

Efficient Search-Space Pruning for Integrated Fusion and Tiling Transformations

Xiaoyang Gao¹, Sriram Krishnamoorthy¹, Swarup Kumar Sahoo¹, Chi-Chung Lam¹,
Gerald Baumgartner², J. Ramanujam³, and P. Sadayappan¹

¹ Department of Computer Science and Engineering
The Ohio State University, Columbus, OH 43210, USA
{gaox, krishnsr, sahoos, clam, saday}@cse.ohio-state.edu

² Department of Computer Science
Louisiana State University, Baton Rouge, LA 70803, USA
gb@csc.lsu.edu

³ Department of Electrical and Computer Engineering
Louisiana State University, Baton Rouge, LA 70803, USA
jxr@ece.lsu.edu

Abstract. Compile-time optimizations involve a number of transformations such as loop permutation, fusion, tiling, array contraction, etc. Determination of the choice of these transformations that minimizes the execution time is a challenging task. We address this problem in the context of tensor contraction expressions involving arrays too large to fit in main memory. Domain-specific features of the computation are exploited to develop an integrated framework that facilitates the exploration of the entire search space of optimizations. In this paper, we discuss the exploration of the space of loop fusion and tiling transformations in order to minimize the disk I/O cost. These two transformations are integrated and pruning strategies are presented that significantly reduce the number of loop structures to be evaluated for subsequent transformations. The evaluation of the framework using representative contraction expressions from quantum chemistry shows a dramatic reduction in the size of the search space using the strategies presented.

1 Introduction

Optimizing compilers incorporate a number of loop transformations such as permutation, tiling, fusion, etc. Considerable work has addressed loop tiling for enhancement of data locality [4, 8, 12, 17, 21, 22, 25–28]. Much work has also been done on improving locality and/or parallelism by loop fusion [10, 11, 13, 14, 24]. Fusion often creates imperfectly nested loops, which are more complex to tile effectively than perfectly nested loops. Several works have addressed the tiling of imperfectly nested loops [2, 25]. Although there has been much progress in developing unified frameworks for modeling a variety of loop transformations [1, 2, 19, 28], their use has so far been restricted to optimization of indirect performance metrics such as reuse distance, degree of parallelism, etc.

The development of model-driven optimization strategies that target direct performance metrics, remains a difficult task. In this paper, we address the problem in the specific domain of tensor contractions (generalized matrix products) involving tensors too large to fit into physical memory. We use special properties of the computations in this domain to integrate the various transformations and investigate pruning strategies to reduce the search space to be explored.

The large sizes of the tensors involved require the development of *out-of-core* implementations that orchestrate the movement of data between disk and main memory. We

have developed an integrated approach to determine the appropriate combination of loop permutation, fusion, and tiling, and the position and ordering of I/O placements. In this paper, we discuss the integration of loop fusion and tiling transformations with the objective of minimizing disk I/O cost. We first evaluate the set of all fusions to be explored. For each fusion structure, all loop permutations and I/O placements would be evaluated. A generalized tiling approach is presented that significantly reduces the number of loop structures to be explored. It also enables subsequent optimizations of I/O placements and loop permutations. This approach enables an exploration of the entire search space using a realistic performance model, without the need to resort to heuristics and search of a limited subspace of the search space to limit search time.

The rest of this paper is organized as follows. In the next section, we elaborate on the computational context of interest and introduce some preliminary concepts. An overview of the program synthesis system, of which the presented framework is a part, is given in Section 3. Section 4 describes a tree partitioning algorithm. In Section 5, we propose a loop structure enumeration algorithm and prove its completeness. The reductions in the space of loop structures to be explored is shown for representative computations in Section 6. Conclusions are provided in Section 7.

2 Computational Context

The work presented in this paper is being developed in the context of the Tensor Contraction Engine (TCE) program synthesis tool [3, 5–7, 16]. The TCE targets a class of electronic structure calculations which involve many computationally intensive components expressed as tensor contraction expressions.

The TCE takes as input a high-level specification of a computation expressed as a set of tensor contraction expressions, and transforms it into efficient parallel code. The current prototype of the TCE incorporates several compile-time optimizations which are treated in a decoupled manner, with the transformations being performed in a pre-determined sequence. In [15], we presented an integrated approach to determine the tile sizes and I/O placements for a fixed structure of the computational loops after fusion and permutation. Techniques to prune the search space of possible I/O placements, orderings, loop permutations and tiling for given a choice of fusion of tensor contractions were presented in [23]. In this paper, we present a technique to enumerate the various fusion structures and develop an algorithm to significantly reduce the number of loop nests to be evaluated for each fusion structure.

In the class of computations considered, the final result to be computed can be expressed using a collection of multi-dimensional summations of the product of several input arrays. As an example, we consider a transformation often used in quantum chemistry codes to transform a set of two-electron integrals from an atomic orbital (AO) basis to a molecular orbital (MO) basis:

$$B(a, b, c, d) = \sum_{p, q, r, s} C1(d, s) \times C2(c, r) \times C3(b, q) \times C4(a, p) \times A(p, q, r, s)$$

Here, all arrays would be initially stored on disk. The indices p , q , r , and s have the same range N . The indices a , b , c , and d have the same range V . Typical values for N range from 60 to 1300; the value for V is usually between 50 and 1000.

The calculation of B is done in four steps to reduce the number of floating point operations,

$$T1(a, q, r, s) = \sum_p C4(a, p) \times A(p, q, r, s)$$

$$\begin{aligned}
T2(a,b,r,s) &= \sum_q C3(b,q) \times T1(a,q,r,s) \\
T3(a,b,c,s) &= \sum_r C2(c,r) \times T2(a,b,r,s) \\
B(a,b,c,d) &= \sum_s C1(d,s) \times T3(a,b,c,s)
\end{aligned}$$

The sequence of contractions in this form can be represented by an operation tree, as shown in Fig. 1(a). The leaves correspond to the input arrays and the root of the operation tree to the output array. The interior nodes, intermediate or output arrays, are produced by the tensor contraction of their immediate children. The edges in the operation tree represent the *producer-consumer* relationship between the different tensor contraction expressions. Note that the operation tree is a binary tree in which each node has either zero or two children.

Assuming that the available memory limit on the machine running this calculation is less than V^4 (which is 3TB for $V = 800$), any of the logical arrays A , $T1$, $T2$, $T3$, and B is too large to entirely fit in memory. Therefore, if the computation is implemented as a succession of four independent steps, the intermediates $T1$, $T2$, and $T3$ have to be written to disk after they are produced, and read from disk before they are used in the next step. Furthermore, the amount of disk access volume could be much larger than the total volume of the data on disk. Since none of these arrays can be fully stored in memory, it may not be possible to read each element only once from disk.

Suitable fusion of the common loops involved in the contractions that produce and consume an intermediate can reduce the size of the intermediate array, making it feasible to retain it in memory. Henceforth, the term intermediate node will be used to refer to both the intermediate array produced in the corresponding interior node of the operation tree, and the contraction that produces it. The reference shall be clear from the context.

Given a choice of fusion, an intermediate node not fused with its parent divides the operation tree into two parts, both of which can be evaluated independently. Such an intermediate node which is not fused, is said to be a *cut-point* in the operation tree. A cut-point node is assumed to be written to disk on production and read back during its consumption. A connected operation tree without any interior cut-points is called a *fused sub-tree*. The divided operation tree for the four-index transform corresponding to $T1$ being a cut-point is shown in Fig. 2(a). The cut-point divides the operation tree into two fused sub-trees, one of which produces $T1$, and the other which consumes it.

The *loop nesting tree* (LNT) represents the loop structure corresponding to a fused sub-tree. Each node in a LNT is labeled by the indices of a set of fully permutable loops that appear together at some level in the overall imperfectly nested loop structure that results after applying loop fusion to the contraction computations in the sub-tree. The leaves represent the innermost loops, while the root represents the outermost loops. Fig. 2(b) shows two possible LNT's for the two fused subtrees in Fig. 2(a), respectively. The corresponding code structure is shown in Fig. 2(c).

3 Integrated Framework

The program synthesis system takes an operation tree representing a set of tensor contractions as input, and generates an efficient loop structure with explicit disk I/O statements to implement the computation. The optimization process may be viewed in terms of the following steps:

1. **Operation Tree Partitioning:** In this step, we divide the original operation tree into several fused subtrees by identifying cut-points. The optimal loop structures for the subtrees are independent of each other, and are determined separately.

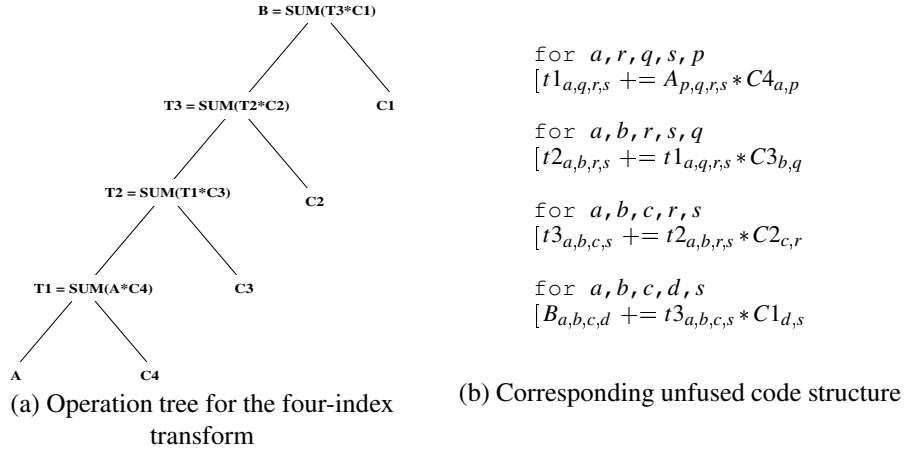


Fig. 1. Operation tree and unfused code structure for the four-index transform

2. Loop Structures Enumeration: For each fused subtree, we enumerate candidate loop structures to be evaluated, as a set of LNT's.
3. Intra-Tile Loop Placements: For a given LNT, we tile all loops at each node and propagate intra-tile loops to all the nodes below it.
4. Disk I/O Placements and Orderings: We then explore various possible placements and orderings of disk I/O statements for each disk array in a tiled loop structure with a pruning strategy to determine the best placement and ordering.
5. Tile Size Selection: For each combination of loop transformations and I/O placements, the I/O cost is formulated as a non-linear optimization problem in terms of the tile sizes. The tile sizes that minimize the disk I/O cost are determined using a general-purpose non-linear optimization solver.
6. Code Generation: We calculate the disk access cost for each solution obtained, and generate code for the one with the minimal disk I/O cost.

Algorithm 1 shows the procedure to evaluate the optimal loop structure with I/O placements for a sub-tree of the operation tree rooted at a given node t . The fused subtree rooted at that node is called a top subtree. Let $t.FS$ denote the optimal loop structure of the sub-tree rooted at node t , which includes three fields: TCS , FFT and $Cost$. TCS represents the set of cut-points that are leaves of the top subtree; FFT represents the fused loop structure of the top subtree; and $Cost$ represents the disk cost incurred by the optimal loop structure.

The top subtrees are first enumerated. This is explained in Section 4. Each of these sub-trees is evaluated in turn to determine the optimal loop structure. Given a fused subtree, its initial cost is evaluated to be sum of the costs of its cot-point nodes. The optimal loop structure and the corresponding cost for a given fused sub-tree are determined by first enumerating all candidate loop nesting trees and determining the disk I/O placements, orderings, loop permutations, and tile sizes that minimize the disk I/O cost. The enumeration of candidate loop nesting trees is discussed in Section 5. The search space of disk I/O placements and orderings, loop permutations, and tile sizes is pruned and modeled as a non-linear optimization problem, which is then solved to determine the disk I/O cost. This process is encapsulated in the procedure $dataLocality()$.

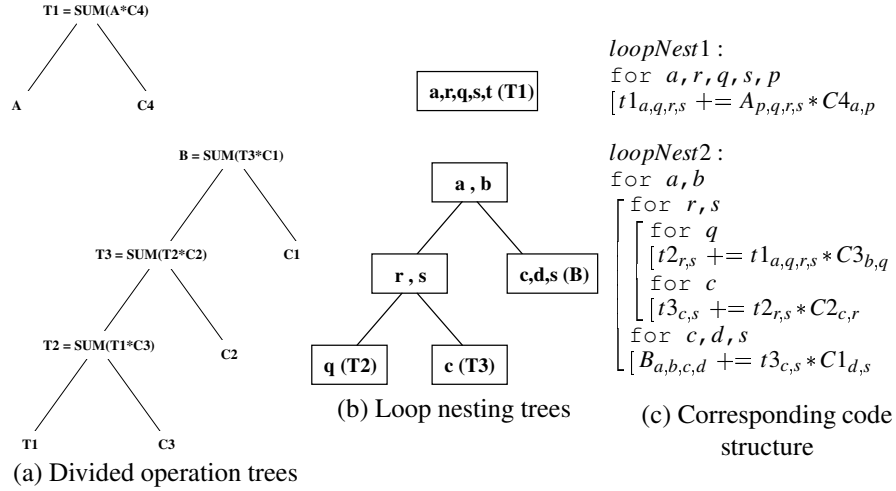


Fig. 2. Representations involved in generation of a fused code structure.

In this paper, we focus on determination of the fused sub-trees and the enumeration of candidate loop nesting trees to be evaluated. Details on subsequent steps can be found in [23].

4 Tree Partitioning

In this section, we discuss the procedure to enumerate the set of all fused subtrees to be evaluated. In general, fusing a loop between the producer of an intermediate array and its consumer eliminates the corresponding dimension of the array and reduces the array size. If the array fits in memory after fusion, no disk I/O is required for that array. On the other hand, if the array does not fit in the physical memory even after fusion, the disk I/O cost is not reduced, and the fusion does not result in any improvement.

Therefore, fusion of any loops corresponding to an intermediate node is assumed to cause the resulting intermediate to reside in memory. It potentially resides in disk if the intermediate node is not fused with its parent.

For an arbitrary operation tree with M intermediate nodes, it has at most $O(2^M)$ possible fused sub-trees, but not all of them can be fused. Consider an intermediate node t . If both its children are fused with it, then the loops corresponding to the summation indices in the given node must be the outermost loops; and it can not be fused with its parent anymore. Thus, either t or one of its children must be a cut-point.

Based on this property, we can restrict the number of top subtrees to $O(M^2)$. The algorithm to enumerate the fused sub-trees rooted at a given node is shown in Algorithm 2. It proceeds in a bottom-up fashion, constructing the fused sub-trees rooted at a given node from those of its children. A node consuming the arrays produced by its children extends the fused sub-trees from each of its children. These sub-trees can further be extended to include the given node's parent. In turn, these sub-trees form a "chain" starting from the given node and terminating at a cut-point. In addition, the given node can be considered as a cut-point. In this scenario, all possible pairs of left and right fused sub-trees form a valid fused sub-tree for the given node.

The field $t.TreeSet$ represents the set of fused sub-trees which can be extended to include the parent of t .

Algorithm 1 SearchOptimalLoopStructure(t : the root of a subtree)

//Given a subtree rooted at t , the algorithm will find the optimal loop structure with disk I/O

```
TreeSet = EnumerateFusedSubtrees( $t$ )
for each subtree  $T_i$  in TreeSet do
  TCS =  $T_i$ .CutpointSet
  LeafCost = 0
  for each cut-point  $ct$  in TCS do
    LeafCost = LeafCost +  $ct$ .FS.Cost
  end for
  //Enumerate all loop structures of fused subtree  $T_i$ 
  LoopSet = EnumerateLoop( $T_i$ )
  OptCost =  $\infty$ 
  //Compute the minimal disk I/O cost of subtree  $T_i$ 
  for each loop structure FFS in LoopSet do
    Cost = dataLocality(FFS)
    if Cost < OptCost then
      OptCost = Cost
      OptFFS = FFS
    end if
  end for
  Cost = OptCost + leafCost
  if Cost <  $t$ .FS.Cost or  $t$ .FS = null then
     $t$ .FS.Cost = Cost
     $t$ .FS.TCS = TCS
     $t$ .FS.FFS = OptFFS
  end if
end for
```

5 Loop Structure Enumeration

In this section, we first present an algorithm that can generate a set of loop structures of a fused subtree. Then, we prove that for any loop structure S of the fused subtree, we can find a corresponding loop structure S' in the generated set, so that S' can be transformed to S by some proper multi-level tiling strategies.

5.1 Enumeration Algorithm

In the previous section, we showed that a fused subtree must be in one of these two forms:

- All contractions form a chain. We call it a *contraction chain*. For instance, Fig. 1 is such an operation tree, in which the contraction chain is $T1, T2, T3, B$.
- The contractions form two chains joining at the root node. In this case, the *contraction chain* is connected by these two chains. An example of such an operation tree is shown in Fig. 3, in which the contraction chain is $T1, T2, B, T3, T4$

Algorithm 2 EnumerateFusedSubtrees(t : the root of a subtree) returns *TreeSet*

t_1 = the left child of t ; t_2 = the right child of t ; *TreeSet* = empty
//Only one subtree
if both t_1 and t_2 are input nodes **then**
 Create a new Tree Tr with $Tr.Cut\ point\ Set = \emptyset$
 Insert Tr into *TreeSet*
end if
//Extending subtrees from the child not an input
if t_1 is an input node and t_2 is an intermediate node **then**
 $childSet = t_2.TreeSet$
 Create a new Tree Tr with $Tr.Cut\ point\ Set = \{t_2\}$
 Insert Tr into *TreeSet*
end if
if t_2 is an input node, and t_1 is an intermediate node **then**
 $childSet = t_1.TreeSet$
 Create a new Tree Tr with $Tr.Cut\ point\ Set = \{t_1\}$
 Insert Tr into *TreeSet*
end if
for each subtree st in $childSet$ **do**
 Create a new Tree Tr with $Tr.Cut\ point\ Set = st.Cut\ point\ Set$
 Insert Tr into *TreeSet*
end for
 $t.TreeSet = TreeSet$
//Extending subtrees from either child, and cutting another child off
if both t and t_2 are intermediate nodes **then**
 $childSet1 = t_1.TreeSet$
 for each subtree st in $childSet1$ **do**
 Create a new Tree Tr with $Tr.Cut\ point\ Set = \{st.Cut\ point\ Set, t_2\}$
 Insert Tr into *TreeSet*
 end for
 $childSet2 = t_2.TreeSet$
 for each subtree st in $childSet2$ **do**
 Create a new Tree Tr with $Tr.Cut\ point\ Set = \{st.Cut\ point\ Set, t_1\}$
 Insert Tr into *TreeSet*
 end for
 Create a new Tree Tr with $Tr.Cut\ point\ Set = \{t_1, t_2\}$
 Insert Tr into *TreeSet*
 $t.TreeSet = TreeSet$
 //Merging subtrees from both children, and extending the result
 for each pair of subtrees $st1$ in $childSet1$ and $st2$ in $childSet2$ **do**
 Create a new Tree Tr
 $Tr.Cut\ point\ Set = \{st1.Cut\ point\ Set, st2.Cut\ point\ Set\}$
 Insert Tr into *TreeSet*
 end for
end if

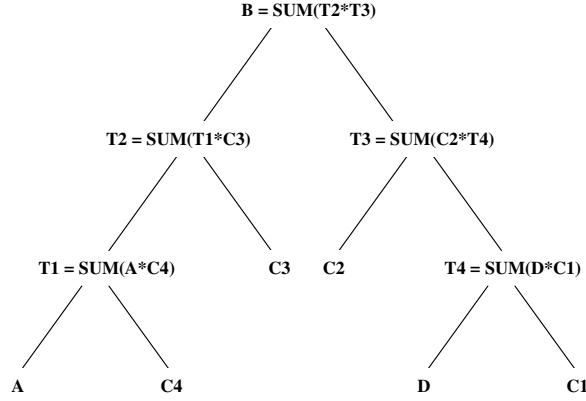


Fig. 3. An operation tree with two chains

Given an operation tree that has n contraction nodes t_1, t_2, \dots, t_n , let t_i .indices denote all loop indices surrounding the contraction node t_i . First, we create a contraction chain of the operation tree. It corresponds to a sequence of perfectly nested loops. Many different choices exist in the ordering of the fusions within this sequence of perfectly nested loop nests. Each of the perfectly nested loops, corresponding to a contraction, can be considered an independent loop nesting tree. The fusion of subtrees producing and consuming an intermediate creates an imperfectly nested loop nests, in which some of the common loops are merged. The process of construction of the loop nesting tree of the fused subtree can be modeled as a parenthesization problem. Consider the sequence of contraction nodes T1, T2, T3, and B in the operation tree shown in Fig. 1. $((T1(T2 T3))B)$ corresponds to a parenthesization in which the contractions producing T3 and consuming T3 are fused first and the resulting loop nest is fused with the contractions producing T1 and B, in that order. Fig. 4 shows one possible parenthesization for the four-index transform and the corresponding loop nesting tree.

We enumerate all possible parenthesizations of the contraction chain. For each parenthesization, a maximally fused loop structure is created by a recursive construction procedure. We call it *maximally fused* since, in the construction procedure, each intermediate node will have its indices fused as much as possible with its parent. The construction procedure is shown in Algorithm 3. It takes a parenthesization P as input, and generate a corresponding LNT. A parenthesization of a contraction chain with n nodes has $n - 1$ pairs of parentheses. Each pair of parentheses includes two elements, left and right element. Each element is either a single contraction node, or a parenthesization of a sub-chain within a pair of parentheses.

For easy understanding, we use an example to explain how the algorithm works. Consider a parenthesization $((T1(T2 T3))B)$ of four-index transform. Fig. 4 shows how the construction procedure creates the corresponding LNT step by step.

5.2 Completeness

In this section, we prove that the set of *maximally fused* loop structures generated by the enumeration algorithm above can represent all loop structures of a fused subtree. The following definitions are provided to clarify terms used in the proof.

Definition 1. Each leaf in a LNT includes a contraction node. The set of contraction nodes from all the leaves in a LNT is called *leafcontractions* of the LNT.

Algorithm 3 Construction(P)

//Given a parenthesization, the algorithm map it to a maximally fused loop structure in LNT

```
 $l = P.left$ 
 $r = P.right$ 
if  $l$  is a parenthesization then
   $lt = \text{Construction}(left)$ 
else if  $l$  is a contraction then
   $lt = \text{Create a new LNT node}$ 
   $lt.indices = l.indices$ 
   $lt.children = null$ 
   $lt.contraction = l$  { $lt$  is a leaf, which includes a contraction node in it}
end if
if  $r$  is a parenthesization then
   $rt = \text{Construction}(right)$ 
else if  $r$  is a contraction then
   $rt = \text{Create a new LNT node}$ 
   $rt.indices = r.indices$ 
   $rt.children = null$ 
   $rt.contraction = r$  { $rt$  is a leaf, which includes a contraction node in it}
end if
 $comindices = lt.indices \cap rt.indices$ 
 $lt.indices = lt.indices - comindices$ 
 $rt.indices = rt.indices - comindices$ 
 $lnt = \text{Create a new LNT node}$ 
 $lnt.indices = comindices$ 
 $lnt.children = \{lt, rt\}$ 
return  $lnt$ 
```

Definition 2. In a LNT, each node t has exactly one path to the root. Let $t.upperindices$ denotes the union of all indices belonging to nodes on the path from t to the root. If a subtree $slnt$ is rooted at t , we also define $slnt.upperindices$ to equal to $t.upperindices$.

Definition 3. In a LNT, suppose two leaves t_i and t_j belong to one subtree $slnt$. If there is no other subtree that contains both t_i and t_j and is a subtree of $slnt$, then we say that $slnt$ is the *minimal common subtree* of t_i and t_j , denoted as $MCS(t_i, t_j)$.

Given an arbitrary loop nesting tree lnt , we can map it to a maximal fused loop nesting tree lnt' , which is in the set of *maximally fused* loop structures generated by the enumeration algorithm above, and can be translated to lnt with some proper multi-level tiling strategy. The mapping algorithm consists of two steps:

1. Take lnt as input, and generate a parenthesization P of the contraction chain using the generation routine provided in Algorithm 4.
2. Apply the construction procedure in Algorithm 3 on P to generate a maximally fused loop structure lnt' .

Parenthesization	LNT
(T2 T3)	<pre> graph TD A["a,b,r,s"] --- B["q (T2)"] A --- C["c (T3)"] </pre>
(T1 (T2 T3))	<pre> graph TD A["a,r,s"] --- B["p,q (T1)"] A --- C["b"] C --- D["q (T2)"] C --- E["c (T3)"] </pre>
((T1 (T2 T3)) B)	<pre> graph TD A["a,s"] --- B["r"] A --- C["b,c,d (B)"] B --- D["p,q (T1)"] B --- E["b"] E --- F["q (T2)"] E --- G["c (T3)"] </pre>

Fig. 4. Construction of a maximally fused loop structure for a particular parenthesization of the four-index transform.

Obviously, lnt' is the set of *maximally fused* loop structures generated by the enumeration algorithm. Afterward, we show that lnt' can be translated to lnt by sinking indices at upper levels down.

Lemma 1. For any pair of contraction nodes t_i and t_j , let $common(lnt, t_i, t_j)$ be the loops shared by t_i and t_j in lnt . We have $common(lnt, t_i, t_j) \subseteq common(lnt', t_i, t_j)$.

Proof. Given a subtree $slnt$, $slnt.upperindices$ represents all common loops shared by $slnt.leafcontractions$.

There is an interesting property of *maximally fused* loop structures in the way they are constructed. For any subtree $slnt$ in the LNT of a *maximally fused* loop structure, $slnt.upperindices$ includes all common loops among $slnt.leafcontractions$. In other words, it includes all possibly shared loops among $slnt.leafcontractions$. In addition, from the mapping method, we can see that if lnt has a subtree $slnt$, then there exist a twin subtree $slnt'$ in lnt' , which satisfied the following conditions:

$$\begin{aligned}
 slnt.leafcontractions &= slnt'.leafcontractions \\
 slnt.upperindices &\subseteq slnt'.upperindices
 \end{aligned}$$

Algorithm 4 Parenthesize(lnt)

//Given an LNT, the algorithm map it to a corresponding parenthesization

```
if  $lnt.children \neq null$  then
   $P = null$ 
  for each child  $c$  in  $lnt.children$  do
     $P' = \text{Parenthesize}(c)$ 
    if  $P$  is null then
       $P = P'$ 
    else
       $P = \text{new Parenthesization}(P, P')$ 
    end if
  end for
else
   $P = c.contraction$  { $c$  is a leaf and includes a contraction node}
end if
return  $P$ 
```

Given any pair of leaf nodes t_i and t_j , we define $mlnt = MCS(t_i, t_j)$ in lnt , where $mlnt.upperindices = common(lnt, t_i, t_j)$. Hence, we can find the corresponding subtree $mlnt'$ in lnt' , where

$$mlnt.upperindices \subseteq mlnt'.upperindices \subseteq common(lnt, t_i, t_j)$$

Thus, we have $common(lnt, t_i, t_j) \subseteq common(lnt', t_i, t_j)$. \square

Lemma 2. If $common(lnt, t_i, t_j) \subset common(lnt', t_i, t_j)$, then we can transform lnt' to form lnt'' by sinking indices down, so that $common(lnt, t_i, t_j) = common(lnt'', t_i, t_j)$

Proof. We define $mlnt$ and $mlnt'$ as $MCS(t_i, t_j)$ in lnt and lnt' respectively. Any loop in $common(lnt', t_i, t_j)$ belongs to the root or an ancestor of $mlnt'$. Assuming loop l is in the difference of $common(lnt, t_i, t_j)$ and $common(lnt', t_i, t_j)$. We remove l from the original node r , and insert it to all children of r . After that, if l still belongs to the root or an ancestor of $mlnt'$, we repeat the sinking operation described above, until l is not in $mlnt'.upperIndices$ any more. The same method is applied for all indices in the difference of $common(lnt, t_i, t_j)$ and $common(lnt', t_i, t_j)$. The new LNT is denoted as lnt'' . Then, we have $common(lnt, t_i, t_j) = common(lnt'', t_i, t_j)$. \square

Applying the sinking operation in Lemma 2 for each pair of contraction nodes (t_i, t_j) , we can transform lnt' to lnt'' , which satisfies the condition: $\forall (t_i, t_j), common(lnt, t_i, t_j) = common(lnt'', t_i, t_j)$. After that, if a node r has no indices in $r.indices$, we remove r from lnt'' , and put all children of r to its parent. Then, lnt'' is same as lnt .

Using *multi-level tiling strategy*, a maximally fused loop structure can be transformed into an arbitrarily fused loop structure by appropriate choice of tile sizes. *Multi-level tiling* can transform the LNT of a loop structure as follows. Each loop present in the root is split into two components, inter-tile loop and intra-tile loop. The intra-tile loop is placed on child nodes of the root. Then the loops present in each of the child nodes including the intra-tile loops from the root, are again split and intra-tile loops are placed on

their respective child nodes. This process is performed recursively till the leaf nodes are encountered. The loop structure corresponding to the LNT can also be transformed accordingly. Figure 5 shows the way to tile loop a in the LNT in Fig. 4 and the relationship between different tiles, where $a.range$ represents the range of loop a .

The sinking operation in LNT can be modeled as the *multi-level tiling* in the loop structure. Given a loop structure, if we tile a fused loop with a tile size equal to its loop range, it leads to the same result as we sink the loop index from original node to all its children. Let S and S' be loop structures represented by lnt and lnt' respectively. Since we can transform lnt' to lnt by sinking operations, we can also transform S' to S by proper multi-level tiling strategies. Next, we will use an example to show the details of the transformation procedure.

An arbitrary fully fused loop structure S of four-index transform is shown in Figure 6(a), and the corresponding maximally fused loop structure S' is in Figure 6(b). After we apply multi-level tiling strategies, S' is translated to the format shown in Figure 7(a). In addition, if we set ranges of inter-tile loops according to the following formulas, and remove all loops with $range = 1$, then S' can be rewritten as the format shown in Figure 7(b), which is exactly the same as S . It should be noted that the indexing of the intermediate arrays has been shown in a more generic way.

$$aT_2 = aT_3 = sT_1 = sT_2 = sT_3 = rT_2 = qT_1 = 1 \quad aT_1 = a.range \quad rI_1 = r.range$$

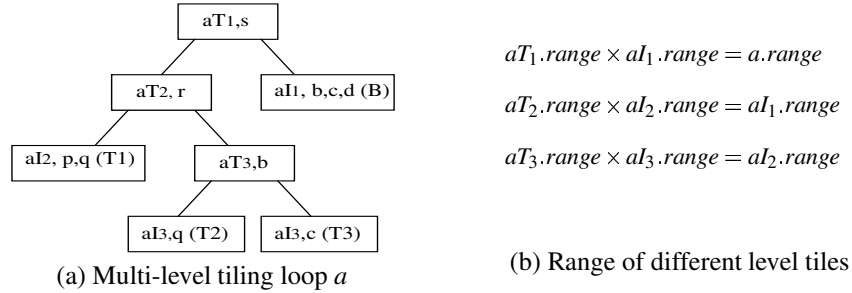


Fig. 5. An example of multi-level tiling in LNT

5.3 Complexity

The total number of loop structures generated by the enumeration algorithm is the same as the number of parenthesizations of the contraction chain. For a contraction chain with n nodes, the number of all possible parenthesizations is called the n^{th} *Catalan Number*. It is exponential in n , and the upper bound is $O(4^n/n^{3/2})$. In contrast, the number of possible loop structures is potentially exponential in the total number of distinct loop indices in the n intermediate nodes, a considerably larger number. The fused operation tree is not very long for most representative computations. In most practical applications, a fused subtree usually has no more than 5 contractions in a single chain. Note that the n^{th} Catalan Number is not very large when n is small. The first six Catalan Numbers are listed here: 1, 1, 2, 5, 14, 42....

```

for a
  for r
    for q, s, p
      [t1s,q += Ap,q,r,s * C4a,p
      for b, s, q
        [t2b,r,s += t1s,q * C3b,q
        for b, c, r, s
          [t3b,c,s += t2b,r,s * C2c,r
          for b, c, d, s
            [Ba,b,c,d += t3b,c,s * C1d,s

```

```

for a, s
  for r
    for q
      for p
        [t1 += Ap,q,r,s * C4a,p
        for b
          [t2b += t1 * C3b,q
          for b, c
            [t3b,c += t2b * C2c,r
            for b, c, d
              [Ba,b,c,d += t3b,c * C1d,s

```

(a) Arbitrary fused loop structure: S (b) Maximally fused loop structure: S'

Fig. 6. An arbitrary loop structure and the corresponding maximally fused structure

```

for aT1, sT1
  for rT1, aT2, sT2
    for qT1, rT2, aT3, sT3
      for p, ql1, rl2, al3, sl3
        [t1al,ql,rl,sl += Ap,q,r,s * C4a,p
        for b, ql1, rl2, al3, sl3
          [t2al,b,rl,sl += t1al,ql,rl,sl * C3b,q
          for b, c, rl1, al2, sl2
            [t3al,b,c,sl += t2al,b,rl,sl * C2c,r
            for al1, b, c, d, sl1
              [Ba,b,c,d += t3al,b,c,sl * C1d,s

```

```

for aT1
  for rT1
    for p, ql1, sl3
      [t1al,ql,rl,sl += Ap,q,r,s * C4a,p
      for b, ql1, sl3
        [t2al,b,rl,sl += t1al,ql,rl,sl * C3b,q
        for b, c, al2, sl2
          [t3al,b,c,sl += t2al,b,rl,sl * C2c,r
          for b, c, d, sl1
            [Ba,b,c,d += t3al,b,c,sl * C1d,s

```

(a) After inserting intra-tile loops

(b) After selecting proper tile counts

Fig. 7. Translate S' to S by multi-level tiling strategy

6 Results

The enumeration algorithm discussed in Section 5.1 generates a set of candidates loop structures to be considered for data locality optimization. Without this algorithm, and generalized tiling, the set of loop structures to be evaluated might be too large, precluding their complete evaluation and necessitating the use of heuristics.

We evaluate the effectiveness of our approach using the following tensor contractions from representative computations from the quantum chemistry domain.

1. **Four-index transform (4index):** This is the sequence of contractions introduced in Section 2.
2. **CCSD:** The second and the third computations are from the class of Coupled Cluster (CC) equations [9, 18, 20] for ab initio electronic structure modeling. The sequence of tensor contraction expressions extracted from this computation is shown as follows:

$$S(j, i, b, a) = \sum_{l,k} (A(l, k, b, a) \times (\sum_d (\sum_c (B(d, c, l, k) \times C(i, c)) \times D(j, d)))$$

3. **CCSDT:** This is a more accurate CC model. A sub-expression from the CCSDT theory is:

Table 1. Effectiveness of pruning of loop structures.

	#Contractions	#Loop structures		Reduction
		Total	Pruned	
4index	4	241	5	98%
CCSD	3	69	2	97%
CCSDT	4	182	5	98%

$$S(h3, h4, p1, p2) = \sum_{p9, h6, h8} (y_{\text{ooovv}}(h8, h6, h4, p9, p1, p2) \times \sum_{h10} (t_{\text{vo}}(p9, h10) \times \sum_{p7} (t_{\text{vo}}(p7, h8) \times \sum_{p5} (t_{\text{vo}}(p5, h6) \times v_{\text{ovv}}(h10, h3, p7, p5))))))$$

We evaluated the fused subtree corresponding to the entire operation tree without any cut-points. The number of all possible loop structures and the number of candidate loop structures enumerated by our approach are shown in Table 1. It can be seen that a very large fraction of the set of possible loop structures, up to 98%, is pruned away using the approach developed in this paper.

7 Conclusions

In this paper we addressed the problem of optimizing the disk access cost of tensor contraction expressions by applying loop transformations. We discussed approaches to partitioning of the operation tree into fused sub-trees and generating a small set of “maximally-fused” loop structures that “cover” all possible imperfectly nested fused loop structures. The approach was evaluated on a set of computations representative of the targeted quantum chemistry domain and a significant reduction was demonstrated in the number of loop structures to be evaluated.

Acknowledgments This work is supported in part by the National Science Foundation through awards 0121676, 0121706, 0403342, 0508245, 0509442, and 0509467.

References

1. N. Ahmed, N. Mateev, and K. Pingali. Synthesizing transformations for locality enhancement of imperfectly nested loops. In *Proc. of ACM Intl. Conf. on Supercomputing*, 2000.
2. N. Ahmed, N. Mateev, and K. Pingali. Tiling imperfectly-nested loops nests. In *Proc. of SC 2000*, 2000.
3. G. Baumgartner, D.E. Bernholdt, D. Cociorva, R. Harrison, S. Hirata, C. Lam, M. Nooijen, R. Pitzer, J. Ramanujam, and P. Sadayappan. A High-Level Approach to Synthesis of High-Performance Codes for Quantum Chemistry. In *Proc. of SC 2002*, November 2002.
4. J. Chame and S. Moon. A tile selection algorithm for data locality and cache interference. In *Proc. of ACM Intl. Conf. on Supercomputing*, pages 492–499, 1999.
5. D. Cociorva, G. Baumgartner, C. Lam, J. Ramanujam P. Sadayappan, M. Nooijen, D. Bernholdt, and R. Harrison. Space-Time Trade-Off Optimization for a Class of Electronic Structure Calculations. In *Proc. of ACM SIGPLAN PLDI 2002*, pages 177–186, 2002.
6. D. Cociorva, X. Gao, S. Krishnan, G. Baumgartner, C. Lam, P. Sadayappan, and J. Ramanujam. Global Communication Optimization for Tensor Contraction Expressions under Memory Constraints. In *Proc. of IPDPS*, 2003.

7. D. Cociorva, J. Wilkins, G. Baumgartner, P. Sadayappan, J. Ramanujam, M. Nooijen, D. E. Bernholdt, and R. Harrison. Towards Automatic Synthesis of High-Performance Codes for Electronic Structure Calculations: Data Locality Optimization. In *Proc. of the Intl. Conf. on High Performance Computing*, volume 2228, pages 237–248. Springer-Verlag, 2001.
8. S. Coleman and K. S. McKinley. Tile Size Selection Using Cache Organization and Data Layout. In *Proc. of the SIGPLAN '95 Conference on Programming Languages Design and Implementation*, 1995.
9. T. Crawford and H. F. Schaefer III. An Introduction to Coupled Cluster Theory for Computational Chemists. In K. Lipkowitz and D. Boyd, editor, *Reviews in Computational Chemistry*, volume 14, pages 33–136. John Wiley, 2000.
10. C. Ding and K. Kennedy. Improving effective bandwidth through compiler enhancement of global cache reuse. *J. Parallel Distrib. Comput.*, 64(1):108–134, 2004.
11. G. Gao, R. Olsen, V. Sarkar, and R. Thekkath. Collective Loop Fusion for Array Contraction. In *Proc. of the Fifth LCPC Workshop*, 1992.
12. S. Ghosh, M. Martonosi, and S. Malik. Precise Miss Analysis for Program Transformations with Caches of Arbitrary Associativity. In *Proc. of the Eighth ACM Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*, 1998.
13. K. Kennedy. Fast greedy weighted fusion. In *Proc. of ACM Intl. Conf. on Supercomputing*, 2000.
14. K. Kennedy and K. S. McKinley. Maximizing loop parallelism and improving data locality via loop fusion and distribution. In *Proc. of Languages and Compilers for Parallel Computing*, pages 301–320. Springer-Verlag, 1993.
15. S. Krishnan, S. Krishnamoorthy, G. Baumgartner, C. Lam, J. Ramanujam, P. Sadayappan, and V. Choppella. Efficient synthesis of out-of-core algorithms using a nonlinear optimization solver. In *Proc. of IPDPS*, page 34b, 2004.
16. C. Lam. *Performance Optimization of a Class of Loops Implementing Multi-Dimensional Integrals*. PhD thesis, The Ohio State University, Columbus, OH, August 1999.
17. M. S. Lam, E. E. Rothberg, and M. E. Wolf. The cache performance and optimizations of blocked algorithms. In *Proc. of Fourth Intl. Conf. on Architectural Support for Programming Languages and Operating Systems*, 1991.
18. T. J. Lee and G. E. Scuseria. Achieving chemical accuracy with coupled cluster theory. In S. R. Langhoff, editor, *Quantum Mechanical Electronic Structure Calculations with Chemical Accuracy*, pages 47–109. Kluwer Academic, 1997.
19. A. W. Lim and M. S. Lam. Maximizing Parallelism and Minimizing Synchronization with Affine Partitions. *Parallel Computing*, 24(3-4):445–475, May 1998.
20. J. M. L. Martin. Benchmark Studies on Small Molecules. In P. v. R. Schleyer, P. R. Schreiner, N. L. Allinger, T. Clark, J. Gasteiger, P. Kollman, and H. F. Schaefer III, editors, *Encyclopedia of Computational Chemistry*, volume 4, pages 115–128. John Wiley, 1998.
21. G. Rivera and C.-W. Tseng. A Comparison of Compiler Tiling Algorithms. In *CC '99: Proc. 8th Intl. Conf. Compiler Construction*, pages 168–182. Springer-Verlag, 1999.
22. G. Rivera and C.-W. Tseng. Tiling optimizations for 3D scientific computations. In *Supercomputing '00: Proc. 2000 ACM/IEEE conference on Supercomputing (CDROM)*, 2000.
23. S. K. Sahoo, S. Krishnamoorthy, R. Panuganti, and P. Sadayappan. Integrated loop optimizations for data locality enhancement of tensor contraction expressions. In *Proc. of Supercomputing (SC 2005)*, 2005.
24. S. Singhai and K. S. McKinley. Loop Fusion for Parallelism and Locality. In *Proc. of Mid-Atlantic States Student Workshop on Programming Languages and Systems*, 1996.
25. Y. Song and Z. Li. New Tiling Techniques to Improve Cache Temporal Locality. In *Proc. of ACM SIGPLAN PLDI*, 1999.
26. M. E. Wolf and M. S. Lam. A Data Locality Algorithm. In *Proc. of ACM SIGPLAN PLDI*, 1991.
27. M. E. Wolf, D. E. Maydan, and D. J. Chen. Combining loop transformations considering caches and scheduling. In *Proc. of the Twenty Ninth Annual International Symposium on Microarchitecture*, pages 274–286, 1996.
28. Q. Yi, V. Adve, and K. Kennedy. Transforming loops to recursion for multi-level memory hierarchies. In *PLDI '00: Proceedings of the ACM SIGPLAN 2000 conference on Programming language design and implementation*, pages 169–181, 2000.