used $O(V^2)$ times). However, it is impractical due to the huge memory requirement. With $O = 100$ and $V = 5000$, the size of $T1, T2$ is $O(10^{14})$ bytes and the size of X, Y is $O(10^{15})$ bytes. By fusing together pairs of producer-consumer loops in the computation, reductions in the needed array sizes may be sought, since the fusion of a loop with common index in the pair of loops allows the elimination of that dimension of the intermediate array. It can be seen that the loop that produces X (with indices a, e, c, f), the loop that produces Y (with indices c, e, a, f) and the loop that consumes X and Y to produce E (with indices c, e, a, f) can all be fully fused together, permitting the elimination of all explicit indices in X and Y to reduce them to scalars. However, the loops producing $T1$ (with indices c, e, b, k) and $T2$ (with indices a, f, b, k) cannot also be directly fused with the other three loops because their indices do not match.

Fig. 3 shows how a reduction of space for $T1$ and $T2$ can be achieved by introducing redundant loops around their producer loops — add loops with the missing indices a, f for $T1$ and c, e for $T2$. Now all five loops have common indices a, e, c, f that can be fused, permitting elimination of those indices from all temporaries. Further, by fusing the producer loops for $T1$ and $T2$ with their consumer loop, which produces Y , the b, k indices can also be eliminated from $T1$ and $T2$. A dramatic reduction of memory space is achieved, reducing all temporaries $T1, T2, X$ and Y to scalars, but the space savings come at the price of a significant increase in computation. No reuse is achieved of the quantities derived from the expensive integral calculations f_1 and f_2 . Since C_f is of the order of 1000 in practice, the integral calculations now dominate the total compute time, increasing the operation count by three orders of magnitude over the unfused form in Fig. 2.

A desirable solution would be somewhere in between the unfused structure of Fig. 2 (with maximal memory requirement and maximal reuse) and the fully fused structure of Fig. 3 (with minimal memory requirement and minimal reuse). This is shown in Fig. 4, where tiling and partial fusion of the loops is employed. The loops with indices

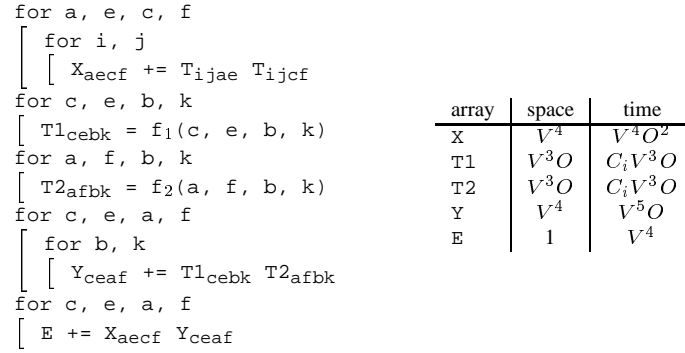


Figure 2. Unfused operation-minimal form.

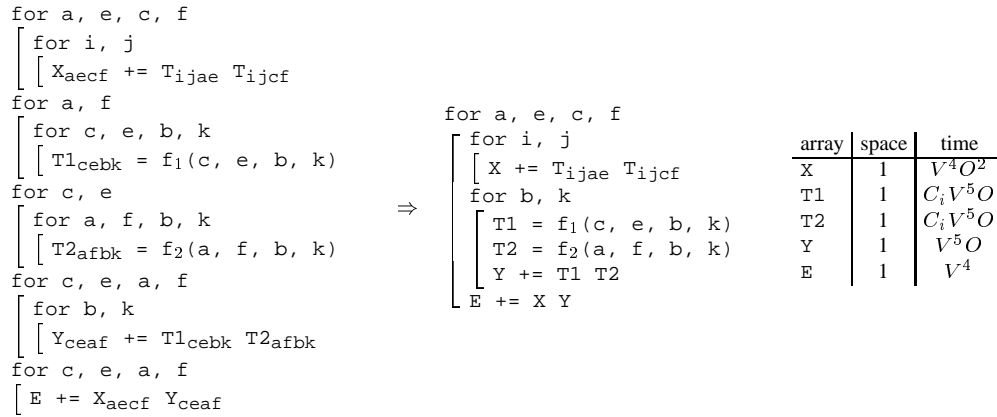


Figure 3. Use of redundant computation to allow full fusion.

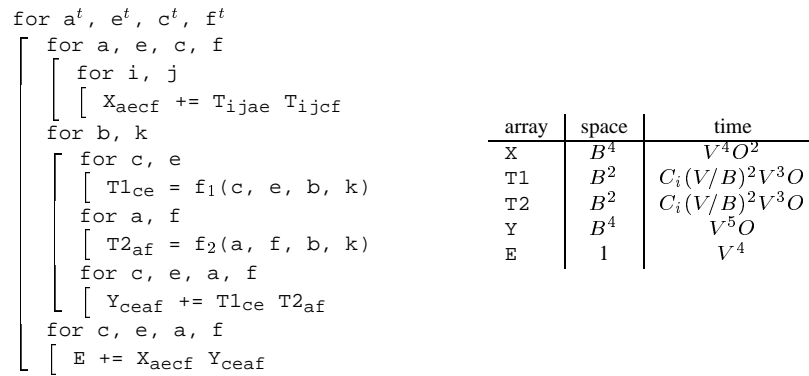


Figure 4. Use of tiling and partial fusion to reduce recomputation cost.

a, e, c, f are tiled by splitting each of those indices into a pair of indices. The indices with a superscript t represent the tiling loops and the unsuperscripted indices now stand for intra-tile loops with a range of B , the block size used for tiling. For each tile (a^t, e^t, c^t, f^t) , blocks of $T1$ and $T2$ of size B^2 are computed and used to form B^4 product contributions to the components of Y , which are stored in an array of size B^4 .

As the tile size B is increased, the cost of function computation for f_1, f_2 decreases by a factor of B^2 , due to the reuse enabled. However, the size of the needed temporary array for Y increases as B^4 (the space needed for X can be reduced back to a scalar by fusing its producer loop with the loop producing E , but Y 's space requirement cannot be decreased). When B^4 becomes larger than the size of physical memory, expensive paging in and out of disk will be required for Y . Further, there is diminishing returns on reuse of $T1$ and $T2$ after B^2 becomes comparable to C_f , since the loop producing Y now becomes the dominant one. So we can expect that as B is increased, performance will improve and then level off and then deteriorate. The optimum value of B will clearly depend on the cost of access at the various levels of the memory hierarchy.

The computation considered here is just one component of the $A3A$ term, which in turn is only one of very many terms that must be computed. Although developers of quantum chemistry codes naturally recognize and perform some of these optimizations, a collective analysis of all these computations to determine their optimal implementation is beyond the scope of manual effort. Further, the time required to develop codes to implement such computational models is quite large, especially since the tensor expressions can get quite complex - Fig. 5 shows an example of the kind of tensor expressions encountered when developing accurate computational models.

In the next section, we provide an overview of a transformation system that we are developing to aid quantum chemists in rapidly developing high-performance parallel codes for computations that they specify as a set of high-level tensor contractions.

4. Overview of the synthesis system

We present in this section a brief description of the basic components of the system being developed. Some of these components are tightly coupled (for example, memory minimization and data distribution), and they are treated together as one combined module in the synthesis system.

High-level language: The input to the synthesis system is a sequence of tensor contraction expressions (essentially sum-of-products array expressions) together with declarations of index ranges and symmetry and sparsity of matrices. This high-level notation provides essential information to the optimization components that would be difficult or impossible to extract out of low-level code.

Algebraic transformations: It takes input from the user as tensor expressions and synthesizes an output computation sequence. The Algebraic Transformations module uses the properties of commutativity and associativity of addition and multiplication and the distributivity of multiplication over addition. It searches for all possible ways of applying these properties to an input sum-of-products expression, and determines a combination that results in an equivalent form of the computation with minimal operation cost.

Memory minimization: The operation-minimal computation sequence synthesized by the Algebraic Transformation module might require an excessive amount of memory due to the large temporary intermediate arrays involved. The Memory Minimization module attempts to perform loop fusion transformations to reduce the memory requirements. This is done without any change to the number of arithmetic operations.

Data distribution and partitioning: This component determines how best to partition the arrays among the processors of a parallel system. We assume a data-parallel model, where each operation in the operation sequence is distributed across the entire parallel machine. The arrays are to be disjointly partitioned between the physical memories of the processors. Since the data distribution pattern affects the memory usage on the parallel machine, this component is closely coupled with the memory minimization component.

Space-time transformation: If the memory minimization module is unable to reduce memory requirements of the computation sequence below the available disk capacity on the system, the computation is infeasible unless a space-time trade-off is performed. If no satisfactory transformation is found, feedback is provided to the Memory Minimization module, causing it to seek a different solution. If the Space-Time Transformation module is successful in bringing down the memory requirement below the disk capacity, the Data Locality Optimization module is invoked.

Data locality optimization: If the space requirement exceeds physical memory capacity, portions of the arrays must be moved between disk and main memory as needed, in a way that maximizes reuse of elements in memory. The same considerations are involved in minimizing cache misses — blocks of data are moved between physical memory and the space available in the cache.

Code generation: The back end of the synthesis system provides the output as pseudocode, Fortran or C code. The generated code can be either serial or parallel, using the MPI libraries. Depending on the circumstances, the synthesised code could also call highly-tuned, machine-specific Basic Linear Algebra Subprograms (BLAS) libraries, or optimized low-level functions from the existing quantum chemistry packages.

In the next sections, we provide some details about the optimizations implemented in some of these modules. For

$$\begin{aligned}
\text{hbar}[a,b,i,j] = & \text{sum}[f[b,c] * t[i,j,a,c], c] - \text{sum}[f[k,c] * t[k,b] * t[i,j,a,c], k,c] + \text{sum}[f[a,c] * t[i,j,c,b], c] - \text{sum}[f[k,c] * t[k,a] * t[i,j,c,b], k,c] - \text{sum}[f[k,j] * t[i,k,a,b], k] - \text{sum}[f[k,c] * \\
& t[j,c] * t[i,k,a,b], k,c] - \text{sum}[f[k,i] * t[j,k,b,a], k] - \text{sum}[f[k,c] * t[i,c] * t[j,k,b,a], k,c] + \text{sum}[t[i,c] * t[j,d] * v[a,b,c,d], c,d] + \text{sum}[t[i,j,c,d] * v[a,b,c,d], c,d] + \text{sum}[t[j,c] * v[a,b,i,c], \\
& c] - \text{sum}[t[k,b] * v[a,k,i,j], k] + \text{sum}[t[i,c] * v[b,a,j,c], c] - \text{sum}[t[k,a] * v[b,k,j,i], k] - \text{sum}[t[k,d] * t[i,j,c,b] * v[k,a,c,d], k,c,d] - \text{sum}[t[i,c] * t[j,k,b,d] * v[k,a,c,d], k,c,d] - \text{sum}[t[j,c] \\
& * t[k,b] * v[k,a,c,i], k,c] + 2 * \text{sum}[t[j,k,b,c] * v[k,a,c,i], k,c] - \text{sum}[t[j,k,c,b] * v[k,a,c,i], k,c] - \text{sum}[t[i,c] * t[j,d] * t[k,b] * v[k,a,d,c], k,c,d] + 2 * \text{sum}[t[k,d] * t[i,j,c,b] * v[k,a,d,c], \\
& k,c,d] - \text{sum}[t[k,b] * t[i,j,c,d] * v[k,a,d,c], k,c,d] - \text{sum}[t[j,d] * t[i,k,c,b] * v[k,a,d,c], k,c,d] + 2 * \text{sum}[t[i,c] * t[j,k,b,d] * v[k,a,d,c], k,c,d] - \text{sum}[t[i,c] * t[j,k,d,b] * v[k,a,d,c], \\
& k,c,d] - \text{sum}[t[j,k,b,c] * v[k,a,i,c], k,c] - \text{sum}[t[i,c] * t[k,b] * v[k,a,j,c], k,c] - \text{sum}[t[i,k,c,b] * v[k,a,j,c], k,c] - \text{sum}[t[i,c] * t[j,d] * t[k,a] * v[k,b,c,d], k,c,d] - \text{sum}[t[k,d] * t[i,j,a,c] \\
& * v[k,b,c,d], k,c,d] - \text{sum}[t[k,a] * t[i,j,c,d] * v[k,b,c,d], k,c,d] + 2 * \text{sum}[t[j,d] * t[i,k,a,c] * v[k,b,c,d], k,c,d] - \text{sum}[t[j,d] * t[i,k,c,a] * v[k,b,c,d], k,c,d] - \text{sum}[t[i,c] * t[j,k,d,a] \\
& * v[k,b,c,d], k,c,d] - \text{sum}[t[i,c] * t[k,a] * v[k,b,c,j], k,c] + 2 * \text{sum}[t[i,k,a,c] * v[k,b,c,j], k,c] - \text{sum}[t[i,k,c,a] * v[k,b,c,j], k,c] + 2 * \text{sum}[t[k,d] * t[i,j,a,c] * v[k,b,d,c], k,c,d] - \\
& \text{sum}[t[j,d] * t[i,k,a,c] * v[k,b,d,c], k,c,d] - \text{sum}[t[j,c] * t[k,a] * v[k,b,i,c], k,c] - \text{sum}[t[j,k,c,a] * v[k,b,i,c], k,c] - \text{sum}[t[i,k,a,c] * v[k,b,j,c], k,c] + \text{sum}[t[i,c] * t[j,d] * t[k,a] * t[l,b] * \\
& v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[k,b] * t[l,d] * t[i,j,a,c] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[k,a] * t[l,d] * t[i,j,c,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[k,a] * t[l,b] * t[i,j,c,d] * v[k,l,c,d], k,l,c,d] \\
& - 2 * \text{sum}[t[j,c] * t[l,d] * t[i,k,a,b] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[j,d] * t[l,b] * t[i,k,a,c] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[j,d] * t[l,b] * t[i,k,c,a] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,c] * \\
& t[l,d] * t[j,k,b,a] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,a] * t[j,k,b,d] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,b] * t[j,k,d,a] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,k,c,d] * t[j,l,b,a] * v[k,l,c,d], \\
& k,l,c,d] + 4 * \text{sum}[t[i,k,a,c] * t[j,l,b,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,k,c,a] * t[j,l,b,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,k,a,b] * t[j,l,c,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,k,a,c] \\
& * t[j,l,d,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,k,c,a] * t[j,l,d,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[j,d] * t[k,l,a,b] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[i,j,c,d] * t[k,l,a,b] * v[k,l,c,d], k,l,c,d] - 2 \\
& * \text{sum}[t[i,j,c,b] * t[k,l,a,d] * v[k,l,c,d], k,l,c,d] - 2 * \text{sum}[t[i,j,a,c] * t[k,l,b,d] * v[k,l,c,d], k,l,c,d] + \text{sum}[t[j,c] * t[k,b] * t[l,a] * v[k,l,c,i], k,l,c] + \text{sum}[t[l,c] * t[j,k,b,a] * v[k,l,c,i], \\
& k,l,c] - 2 * \text{sum}[t[l,a] * t[j,k,b,c] * v[k,l,c,i], k,l,c] + \text{sum}[t[l,a] * t[j,k,c,b] * v[k,l,c,i], k,l,c] - 2 * \text{sum}[t[k,c] * t[j,l,b,a] * v[k,l,c,i], k,l,c] + \text{sum}[t[k,a] * t[j,l,b,c] * v[k,l,c,i], k,l,c] \\
& + \text{sum}[t[k,b] * t[j,l,c,a] * v[k,l,c,i], k,l,c] + \text{sum}[t[j,c] * t[l,k,a,b] * v[k,l,c,i], k,l,c] + \text{sum}[t[i,c] * t[k,a] * t[l,b] * v[k,l,c,j], k,l,c] + \text{sum}[t[l,c] * t[i,k,a,b] * v[k,l,c,j], k,l,c] - 2 * \\
& \text{sum}[t[l,b] * t[i,k,a,c] * v[k,l,c,j], k,l,c] + \text{sum}[t[l,b] * t[i,k,c,a] * v[k,l,c,j], k,l,c] + \text{sum}[t[i,c] * t[k,l,a,b] * v[k,l,c,j], k,l,c] + \text{sum}[t[j,c] * t[l,d] * t[i,k,a,b] * v[k,l,d,c], k,l,c,d] + \\
& \text{sum}[t[j,d] * t[l,b] * t[i,k,a,c] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[j,d] * t[l,a] * t[i,k,c,b] * v[k,l,d,c], k,l,c,d] - 2 * \text{sum}[t[i,k,c,d] * t[j,l,b,a] * v[k,l,d,c], k,l,c,d] - 2 * \text{sum}[t[i,k,a,c] * t[j,l,b,d] \\
& * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,c,a] * t[j,l,b,d] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,a,b] * t[j,l,c,d] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,c,b] * t[j,l,d,a] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[i,k,a,c] * \\
& t[j,l,d,b] * v[k,l,d,c], k,l,c,d] + \text{sum}[t[k,a] * t[l,b] * v[k,l,i,j], k,l] + \text{sum}[t[k,l,a,b] * v[k,l,i,j], k,l] + \text{sum}[t[k,b] * t[l,d] * t[i,j,a,c] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[k,a] * t[l,d] * t[i,j,c,b] \\
& * v[l,k,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,d] * t[j,k,b,a] * v[l,k,c,d], k,l,c,d] - 2 * \text{sum}[t[i,c] * t[l,a] * t[j,k,b,d] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[i,c] * t[l,a] * t[j,k,d,b] * v[l,k,c,d], k,l,c,d] + \\
& \text{sum}[t[i,j,c,b] * t[k,l,a,d] * v[l,k,c,d], k,l,c,d] + \text{sum}[t[i,j,a,c] * t[k,l,b,d] * v[l,k,c,d], k,l,c,d] - 2 * \text{sum}[t[l,c] * t[i,k,a,b] * v[l,k,c,j], k,l,c] + \text{sum}[t[l,b] * t[i,k,a,c] * v[l,k,c,j], k,l,c] + \\
& \text{sum}[t[l,a] * t[i,k,c,b] * v[l,k,c,j], k,l,c] + v[a,b,i,j]
\end{aligned}$$

Figure 5. A contraction expression from quantum chemistry.

details of the operation minimization algorithm see [16, 17].

5. Memory minimization and space-time trade-offs

As discussed in Section 2, the operation minimization procedure results in the creation of intermediate temporary arrays. For typical computations in computational chemistry, the space required for storing these temporary arrays can be several tera bytes, which makes the computation impractical. As shown in Fig. 1(c), the problem with memory requirements of large intermediate arrays can be mitigated through loop fusion. Loop fusion merges loop nests with common outer loops into larger imperfectly nested loops. When one loop nest produces an intermediate array that is consumed by another loop nest, fusing the two loop nests allows the dimension corresponding to the fused loop to be eliminated in the array. This results in a smaller intermediate array and thus lowers the memory requirement. The use of loop fusion can be seen to result in significant potential reduction to the total memory requirement. For a computation comprising of a number of nested loops, there will generally be a number of fusion choices, that are not all mutually compatible. This is because different fusion choices could require different loops to be made the outermost. A data structure that we call a *fusion graph* can be used to facilitate enumeration of all possible compatible fu-

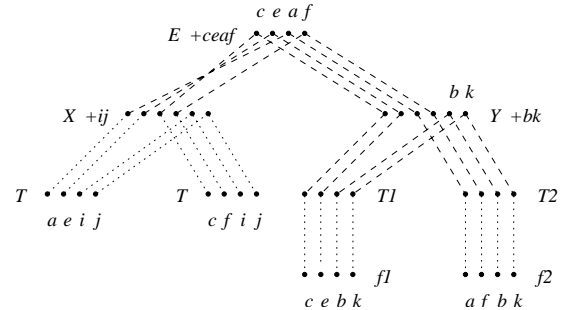


Figure 6. Fusion graph for unfused operation-minimal form of loop in Fig. 2

sion configurations for a given computation tree.

Fig. 6 shows the fusion graph for the unfused form of the computation from Fig. 2. Corresponding to each node in a computation tree, the fusion graph has a set of vertices corresponding to the loop indices of the node of the computation tree. In Fig. 6, we do not show the operator tree corresponding to the computation, but directly illustrate the fusion graph. The potential for fusion of a common loop among a producer-consumer pair of loop nests is indicated in the fusion graph through a dashed *potential fusion* edge connecting the corresponding vertices. Leaf nodes in the

fusion graph correspond to input arrays or primitive function evaluations and do not represent a loop nest. The edges from the leaves to their parents are shown as dotted edges and do not affect the fusion possibilities. If a pair of loop nests is fused using one or more common loops, it is captured in the fusion graph by changing the dashed potential-fusion edges to continuous fusion edges. If more than two loop nests are fused together, a chain of fusion edges results, called a *fusion chain*. The *scope of a fusion chain* is the set of nodes it spans. The fusion graph allows us to characterize the condition for feasibility of a particular combination of fusions: the scope of any two fusion chains in a fusion graph must either be disjoint or a subset/superset of each other. Scopes of fusion chains do not partially overlap because loops do not (i.e., loops must be either separate or nested).

The fusion graph in Fig. 6 can be used to determine the fusion possibilities. On the left side of the graph, the edges corresponding to (a, e, c, f) can all be made fusion edges, suggesting that complete fusion is possible for the loop nests producing and consuming X , reducing it to a scalar. Similarly, on the right side of the graph, the edges corresponding to (c, e, a, f) can also be made fusion edges, reducing Y to a scalar. Further, by creating fusion edges for indices (c, e) , the producer loop for $T1$ can be fully fused with the Y loop that consumes it. However, now the producer loop for $T2$ cannot be fused since the addition of any fusion edge (say for index a) will result in partially overlapping fusion chains for a and (c, e) . The fully fused version from Fig. 3 can be represented graphically as shown in Fig. 7(a). Additional vertices have been added for indices (c, e) and (a, f) , respectively, at the nodes corresponding to the producer loops for $T1$ and $T2$. Now, complete fusion chains can be created without any partial overlap in the scopes of the fusion chains. From the figure, it can be seen that in fact the redundant computation need only be added to one of $T1$ or $T2$ to achieve complete fusion. For example, removing the additional vertices for (a, f) at $T2$ does not violate the non-partial-overlap condition for fusion.

The fusion graph was used to develop an algorithm [15, 13] to determine the combination of fusions that minimizes the total storage required for all the temporary intermediate arrays. A bottom-up dynamic programming approach was used that maintains a set of pareto-optimal fusion configurations at each node, merging solutions from children nodes to generate the optimal configurations at a parent. The two metrics used are the total memory required under the subtree rooted at the node, and the constraints imposed by a configuration on fusion further up the tree. A configuration is inferior to another if it is *more or equally constraining* with respect to further fusions than the other, and uses no less memory. At the root of the tree, the configuration with the lowest memory requirement is chosen.

Although the complexity of the algorithm is exponential

in the number of index variables and the number of solutions could in theory grow exponentially with the size of the expression tree, the number of index variables in practical applications is small enough and there is indication that the pruning is effective in keeping the size of the solution set at each node small.

If after memory minimization the storage requirements still exceed the disk capacity, we need to recompute some (parts of) temporary arrays in order to further reduce the space requirements. The space-time trade-off algorithm we have developed [4] employs a combination of fusion and tiling to achieve a good balance between recomputation and memory usage. The first step of the space-time trade-off algorithm uses a dynamic programming approach similar to the memory minimization algorithm that maintains a set of pareto-optimal fusion/recomputation configurations, in which the recomputation cost is used as a third metric. Solutions exceeding the memory limit are pruned out. The result of the search is a set of loop structures with different combinations of space requirements and recomputation cost.

In the second step of the algorithm, recomputation indices are split into tiling and intra-tile loop pairs. By making intra-tile loops the inner-most loops, any recomputation only needs to be performed once per iteration of the tiling loop in exchange for increasing the storage requirements for temporaries in which the dimension corresponding to the tiled loop had been eliminated. For each solution from the first step of the algorithm, we then search for tile sizes that minimize the recomputation cost, and take the solution that results in the lowest recomputation cost.

6. Data locality optimization

Once a solution is found that fits onto disk, we optimize the data locality to reduce memory and disk access times. We developed algorithms [2, 3] that, given a memory-reduced (fused) version of the code, find the appropriate blocking of the loops in order to maximize data reuse. These algorithms can be applied at different levels of the memory hierarchy, for example, to minimize data transfer between main memory and disk (disk access minimization), or to minimize data transfer between main memory and the cache (cache optimization). In this section, we briefly describe the main points of our algorithm [3], focusing mostly on the cache management problem. For the disk access minimization problem, the same approach is used, replacing the cache size by the physical memory size.

We introduce a memory access cost model (*Cost*), an estimate on the number of cache misses, as a function of tile sizes and loop bounds. In a bottom-up traversal of the abstract syntax tree, we count for each loop the number (*Accesses*) of distinct array elements accessed in its scope. If this number is smaller than the number of elements that fit into the cache, then $Cost = Accesses$. Otherwise, it means

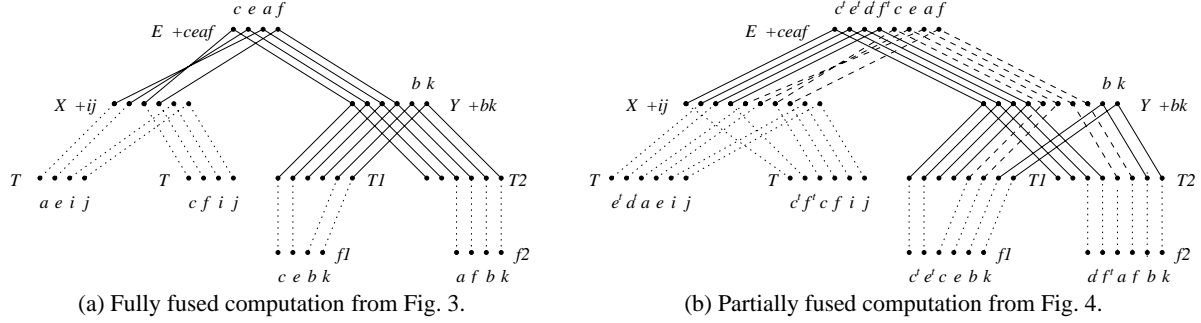


Figure 7. Fusion graphs showing redundant computation and tiling.

that the elements in the cache are not reused from one loop iteration to the next, and the cost is obtained by multiplying the loop range by the cost of its inner loop(s).

Using this cost model, we can compute the total memory access cost for given tile sizes. The procedure is repeated for different sets of tile sizes, and new costs are computed. In the end the lowest possible cost is chosen, thus determining the optimal tile sizes. We define our tile size search space in the following way: if N_i is a loop range, we use a tile size starting from $T_i = 1$ (no tiling), and successively increasing T_i by doubling it until it reaches N_i . This approach ensures a slow (logarithmic) growth of the search space with increasing array dimension for large N_i . If N_i is small enough, an exhaustive search is performed instead.

7. Data partitioning and communication minimization

Given a sequence of formulae, we need to find an effective partitioning of arrays and operations among the processors and a choice of loop fusions in order to minimize inter-processor communication, while staying within the available memory in implementing the computation on a message-passing parallel computer.

Since primitive tensor contractions are essentially generalized multi-dimensional matrix multiplications, we use a generalized form of the memory efficient Cannon algorithm [12]. A logical view of the P processors as a two-dimensional $\sqrt{P} \times \sqrt{P}$ grid is used, and each array is fully distributed along the two processor dimensions. We use a pair of indices to denote the partitioning or *distribution* of the elements of a data array on a two-dimensional processor array. The d -th position in a pair α , denoted $\alpha[d]$, where d can be either 1 or 2, corresponds to the d -th processor dimension. Each position is an index variable distributed along that processor dimension. As an example, suppose 16 processors form a two-dimensional 4×4 logical array. For the array $B(b, e, f, l)$ in Fig. 1(a), the pair $\alpha = \langle b, f \rangle$ specifies that the first (b) and the third (f) dimensions of B are distributed along the first and second processor dimensions respectively, and that the sec-

ond (e) and fourth (l) dimensions of B are not distributed. Thus, a processor whose id is P_{z_1, z_2} , with z_1 and z_2 between 1 and 4, will be assigned a portion of B specified by $B(\text{myrange}(z_1, N_b, 4), 1 : N_e, \text{myrange}(z_2, N_e, 4), 1 : N_l)$, where $\text{myrange}(z, N, p)$ is the range $(z - 1) \times N/p + 1$ to $z \times N/p$.

A tensor contraction formula can be expressed as a generalized matrix multiplication $C(I, J) += A(I, K) * B(K, J)$, where I, J , and K represent index collections, or index sets. This observation follows from a special property of tensor contractions: all the indices appearing on the left-hand side must appear on the right-hand side only *once* (index sets I and J , for A and B , respectively), and all summation indices must appear on both right-hand side arrays (index set K). For example, the tensor contraction $T1(b, c, d, f) = \sum_{e, l} B(b, e, f, l) \times D(c, d, e, l)$ is characterized by the index sets $I = \{b, f\}$, $J = \{c, d\}$, and $K = \{e, l\}$.

We generalize Cannon’s algorithm for multi-dimensional arrays as follows: a triplet $\{i, j, k\}$ formed by one index from each index set I, J , and K defines a distribution $\langle i, j \rangle$ for the result array C , and distributions $\langle i, k \rangle$ and $\langle k, j \rangle$ for the input arrays A and B , respectively. In addition, one of the 3 indices $\{i, j, k\}$ is chosen as the “rotation index,” along which the processor communication takes place. For example, in the traditional Cannon algorithm for matrix multiplication, the summation index k plays that role; blocks of the input arrays A and B are rotated among processors, and each processor holds a different block of A and B and the same block of C after each rotation step. At every step, processors multiply their local blocks of A and B , and add the result to the block of C .

Due to the symmetry of the problem, any of the 3 indices $\{i, j, k\}$ can be chosen as the rotation index, so it is always possible to keep one of the arrays in a fixed distribution and communicate (“rotate”) the other two arrays. Therefore, the number of distinct communication patterns within the generalized Cannon’s algorithm framework is given by $3 \times NI \times NJ \times NK$, where NI is defined as the number of indices in the index set I . The communication costs of the tensor con-

traction depend on the distribution choice $\{i, j, k\}$ and the choice of rotation index.

In addition to the communication of array blocks during the rotation phase of the Cannon algorithm, array re-distribution may be necessary between the Cannon steps. For instance, suppose the arrays $B(b, e, f, l)$ and $D(c, d, e, l)$ have initial distributions $\langle b, f \rangle$ and $\langle e, c \rangle$ respectively. If we want $T1$ to have the distribution $\langle b, c \rangle$ when evaluating $T1(b, c, d, f) = \sum_{e,l} B(b, e, f, l) \times D(c, d, e, l)$, B would have, for example, to be re-distributed from $\langle b, f \rangle$ to $\langle b, e \rangle$ for the generalized Cannon algorithm to be possible. But since the initial distribution $\langle e, c \rangle$ of $D(c, d, e, l)$ is the same as the distribution required to perform the Cannon rotations, no re-distribution is necessary for array D .

The partitioning of data arrays among the processors and the fusions of loops both affect the total inter-processor communication cost. Fusion generally results in an increase of communication cost, but can significantly reduce the per-processor memory requirement. We use a dynamic programming algorithm to search among all combinations of loop fusions and array distributions to find the one with minimal total communication cost, that also fits within the available memory. We omit details here and refer the reader to [5, 6].

8. Related work and discussion

The approach undertaken in this project bears similarities to some projects in other domains, such as the SPIRAL project which is aimed at the design of a system to generate efficient libraries for digital signal processing algorithms [23, 10, 26]. SPIRAL generates efficient implementations of algorithms expressed in a domain-specific language called SPL by a systematic search through the space of possible implementations.

Other efforts in automatically generating efficient implementations of programs include FFTW [8], the telescoping languages project [11], ATLAS [25] for deriving efficient implementation of BLAS routines, the PHiPAC [1] project, and the TUNE project [24]. In addition, motivated by the difficulty of detecting and optimizing matrix operations hidden in array subscript expressions within loop nests, several projects have worked on efficient code generation from high-level languages such as MATLAB and Maple [7, 20, 21, 22].

While our effort shares some common goals with several of the projects mentioned above, there are also significant differences. Some of the optimizations we consider, such as the algebraic optimizations, memory minimization, and space-time trade-offs, do not appear to have been previously explored, to the best of our knowledge. We also take advantage of certain domain-specific properties of the computations; for example, since all expressions considered

in this framework are tensor contractions, the loops of the resulting code are fully permutable, and there are no dependencies preventing fusion. This observation is crucial for the optimization algorithms of several components (memory minimization, space-time transformation, data locality). Also, some of the multi-dimensional arrays involved in the computation have certain domain-specific symmetry properties that can be exploited in order to lower the number of arithmetic operations, and thus total execution time.

While optimization of performance is a significant goal, more important in our context is the potential for dramatically reducing the developmental effort required of a quantum chemist to develop a new ab initio computational model. Currently, the manual development and testing of a reasonably efficient parallel code for a computational model such as the coupled cluster model typically takes several weeks to many months for a computational chemist. We aim to reduce the time to prototype a new model to under a day, through use of the synthesis system being developed.

References

- [1] J. Bilmes, K. Asanovic, C. Chin, and J. Demmel. Optimizing matrix multiply using PHiPAC. In *Proc. ACM Intl. Conf. on Supercomputing*, pages 340–347, 1997.
- [2] D. Cociorva, J. Wilkins, C. Lam, G. Baumgartner, P. Sadayappan, J. Ramanujam. Loop Optimizations for a Class of Memory-Constrained Computations. In *Proc. 15th ACM Intl. Conf. on Supercomputing*, 2001.
- [3] D. Cociorva, J. Wilkins, G. Baumgartner, P. Sadayappan, J. Ramanujam, M. Nooijen, D. Bernholdt, and R. Harrison. Towards Automatic Synthesis of High-Performance Codes for Electronic Structure Calculations: Data Locality Optimization. *Proc. of the Intl. Conf. on High Performance Computing*, Lecture Notes in Computer Science, Vol. 2228, pp. 237–248, Springer-Verlag, 2001.
- [4] D. Cociorva, J. Wilkins, G. Baumgartner, P. Sadayappan, J. Ramanujam, M. Nooijen, D. Bernholdt, and R. Harrison. Space-Time Trade-Off Optimization for a Class of Electronic Structure Calculations. *Proc. of the 2002 Conference on Programming Language Design and Implementation*, Berlin, Germany, June 2002.
- [5] D. Cociorva, G. Baumgartner, C.-C Lam, P. Sadayappan, and J. Ramanujam. Memory-Constrained Communication Minimization for a Class of Array Computations. *Languages and Compilers for Parallel Computing*, 2002.
- [6] D. Cociorva, X. Gao, S. Krishnan, G. Baumgartner, C.-C Lam, P. Sadayappan, and J. Ramanujam. Global Communication Optimization for Tensor Contraction Expressions under Memory Constraints. *Tech. Report, Dept. of Computer & Information Science, Ohio State University*, 2002.
- [7] L. De Rose and D. Padua. A MATLAB to Fortran 90 translator and its effectiveness. In *Proc. 10th ACM Intl. Conf. on Supercomputing*, 1996.
- [8] M. Frigo and S. Johnson. FFTW: An adaptive software architecture for the FFT. In *Proc. ICASSP 98*, Volume 3, pages 1381–1384, 1998, <http://www.fftw.org>.
- [9] High Performance Computational Chemistry Group. NWChem, A computational chemistry package for parallel computers, Version 3.3, 1999. Pacific Northwest National Laboratory, Richland, WA 99352.

- [10] J. Johnson, R. Johnson, D. Padua, and J. Xiong. Searching for the best FFT formulas with the SPL compiler. In *Languages and Compilers for High-Performance Computing*, Springer-Verlag, 2001.
- [11] K. Kennedy, B. Broo, K. Cooper, J. Dongarra, R. Fowler, D. Gannon, L. Johnsson, J. Mellor-Crummey, and L. Torczon. Telescoping Languages: A Strategy for Automatic Generation of Scientific Problem-Solving Systems from Annotated Libraries. *J. Parallel and Distributed Computing*, 2001.
- [12] For a brief description, see V. Kumar, A. Grama, A. Gupta, and G. Karypis, *Introduction to Parallel Computing*, The Benjamin/Cummings Publishing Company, pp. 171, 1994.
- [13] C.-C Lam. *Performance Optimization of a Class of Loops Implementing Multi-Dimensional Integrals*. PhD thesis, The Ohio State University, 1999.
- [14] C.-C Lam, D. Cociorva, G. Baumgartner, and P. Sadayappan. Memory-optimal evaluation of expression trees involving large objects. In *Intl. Conf. on High Performance Computing*, 1999.
- [15] C.-C Lam, D. Cociorva, G. Baumgartner, and P. Sadayappan. Optimization of memory usage for a class of loops implementing multi-dimensional integrals. In *Languages and Compilers for Parallel Computing*, 1999.
- [16] C.-C Lam, P. Sadayappan, and R. Wenger. On optimizing a class of multi-dimensional loops with reductions for parallel execution. *Parallel Processing Letters*, 7(2):157–168, 1997.
- [17] C.-C Lam, P. Sadayappan, and R. Wenger. Optimization of a class of multi-dimensional integrals on parallel machines. In *8th SIAM Conf. on Parallel Processing for Scientific Computing*, 1997.
- [18] T. J. Lee and G. E. Scuseria. Achieving chemical accuracy with coupled cluster theory. In S. R. Langhoff (Ed.), *Quantum Mechanical Electronic Structure Calculations with Chemical Accuracy*, pages 47–109, Kluwer Academic, 1997.
- [19] J. M. L. Martin. In P. v. R. Schleyer, P. R. Schreiner, N. L. Allinger, T. Clark, J. Gasteiger, P. Kollman, H. F. Schaefer III (Eds.), *Encyclopedia of Computational Chemistry*. Wiley & Sons, Berne (Switzerland). Vol. 1, pp. 115–128, 1998.
- [20] The Match Project. A MATLAB compilation environment for distributed heterogeneous adaptive computing systems. www.ece.nwu.edu/cpdc/Match/Match.html.
- [21] Vijay Menon, Keshav Pingali. High-Level Semantic Optimization of Numerical Codes. In *Proc. ACM Intl. Conf. on Supercomputing*, 1999.
- [22] Vijay Menon, Keshav Pingali. A Case for Source-Level Transformations in MATLAB. In *Proc. 2nd Conf. on Domain-Specific Languages*, 1999.
- [23] J. Moura, J. Johnson, R. Johnson, D. Padua, V. Prasanna, M. Puschel, and M. Veloso. SPIRAL: Portable library of optimized signal processing algorithms, 1998. <http://www.ece.cmu.edu/~spiral>.
- [24] M. Thottethodi, S. Chatterjee, and A. Lebeck. Tuning Strassen's matrix multiplication for memory hierarchies. In *Proc. SuperComputing '98*.
- [25] R. Whaley and J. Dongarra. Automatically Tuned Linear Algebra Software (ATLAS). In *Proc. SC '98 (Electronic Publication)*, IEEE Publication, 1998.
- [26] J. Xiong, D. Padua, and J. Johnson. SPL: A language and compiler for DSP algorithms. In *Proc. ACM SIGPLAN Conf. on Programming Language Design and Implementation*, 2001.