

Shared-Memory Parallelization of MTTKRP for Dense

Koby Hayashi, Grey Ballard, Yujie Jiang, Michael J. Tobia

Wake Forest University

P.O. Box 1234

Winston Salem, NC 27109

{hayakb13,ballard,jiany14,tobiamj}@wfu.edu

ABSTRACT

The matricized-tensor times Khatri-Rao product (MTTKRP) is the computational bottleneck for algorithms computing CP decompositions of tensors. In this paper, we develop shared-memory parallel algorithms for MTTKRP involving dense tensors. The algorithms cast nearly all of the computation as matrix operations in order to use optimized BLAS subroutines, and they avoid reordering tensor entries in memory. We benchmark sequential and parallel performance of our implementations, demonstrating high sequential performance and efficient parallel scaling. We use our parallel implementation to compute a CP decomposition of a neuroimaging data set and achieve a speedup of up to 7.4 \times over existing parallel software.

ACM Reference format:

Koby Hayashi, Grey Ballard, Yujie Jiang, Michael J. Tobia. 2018. Shared-Memory Parallelization of MTTKRP for Dense. In *Proceedings of*, , , 12 pages.

DOI: 10.1145/nnnnnnnn.nnnnnnnn

1 INTRODUCTION

Tensor decompositions provide a means of data analysis for multi-dimensional data. In particular, the CP decomposition is a generalization of the matrix singular value decomposition (or principal component analysis), providing a low-rank approximation of data. This model representation of the data can be used in applications such as blind source separation (interpreting each component as a source signal [8]), in anomaly detection (identifying data points that are not explained by the model [24]), and for predicting missing or future data [1]. Interest in tensor analysis and the use of the CP decomposition has been growing recently; we refer the reader to survey papers for a more exhaustive list of references [3, 13, 20].

In addition to the growing interest, the increasing size of today's data sets has brought a higher demand for high-performance implementations of the fundamental computational kernels. For example, nearly all of the time computing CP decompositions occurs in an operation known as matricized-tensor times Khatri-Rao product (MTTKRP). Most of the available tensor analysis software packages [7, 26] are written in Matlab, yielding limitations on performance and utility of multicore and other high-performance architectures. While there have been many recent developments in efficient software for sparse tensor decompositions [15, 22], there remain few options in the case of dense tensors, which is the subject of this work.

Our motivating application is a neuroimaging data analysis problem involving functional MRI (fMRI) data. We are given an input tensor representing correlations between pairs of regions of interest in the brain over time and for various human subjects. We wish to extract functional brain networks to study how they behave over time relative to a cognitive task and how they relate to and differentiate among subjects. We discuss this and related problems in more detail in Section 3. Because of limitations in memory and computational time, the regions of interest are highly coarsened versions of the data. The existing approach uses the Matlab Tensor Toolbox [7], but the computational time is a bottleneck in the analysis process. In order to decrease the time and allow for analysis of larger data sets with finer granularity, our goal is to develop shared-memory parallelizations of the MTTKRP computation in order to utilize multi-core servers.

One advantage of tensor computations is that they can often be cast as matrix operations, which have been well-optimized via the BLAS interface for today's architectures. In particular, the bulk of MTTKRP corresponds to a single matrix-matrix multiplication. Unfortunately, using BLAS requires that matrices be stored in regular layouts in memory (e.g., column-major ordering), and it is impossible to choose a dense tensor data layout in memory that is conducive to direct BLAS calls in all cases. Thus, using BLAS directly requires reordering tensor entries in memory, which is usually too expensive. The main task in optimizing dense MTTKRP is to employ BLAS in a way that respects a single tensor data layout and avoids tensor reordering. We discuss MTTKRP in context of the CP decomposition, along with related work, in Section 2.

We consider two MTTKRP algorithms, which we refer to as 1-step and 2-step, that cast the computation as calls to BLAS and never reorder the tensor. The 1-step algorithm is novel for MTTKRP, using ideas from optimization of a related tensor computation [5, 14]; the 2-step algorithm was developed by Phan *et al.* [19]. We also develop a parallel algorithm for computing the Khatri-Rao product of matrices, which is needed for the 1-step and 2-step algorithms. These sequential and parallel algorithms are presented in Section 4. We benchmark the algorithms in Section 5, comparing their performance to baselines, and demonstrating high sequential performance and efficient parallel scaling on a multicore server.

To summarize, the primary contributions of this work are as follows:

- we develop a parallel row-wise algorithm for computing a Khatri-Rao product of multiple matrices;
- we implement a new 1-step and an existing 2-step MTTKRP algorithm and parallelize the algorithms using a combination of OpenMP and multithreaded BLAS;

- we demonstrate performance improvement over a baseline approach and achieve parallel speedups of up to 12× and 8× over 12 threads; and
- we obtain up to a 7.4× speedup over existing software for computing the CP decomposition of fMRI tensors.

2 BACKGROUND

2.1 Notation

Tensors will be denoted using Euler script (e.g., \mathcal{X}), matrices will be denoted using upper case bold face type (e.g., \mathbf{M}), and vectors will be denoted as lower case bold face type (e.g., \mathbf{v}). We use Matlab-style notation to index into tensors, matrices, and vectors. For example, $\mathbf{M}(:, c)$ is the c th column of matrix \mathbf{M} . Scalar integer values will not be bold-faced, and we use brackets to indicate sets of integers: $[N] = \{0, 1, \dots, N-1\}$. Note that we use 0-indexing throughout. An N -dimensional tensor will be referred to as N -way or order N . An N -way tensor is rank-1 if it can be represented by an outer product of N vectors, one vector in each mode.

We use the notation $I_0 \times \dots \times I_{N-1}$ to specify the dimensions of an N -way tensor. For shorthand, we let $I = \prod_{k \in [N]} I_k$ be the total number of entries in the tensor. We also define $I_{\neq n} = \prod_{n \neq k \in [N]} I_k$ to be the product of all modes but n , $I_n^L = \prod_{n > k \in [N]} I_k$ to be the product of all modes to the left of n , and $I_n^R = \prod_{n < k \in [N]} I_k$ to be the product of all modes to the right of n .

The *Hadamard* product, or element-wise product, is denoted by $*$. For example, $\mathbf{C} = \mathbf{A} * \mathbf{B}$ implies $\mathbf{C}(i, j) = \mathbf{A}(i, j) \cdot \mathbf{B}(i, j)$. The Kronecker product, a generalization of an outer product of vectors, is denoted by \otimes . The *Khatri-Rao* product is denoted by \odot and will be central to this work. It can be considered a column-wise Kronecker product, or it can be defined row-wise using the Hadamard product. Given an $I_A \times C$ matrix \mathbf{A} and an $I_B \times C$ matrix \mathbf{B} , the Khatri-Rao product $\mathbf{K} = \mathbf{A} \odot \mathbf{B}$ has dimension $I_A I_B \times C$ (note that the input matrices must have the same number of columns). Defined column-wise, we have $\mathbf{K}(:, c) = \mathbf{A}(:, c) \otimes \mathbf{B}(:, c)$ for $c \in [C]$. Defined row-wise, we have $\mathbf{K}(r_B + r_A I_B, :) = \mathbf{A}(r_A, :) * \mathbf{B}(r_B, :)$.

To describe how tensors are stored in memory, we define the standard linearization of tensor entries that generalizes column-major order of matrix entries. Given a tensor entry (i_0, \dots, i_{N-1}) , its index in the linearization is given by $\ell = \sum_{n \in [N]} i_n \cdot I_n^L$.

We also *matricize* or *unfold* a tensor into a matrix. A mode- n fiber of a tensor is a vector of entries that share all indices but one; for example, $\mathcal{X}(i, :, k)$ is a mode-1 fiber of \mathcal{X} . Arranging all of the mode- n fibers into the columns of a matrix, we obtain the mode- n matricization $\mathbf{X}_{(n)}$, which is an $I_n \times I_{\neq n}$ matrix. The order of the columns corresponds to a linearization of the remaining modes (excluding mode n). We also use a generalization of this concept, assigning multiple modes to the rows of the matrix and the remaining modes to the columns. In this matricization, an entry's row index corresponds to a linearization of the row modes, and the column index corresponds to a linearization of the column modes. We use the notation $\mathbf{X}_{(m:n)}$ to denote such a matricization with contiguous row modes, where modes $\{m, m+1, \dots, n\}$ are the row modes.

Finally, *tensor-times-matrix* (TTM) is denoted by \times_n for mode n and is defined such that $\mathcal{Y} = \mathcal{X} \times_n \mathbf{M}$ is equivalent to $\mathbf{Y}_{(n)} = \mathbf{M}^T \mathbf{X}_{(n)}$.

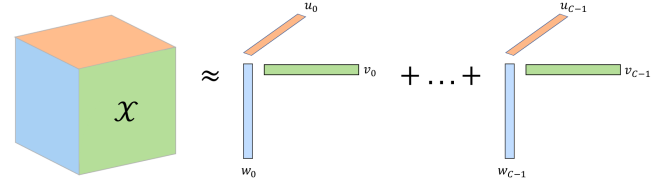


Figure 1: Rank- C CP decomposition of a 3-way tensor.

When \mathbf{M} is a column vector, we refer to the operation as tensor-times-vector (TTV).

2.2 CP Decomposition

A *CP decomposition* is an approximation of a N -way tensor \mathcal{X} by a model tensor \mathcal{Y} that is a sum of C rank-1 tensors, as shown in Figure 1. In this section we focus on a 3-way example, but the CP decomposition extends to any $N > 3$ (for more details, see [13], for example). A CP model is an N -way, rank- C tensor. This tensor is represented as a set of N matrices called *factor matrices*. In general, the n^{th} factor matrix, denoted by \mathbf{U}_n , has I_n rows and C columns. In our 3-way example, the model has factor matrices \mathbf{U} , \mathbf{V} , and \mathbf{W} , and its entries are given by

$$\mathcal{Y}(i, j, k) = \sum_{c \in [C]} \mathbf{U}(i, c) \cdot \mathbf{V}(j, c) \cdot \mathbf{W}(k, c).$$

We also use the notation $\mathcal{Y} = \llbracket \mathbf{U}, \mathbf{V}, \mathbf{W} \rrbracket$.

One commonly used method of computing the CP decomposition is the Alternating Least Squares (CP-ALS) method. In CP-ALS, one factor matrix is updated at a time and the rest are kept fixed. This update involves solving a linear least squares problem that can be expressed in matrix notation; for example, the update of \mathbf{V} has the form $\mathbf{V} = \mathbf{X}_{(1)}(\mathbf{U} \odot \mathbf{W})(\mathbf{U}^T \mathbf{U} * \mathbf{W}^T \mathbf{W})^\dagger$.

In general, each factor matrix update consists of 3 operations:

- Matricized Tensor Times Khatri Rao Product (MTTKRP),

$$\mathbf{M} = \mathbf{X}_{(n)}(\mathbf{U}_{N-1} \odot \dots \odot \mathbf{U}_{n+1} \odot \mathbf{U}_{n-1} \odot \dots \odot \mathbf{U}_0),$$
- Gram matrix and Hadamard product computation,

$$\mathbf{H} = \bigotimes_{n \neq k \in [N]} (\mathbf{U}_k)^T (\mathbf{U}_k)$$

- and solving a linear system: $\mathbf{U}_n = \mathbf{M} \mathbf{H}^\dagger$,

where \dagger denotes the pseudoinverse. Ignoring the cost of forming the Khatri-Rao product matrix (KRP), the number of flops required to multiply $\mathbf{X}_{(n)}$ by the KRP is $O(IC)$. The number of flops required to compute \mathbf{H} is $O(C^2 N + C \sum_{n \neq k \in [N]} I_k^2)$, and the number of flops involved in solving the linear system is $O(C^3 + I_n C^2)$. Thus, as I is the product of all the tensor dimensions, in the typical case nearly all of the computation is spent in MTTKRP. We note that there are alternative optimization schemes to CP-ALS, but because MTTKRP is part of the gradient, nearly all of them require computing and are bottlenecked by MTTKRP.

2.3 MTTKRP

The most straightforward way to compute the MTTKRP is to form the matricized tensor explicitly (as a column- or row-major matrix), form the KRP explicitly, and then use a BLAS call to perform the

matrix multiplication [6]. This approach benefits from an efficient matrix multiplication. However, forming an explicit matricized tensor involves reordering the tensor entries (for most modes), which is a completely memory-bound operation and can become the bottleneck. Likewise, the KRP computation involves $O(I_{\neq n}C)$ flops to produce a matrix of size $I_{\neq n} \times C$, which has a very low arithmetic intensity and will also be memory bound. The goal in this paper is to avoid reordering tensor entries and perform the KRP computation and matrix multiplication as fast as possible.

2.4 Related Work

There are several Matlab software packages that implement optimization techniques for computing CP decompositions and include functions for computing MTTKRP, including N-way Toolbox [4], Tensor Toolbox [7], and Tensorlab [26]. Recently, there have been efforts to develop more efficient implementations of MTTKRP to compute CP decompositions of sparse tensors, involving parallelizations for multi-core, many-core, and distributed-memory systems. Kaya and Uçar [12] develop a distributed-memory implementation of CP (and MTTKRP) for sparse tensors using hypergraph partitioning techniques to optimize performance. Smith *et al.* [23] develop SPLATT, a shared-memory parallel library for sparse CP, which has been extended to distributed-memory [21] and many-core platforms [22]. Li *et al.* [15] present AdaTM, a shared-memory parallel framework for sparse CP that reuses intermediate quantities across the MTTKRPs in different modes to save computation.

In the dense case, Bader and Kolda [6] proposed the straightforward approach described in the previous section. Phan *et al.* [19] introduce an alternative approach that avoids reordering tensor entries but still casts most of the computation in terms of matrix multiplication. We implement sequential and parallel versions of their algorithm in Section 4. They also show how to avoid redundant computation across MTTKRPs for different modes, but we do not consider that optimization in this paper. Vannieuwenhoven *et al.* [25] also implement the algorithm of Phan *et al.* and combine it with a blocking approach to minimize the temporary memory footprint of a dense MTTKRP. They show that minimizing the memory footprint does not have an adverse effect on performance, though they do not parallelize the algorithm.

Other related work exploits the data layout of matricized tensors and avoid reordering tensor entries using similar ideas to ours for a different tensor computation, known as tensor-times-matrix (TTM). Li *et al.* [14] develop a parallelization framework for computing TTMs with dense tensors on multicore platforms. Austin *et al.* [5] present distributed-memory parallel algorithms for computing the Tucker decomposition of dense tensors, which includes the sequential implementation of TTM that avoids reordering entries.

Other parallelizations of the CP decomposition and MTTKRP for dense tensors include those of Liavas *et al.* [16] and Aggour and Yenner [2]. Liavas *et al.* consider the nonnegative case and implement a distributed-memory parallel MTTKRP, presenting performance results for 3-way tensors of equal-sized dimensions. Aggour and Yenner also use distributed-memory parallelization and focus on 3-way tensors that have a single long dimension. We are unaware of any work that parallelizes MTTKRP for dense tensors on shared-memory platforms.

3 NEUROIMAGING APPLICATION

Our motivation for this work is a need for more efficient software to extract brain connectivity information from functional Magnetic Resonance Imaging (fMRI) data. The CP decomposition is a useful tool for neuroimaging research in general because it affords a multidimensional approach to the analysis of large and potentially heterogeneous data sets. It allows researchers to extract commonalities from diverse groups of human data and generate dynamic brain connectivity maps [11].

In our case, for example, subjects are given a cognitive task that lasts several minutes, and fMRI information is gathered for a discrete set of voxels in the brain. The voxel information is aggregated into regions of interest, and then for each subject and at each time point, the instantaneous correlations between all pairs of regions of interest are computed. This produces a time-by-subject-by-region-by-region dense tensor, and we use a CP decomposition to extract components that represent brain networks that vary over time and have varying representation over subjects. Analysis of these components yields a more complete picture of how brain connectivity relates to tasks and individual performance, and it can be used to differentiate among individuals. Such advanced data analyses hold promise to reveal important early onset symptoms of neurodegenerative disease so that prophylactic treatments may be devised.

The spatial and temporal resolution of human functional neuroimaging data is always improving due to developments in MRI technology. MRI techniques that push the limits of achievable spatial and temporal resolution lead to larger and richer brain imaging data sets, which will rely on efficient algorithms for analysis. In addition to the ever-improving spatial and temporal resolution of human functional MRI data collection, research studies are employing larger and larger group sizes as well—referred to as population imaging studies—with the aim of building a data driven discovery science of the human mind and brain (e.g., Human Connectome Project [10], UK BioBank [18]). The increasing sample sizes in combination with the improving MRI data collection requires efficient and scalable methods for analysis. The need to discover the optimal rank of multi-way and multi-modal data, and employ multiple random starts to ensure uniqueness, reliability, and reproducibility of the solutions for the massive data sets, all implicate large-scale computing as crucial to the success and advancement of these and similar projects.

4 ALGORITHMS

4.1 Khatri-Rao Product

We consider the Khatri-Rao product (KRP) of Z matrices. In the context of MTTKRP, the output for the n th mode mathematically depends on the *full KRP*:

$$\mathbf{K} = \mathbf{U}_{N-1} \circ \cdots \circ \mathbf{U}_{n+1} \circ \mathbf{U}_{n-1} \circ \cdots \circ \mathbf{U}_0.$$

However, we consider MTTKRP algorithms that form the full KRP as well as those that do not form it explicitly and instead compute partial KRPs, such as the *left KRP*:

$$\mathbf{K}_L = \mathbf{U}_{(0)} \circ \cdots \circ \mathbf{U}_{(n-1)},$$

and the *right KRP*:

$$\mathbf{K}_R = \mathbf{U}_{(n+1)} \odot \dots \odot \mathbf{U}_{(N-1)}.$$

Here, we consider the generic case of Z input matrices.

The KRP is often defined as column-wise Kronecker product, and it can be computed that way. However, we consider a row-wise approach that is more easily parallelized. Recall that each row of a KRP is the Hadamard product of a set of rows, one for each input factor matrix. For example, the j th row of $\mathbf{K} = \mathbf{A} \odot \mathbf{B} \odot \mathbf{C}$ can be expressed as $\mathbf{K}(j, :) = \mathbf{A}(a, :) * \mathbf{B}(b, :) * \mathbf{C}(c, :)$, where $j = aI_B I_C + bI_C + c$ and I_B and I_C are the number of rows of \mathbf{B} and \mathbf{C} , respectively.

Naively, a KRP of Z matrices requires $Z-1$ Hadamard products per row of the output matrix. This number can be reduced by storing and re-using partially computed Hadamard products. In the example above, $\mathbf{A}(a, :) * \mathbf{B}(b, :)$ will be used I_C times in computing \mathbf{K} . Saving this partial Hadamard product when it is first computed allows for reuse in computing subsequent rows of \mathbf{K} . For a KRP of $Z \geq 3$ matrices, one can store $Z-2$ Hadamard products to save computation and perform roughly one Hadamard product per row of the output matrix.

4.1.1 Sequential. Algorithm 1 presents the pseudocode that implements this technique, computing the output matrix \mathbf{K} one row at a time and re-using partial Hadamard products. The vector ℓ is a multi-index that stores the row indices of the input matrices that correspond to a row of the output. The matrix \mathbf{P} is a temporary matrix that stores the $Z-2$ intermediate Hadamard products, which are each of length C , the number of columns of the input matrices. For the sequential algorithm, ℓ is initialized to $\mathbf{0}$ in line 2, $\mathbf{P}(0, :)$ is initialized to $\mathbf{U}_0(0, :) * \mathbf{U}_1(0, :)$ and $\mathbf{P}(z, :)$ is initialized to $\mathbf{P}(z-1, :) * \mathbf{U}_{z+1}(0, :)$ for $0 < z \leq Z-3$ in line 3. Inside the for loop, the multi-index is incremented at line 6. For every index that changes (except for the last one), the corresponding temporary Hadamard products in \mathbf{P} must be re-computed (at line 7). However, this update occurs infrequently – extra computation is required only one out of every I_{Z-1} iterations. Thus, the dominant cost is that of the single Hadamard product of line 5, which occurs once per row.

4.1.2 Parallel. Parallelizing Alg. 1 is straightforward, so we do not provide pseudocode. The parallel variant works as follows. We assign the rows of the output matrix to threads in contiguous blocks. Each thread initializes ℓ and \mathbf{P} according to its starting row index (rather than starting with row 0). Then the algorithm proceeds exactly as in the sequential case, except that the thread stops iterating after it computes its last assigned row.

4.2 1-Step MTTKRP

We now consider a 1-step approach to compute the MTTKRP. Given a matricized tensor and an explicit KRP, the idea is to perform the matrix multiplication efficiently. The benefit of this approach is that most of the computation is cast as matrix multiplication, for which high-performance implementations exist (via the BLAS interface). However, the principal complication is that the matricizations for all internal modes ($0 < n < N-1$) of a tensor whose elements are linearized in a natural way are not column- or row-major in memory, which is a requirement for the BLAS interface. The time

Algorithm 1 Row-wise Khatri-Rao Product with Reuse

Require: \mathbf{U}_z is an $J_z \times C$ matrix, for $z \in [Z]$, $Z \geq 3$

```

1: function  $\mathbf{K} = \text{KRP}(\mathbf{U}_0, \dots, \mathbf{U}_{Z-1})$ 
2:   initialize( $\ell$ )           ▶ Initialize multi-index of length  $Z$ 
3:   initialize( $\mathbf{P}$ )           ▶  $Z-2 \times C$  matrix for intermediate products
4:   for  $j \in [\llbracket J_Z \rrbracket]$  do           ▶ Loop over rows of output  $\mathbf{K}$ 
5:      $\mathbf{K}(j, :) \leftarrow \mathbf{P}(Z-3, :) * \mathbf{U}_{Z-1}(\ell_{Z-1}, :)$ 
6:     increment( $\ell$ )
7:     update( $\mathbf{P}$ )           ▶ Update intermediate products if needed
8:   end for
9: end function

```

Ensure: $\mathbf{K} = \mathbf{U}_0 \odot \dots \odot \mathbf{U}_{Z-1}$ is $\llbracket J_Z \rrbracket \times C$ matrix

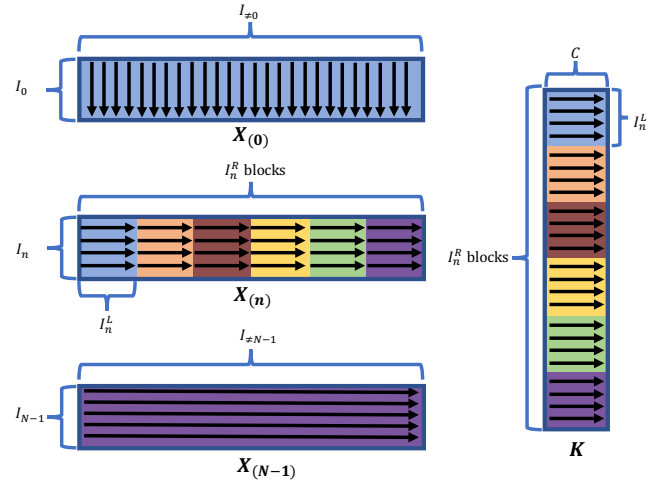


Figure 2: Data layout of the matricizations of an N -way tensor \mathcal{X} . Note that $\mathbf{X}_{(0)}$ is column-major, $\mathbf{X}_{(N-1)}$ is row-major, and $\mathbf{X}_{(n)}$ consists of I_n^R contiguous row-major submatrices of dimension $I_n \times I_n^L$. A conformal partitioning of the n th-mode KRP for the n th mode MTTKRP is depicted on the right.

required to re-order tensor elements to obtain a column- or row-major matricization is usually prohibitive, negating the benefit of BLAS performance.

Our main idea of 1-Step MTTKRP is to perform the matrix multiplication without reordering tensor entries, using multiple BLAS calls. Our algorithm is based on the observation that given the natural linearization of tensor entries, the n th mode matricization can be seen as a contiguous set of submatrices, each of which is stored row-major in memory [5, 14]. Figure 2 shows how $\mathbf{X}_{(n)}$ is ordered in memory, and it also shows how the KRP matrix \mathbf{K} can be conformally partitioned to perform the matrix multiplication as a block inner product.

4.2.1 Sequential. Algorithm 2 shows the pseudocode for the sequential 1-step algorithm. The first step is to compute the full KRP using Algorithm 1. In the case of mode 0, $\mathbf{X}_{(0)}$ is column-major so MTTKRP can be performed with a single BLAS call (line 4). For other modes, to avoid reordering tensor entries, the algorithm may have to make multiple BLAS calls. As shown in Fig. 2, lines 6 and

7 define a conformal partitioning of the two matrices so that the MTTKRP can be performed as a block inner product (the sum of submatrix multiplies). This partitioning is such that each submatrix of the matricization (of size $I_n \times I_n^L$) and each submatrix of the KRP (of size $I_n^L \times C$) is row-major in memory. Thus, each multiplication in line 9 is a BLAS call.

Algorithm 2 Sequential 1-step MTTKRP

Require: \mathcal{X} is $I_0 \times \dots \times I_{N-1}$, $\mathcal{Y} = \llbracket \mathbf{U}_0, \dots, \mathbf{U}_{N-1} \rrbracket$, $n \in [N]$

- 1: **function** $\mathbf{M} = \text{MTTKRP-1STEP-SEQ}(\mathcal{X}, \mathcal{Y}, n)$
- 2: $\mathbf{K} \leftarrow \text{KRP}(\mathbf{U}_{N-1}, \dots, \mathbf{U}_{n+1}, \mathbf{U}_{n-1}, \dots, \mathbf{U}_0)$ ▷ Alg. 1
- 3: **if** $n = 0$ **then**
- 4: $\mathbf{M} \leftarrow \mathbf{X}_{(0)} \cdot \mathbf{K}$ ▷ $\mathbf{X}_{(0)}$ is column-major
- 5: **else**
- 6: Partition $\mathbf{X}_{(n)}$ into I_n^R column blocks of size $I_n \times I_n^L$
- 7: Partition \mathbf{K} into I_n^R row blocks of size $I_n^L \times C$
- 8: **for** $j \in [I_n^R]$ **do** ▷ Loop over column blocks of $\mathbf{X}_{(n)}$
- 9: $\mathbf{M} \leftarrow \mathbf{M} + \mathbf{X}_{(n)}[j] \cdot \mathbf{K}[j]$ ▷ $\mathbf{X}_{(n)}[j]$ is row-major
- 10: **end for**
- 11: **end if**
- 12: **end function**

Ensure: $\mathbf{M} = \mathbf{X}_{(n)} \cdot (\mathbf{U}_{N-1} \odot \dots \odot \mathbf{U}_{n+1} \odot \mathbf{U}_{n-1} \odot \dots \odot \mathbf{U}_0)$ is $I_n \times C$

4.2.2 *Parallel.* We use two different techniques to parallelize 1-step MTTKRP, depending on the mode. We distinguish between external ($n = 0$ or $n = N-1$) and internal ($0 < n < N-1$) modes.

For external modes, we parallelize over the *columns* of the matricization. Each thread is assigned a contiguous set of columns of $\mathbf{X}_{(n)}$. In line 7, the thread computes the corresponding rows of the KRP \mathbf{K} (using a variant of Algorithm 1), and in line 8 it performs the multiplication with a BLAS call. Each thread computes a contribution to the output matrix, so a parallel reduction is performed at the end of the algorithm (line 23).

For internal modes, we parallelize over the *blocks* of the matricization. In this case, each thread is assigned a set of $I_n \times I_n^L$ blocks. The row block of \mathbf{K} that corresponds to the j th column block of $\mathbf{X}_{(n)}$ is the Khatri-Rao product of the left KRP matrix and the j th row of the right KRP matrix. Thus, \mathbf{K}_L is pre-computed in parallel in line 11, and each thread computes the corresponding row of \mathbf{K}_R (line 14) and Khatri-Rao product (line 15) to obtain the necessary row blocks of \mathbf{K} . The matrix multiplication of line 16 is performed with a BLAS call because each block is row-major. Again, a parallel reduction is necessary at the end of the algorithm.

Note that the internal-mode parallelization scheme assumes that the number of threads is much less than I_n^R in order to achieve load balance. We expect this to hold in nearly all cases because I_n^R is a product of tensor dimensions. If this is not the case (say I_n^R corresponds to only one mode of very small dimension), then an alternative approach would be to use the sequential algorithm with multithreaded BLAS.

4.3 2-Step MTTKRP

The 2-step algorithm first performs a *partial MTTKRP* and then finishes the computation with multiple tensor-times-vector operations (*multi-TTV*). The algorithm was developed by Phan *et al.*

Algorithm 3 Parallel 1-Step MTTKRP

Require: \mathcal{X} is $I_0 \times \dots \times I_{N-1}$, $\mathcal{Y} = \llbracket \mathbf{U}_0, \dots, \mathbf{U}_{N-1} \rrbracket$, $n \in [N]$

Require: T is the number of threads

- 1: **function** $\mathbf{M} = \text{MTTKRP-1STEP-PAR}(\mathcal{X}, \mathcal{Y}, n, T)$
- 2: **if** $n = 0$ or $n = N-1$ **then**
- 3: $b \leftarrow \lceil I_n^R / T \rceil$
- 4: Partition $\mathbf{X}_{(n)}$ into T column blocks of size $I_n \times b$
- 5: Partition KRP \mathbf{K} into T row blocks of size $b \times C$
- 6: **parallel for** $t \in [T]$, private($\bar{\mathbf{M}}_t, \bar{\mathbf{K}}[t]$) **do**
- 7: Compute $\bar{\mathbf{K}}[t]$ ▷ Variant of Alg. 1
- 8: $\bar{\mathbf{M}}_t \leftarrow \mathbf{X}_{(n)}[t] \cdot \bar{\mathbf{K}}[t]$
- 9: $\bar{\mathbf{X}}_{(n)}[t]$ is submatrix of column- or row-major matrix **end for**
- 10: **else**
- 11: $\mathbf{K}_L \leftarrow \text{KRP}(\mathbf{U}_{N-1}, \dots, \mathbf{U}_0)$ ▷ Parallel variant of Alg. 1
- 12: Partition $\mathbf{X}_{(n)}$ into I_n^R column blocks of size $I_n \times I_n^L$
- 13: **parallel for** $j \in [I_n^R]$, private($\bar{\mathbf{M}}_t, \bar{\mathbf{K}}_t$) **do**
- 14: Compute $\mathbf{K}_R(j, :)$ ▷ j th row of \mathbf{K}_R
- 15: $\bar{\mathbf{K}}_t \leftarrow \mathbf{K}_R(j, :) \odot \mathbf{K}_L$
- 16: $\bar{\mathbf{M}}_t \leftarrow \bar{\mathbf{M}}_t + \mathbf{X}_{(n)}[j] \cdot \bar{\mathbf{K}}_t$ ▷ $\mathbf{X}_{(n)}[j]$ is row-major
- 17: **end for**
- 18: **end if**
- 19: $\mathbf{M} \leftarrow \sum_t \bar{\mathbf{M}}_t$ ▷ Parallel reduction
- 20: **end function**

Ensure: $\mathbf{M} = \mathbf{X}_{(n)} \cdot (\mathbf{U}_{N-1} \odot \dots \odot \mathbf{U}_{n+1} \odot \mathbf{U}_{n-1} \odot \dots \odot \mathbf{U}_0)$ is $I_n \times C$

[19, Section III.B], but we explain it again here using our notation. For external modes, the 2-step algorithm degenerates to the 1-step algorithm. The pseudocode appears in Algorithm 4, and the data layouts of each computation are show in Figure 3.

As in the case of the 1-step algorithm, our goal will be to perform the computation without reordering tensor entries. The first observation is that more general matricizations of the tensor are column-major in memory. That is, using the notation defined in Section 2.1, $\mathbf{X}_{(0:n)}$ is column major in memory for all n . This implies that we can compute the matrix product of $\mathbf{X}_{(0:n)}$ and \mathbf{K}_R (the right KRP) with a single BLAS call, as shown in Figure 3a. Because this matrix multiplication involves all the tensor entries but only a subset of the input matrices, we refer to this operation as a partial MTTKRP. (Note that when $n = 0$, it is a full MTTKRP.)

The output of a partial MTTKRP is an intermediate quantity which must be combined with the remaining input matrices to obtain the final MTTKRP output. We interpret the output of this partial MTTKRP as a tensor of dimension $n + 2$, defining

$$\mathcal{R}_{(0:n)} = \mathbf{X}_{(0:n)} \cdot \mathbf{K}_R,$$

so that \mathcal{R} has dimensions $I_0 \times \dots \times I_n \times C$. Given \mathcal{R} , the j th column of the MTTKRP output matrix \mathbf{M} is a tensor-times-vector (TTV) operation involving the j th subtensor of \mathcal{R} and the j th columns of the remaining input matrices:

$$\mathbf{M}(:, j) = \mathcal{R}(:, \dots, :, j) \times_0 \mathbf{U}_0(:, j) \cdots \times_{n-1} \mathbf{U}_{n-1}(:, j).$$

Because the operation must be performed for each column of \mathbf{M} , we refer to the overall 2nd step as a multi-TTV.

We note that the operation can also be interpreted as an MTTKRP involving subtensor $\mathcal{R}(:, \dots, :, j)$ and the set of columns. Because

the subtensor involves a leading set of modes, it is stored as a tensor in natural order, and we can apply the same computational techniques. In particular, this MTTKRP is done with respect to the last mode of the subtensor, so it requires only one BLAS call. However, the KRP has only one column in this case, so the BLAS call is for matrix-vector rather than matrix-matrix multiplication. The matrix-vector product must be performed for each of the C output columns, as shown in Figure 3b.

The 2 steps described above incorporate the modes to the right in the 1st step and the modes to the left in the 2nd step, but that order can be reversed. To compute the partial MTTKRP involving the left modes, we observe that $\mathbf{X}_{(0:n-1)}^\top$ is row major in memory, and we can compute a different temporary quantity

$$\mathbf{L}_{(0:N-n-1)} = \mathbf{X}_{(0:n-1)}^\top \cdot \mathbf{K}_L,$$

where \mathcal{L} is $I_n \times \dots \times I_{N-1} \times C$ (see Figure 3c). The second step multi-TTV involves subtensors of \mathcal{L} and columns of \mathbf{K}_R . The j th column of the output is given by

$$\mathbf{M}(:, j) = \mathcal{L}(:, \dots, :, j) \times_{n+1} \mathbf{U}_{n+1}(:, j) \cdots \times_{N-1} \mathbf{U}_{N-1}(:, j).$$

Computationally, we can interpret each TTV as an MTTKRP with respect to the first mode of the subtensor, so again it involves only one BLAS call for each matrix-vector multiplication (see Figure 3d).

The pseudocode for the 2-step algorithm is given in Algorithm 4. It starts by computing both left and right partial KRPs. Given that either ordering of the steps is correct, the algorithm chooses the ordering that minimizes computation in the 2nd step (the number of flops in the 1st step is the same). If it is more efficient to do the left side first, it computes the left partial MTTKRP in line 5 and the multi-TTV in lines 6–9. Otherwise, it computes the right partial MTTKRP in line 11 and the multi-TTV in lines 12–15.

The bulk of the computation occurs in the partial MTTKRP, which is a matrix multiplication requiring a single BLAS call. We note that the dimensions of this matrix multiplication are more balanced than in the case of the full MTTKRP. Parallelization of this algorithm is all within the BLAS calls, so Algorithm 4 applies for both sequential and parallel variants.

5 EXPERIMENTAL RESULTS

5.1 Experimental Setup

All experiments are benchmarked on a dual-socket Intel Xeon E5-2620 (Sandy Bridge) server with a total of 12 cores. Each socket has a 15 MB L3 cache, and each core has a 256 KB L2 cache and 32 KB L1 data cache. Each core has a clock rate of 2.00 GHz (with turbo boost disabled) and peak flop rate of 16 GFLOPS. Our code is written in C and compiled with GCC version 5.4.0. We use Intel’s Math Kernel Library (MKL) version 2017.2.174. The MATLAB benchmarks use version 9.0.0.341360 (R2016a) and Tensor Toolbox version 2.6. All experiments are performed in double precision.

5.2 KRP

We first consider the performance of Algorithm 1, which computes the Khatri-Rao product of Z input matrices. Figure 4 presents performance results for Algorithm 1, which exploits re-use of intermediate quantities, in comparison with a naive version of the algorithm and the STREAM benchmark [17]. We consider $C = \{25, 50\}$ (in

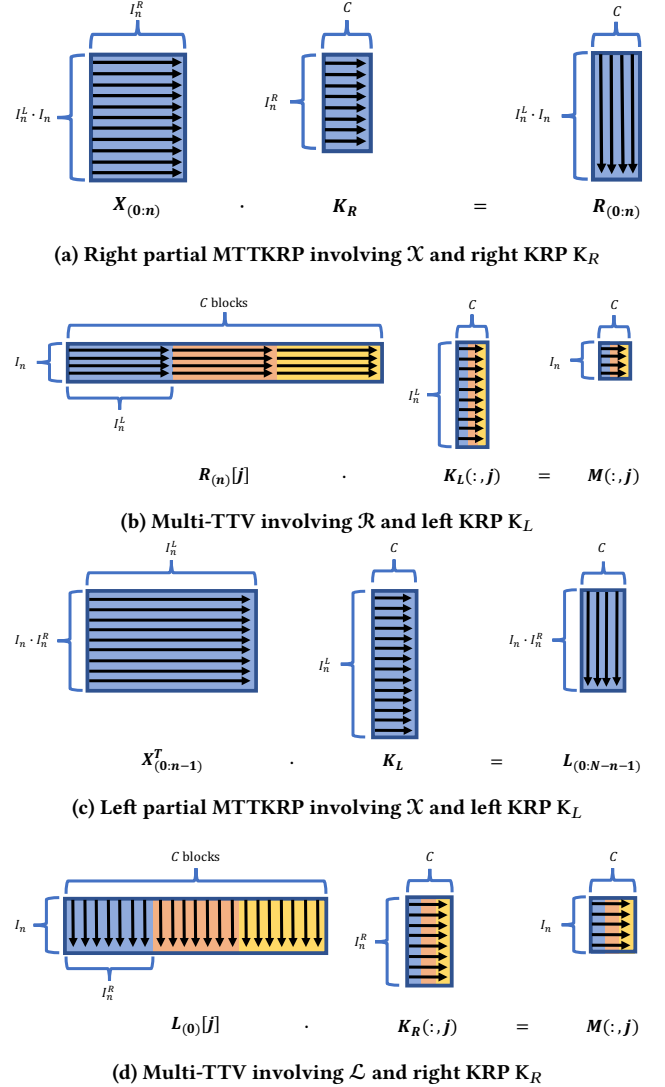


Figure 3: Data layouts of tensors and matrices involved in computations of 2-step MTTKRP. Figs. 3a and 3b correspond to the 2 steps with right-left ordering; Figs. 3c and 3d correspond to the 2 steps with left-right ordering.

Figs. 4a and 4b) and in each case experiment with $Z = \{2, 3, 4\}$. For each value of C we choose the input matrix row dimensions to be all equivalent and such that their product is approximately 20 million. This implies that for all experiments shown in Fig. 4a, the output matrix has a size of roughly 500 million entries; in Fig. 4b the output matrix has approximately 1 billion entries.

We measure the running time of three algorithms over various numbers of threads. The results labeled “Reuse” correspond to Alg. 1, which is the algorithm we use in the MTTKRP experiments. The results labeled “Naive” correspond to a row-wise algorithm that does not store and re-use intermediate Hadamard products. The STREAM benchmark we report is based on reading, scaling,

Algorithm 4 Sequential/Parallel 2-step MTTKRP

Require: \mathcal{X} is $I_0 \times \dots \times I_{N-1}$, $\mathcal{Y} = [\mathbf{U}_0, \dots, \mathbf{U}_{N-1}]$, $n \in [N]$

- 1: **function** $\mathbf{M} = \text{MTTKRP-2STEP}(\mathcal{X}, \mathcal{Y}, n)$
- 2: $\mathbf{K}_L \leftarrow \text{KRP}(\mathbf{U}_{n-1}, \dots, \mathbf{U}_0)$ \triangleright Alg. 1, Left Partial KRP
- 3: $\mathbf{K}_R \leftarrow \text{KRP}(\mathbf{U}_{N-1}, \dots, \mathbf{U}_{n+1})$ \triangleright Alg. 1, Right Partial KRP
- 4: **if** $I_n^L > I_n^R$ **then** \triangleright Use Left Partial MTTKRP
- 5: $\mathbf{L}_{(0:N-n-1)} = \mathbf{X}_{(0:n-1)}^T \cdot \mathbf{K}_L$
- 6: Partition $\mathbf{L}_{(0)}$ into C column blocks of size $I_n \times I_n^R$
- 7: **for** $j \in [C]$ **do**
- 8: $\mathbf{M}(:, j) \leftarrow \mathbf{L}_{(0)}[j] \cdot \mathbf{K}_R(:, j)$ $\triangleright \mathbf{L}_{(0)}[j]$ is column-major
- 9: **end for**
- 10: **else** \triangleright Use Right Partial MTTKRP
- 11: $\mathbf{R}_{(0:n)} = \mathbf{X}_{(0:n)} \cdot \mathbf{K}_R$
- 12: Partition $\mathbf{R}_{(n)}$ into C column blocks of size $I_n \times I_n^L$
- 13: **for** $j \in [C]$ **do**
- 14: $\mathbf{M}(:, j) \leftarrow \mathbf{R}_{(n)}[j] \cdot \mathbf{K}_L(:, j)$ $\triangleright \mathbf{R}_{(n)}[j]$ is row-major
- 15: **end for**
- 16: **end if**
- 17: **end function**

Ensure: $\mathbf{M} = \mathbf{X}_{(n)} \cdot (\mathbf{U}_{N-1} \odot \dots \odot \mathbf{U}_{n+1} \odot \mathbf{U}_{n-1} \odot \dots \odot \mathbf{U}_0)$ is $I_n \times C$

and writing a matrix the same size as the output KRP matrix. Each reported time is the average of 100 trials.

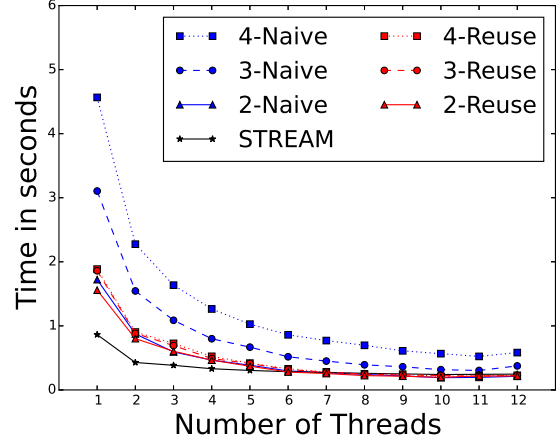
Our first conclusion from the data is that exploiting reuse is an important optimization for KRP. Algorithm 1 outperforms its naive alternative, and the difference increases with Z (note that for $Z = 2$ there is no difference in algorithm). For $Z = \{3, 4\}$, the speedups of Reuse over Naive range from 1.5 \times to 2.5 \times .

Our second conclusion is that Algorithm 1 is essentially a memory-bound operation, achieving competitive performance with the STREAM benchmark. This is expected, as the number of flops in the optimized Algorithm 1 is the same as the number of output matrix entries. However, because the input matrices are relatively small, we see that KRP can take even less time than STREAM (as in the case of $C = 50$), which involves both a read and a write of the large matrix.

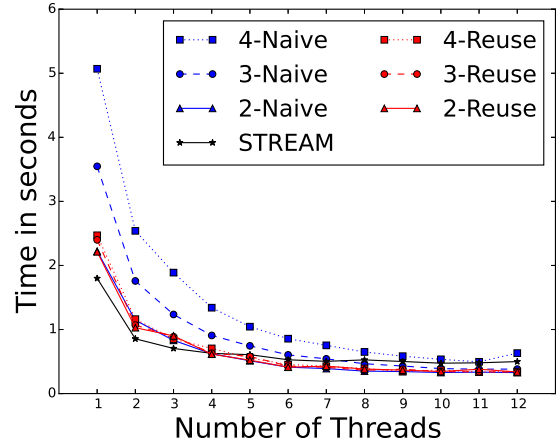
Finally, we see efficient scaling of our parallel variant of Algorithm 1. For $C = 25$, we observe a parallel speedup range of 6.6–7.4 \times for 12 threads; for $C = 50$ the speedup range is 7.9–8.3 \times .

5.3 MTTKRP

In this section we discuss performance results for our proposed MTTKRP algorithms. We compare the performance of 1-step (Algorithm 3) and 2-step (Algorithm 4) algorithms over various numbers of threads, noting that the two algorithms are equivalent for external modes ($n = 0$ and $n = N-1$). For a baseline, we also compare against the performance of a single BLAS call (MKL’s implementation of DGEMM). This benchmark is run on a single matrix multiplication between two column-major matrices that are the same size as the matricized tensor and the KRP, respectively. It can be viewed as a lower bound on the performance of the most straightforward approach to MTTKRP (that reorders tensor entries) because it does not include the time required to reorder entries or



(a) Output matrix with dimension $J \times 25$



(b) Output matrix with dimension $J \times 50$

Figure 4: Time comparison of Algorithm 1 against a naive algorithm and STREAM benchmark over varying numbers of threads. Each experiment involves $J \approx 2e7$ output matrix rows and either 2, 3, or 4 input matrices and either 25 or 50 columns. Both KRP algorithms compute a row of the output at time; Algorithm 1 avoids the redundant computation performed by the naive algorithm.

form the explicit KRP. Each reported result in this section is the median of 10 runs.

We note that in the case of the 1-step approach, the parallel algorithm (Alg. 3) run with 1 thread is slightly different than the sequential algorithm (Alg. 2) for internal modes. Instead of forming the full KRP \mathbf{K} explicitly, the parallel algorithm forms the left partial KRP and computes blocks of \mathbf{K} as needed. Because we observed the parallel approach (when run with 1 thread) is slightly more efficient and uses less memory than the sequential approach, we use the parallel approach for all sequential benchmarks.

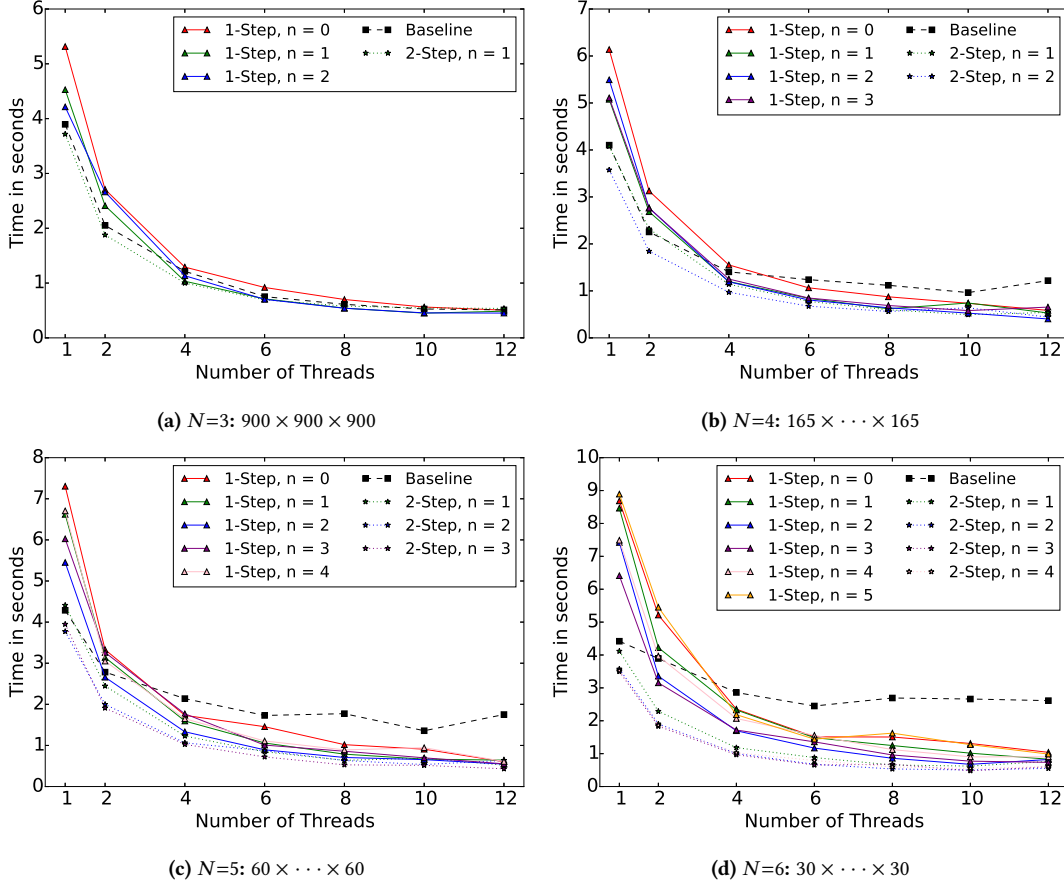


Figure 5: Time comparison of 1-step and 2-step MTTKRP algorithms for different modes over varying numbers of threads. Each subfigure corresponds to a number of modes; all experiments involve approximately 750 million tensor entries. The 2-step algorithm is defined only for inner modes. The baseline DGEMM benchmark is the time to multiply column-major matrices of the same dimensions as the MTTKRP.

5.3.1 Parallel Scaling. We first consider the overall time and parallel scaling of each algorithm for different tensors. Figure 5 presents results for four different tensors, with $N = \{3, 4, 5, 6\}$. For each tensor, each dimension is the same and chosen so that the total number of tensor entries is approximately 750 million. In each experiment, the number of columns in the output matrix is $C = 25$.

Focusing first on sequential performance, Figure 5 shows that the 2-step algorithm is faster than the baseline and that the 1-step algorithm is slower than the baseline, relationships that are consistent across both tensors and modes. We note that the 2-step algorithm is not available for external modes. In the worst case, the 1-step algorithm takes about $2\times$ as long as the baseline; the baseline is never slower than the 2-step algorithm by more than 25% and never faster by more than 3%. Section 5.3.2 further investigates the reasons for these relative performance differences.

Our next observation is that both 1-step and 2-step algorithms scale more efficiently than the baseline, particularly for larger N . In fact, even at 4 threads, all of the proposed implementations are comparable or better than the single BLAS call, and they continue to improve up to 12 threads. At 12 threads and for $N > 3$, the

speedup of 1-step and 2-step algorithms over the baseline range from $2\times$ to $4.7\times$, and the baseline still does not include time for reordering tensor entries or computing the KRP.

We believe part of the explanation for the poor scaling of the baseline implementation is that MKL has not fully optimized matrix multiplication of this shape, as has been observed in previous work [9]. As N increases in our benchmarks, the shape of the MTTKRP matrix multiplication approaches an inner product, with a long inner matrix multiplication dimension and a small output matrix. The optimal parallelization of this computation involves write conflicts, for which we use temporary private memory and a parallel reduction, but MKL’s implementation may be avoiding the memory footprint overhead of such an approach.

Unlike the baseline implementation, the 1-step and 2-step algorithms scale well to 12 threads. The parallel speedup of the 1-step algorithm ranges from $8 - 12\times$ on 12 threads, and the 2-step parallel speedup ranges from $6 - 8\times$. We note that the 2-step algorithm relies on the parallel performance of MKL, but it sees better parallel scaling than the baseline because the matricization within the partial MTTKRP is more square than that of the baseline approach

(though it involves the same number of flops). The fact that the 1-step algorithm scales slightly better than the 2-step algorithm implies that the parallel running times of the two approaches are fairly comparable at 12 threads. We explore this further in Section 5.3.2.

5.3.2 Time Breakdown. Figure 6 gives a detailed breakdown of the computation time of MTTKRP algorithms. We consider the same four tensors as in Section 5.3.1, with $N = \{3, 4, 5, 6\}$, and we benchmark both sequential ($T = 1$) and parallel ($T = 12$) cases. Each experiment uses $C = 25$.

The baseline implementation has only one category: matrix multiplication (labeled Baseline in the legend). The time for the 1-step algorithm (Alg. 3) is broken down into matrix multiplication (DGEMM, line 8 or 16), forming the full KRP (Full KRP, line 7), forming the left KRP for internal modes and multiplying it by a row of the right KRP (Left & Right KRP, lines 11 and 15), and performing the final parallel reduction (REDUCE, line 19). The time for the 2-step algorithm (Alg. 4) is broken down into matrix multiplication (DGEMM, line 5 or 11), matrix-vector multiplication (DGEMV, line 8 or 14), and forming the left and right KRPs (Left & Right KRP, line 2 and 3).

Our first observation is that a considerable amount of time in the 1-step algorithm is spent computing the KRP, particularly for external modes. For internal modes, only the left KRP is computed explicitly (which requires negligible time), and the rest of the KRP time is spent in computing blocks of the full KRP using a row of the right KRP. In fact, this extra cost is the main reason the 1-step algorithm is slower than the baseline in the sequential case; the matrix multiplication time is generally comparable to the baseline even though it involves multiple BLAS calls for smaller matrix dimensions. (Recall that the baseline ignores the cost of forming the KRP.) Comparing the sequential 1-step performance to the parallel 1-step performance, we see that each category scales similarly and the proportions are generally preserved.

For internal modes, we observe that the 2-step algorithm spends almost all of its time in matrix multiplication. The time spent in other computations is generally negligible. Comparing the matrix multiplication time to the baseline, we see that the 2-step algorithm demonstrates slightly better performance because the matrix dimensions are more amenable for MKL. In the sequential case, the 2-step algorithm is generally faster than the 1-step algorithm (for internal modes). In the parallel case, the two algorithms are comparable, though the 2-step algorithm is usually slightly faster.

5.3.3 Neuroscience Data. The underlying motivation for this work is to speed up CP-ALS in order to analyze neuroscience (fMRI) data. Our data is a 4-way tensor of size $225 \times 59 \times 200 \times 200$, representing for 225 time steps and for 59 subjects the correlation between fMRI signals measured at 200 different brain regions. For a different type of analysis, we linearize the last two modes; because the tensor is symmetric in these two modes this linearization also reduces the number of tensor entries by a factor of 2. The corresponding 3-way tensor is $225 \times 59 \times 19900$.

In this section, we compare the performance of Matlab code using the Tensor Toolbox with our implementation of CP-ALS using efficient MTTKRP kernels. For dense tensors, the current available software packages (e.g., N-way Toolbox [4], Tensor Toolbox [7],

Tensorlab [26]) are all written in Matlab. Because the Matlab software packages cast computations as matrix operations, and Matlab uses efficient BLAS implementations like MKL, we can expect reasonable performance from Matlab. However, on multicore servers, the only opportunity for parallelization in the packages is within BLAS calls.

Figure 7 shows the per-iteration run times for computing CP decompositions on our 3D and 4D application tensors with $C = \{10, 15, 20, 25, 30\}$, using both sequential and parallel, MATLAB and C implementations. Our C implementation of CP-ALS employs Algorithm 3 (1-step) for both outer modes and Algorithm 4 (2-step) for all inner modes. We observe up to a $2\times$ speedup of our sequential implementation over Matlab, running on only 1 core. In the parallel case, the highest speedup of our implementation over Matlab (running with all 12 cores available) comes for the largest C : $6.7\times$ for the 3D tensor and $7.4\times$ for the 4D tensor.

Figure 8 gives the time breakdown for our implementation of MTTKRP for the application tensors, which have varying dimensions across modes. This plot can be contrasted with Figure 6, which depicts data for synthetic tensors that have all the same dimensions across modes. In particular, note that the KRP cost is relatively more significant in small modes ($n = 1, I_1 = 59$), which agrees with the larger ratio of flops. We observe that for large modes, the 1-step algorithm is competitive with the baseline in the sequential case and again outperforms the baseline in the parallel case. The 2-step algorithm is consistently better than the baseline, and significantly better in the parallel case. For mode $n = 1$ the parallel MTTKRP algorithms are $2.8\times$ and $3.5\times$ faster than the baseline for 3D and 4D, respectively.

6 CONCLUSION

In this paper, we introduce a parallel algorithm for KRP and two parallel algorithms for MTTKRP. Our performance results indicate that our implementations perform well sequentially and scale efficiently up to 12 threads. We also show that improving the performance of these kernels results in faster CP-ALS iterations for application problems.

One conclusion that we wish to highlight from the performance results is the high relative cost of the KRP computation in the 1-step algorithm. For example, in the case of the external modes of synthetic 6-way tensor where each mode has dimension 30 (see Fig. 6d), the KRP takes about a third to a half of the time even though the number of flops is 1/30th the number of flops involved in the matrix multiplication. This is due in large part to the memory boundedness of the KRP computation. Just as tensor reordering should be avoided, future optimization of MTTKRP should avoid computing large KRPs.

The natural next step for this work is to implement the algorithm proposed by Phan *et al.* [19, Section III.C] for avoiding recomputation across MTTKRPs of different modes, which can be done for CP-ALS and other gradient-based optimization methods. In particular, the computational kernels of the full algorithm are the same as the single-mode computation. Using this algorithm, we could expect a further reduction in per-iteration CP-ALS time of around 50% in the 3D case and $2\times$ in the 4D case (and higher for larger N).

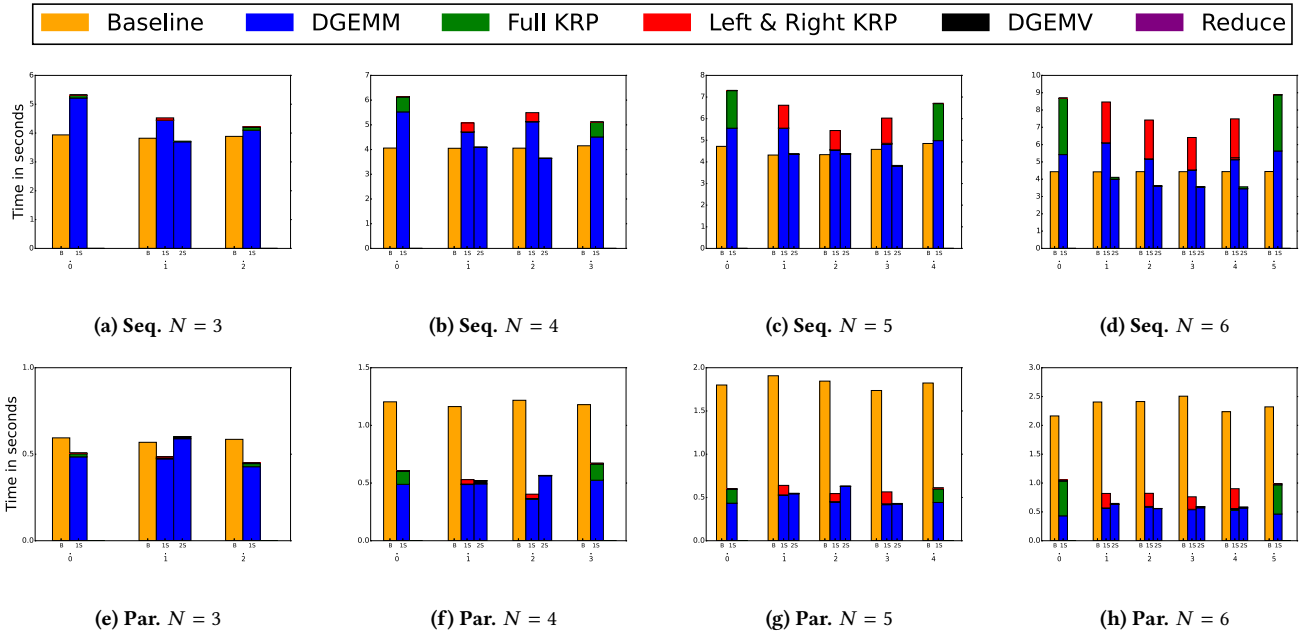


Figure 6: Time breakdown of 1-step and 2-step MTTKRP (and baseline DGEMM) across modes for varying numbers of modes. The top row corresponds to sequential time ($T = 1$), and the bottom row corresponds to parallel time ($T = 12$). Each column corresponds to a number of modes; all experiments involve approximately 750 million tensor entries and $C = 25$ matrix columns. The baseline DGEMM benchmark is the time to multiply column-major matrices of the same dimensions as the MTTKRP.

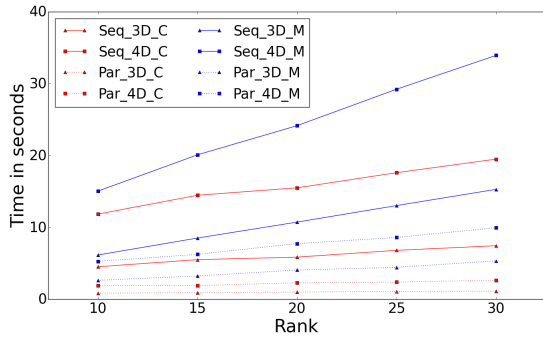


Figure 7: Per-iteration CP-ALS time for Matlab and C implementations over different ranks. The Matlab implementation is Tensor Toolbox’s `cp_als` function, and the C implementation is ours, using the most efficient 1-step and 2-step MTTKRP algorithm for each mode. Parallel runs are given all 12 cores on the machine. Tensor sizes are $225 \times 59 \times 200 \times 200$ (4D) and $225 \times 59 \times 19900$ (3D).

We note that this algorithm also benefits from avoiding computing large KRPs.

We have also noticed that with improvements of MTTKRP performance, other computations within CP-ALS, such as the residual error computation, have become relatively more costly. Improving

the efficiency of these new bottlenecks could yield overall performance increases.

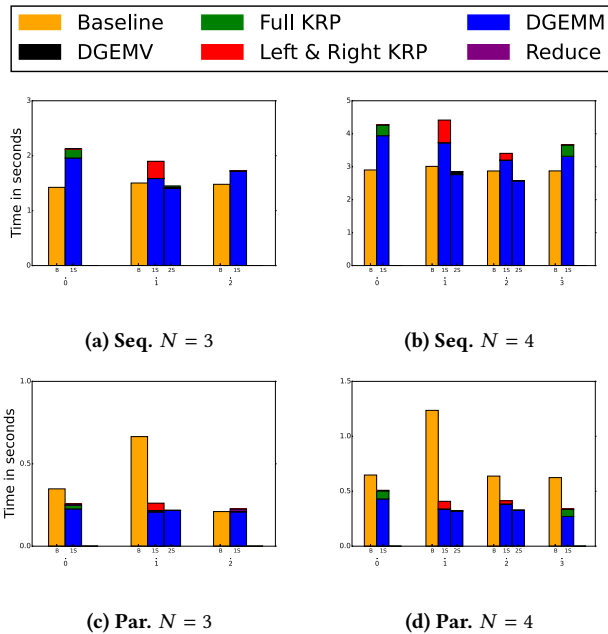


Figure 8: Time breakdown of 1-step and 2-step MTTKRP (and baseline DGEMM) across modes for 3D and 4D fMRI tensors. The top row corresponds to sequential time ($T = 1$), while the bottom row corresponds to parallel time ($T = 12$). The baseline DGEMM benchmark is the time to multiply column-major matrices of the same dimensions as the MTTKRP.

REFERENCES

- [1] Evrim Acar, Daniel M. Dunlavy, Tamara G. Kolda, and Morten Mørup. 2011. Scalable Tensor Factorizations for Incomplete Data. *Chemometrics and Intelligent Laboratory Systems* 106, 1 (March 2011), 41–56. <https://doi.org/10.1016/j.chemolab.2010.08.004>
- [2] Kareem S. Aggour and Bülent Yener. 2016. *A Parallel PARAFAC Implementation & Scalability Testing for Large-Scale Dense Tensor Decomposition*. Technical Report 16-02. Rensselaer Polytechnic Institute. <https://www.cs.rpi.edu/research/pdf/16-02.pdf>
- [3] Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. 2014. Tensor Decompositions for Learning Latent Variable Models. *Journal of Machine Learning Research* 15 (2014), 2773–2832. <https://doi.org/10.21236/ada604494>
- [4] Claus A Andersson and Rasmus Bro. 2000. The N-way Toolbox for MATLAB. *Chemometrics and Intelligent Laboratory Systems* 52, 1 (2000), 1–4. [https://doi.org/10.1016/S0169-7439\(00\)00071-X](https://doi.org/10.1016/S0169-7439(00)00071-X)
- [5] Woody Austin, Grey Ballard, and Tamara G. Kolda. 2016. Parallel Tensor Compression for Large-Scale Scientific Data. In *Proceedings of the 30th IEEE International Parallel and Distributed Processing Symposium*. 912–922. <https://doi.org/10.1109/IPDPS.2016.67>
- [6] Brett W. Bader and Tamara G. Kolda. 2007. Efficient MATLAB computations with sparse and factored tensors. *SIAM Journal on Scientific Computing* 30, 1 (December 2007), 205–231. <https://doi.org/10.1137/060676489>
- [7] Brett W. Bader, Tamara G. Kolda, et al. 2015. MATLAB Tensor Toolbox Version 2.6. Available online. (February 2015). <http://www.sandia.gov/~tgkolda/TensorToolbox/>
- [8] Andrzej Cichocki, Danilo Mandic, Lieven De Lathauwer, Guoxu Zhou, Qibin Zhao, Cesar Caiafa, and Anh-Huy Phan. 2015. Tensor Decompositions for Signal Processing Applications: From two-way to multiway component analysis. *IEEE Signal Processing Magazine* 32, 2 (March 2015), 145–163. <https://doi.org/10.1109/MSP.2013.2297439>
- [9] James Demmel, David Eiahu, Armando Fox, Shoaib Kamil, Benjamin Lipshitz, Oded Schwartz, and Omer Spillinger. 2013. Communication-Optimal Parallel Recursive Rectangular Matrix Multiplication. In *Proceedings of the 27th IEEE International Symposium on Parallel and Distributed Processing (IPDPS '13)*. 261–272. <https://doi.org/10.1109/IPDPS.2013.80>
- [10] D.C. Van Essen, K. Ugurbil, E. Auerbach, D. Barch, T.E.J. Behrens, R. Bucholz, A. Chang, L. Chen, M. Corbetta, S.W. Curtiss, S. Della Penna, D. Feinberg, M.F. Glasser, N. Harel, A.C. Heath, L. Larson-Prior, D. Marcus, G. Michalareas, S. Moeller, R. Oostenveld, S.E. Petersen, F. Prior, B.L. Schlaggar, S.M. Smith, A.Z. Snyder, J. Xu, and E. Yacoub. 2012. The Human Connectome Project: a data acquisition perspective. *Neuroimage* 62, 4 (2012), 2222–2231. <https://doi.org/10.1016/j.neuroimage.2012.02.018>
- [11] R. Matthew Hutchison, Thilo Womelsdorf, Elena A. Allen, Peter A. Bandettini, Vince D. Calhoun, Maurizio Corbetta, Stefania Della Penna, Jeff H. Duyn, Gary H. Glover, Javier Gonzalez-Castillo, Daniel A. Handwerker, Shella Keilholz, Vesa Kiviniemi, David A. Leopold, Francesco de Pasquale, Olaf Sporns, Martin Walter, and Catie Chang. 2013. Dynamic functional connectivity: Promise, issues, and interpretations. *NeuroImage* 80 (2013), 360–378. <https://doi.org/10.1016/j.neuroimage.2013.05.079> Mapping the Connectome.
- [12] Oguz Kaya and Bora Uçar. 2015. Scalable Sparse Tensor Decompositions in Distributed Memory Systems. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '15)*. ACM, New York, NY, USA, Article 77, 11 pages. <https://doi.org/10.1145/2807591.2807624>
- [13] Tamara G. Kolda and Brett W. Bader. 2009. Tensor Decompositions and Applications. *SIAM Rev.* 51, 3 (September 2009), 455–500. <https://doi.org/10.1137/07070111X>
- [14] Jiajia Li, Casey Battaglini, Ioakeim Perros, Jimeng Sun, and Richard Vuduc. 2015. An Input-Adaptive and In-Place Approach to Dense Tensor-Times-Matrix Multiply. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis (SC '15)*. ACM, New York, NY, USA, Article 76, 12 pages. <https://doi.org/10.1145/2807591.2807671>
- [15] Jiajia Li, Jee Choi, Ioakeim Perros, Jimeng Sun, and Richard Vuduc. 2017. Model-Driven Sparse CP Decomposition for Higher-Order Tensors. In *IEEE International Parallel and Distributed Processing Symposium (IPDPS)*. 1048–1057. <https://doi.org/10.1109/IPDPS.2017.80>
- [16] Athanasios P. Liavas, Georgios Kostoulas, Georgios Lourakis, Kejun Huang, and Nicholas D. Sidiropoulos. 2017. Nesterov-based parallel algorithm for large-scale nonnegative tensor factorization. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 5895–5899. <https://doi.org/10.1109/ICASSP.2017.7953287>
- [17] John D. McCalpin. 1991–2007. *STREAM: Sustainable Memory Bandwidth in High Performance Computers*. Technical Report. University of Virginia, Charlottesville, Virginia. <http://www.cs.virginia.edu/stream/>
- [18] Karla L. Miller, Fidel Alfaro-Almagro, Neal K. Bangerter, David L. Thomas, Essa Yacoub, Junqian Xu, Andreas J. Bartsch, Saad Jbabdi, Stamatios N. Sotiropoulos, Jesper L.R. Andersson Andersson, Ludovica Griffanti, Gwenaëlle Douaud, Thomas W. Okell, Peter Weale, Julius Dragonu, Steve Garratt, Sarah Hudson, Rory Collins, Mark Jenkinson, Paul M. Matthews, and Stephen M. Smith. 2016. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nature Neuroscience* 19, 11 (2016), 1523–1536. <https://doi.org/10.1038/nn.4393>
- [19] Anh-Huy Phan, Petr Tichavsky, and Andrzej Cichocki. 2013. Fast Alternating LS Algorithms for High Order CANDECOMP/PARAFAC Tensor Factorizations. *IEEE Transactions on Signal Processing* 61, 19 (Oct 2013), 4834–4846. <https://doi.org/10.1109/TSP.2013.2269903>
- [20] Nicholas D. Sidiropoulos, Lieven De Lathauwer, Xiao Fu, Kejun Huang, Evangelos E. Papalexakis, and Christos Faloutsos. 2017. Tensor Decomposition for Signal Processing and Machine Learning. *IEEE Transactions on Signal Processing* 65, 13 (July 2017), 3551–3582. <https://doi.org/10.1109/TSP.2017.2690524>
- [21] Shaden Smith and George Karypis. 2016. A Medium-Grained Algorithm for Distributed Sparse Tensor Factorization. In *30th IEEE International Parallel and Distributed Processing Symposium*. 902–911. <https://doi.org/10.1109/IPDPS.2016.113>
- [22] Shaden Smith, Jongsoo Park, and George Karypis. 2017. Sparse Tensor Factorization on Many-Core Processors with High-Bandwidth Memory. *31st IEEE International Parallel & Distributed Processing Symposium (IPDPS'17)* (May 2017), 1058–1067. <https://doi.org/10.1109/IPDPS.2017.84>
- [23] Shaden Smith, Niranjay Ravindran, Nicholas D. Sidiropoulos, and George Karypis. 2015. SPLATT: Efficient and Parallel Sparse Tensor-Matrix Multiplication. In *Proceedings of the 2015 IEEE International Parallel and Distributed Processing Symposium (IPDPS '15)*. IEEE Computer Society, Washington, DC, USA, 61–70. <https://doi.org/10.1109/IPDPS.2015.27>
- [24] Jimeng Sun, Dacheng Tao, and Christos Faloutsos. 2006. Beyond Streams and Graphs: Dynamic Tensor Analysis. In *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '06)*. ACM, New York, NY, USA, 374–383. <https://doi.org/10.1145/1150402.1150445>
- [25] Nick Vannieuwenhoven, Karl Meerbergen, and Raf Vandebril. 2015. Computing the Gradient in Optimization Algorithms for the CP Decomposition in Constant Memory through Tensor Blocking. *SIAM Journal on Scientific Computing* 37, 3 (2015), C415–C438. <https://doi.org/10.1137/14097968X>

- [26] Nico Vervliet, Otto Debaels, Laurent Sorber, Marc Van Barel, and Lieven De Lathauwer. 2016. Tensorlab 3.0. (Mar. 2016). <http://www.tensorlab.net> Available online.