

Parallel Candecomp/Parafac Decomposition of Sparse Tensors Using Dimension Trees

Oguz Kaya, Bora Uçar

▶ To cite this version:

Oguz Kaya, Bora Uçar. Parallel Candecomp/Parafac Decomposition of Sparse Tensors Using Dimension Trees. SIAM Journal on Scientific Computing, 2018, 40 (1), pp.C99 - C130. 10.1137/16M1102744 . hal-01397464v2

HAL Id: hal-01397464 https://inria.hal.science/hal-01397464v2

Submitted on 17 Oct 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés. $\frac{1}{2}$

PARALLEL CANDECOMP/PARAFAC DECOMPOSITION OF SPARSE TENSORS USING DIMENSION TREES*

3

OGUZ KAYA[†] AND BORA UÇAR[‡]

Abstract. CANDECOMP/PARAFAC (CP) decomposition of sparse tensors has been success-4 5 fully applied to many problems in web search, graph analytics, recommender systems, health care 6 data analytics, and many other domains. In these applications, efficiently computing the CP decomposition of sparse tensors is essential in order to be able to process and analyze data of massive scale. For this purpose, we investigate an efficient computation of the CP decomposition of sparse 8 9 tensors and its parallelization. We propose a novel computational scheme for reducing the cost of a core operation in computing the CP decomposition with the traditional alternating least squares 10 (CP-ALS) based algorithm. We then effectively parallelize this computational scheme in the context 11 12 of CP-ALS in shared and distributed memory environments, and propose data and task distribution 13 models for better scalability. We implement parallel CP-ALS algorithms and compare our imple-14 mentations with an efficient tensor factorization library using tensors formed from real-world and synthetic datasets. With our algorithmic contributions and implementations, we report up to 5.96x, 15 5.65x, and 3.9x speedup in sequential, shared memory parallel, and distributed memory parallel 16executions over the state of the art, and achieve strong scalability up to 4096 cores on an IBM 17 18 BlueGene/Q supercomputer.

19 Key words. sparse tensors, CP decomposition, dimension tree, parallel algorithms

20 AMS subject classifications. 15-04, 05C70, 15A69, 15A83

1. Introduction. With growing features and dimensionality of data, tensors, or 21 multi-dimensional arrays, have been increasingly used in many fields including the analysis of Web graphs [28], knowledge bases [10], recommender systems [36, 37, 43], 23 signal processing [30], computer vision [46], health care [34], and many others [29]. 24 25 Tensor decomposition algorithms are used as an effective tool for analyzing data in order to extract latent information within the data, or predict missing data elements. 26There have been considerable efforts in designing numerical algorithms for different 27 tensor decomposition problems (see the survey [29]), and algorithmic and software 28contributions go hand in hand with these efforts [2, 5, 16, 22, 26, 27, 42, 40]. 29

30 One of the well known tensor decompositions is the CANDECOMP/PARAFAC (CP) formulation, which approximates a given tensor as a sum of rank-one tensors. 31 Among the commonly used algorithms for computing a CP decomposition is CP-ALS [11, 19], which is based on the alternating least squares method, though other 33 variants also exist [1, 44]. These algorithms are iterative, in which the computa-34 35 tional core of each iteration involves a special operation called matricized tensor-times Khatri-Rao product (MTTKRP). When the input tensor is sparse and N dimensional, 36 MTTKRP operation amounts to element-wise multiplication of N-1 row vectors from 37 N-1 matrices and their scaled sum reduction according to the nonzero structure of the 38 tensor. As the dimensionality of the tensor increases, this operation gets computation-39 40 ally more expensive; hence, efficiently carrying out MTTKRP for higher dimensional tensors is of our particular interest in emerging applications [34]. This operation 41 has received recent interest for efficient execution in different settings such as MAT-42 LAB [2, 5], MapReduce [22], shared memory [42], and distributed memory [16, 26, 40]. 43

^{*}Submitted to the editors November 8, 2016. A preliminary version appeared in SC'15 [26].

[†]INRIA and LIP, UMR5668 (CNRS - ENS Lyon - UCBL - Université de Lyon - INRIA), Lyon, France (oguz.kaya@ens-lyon.fr).

[‡]CNRS and LIP, UMR5668 (CNRS - ENS Lyon - UCBL - Université de Lyon - INRIA), Lyon, France (bora.ucar@ens-lyon.fr)

O. KAYA AND B. UÇAR

44 We are interested in a fast computation of MTTKRP as well as CP-ALS for sparse 45 tensors using efficient computational schemes and effective parallelization in shared 46 and distributed memory environments.

Our contributions in this paper are as follows. We investigate the parallelization 47 of CP-ALS algorithm for sparse tensors in shared and distributed memory systems. 48 For the shared-memory computations, we propose a novel computational scheme that 49 significantly reduces the computational cost while offering an effective parallelism. We 50then perform theoretical analyses corresponding to the computational gains and the 51memory utilization, which are also validated with the experiments. We propose a finegrain distributed memory parallel algorithm, and compare it against a medium-grain 53 variant [40]. Finally, we discuss effective partitioning routines for these algorithms. 5455Even though the discussion is confined to CP-ALS in the paper, the contributions apply to any other algorithm involving MTTKRP in its core [1]. 56

The organization of the rest of the paper is as follows. In the next section, we introduce our notation, describe the CP-ALS method, and present a data structure 58 called dimension tree which enables efficient CP-ALS computations. Next, in section 3, we explain how to use dimension trees to carry out MTTKRPs within CP-ALS. 60 Afterwards, we discuss an effective shared and distributed memory parallelization of 61 CP-ALS iterations in section 4. We discuss partitioning methods pertaining to dis-62 tributed memory parallel performance, and use these partitioning methods in our 63 experiments using real-world tensors. In section 5, we give an overview of the existing 64 literature. Finally, we present experimental results in section 6 to demonstrate per-66 formance gains using our algorithms with shared and distributed memory parallelism over an efficient state of the art implementation, and then conclude the paper.

68 2. Background and notation.

2.1. Tensors and CP-ALS. We denote the set $\{1, \ldots, M\}$ of integers as \mathbb{N}_M 69 for $M \in \mathbb{Z}^+$. For vectors, we use bold lowercase Roman letters, as in **x**. For matrices, 70 we use bold uppercase Roman letters, e.g., \mathbf{X} . For tensors, we generally follow the 7172notation in Kolda and Bader's survey [29]; particularly, we use bold calligraphic fonts, e.g., \mathcal{X} , to represent tensors. The *order* of a tensor is defined as the number of its 73 dimensions, or equivalently, modes, which we denote by N. A slice of a tensor in the 74 nth mode is a set of tensor elements obtained by fixing the index only along the nth mode. We use the MATLAB notation to refer to matrix rows and columns as well 77 as tensors slices, e.g., $\mathbf{X}(i, :)$ and $\mathbf{X}(:, j)$ are the *i*th row and the *j*th column of \mathbf{X} , whereas $\mathcal{X}(:,:,k)$ represents the kth slice of \mathcal{X} in the third dimension. We use italic 78 lowercase letters with subscripts to represent vector, matrix, and tensor elements, 79 e.g., x_i for a vector \mathbf{x} , $x_{i,j}$ for a matrix \mathbf{X} , and $x_{i,j,k}$ for a 3-dimensional tensor \mathcal{X} . 80 For the column vectors of a matrix, we use the same letter in lowercase and with a 81 subscript corresponding to the column index, e.g., \mathbf{x}_i to denote $\mathbf{X}(:,i)$, whereas a row 82 of a matrix is always expressed in MATLAB notation, as in $\mathbf{X}(i, :)$. Sets, lists, trees, 83 and hypergraphs are expressed in non-bold calligraphic fonts. 84

Let $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ be an *N*-mode tensor whose size in mode *n* is I_n for $n \in \mathbb{N}_N$. The multiplication of \mathcal{X} along the mode *n* with a vector $\mathbf{v} \in \mathbb{R}^{I_n}$ is a tensor $\mathcal{Y} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times 1 \times I_{n+1} \times \cdots \times I_N}$ with elements

88 (1)
$$y_{i_1,\dots,i_{n-1},1,i_{n+1},\dots,i_N} = \sum_{j=1}^{I_n} v_j x_{i_1,\dots,i_{n-1},j,i_{n+1},\dots,i_N}$$

89 This operation is called tensor-times-vector multiply (TTV) and is denoted by $\boldsymbol{\mathcal{Y}}=$

90 $\mathcal{X} \times_n \mathbf{v}$. The order of a series of TTVs is irrelevant, i.e., $\mathcal{X} \times_i \mathbf{u} \times_j \mathbf{v} = \mathcal{X} \times_j \mathbf{v} \times_i \mathbf{u}$ 91 for $\mathbf{u} \in \mathbb{R}^{I_i}$, $\mathbf{v} \in \mathbb{R}^{I_j}$, $i \neq j$, and $i, j \in \mathbb{N}_N$.

A tensor \mathcal{X} can be *matricized*, meaning that a matrix **X** can be associated with 92 \mathcal{X} by identifying a subset of its modes with the rows of **X**, and the rest of the modes with the columns of **X**. This involves a mapping of the elements of \mathcal{X} to those of 94 **X**. We will be exclusively dealing with the matricizations of tensors along a single 95 mode, meaning that a single mode is mapped to the rows of the resulting matrix, 96 and the rest of the modes correspond to the columns of the resulting matrix. We use 97 $\mathbf{X}_{(d)}$ to denote the matricization along mode d, e.g., for $\boldsymbol{\mathcal{X}} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, the matrix 98 $\mathbf{X}_{(1)}$ denotes the mode-1 matricization of $\boldsymbol{\mathcal{X}}$. Specifically in this matricization, the 99 tensor element x_{i_1,\ldots,i_N} corresponds to the element $\left(i_1, i_2 + \sum_{j=3}^N \left[(i_j - 1) \prod_{k=2}^{j-1} I_k\right]\right)$ 100 of $\mathbf{X}_{(1)}$. Matricizations in other modes are defined similarly. 101 The Hadamard product of two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{I}$ is a vector $\mathbf{w} = \mathbf{u} * \mathbf{v}, \mathbf{w} \in \mathbb{R}^{I}$, 102where $w_i = u_i \cdot v_i$. The outer product of K > 1 vectors $\mathbf{u}^{(1)}, \ldots, \mathbf{u}^{(K)}$ of corresponding 103 sizes I_1, \ldots, I_K is denoted by $\mathcal{X} = \mathbf{u}^{(1)} \circ \cdots \circ \mathbf{u}^{(K)}$ where $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_K}$ is a K-104 105

dimensional tensor with elements $x_{i_1,...,i_K} = \prod_{k \in \mathbb{N}_K} u_{i_k}^{(k)}$. The Kronecker product of vectors $\mathbf{u} \in \mathbb{R}^I$ and $\mathbf{v} \in \mathbb{R}^J$ results in a vector $\mathbf{w} = \mathbf{u} \otimes \mathbf{v}, \mathbf{w} \in \mathbb{R}^{IJ}$ defined as

107
$$\mathbf{w} = \mathbf{u} \otimes \mathbf{v} = \begin{bmatrix} u_1 \mathbf{v} \\ u_2 \mathbf{v} \\ \vdots \\ u_I \mathbf{v} \end{bmatrix}$$

108 For matrices $\mathbf{U} \in \mathbb{R}^{I \times K}$ and $\mathbf{V} \in \mathbb{R}^{J \times K}$, their *Khatri-Rao product* corresponds to

109 (2)
$$\mathbf{W} = \mathbf{U} \odot \mathbf{V} = [\mathbf{u}_1 \otimes \mathbf{v}_1, \dots, \mathbf{u}_K \otimes \mathbf{v}_K],$$

110 where $\mathbf{W} \in \mathbb{R}^{IJ \times K}$.

111 For the operator \circ , we use the shorthand notation $\circ_{i\neq n} \mathbf{U}^{(i)}$ to denote the opera-112 tion $\mathbf{U}^{(1)} \circ \cdots \circ \mathbf{U}^{(n-1)} \circ \mathbf{U}^{(n+1)} \circ \cdots \circ \mathbf{U}^{(N)}$ over a set $\{\mathbf{U}^{(1)}, \ldots, \mathbf{U}^{(N)}\}$ of matrices (and 113 similarly for vectors). Similarly, $\mathcal{X} \times_{i \in \mathcal{I}} \mathbf{u}^{(i)}$ denotes the operation $\mathcal{X} \times_{i_1} \mathbf{u}^{(i_1)} \times_{i_2}$ 114 $\cdots \times_{i_{|\mathcal{I}|}} \mathbf{u}^{(i_{|\mathcal{I}|})}$ over a set $\mathcal{I} = \{i_1, \ldots, i_{|\mathcal{I}|}\}$ of dimensions using a set $\{\mathbf{u}^{(1)}, \ldots, \mathbf{u}^{(N)}\}$ 115 of vectors.

2.2. CP decomposition. The rank-R CP-decomposition of a tensor \mathcal{X} ex-116 presses or approximates \mathcal{X} as a sum of R rank-1 tensors. For instance, for $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$, we obtain $\mathcal{X} \approx \sum_{r=1}^{R} \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$ where $\mathbf{a}_r \in \mathbb{R}^I$, $\mathbf{b}_r \in \mathbb{R}^J$, and $\mathbf{c}_r \in \mathbb{R}^K$. This decomposition results in the element-wise approximation (or equality) $x_{i,j,k} \approx$ 117 118119 $\sum_{r=1}^{R} a_{ir} b_{jr} c_{kr}$. The minimum R value rendering this approximation an equality is 120 called as the rank (or CP-rank) of the tensor \mathcal{X} , and computing this value is NP-121hard [21]. Here, the matrices $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R], \mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R]$, and $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_R]$ 122are called the *factor matrices*, or *factors*. For *N*-mode tensors, we use $\mathbf{U}^{(1)}, \ldots, \mathbf{U}^{(N)}$ to refer to the factor matrices having I_1, \ldots, I_N rows and *R* columns, and $\mathbf{u}_i^{(i)}$ to refer 123124to the *j*th column of $\mathbf{U}^{(i)}$. The standard algorithm for computing a CP decomposition 125is the alternating least squares (CP-ALS) method, which establishes a good trade-off 126between the number of iterations and the cost per iteration [29]. It is an iterative 127 algorithm, shown in Algorithm 1, that progressively updates the factors $\mathbf{U}^{(n)}$ in an 128 alternating fashion starting from an initial guess. CP-ALS runs until it can no longer 129 improve the solution, or it reaches the allowed maximum number of iterations. The 130initial factor matrices can be randomly set, or computed using the truncated SVD of 131

O. KAYA AND B. UÇAR

132 the matricizations of \mathcal{X} [29]. Each iteration of CP-ALS consists of N subiterations,

where in the *n*th subiteration $\mathbf{U}^{(n)}$ is updated using \mathcal{X} and the current values of all other factor matrices.

Algorithm 1 CP-ALS: ALS algorithm for computing CP decomposition

Input: $\boldsymbol{\mathcal{X}}$: An *N*-mode tensor, $\boldsymbol{\mathcal{X}} \in \mathbb{R}^{I_1,...,I_N}$ R: The rank of CP decomposition $\mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}$: Initial factor matrices **Output:** $[\![\lambda]; \mathbf{U}^{(1)}, \ldots, \mathbf{U}^{(N)}]\!]$: A rank-*R* CP decomposition of $\boldsymbol{\mathcal{X}}$ 1: repeat for $n = 1, \ldots, N$ do 2: $\mathbf{M}^{(n)} \leftarrow \mathbf{X}_{(n)} \odot_{i \neq n} \mathbf{U}^{(i)}$ 3: $\mathbf{H}^{(n)} \leftarrow *_{i \neq n} (\mathbf{U}^{(i)T} \mathbf{U}^{(i)})$ 4: $\mathbf{U}^{(n)} \leftarrow \mathbf{M}^{(n)} \mathbf{H}^{(n)^{\dagger}}$ ▶ $\mathbf{H}^{(n)^{\dagger}}$ is the pseudoinverse of $\mathbf{H}^{(n)}$. 5: $\boldsymbol{\lambda} \leftarrow \text{COLUMN-NORMALIZE}(\mathbf{U}^{(n)})$ 6: 7: until convergence is achieved or the maximum number of iterations is reached. 8: return $[\![\boldsymbol{\lambda}; \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}]\!]$

Computing the matrix $\mathbf{M}^{(n)} \in \mathbb{R}^{I_n \times R}$ at Line 3 of Algorithm 1 is the sole part 135involving the tensor \mathcal{X} , and it is the most expensive computational step, for both 136 sparse and dense tensors. The operation $\mathbf{X}_{(n)} \odot_{i \neq n} \mathbf{U}^{(i)}$ is called *matricized tensor*-137times Khatri-Rao product (MTTKRP). The Khatri-Rao product of the involved $\mathbf{U}^{(n)}$ s 138defines a matrix of size $(\prod_{i \neq n} I_i) \times R$ according to (2), and can get very costly in 139 terms of computational and memory requirements when I_i or N is large—which is 140 the case for many real-world sparse tensors. To alleviate this, various methods are 141 142 proposed in the literature that enable performing MTTKRP without forming Khatri-Rao product. One such formulation [4], also used in Tensor Toolbox [5], expresses 143MTTKRP in terms of a series of TTVs, and computes the resulting matrix $\mathbf{M}^{(n)}$ 144column by column. With this formulation, the rth column of $\mathbf{M}^{(n)}$ can be computed 145using N-1 TTVs as in $\mathbf{M}^{(n)}(:,r) \leftarrow \mathcal{X} \times_{i \neq n} \mathbf{u}_r^{(i)}$, or equivalently, 146

147 (3) $\mathbf{M}^{(n)}(:,r) \leftarrow \mathcal{X} \times_1 \mathbf{u}_r^{(1)} \times_2 \cdots \times_{n-1} \mathbf{u}_r^{(n-1)} \times_{n+1} \mathbf{u}_r^{(n+1)} \times_{n+2} \cdots \times_N \mathbf{u}_r^{(N)}$.

Once $\mathbf{M}^{(n)}$ is obtained, the Hadamard product of matrices $\mathbf{U}^{(i)}^T \mathbf{U}^{(i)}$ of size 148 $R \times R$ is computed for $1 \leq i \leq N, i \neq n$ to form the matrix $\mathbf{H}^{(n)} \in \mathbb{R}^{R \times R}$. Note 149 that within the *n*th subiteration, only $\mathbf{U}^{(n)}$ is updated among all factor matrices. 150Therefore, for efficiency, one can precompute all matrices $\mathbf{U}^{(i)}{}^{T}\mathbf{U}^{(i)}$ of size $R \times R$ 151for $i \in \mathbb{N}_N$, then update $\mathbf{U}^{(n)T}\mathbf{U}^{(n)}$ once $\mathbf{U}^{(n)}$ changes. As in many cases the rank 152R of approximation is chosen as a small constant in practice for sparse tensors (less 153than 50) [47], performing these Hadamard products to compute $\mathbf{H}^{(n)}$ and the matrix-154matrix multiplication to compute $\mathbf{U}^{(n)T}\mathbf{U}^{(n)}$ become relatively cheap compared with 155the TTV step. Once both $\mathbf{M}^{(n)}$ and $\mathbf{H}^{(n)}$ are computed, another matrix-matrix 156multiplication is performed using $\mathbf{M}^{(n)}$ and the pseudoinverse of $\mathbf{H}^{(n)}$ in order to 157update the matrix $\mathbf{U}^{(n)}$, which is not expensive when R is small. Finally, $\mathbf{U}^{(n)}$ is 158normalized column-wise, and the column vector norms are stored in a vector $\boldsymbol{\lambda} \in \mathbb{R}^{R}$. 159The convergence is achieved when the relative reduction in the norm of the error, 160 i.e., $\|\boldsymbol{\mathcal{X}} - \sum_{r=1}^{R} \lambda_r (\mathbf{u}_r^{(1)} \circ \cdots \circ \mathbf{u}_r^{(N)})\|$, is small. The cost of this computation is 161 insignificant. 162

4

163 2.3. Hypergraphs and hypergraph partitioning. We partition the data and 164the computation using the standard hypergraph partitioning tools. Necessary definitions follow. 165

A hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ is a set \mathcal{V} of vertices and a set \mathcal{E} of hyperedges. Each 166 hyperedge is a subset of \mathcal{V} . Weights, denoted with $w[\cdot]$, and costs, denoted with $c[\cdot]$, 167 can be associated with, respectively, the vertices and the hyperedges of \mathcal{H} . For a given 168 integer $K \geq 2$, a K-way vertex partition of a hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ is denoted by 169 $\Pi = \{\mathcal{V}_1, \ldots, \mathcal{V}_K\}$, where the parts are non-empty, i.e., $\mathcal{V}_k \neq \emptyset$ for $k \in \mathbb{N}_K$; mutually 170exclusive, i.e., $\mathcal{V}_k \cap \mathcal{V}_\ell = \emptyset$ for $k \neq \ell$; and collectively exhaustive, i.e., $\mathcal{V} = \bigcup \mathcal{V}_k$. Let $W_k = \sum_{v \in \mathcal{V}_k} w[v]$ be the total vertex weight in \mathcal{V}_k , and $W_{avg} = \sum_{v \in \mathcal{V}} w[v]/K$ 171

172denote the average part weight. If each part $\mathcal{V}_k \in \Pi$ satisfies the balance criterion 173

174 (4)
$$W_k \leq W_{avg}(1+\varepsilon)$$
 for $k \in \mathbb{N}_K$,

we say that Π is *balanced* where ε represents the allowed maximum imbalance ratio. 175In a partition Π , a hyperedge that has at least one vertex in a part is said to 176 connect that part. The number of parts connected by a hyperedge h is called its 177connectivity, and is denoted by κ_h . Given a vertex partition Π of a hypergraph 178 $\mathcal{H} = (\mathcal{V}, \mathcal{E})$, one can measure the *cutsize* metric induced by Π as 179

180 (5)
$$\chi(\Pi) = \sum_{h \in \mathcal{E}} c[h](\kappa_h - 1) .$$

This cut measure is called the *connectivity-1* cutsize metric. 181

Given $\varepsilon > 0$ and an integer K > 1, the standard hypergraph partitioning problem 182 is defined as the task of finding a balanced partition Π with K parts such that $\chi(\Pi)$ 183 is minimized. Hypergraph partitioning problem is NP-hard [31]. 184

185A common variant of the above problem is the *multi-constraint hypergraph par*titioning [15, 24]. In this variant, each vertex has an associated vector of weights. 186 The partitioning objective is the same as above, and the partitioning constraint is to 187 satisfy a balancing constraint for each weight. Let w[v, i] denote the C weights of a 188 vertex v for $i \in \mathbb{N}_C$. In this variant, the balance criterion (4) is rewritten as 189

190 (6)
$$W_{k,i} \leq W_{ava,i} (1+\varepsilon)$$
 for $k \in \mathbb{N}_K$ and $i \in \mathbb{N}_C$,

where the *i*th weight $W_{k,i}$ of a part \mathcal{V}_k is defined as the sum of the *i*th weights of the 191vertices in that part, i.e., $W_{k,i} = \sum_{v \in \mathcal{V}_k} w[v,i]$, and $W_{avg,i}$ represents the average part weight for the *i*th weight of all vertices, i.e., $W_{avg,i} = \sum_{v \in \mathcal{V}} w[v,i]/K$. 192193

2.4. Dimension tree. A *dimension tree* is a data structure that partitions the 194 195mode indices of an N-dimensional tensor in a hierarchical manner for computing tensor decompositions efficiently. It was first used in the hierarchical Tucker format 196representing the hierarchical Tucker decomposition of a tensor [18], which was intro-197 duced as a computationally feasible alternative to the original Tucker decomposition 198 for higher order tensors. We provide the formal definition of a dimension tree along 199200 with some basic properties as follows.

201 DEFINITION 1. A dimension tree \mathcal{T} for N dimensions is a tree with a root, denoted by ROOT(\mathcal{T}), and N leaf nodes, denoted by the set LEAVES(\mathcal{T}). In a dimension 202tree \mathcal{T} , each non-leaf node has at least two children, and each node $t \in \mathcal{T}$ is associated 203 with a mode set $\mu(t) \subseteq \mathbb{N}_N$ satisfying the following properties: 204

205 1. $\mu(\operatorname{ROOT}(\mathcal{T})) = \mathbb{N}_N$.

O. KAYA AND B. UÇAR

206	2.	For	each	nor	n-leaf	node	$t \in T$	Τ,	the	mode	sets of	^{r}its	children	partition	$\mu(t)$
	~	m 1		1 0				-		-	$\langle - \rangle$	-	(1)	()	

3. The nth leaf node, denoted by $l_n \in \text{LEAVES}(\mathcal{T})$, has $\mu(l_n) = \{n\}$.

For the simplicity of the discussion, we assume without loss of generality that the 208 sequence l_1, \ldots, l_N corresponds to a relative ordering of the leaf nodes in a post-order 209traversal of the dimension tree. If this is not the case, we can relabel tensor modes 210accordingly. We define the *inverse mode set* of a node t as $\mu'(t) = \mathbb{N}_N \setminus \mu(t)$. For each 211node t with a parent $\mathcal{P}(t), \mu(t) \subset \mu(\mathcal{P}(t))$ holds due to the second property, which in 212 turn yields $\mu'(t) \supset \mu'(\mathcal{P}(t))$. If a dimension tree has the height $\lceil \log(N) \rceil$ with its first 213 $\lfloor \log(N) \rfloor$ levels forming a complete binary tree, we call it a balanced binary dimension 214tree (BDT). In Figure 1, we show a BDT for 4 dimensions (associated with a sparse 215216 tensor described later).

3. Computing CP decomposition using dimension trees. In this section, 217we propose a novel way of using dimension trees for computing the standard CP 218decomposition of tensors with a formulation that asymptotically reduces the compu-219tational cost. In doing so, we do not alter the original CP decomposition in any way. 220 The reduction in the computational cost is made possible by storing partial TTV re-221 222 sults, and hence by trading off more memory. A similar idea of reusing partial results without the use of a tree framework was moderately explored by Baskaran et al. [6] 223 for computing the Tucker decomposition of sparse tensors, and by Phan et al. for 224 computing the CP decomposition of dense tensors [35]. We generalized the approach 225of Baskaran et al. [6] using dimension trees for better computational gains [25] in the 226227 standard algorithm for sparse Tucker decomposition. Here, we adopt the same data structure for reducing the cost of MTTKRP operations. 228

3.1. Using dimension trees to perform successive tensor-times-vector 229 multiplies. At each subiteration of the CP-ALS algorithm, \mathcal{X} is multiplied with 230the column vectors of matrices in N-1 modes using (3) in performing MTTKRP. 231 232 Some of these TTVs involve the same matrices as the preceding subiterations. As a series of TTVs can be done in any order, this opens up the possibility to factor out 233and reuse TTV results that are common in consecutive subiterations for reducing the 234computational cost. For instance, in the first subiteration of CP-ALS using a 4-mode tensor \mathcal{X} , we compute $\mathcal{X} \times_2 \mathbf{u}_r^{(2)} \times_3 \mathbf{u}_r^{(3)} \times_4 \mathbf{u}_r^{(4)}$ and eventually update $\mathbf{u}_r^{(1)}$ for each 235236 $r \in \mathbb{N}_R$, whereas in the second subiteration we compute $\mathcal{X} \times_1 \mathbf{u}_r^{(1)} \times_3 \mathbf{u}_r^{(3)} \times_4 \mathbf{u}_r^{(4)}$ 237and update $\mathbf{u}_r^{(2)}$. In these two subiterations, the matrices $\mathbf{U}^{(3)}$ and $\mathbf{U}^{(4)}$ remain 238unchanged, and both TTV steps involve the TTV of \mathcal{X} with $\mathbf{u}_r^{(3)}$ and $\mathbf{u}_r^{(4)}$. Hence, 239 we can compute $\mathcal{Y}_r = \mathcal{X} \times_3 \mathbf{u}_r^{(3)} \times_4 \mathbf{u}_r^{(4)}$, then reuse it in the first and the second subiterations as $\mathcal{Y}_r \times_2 \mathbf{u}_r^{(2)}$ and $\mathcal{Y}_r \times_1 \mathbf{u}_r^{(1)}$ to obtain the required TTV results. 240241

We use dimension trees to systematically detect and reuse such partial results by 242 associating a tree \mathcal{T} with an N-dimensional tensor \mathcal{X} as follows. With each node 243 $t \in \mathcal{T}$, we associate R tensors $\mathcal{X}_1^{(t)}, \ldots, \mathcal{X}_R^{(t)}$. $\mathcal{X}_r^{(t)}$ corresponds to the TTV result $\mathcal{X}_r^{(t)} = \mathcal{X} \times_{d \in \mu'(t)} \mathbf{u}_r^{(d)}$. Therefore, the inverse mode set $\mu'(t)$ corresponds to the set 244245of modes in which TTV is performed on \mathcal{X} to form $\mathcal{X}_r^{(t)}$. For the root of the tree, $\mu'(\text{ROOT}(\mathcal{T})) = \emptyset$; thus, all tensors of the root node correspond to the original tensor \mathcal{X} , i.e., $\mathcal{X}_r^{(\text{ROOT}(\mathcal{T}))} = \mathcal{X}$ for $r \in \mathbb{N}_R$. Since $\mu'(t) \supset \mu'(\mathcal{P}(t))$ for a node t and its 246 247 248 parent $\mathcal{P}(t)$, tensors of $\mathcal{P}(t)$ can be used as partial results to update the tensors of t. 249Let $\delta(t) = \mu'(t) \setminus \mu'(\mathcal{P}(t))$. We can then compute each tensor of t from its parent's 250as $\boldsymbol{\mathcal{X}}_{r}^{(t)} = \boldsymbol{\mathcal{X}}_{r}^{(\mathcal{P}(t))} \times_{d \in \delta(t)} \mathbf{u}_{r}^{(d)}$. This procedure is called DTREE-TTV and is shown in Algorithm 2. DTREE-TTV first checks if the tensors of t are already computed, 251252

and immediately returns if so. This happens, for example, when two children of t call 253

254DTREE-TTV on t consecutively, in which case in the second DTREE-TTV call, the

255tensors of t would be already computed. If the tensors of t are not already computed,

DTREE-TTV on $\mathcal{P}(t)$ is called first to make sure that $\mathcal{P}(t)$'s tensors are up-to-date. 256

Then, each $\mathcal{X}_r^{(t)}$ is computed by performing a TTV on the corresponding tensor $\mathcal{X}_r^{(\mathcal{P}(t))}$ of the parent. We use the notation $\mathcal{X}_{\cdot}^{(t)}$ to denote all R tensors of a node t. 257

258

Algorithm 2 DTREE-TTV: Dimension tree-based TTV with R vectors in each mode

Input: *t*: A node of the dimension tree **Output:** Tensors $\boldsymbol{\mathcal{X}}_{:}^{(t)}$ of t are computed. 1: if EXISTS($\boldsymbol{\mathcal{X}}^{(t)}$) then ▶ Tensors $\boldsymbol{\mathcal{X}}_{:}^{(t)}$ are already computed. 2: return • Compute the parent's tensors $\boldsymbol{\mathcal{X}}_{:}^{(\mathcal{P}(t))}$ first. 3: DTREE-TTV($\mathcal{P}(t)$) 4: for r = 1, ..., R do 5: $\boldsymbol{\mathcal{X}}_{r}^{(t)} \leftarrow \boldsymbol{\mathcal{X}}_{r}^{(\mathcal{P}(t))} \times_{d \in \delta(t)} \mathbf{u}_{r}^{(d)}$ \blacktriangleright Now update all tensors of t using parent's.

259**3.2.** Dimension tree-based CP-ALS algorithm. At the *n*th subiteration of Algorithm 1, we need to compute $\mathcal{X} \times_{i \neq n} \mathbf{u}_r^{(i)}$ for all $r \in \mathbb{N}_R$ in order to form 260 $\mathbf{M}^{(n)}$. Using a dimension tree \mathcal{T} with leaves l_1, \ldots, l_N , we can perform this simply by 261executing DTREE-TTV (l_n) , after which the rth tensor of l_n provides the rth column of 262 $\mathbf{M}^{(n)}$. Once $\mathbf{M}^{(n)}$ is formed, the remaining steps follow as before. We show the whole 263CP-ALS using a dimension tree in Algorithm 3. At Line 1, we construct a dimension 264tree \mathcal{T} with the leaf order l_1, \ldots, l_N obtained from a post-order traversal of \mathcal{T} . This 265tree can be constructed in any way that respects the properties of a dimension tree 266267described in section 3: but for our purposes we assume that it is formed as a BDT. At Line 8 within the *n*th subiteration, we destroy all tensors of a node t if its set 268of multiplied modes $\mu'(t)$ involve n, as in this case its tensors involve multiplication 269using the old value of $\mathbf{U}^{(n)}$ which is about to change. Note that this step destroys the 270tensors of all nodes not lying on a path from l_n to the root. Afterwards, DTREE-TTV 271 is called at Line 9 for the leaf node l_n to compute its tensors. This step computes (or 272273reuses) the tensors of all nodes from the path from l_n to the root. Next, the rth column of $\mathbf{M}^{(n)}$ is formed using $\mathcal{X}_r^{(l_n)}$ for $r \in \mathbb{N}_R$. Once $\mathbf{M}^{(n)}$ is ready, $\mathbf{H}^{(n)}$ and $\mathbf{U}^{(n)}$ 274are computed as before, after which $\mathbf{U}^{(n)}$ is normalized. 275

Performing TTVs in CP-ALS using a BDT in this manner provides significant 276computational gains with a moderate increase in the memory cost. We now state two 277theorems pertaining to the computational and memory efficiency of DTREE-CP-ALS. 278

THEOREM 2. Let \mathcal{X} be an N-mode tensor. The total number of TTVs at each 279280 iteration of Algorithm 3 using a BDT is at most $RN \lceil \log N \rceil$.

Proof. As we assume that the sequence l_1, \ldots, l_N is obtained from a post-order 281 traversal of \mathcal{T} , for each internal node t, the subtree rooted at t has the leaves 282 $l_i, l_{i+1}, \ldots, l_{i+k-1}$ corresponding to k consecutive mode indices for some positive 283integers i and k. As we have $\mu(l_i) = i$ for the *i*th leaf node, we obtain $\mu(t) = i$ 284 $\{i, i+1, \ldots, i+k-1\}$ due to the second property of dimension trees. As a result, 285for each leaf node $l_{i+k'}$ for $0 \le k' < k$, we have $i + k' \in \mu(t)$; hence $i + k' \notin \mu'(t)$ 286 as $\mu'(t) = \mathbb{N}_N \setminus \mu(t)$. Therefore, within an iteration of Algorithm 3, the tensors of t 287get computed at the *i*th subiteration, stay valid (not destroyed) and get reused until 288 the i + k - 1th subiteration, and finally get destroyed in the following subiteration. 289290 Once destroyed, the tensors of t are never recomputed in the same iteration, as all

7

Algorithm 3 DTREE-CP-ALS: Dimension tree-based CP-ALS algorithm

Input: $\boldsymbol{\mathcal{X}}$: An *N*-mode tensor R: The rank of CP decomposition **Output:** $[\lambda; \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}]$: Rank-*R* CP decomposition of $\boldsymbol{\mathcal{X}}$ 1: $\mathcal{T} \leftarrow \text{CONSTRUCT-DIMENSION-TREE}(\mathcal{X})$ ▶ The tree has leaves $\{l_1, \ldots, l_N\}$. 2: for n = 2...N do $\mathbf{W}^{(n)} \leftarrow \mathbf{U}^{(n)T}\mathbf{U}^{(n)}$ 3: 4: repeat for $n = 1, \ldots, N$ do 5:for all $t \in \mathcal{T}$ do 6: 7: if $n \in \mu'(t)$ then DESTROY($\boldsymbol{\mathcal{X}}_{\cdot}^{(t)}$) \blacktriangleright Destroy all tensors that are multiplied by $\mathbf{U}^{(n)}$. 8: DTREE-TTV (l_n) ▶ Perform the TTVs for the leaf node tensors. 9: for r = 1, ..., R do Form $\mathbf{M}^{(n)}$ column-by-column (done implicitly). 10: $\mathbf{M}^{(n)}(:,r) \leftarrow \boldsymbol{\mathcal{X}}_{r}^{(l_{n})}$ $\blacktriangleright \mathcal{X}_{r}^{(l_{n})}$ is a vector of size I_{n} . 11: $\mathbf{H}^{(n)} \leftarrow *_{i \neq n} \mathbf{W}^{(i)}$ 12: $\mathbf{U}^{(n)} \leftarrow \mathbf{M}^{(n)} {\mathbf{H}^{(n)}}^{\dagger}$ 13: $\boldsymbol{\lambda} \leftarrow \text{Column-Normalize}(\mathbf{U}^{(n)})$ 14: $\mathbf{W}^{(n)} \leftarrow \mathbf{U}^{(n)}{}^T \mathbf{U}^{(n)}$ 15:16: **until** converge is achieved or the maximum number of iterations is reached. 17: return $[\![\boldsymbol{\lambda}; \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}]\!]$

the modes associated with the leaf tensors in its subtree are processed (which are the 291 only nodes that can reuse the tensors of t). As a result, in one CP-ALS iteration, 292tensors of every tree node (except the root) get to be computed and destroyed exactly 293once. Therefore, the total number of TTVs in one iteration becomes the sum of the 294 295number of TTVs performed to compute the tensors of each node in the tree once. In computing its tensors, every node t has R tensors, and for each tensor it performs 296TTVs for each dimension in the set $\delta(t)$ in Algorithm 2, except the root node, whose 297tensors are all equal to \mathcal{X} and never change. Therefore, we can express the total 298number of TTVs performed within a CP-ALS iteration due to one of these R tensors 299300 as

301 (7)
$$\sum_{t \in \mathcal{T} \setminus \{\operatorname{Root}(\mathcal{T})\}} |\delta(t)| = \sum_{t \in \mathcal{T} \setminus \{\operatorname{Root}(\mathcal{T})\}} |\mu(\mathcal{P}(t)) \setminus \mu(t)|.$$

Since in a BDT every non-leaf node t has exactly two children, say t_1 and t_2 , we obtain $|\mu(t) \setminus \mu(t_1)| + |\mu(t) \setminus \mu(t_2)| = |\mu(t)|$, as $\mu(t)$ is partitioned into two sets $\mu(t_1)$ and $\mu(t_2)$. With this observation, we can reformulate (7) as

305 (8)
$$\sum_{t \in \mathcal{T} \setminus \{\text{Root}(\mathcal{T})\}} |\mu(\mathcal{P}(t)) \setminus \mu(t)| = \sum_{t \in \mathcal{T} \setminus \text{Leaves}(\mathcal{T})} |\mu(t)| .$$

Note that in constructing a BDT, at the root node we start with the mode set $\mu(\text{ROOT}(\mathcal{T})) = \mathbb{N}_N$. Then, at each level k > 0, we form the mode sets of the nodes at level k by partitioning the mode sets of their parents at level k - 1. As a result, at each level k, each dimension $n \in \mathbb{N}_N$ can appear in only one set $\mu(t)$ for a node t belonging to the level k of the BDT. With this observation in mind, rewriting (8) by

iterating over nodes by levels of the BDT yields 311

312
$$\sum_{t \in \mathcal{T} \setminus \text{LEAVES}(\mathcal{T})} |\mu(t)| = \sum_{k=1} \sum_{t \in \mathcal{T} \setminus \text{LEAVES}(\mathcal{T}), \text{LEVEL}(t) = k} |\mu(t)|$$

313
$$\leq \sum^{\lceil \log N \rceil} N = N \lceil \log N \rceil .$$

314

315 As there are R tensors in each BDT tree node, the overall cost becomes $RN[\log N]$ TTVs for a CP-ALS iteration. Π 316

k=1

 $\lceil \log N \rceil$

In comparison, the traditional scheme [42] incurs R(N-1) TTVs in each mode, 317 and RN(N-1) TTVs in total in an iteration. This yields a factor of $(N-1)/\log N$ 318 319 reduction in the number performed of TTVs using dimension trees. We note that in terms of the actual computational cost, this corresponds to a lower bound on the ex-320 pected speedup for the following reason. As the tensor is multiplied in different dimen-321 sions, resulting tensors are expected to have many index overlaps, effectively reducing their number of nonzeros and rendering subsequent TTVs significantly cheaper. This 323 renders using a BDT much more effective as it avoids repeating such expensive TTVs 324 at the higher levels of the dimension tree by reusing partial results. That is, the for-325 mula for the potential gain is a complicated function, depending on the sparsity of the 326 tensor. On one extreme, multiplying the tensor in certain dimensions might create no 327 or few index overlaps, which makes the cost of each TTV approximately equal, yield-328 ing the stated speedup. On the other extreme, the first TTVs performed the original 329 tensor may drastically reduce the number of tensor elements so that the cost of the 330 subsequent TTVs becomes negligible. The traditional scheme multiplies the original 331 tensor N times, once per dimension, whereas a BDT suffices with 2 such TTVs as 332 the root node has only two children, yielding a speedup factor of N/2. Therefore, in 333 practice the actual speedup is expected to be between these two extremes depending 334 335 on the sparsity of the tensor, and having more speedup with higher index overlap after multiplications. 336

For sparse tensors, one key idea we use for obtaining high performance is per-337 forming TTVs for all R tensors $\mathcal{X}^{(t)}$ of a node $t \in \mathcal{T}$ in a vectorized manner. We 338 illustrate this on a 4-dimensional tensor \mathcal{X} and a BDT, and for clarity, we put the mode set $\mu(t)$ of each tree node t in the subscript, as in t_{1234} . Let t_{1234} represent 340 the root of the BDT with $\mathcal{X}_r^{(t_{1234})} = \mathcal{X}$ for all $r \in \mathbb{N}_R$. The two children of t_{1234} 341 are t_{12} and t_{34} with the corresponding tensors $\mathcal{X}_r^{(t_{12})} = \mathcal{X}_r^{(t_{1234})} \times_3 \mathbf{u}_r^{(3)} \times_4 \mathbf{u}_r^{(4)}$ and $\mathcal{X}_r^{(t_{34})} = \mathcal{X}_r^{(t_{1234})} \times_1 \mathbf{u}_r^{(1)} \times_2 \mathbf{u}_r^{(2)}$, respectively. Since $\mathcal{X}_r^{(t_{1234})}$ are identical for all 342 343 $r \in \mathbb{N}_R$, the nonzero pattern of tensors $\mathcal{X}_r^{(t_{12})}$ are also identical. This is also the case 344 for $\mathcal{X}_{r}^{(t_{34})}$, and the same argument applies to the children t_1 and t_2 of t_{12} , as well 345 as t_3 and t_4 of t_{34} . As a result, each node in the tree involves R tensors with identical nonzero patterns. This opens up two possibilities in terms of efficiency. First, 347 it is sufficient to compute only one set of nonzero indices for each node $t \in \mathcal{T}$ to 348 represent the nonzero structure of all of its tensors, which reduces the computational 349 and memory cost by a factor of R. Second, we can perform the TTVs for all tensors 350 at once in a "vectorized" manner by modifying (1) to perform R TTVs of the form 351 $\boldsymbol{\mathcal{Y}}_r \leftarrow \boldsymbol{\mathcal{X}}_r \times_d \mathbf{v}_r$ for $\mathbf{V} = [\mathbf{v}_1 | \cdots | \mathbf{v}_R] \in \mathbb{R}^{I_d \times R}$ as 352

353 (9)
$$\mathbf{y}_{i_1,\dots,i_{d-1},1,i_{d+1},\dots,i_N}^{(:)} = \sum_{j=1}^{I_d} \mathbf{V}(j,:) * \mathbf{x}_{i_1,\dots,i_{d-1},j,i_{d+1},\dots,i_N}^{(:)},$$

where $\mathbf{y}_{i_1,\ldots,i_{d-1},1,i_{d+1},\ldots,i_N}^{(:)}$ and $\mathbf{x}_{i_1,\ldots,i_{d-1},j,i_{d+1},\ldots,i_N}^{(:)}$ are vectors of size R with elements $y_{i_1,\ldots,i_{d-1},1,i_{d+1},\ldots,i_N}^{(r)}$ and $x_{i_1,\ldots,i_{d-1},j,i_{d+1},\ldots,i_N}^{(r)}$ in \mathcal{Y}_r and \mathcal{X}_r , for all $r \in \mathbb{N}_R$. We call this operation tensor-times-multiple-vector multiplication (TTMV) as R column vectors of \mathbf{V} are multiplied simultaneously with R tensors of identical nonzero patterns. We can similarly extend this formula to the multiplication $\mathcal{Z}_r \leftarrow \mathcal{Y}_r \times_e \mathbf{w}_r =$ $(\mathcal{X}_r \times_d \mathbf{v}_r) \times_e \mathbf{w}_r$ in two modes d and e, d < e, with matrices $\mathbf{V} \in \mathbb{R}^{I_d \times R}$ and $\mathbf{W} \in \mathbb{R}^{I_e \times R}$ as

361
$$\mathbf{z}_{i_{1},...,i_{d-1},1,i_{d+1},...,i_{e-1},1,i_{e+1},...,i_{N}}^{(:)} = \sum_{j_{2}=1}^{I_{e}} \mathbf{W}(j,:) * \mathbf{y}_{i_{1},...,i_{d-1},1,i_{d+1},...,i_{e-1},j_{2},i_{e+1},...,i_{N}}^{(:)}$$

362 (10) $= \sum_{(j_{1},j_{2})=(1,1)}^{(I_{d},I_{e})} \mathbf{V}(j_{1},:) * \mathbf{W}(j_{2},:) * \mathbf{x}_{i_{1},...,i_{d-1},j_{1},i_{d+1},...,i_{e-1},j_{2},i_{e+1},...,i_{N}}^{(:)}$

The formula similarly generalizes to any number of dimensions, where each dimension 363 adds another Hadamard product with a corresponding matrix row. This "thick" mode 364 of operation provides a significant performance gain thanks to the increase in locality. 365Also, performing TTVs in this manner in CP-ALS effectively reduces the $RN \log N$ 366 TTVs required in Theorem 2 to $N[\log N]$ TTMV calls within an iteration. In our 367 approach, for each node $t \in \mathcal{T}$, we store a single list \mathcal{I}_t containing, for each $k \in \mathbb{N}_{|\mathcal{I}_t|}$, an 368 index tuple of the form $\mathcal{I}_t(k) = (i_1, \ldots, i_N)$, to represent the nonzeros $x_{i_1, \ldots, i_N}^{(r)} \in \mathcal{X}_r^{(t)}$ for all $r \in \mathbb{N}_R$. Also, for each such index tuple we hold a vector of size R corresponding 369 370 to the values of this nonzero in each one of R tensors, and we denote this value vector 371 as $\mathbf{V}_t(k,:)$, where $\mathbf{V}_t \in \mathbb{R}^{|\mathcal{I}_t| \times R}$ is called the *value matrix* of t. Finally, carrying 372 out the summation (9) requires a "mapping" $\mathcal{R}_t(k)$ indicating the set of elements of 373 $\mathcal{X}_{:}^{(P(t))}$ contributing to this particular element of $\mathcal{X}_{:}^{(t)}$. Altogether, we represent the 374 sparse tensors $\mathcal{X}^{(t)}$ of each tree node t with the tuple $(\mathcal{I}_t, \mathcal{R}_t, \mathbf{V}_t)$ whose computational 375 details follow next. 376

The sparsity of input tensor \mathcal{X} determines the sparsity structure of tensors in the dimension tree, i.e., \mathcal{I}_t and \mathcal{R}_t , for each tree node t. Computing this sparsity structure for each TTV can get very expensive, and is redundant. As \mathcal{X} stays fixed, we can compute the sparsity of each tree tensor once, and reuse it in all CP-ALS iterations. To this end, we need a data structure that can express this sparsity, while exposing parallelism to update the tensor elements in numerical TTV computations. We now describe computing this data structure in what we call the *symbolic TTV* step.

3.2.1. Symbolic TTV. For simplicity, we proceed with describing how to perform the symbolic TTV using the same 4-dimensional tensor \mathcal{X} and the BDT. The approach naturally generalizes to any N-dimensional tensor and dimension tree.

The first information we need is the list of nonzero indices \mathcal{I}_t for each node t in 387 a dimension tree, which we determine as follows. As mentioned, the two children t_{12} 388 and t_{34} of t_{1234} have the corresponding tensors $\boldsymbol{\mathcal{X}}_{r}^{(t_{12})} = \boldsymbol{\mathcal{X}}_{r}^{(t_{1234})} \times_{3} \mathbf{u}_{r}^{(3)} \times_{4} \mathbf{u}_{r}^{(4)}$ and $\boldsymbol{\mathcal{X}}_{r}^{(t_{34})} = \boldsymbol{\mathcal{X}}_{r}^{(t_{1234})} \times_{1} \mathbf{u}_{r}^{(1)} \times_{2} \mathbf{u}_{r}^{(2)}$, respectively. Using (10), the nonzero indices of $\boldsymbol{\mathcal{X}}_{r}^{(t_{12})}$ 389 390 and $\mathcal{X}_{r}^{(t_{34})}$ take the form (i, j, 1, 1) and (1, 1, k, l), respectively (tensor indices in the 391 multiplied dimensions are always 1; hence we omit storing them in practice), and such 392 a nonzero index exists in these tensors only if there exists a nonzero $x_{i,j,k,l} \in \mathcal{X}_r^{(t_{1234})}$. 393 Determining the list $\mathcal{I}_{t_{12}}$ (or $\mathcal{I}_{t_{34}}$) can simply be done by starting with a list $\mathcal{I}_{t_{1234}}$ 394of tuples, then replacing each tuple (i, j, k, l) in the list with the tuple (i, j) (or with 395(k, l) for $\mathcal{I}_{t_{34}}$). This list may contain duplicates, which can be efficiently eliminated 396

Fig. 1: BDT of a 4-dimensional sparse tensor $\mathcal{X} \in \mathbb{R}^{4 \times 4 \times 4}$ having 7 nonzeros. Each closed box refers to a tree node. Within each node, the index array and the mode set corresponding to that node are given. The reduction sets of two nodes in the tree are indicated with the dashed lines.



by sorting. Once the list of nonzeros for t_{12} (or t_{34}) is determined, we proceed to detecting the nonzero patterns of its children t_1 and t_2 (or, t_3 and t_4).

To be able to carry out the numerical calculations (9) and (10) for each nonzero 399 index at a node t, we need to identify the set of nonzeros of the parent node $\mathcal{P}(t)$'s 400tensors that contribute to this nonzero. Specifically, at t_{12} , for each nonzero index 401 $\mathcal{I}_{t_{12}}(m) = (i, j)$ we need to bookmark all nonzero index tuples of t_{1234} of the form $\mathcal{I}_{t_{1234}}(n) = (i, j, k, l)$, as this is this set of nonzeros of $\mathcal{X}_{1234}^{(t_{1234})}$ that contribute to the 402403nonzero of $\mathcal{X}^{(t_{12})}_{\cdot}$ with index (i, j, 1, 1) in (10). Therefore, for each such index tuple of 404 t_{12} we need a *reduction set* $\mathcal{R}_{t_{12}}(m)$ which contains all such index tuples of the parent, 405i.e., $n \in \mathcal{R}_{t_{12}}(m)$. We determine these sets simultaneously with \mathcal{I}_t . In Figure 1, we 406 illustrate a sample BDT for a 4-dimensional sparse tensor with \mathcal{I}_t , shown with arrays, 407 and \mathcal{R}_t , shown using arrows. 408

Next, we provide the following theorem to help us analyze the computational and memory cost of symbolic TTV using a BDT for sparse tensors.

411 THEOREM 3. Let \mathcal{X} be an N-mode sparse tensor. The total number of index 412 arrays in a BDT of \mathcal{X} is at most $N(\lceil \log N \rceil + 1)$.

413 Proof. Each node t in the dimension tree holds an index array for each mode in 414 its mode set $\mu(t)$. As stated in the proof of Theorem 2, the total size of mode sets 415 at each tree level is at most N. Therefore, the total number of index arrays cannot 416 exceed ($\lceil \log N \rceil + 1$)N in a BDT.

Theorem 3 shows that the storage requirement for the tensor indices of a BDT cannot exceed ($\lceil \log N \rceil + 1$)-times the size of the original tensor in the coordinate

format (which has N index arrays of size $nnz(\mathcal{X})$), yielding the overall worst-case 419 memory cost $nnz(\mathcal{X})N(\lceil \log N \rceil + 1)$. Assuming that the tensor does not contain any 420 empty slices (which can otherwise be removed in a preprocessing step), the index 421 array of a leaf node (whose tensors are vectors) corresponding to mode n is simply 422 $[1, \ldots, I_n]$ hence need not be explicitly stored, which effectively reduces this cost to 423 $\operatorname{nnz}(\boldsymbol{\mathcal{X}}) N \log N$. This cost is indeed a pessimistic estimate for real-world tensors, as 424 with significant index overlap in non-root tree tensors, the total memory cost could 425 reduce to as low as $O(\operatorname{nnz}(\boldsymbol{\mathcal{X}})N)$. This renders the approach very suitable for higher 426 dimensional tensors. Another minor cost is the storage of reduction pointers, which 427necessitates one array per non-root tree node, taking 2N-2 arrays in total. The size 428 of these arrays similarly diminishes towards the leaves with a potential index overlap. 429In computing the symbolic TTV, we sort $|\mu(t)|$ index arrays for each node $t \in \mathcal{T}$. 430 In addition, for each non-root node t of a BDT, we sort an extra array to determine 431 the reduction set \mathcal{R}_t . In the worst case, each array can have up to $nnz(\mathcal{X})$ elements. 432 Therefore, combining with the number of index arrays as given in Theorem 3, the over-433 all worst case cost of sorting becomes $O((N \lceil \log N \rceil + 2N - 2) \operatorname{nnz}(\mathcal{X}) \log(\operatorname{nnz}(\mathcal{X}))) =$ 434 $O(N \log N \operatorname{nnz}(\boldsymbol{\mathcal{X}}) \log(\operatorname{nnz}(\boldsymbol{\mathcal{X}})))$. We note, however, that both the total index array 435 size and sorting cost are pessimistic overestimates, since the nonzero structure of real-436 world tensors exhibits significant locality in indices. For example, on two tensors from 437 our experiments (Delicious and Flickr), we observed a reduction factor of 2.57 and 5.5 438 in the number of nonzeros of the children of the root of the BDT. Consequently, the 439number of nonzeros in a node's tensors reduces dramatically as we approach towards 440 441 the leaves. In comparison, existing approaches [39] sort the original tensor once with a cost of $O(Nnnz(\mathcal{X}) \log(nnz(\mathcal{X})))$ at the expense of computing TTVs from scratch 442 in each CP-ALS iteration. 443

Symbolic TTV is a one-time computation whose cost is amortized. Normally, 444 choosing an appropriate rank R for a sparse tensor \mathcal{X} requires several executions of 445CP-ALS. Also, CP-ALS is known to be sensitive to the initialization of factor matrices; 446 447 therefore, it is often executed with multiple initializations [29]. In all of these use cases, the tensor $\boldsymbol{\mathcal{X}}$ is fixed; therefore, the symbolic TTV is required only once. Moreover, 448 CP-ALS usually has a number of iterations which involve many costly numeric TTV 449 calls. As a result, the cost of the subsequent numeric TTV calls over many iterations 450and many CP-ALS executions easily amortizes that of this symbolic preprocessing. 451Nevertheless, in case of need, this step can efficiently be parallelized in multiple ways. 452First, symbolic TTV is essentially a sorting of multiple index arrays; hence, one can 453use parallel sorting methods. Second, the BDT structure naturally exposes a coarser 454level of parallelism; once a node's symbolic TTV is computed, one can proceed with 455those of its children in parallel, and process the whole tree in this way. Finally, in a 456distributed memory setting where we partition the tensor to multiple processes, each 457 process can perform the symbolic TTV on its local tensor in parallel. We benefit only 458 from this parallelism in our implementation. 459

460 After symbolic TTV is performed, index arrays of all nodes in the tree stay fixed 461 and are kept throughout CP-ALS iterations. However, at each subiteration n, only 462 the value matrices \mathbf{V}_t of tree tensors which are necessary to compute DTREE-TTV (l_n) 463 are kept. The following theorem provides an upper bound on the number of such value 464 matrices, which gives an upper bound on the memory usage for tensor values.

465 THEOREM 4. For an N-mode tensor \mathcal{X} , the total number of tree nodes whose 466 value matrices are allocated is at most $\lceil \log N \rceil$ at any instant of Algorithm 3 using a 467 BDT.

Proof. Note that at the beginning of the nth subiteration of Algorithm 3, the 468 tensors of each node $t \in \mathcal{T}$ involving n in $\mu'(t)$ are destroyed at Line 8. These are 469exactly the tensors that do not lie in the path from the leaf l_n to the root, as they 470 do not involve n in their mode set μ . The TTMV result for l_n depends only on the 471nodes on this path from l_n to the root; therefore, at the end of the *n*th subiteration, 472 only the tensors of the nodes on this path will be computed using Algorithm 2. As 473 this path length cannot exceed $\left[\log N\right]$ in a BDT, the number of nodes whose value 474matrices are not destroyed cannot exceed $\lceil \log N \rceil$ at any instant of Algorithm 3. 475

Theorem 4 puts an upper bound of $\lceil \log N \rceil$ on the maximum number of allo-476cated value matrices, thus on the maximum memory utilization due to tensor values, 477in DTREE-CP-ALS. In the worst case, each value matrix may have up to $nnz(\boldsymbol{\mathcal{X}})$ 478 elements, requiring $nnz(\mathcal{X})R[\log N]$ memory in total to store the values. Combin-479ing with Theorems 2 and 3, in the worst case, the intermediate results increase the 480 memory requirement by a factor of $O(\log N(1 + R/N))$ (with respect to the storage 481 of \mathcal{X} in coordinate format), while the computational cost is reduced by a factor of 482 $O(N/\log N)$. These are the largest increase in the memory and the smallest decrease 483in the computational cost. In practice, we expect less increase in the memory and 484 485 more gains in the computational cost thanks to the overlap of indices towards the leaves. 486

4. Parallel CP-ALS for sparse tensors using dimension trees. We first present shared memory parallel algorithms involving efficient parallelization of the dimension tree-based TTMVs in subsection 4.1. Later, in subsection 4.2, we present distributed memory parallel algorithms that use this shared memory parallelization.

Algorithm 4 SMP-DTREE-TTV Input: t: A dimension tree node/tensor

	auto di la annononon di de nou	ic/ tensor
Ou	tput: Numerical values \mathbf{V}_t a	are computed.
1:	if $EXISTS(\mathbf{V}_t)$ then	
2:	return	▶ Numerical values \mathbf{V}_t are already computed.
3:	SMP-DTREE-TTV $(\mathcal{P}(t))$	► Compute the parent's values $\mathbf{V}_{\mathcal{P}(t)}$ first.
4:	parallel for $1 \leq i \leq \mathcal{I}_t $ do	▶ Process each $\mathcal{I}_t(i) = (i_1, \ldots, i_N)$ in parallel.
5:	$\mathbf{V}_t(i,:) \leftarrow 0$	▶ Initialize with a zero vector of size $1 \times R$.
6:	for all $j \in \mathcal{R}_t(i)$ do \blacktriangleright R	educe from the element with index $\mathcal{I}_{\mathcal{P}(t)}(j) = (j_1, \ldots, j_N).$
7:	$\mathbf{r} \leftarrow \mathbf{V}_{\mathcal{P}(t)}(j,:)$	
8:	for all $d \in \delta(t)$ do	\blacktriangleright Multiply the vector r with corresponding matrix rows.
9:	$\mathbf{r} \leftarrow \mathbf{r} \ast \mathbf{U}^{(d)}(j_d,:)$	
10:	$\mathbf{V}_t(i,:) \leftarrow \mathbf{V}_t(i,:) + \mathbf{r}$	• Add the update due to the parent's j th element.

4.1. Shared memory parallelism. For the given tensor \mathcal{X} , after forming the 491 dimension tree \mathcal{T} with symbolic structures \mathcal{I}_t and \mathcal{R}_t for all tree nodes, we can 492 perform numeric TTMV computations in parallel. In Algorithm 4, we provide the 493 shared memory parallel TTMV algorithm, called SMP-DTREE-TTV, for a node t of a 494dimension tree. The goal of SMP-DTREE-TTV is to compute the tensor values \mathbf{V}_t for 495a given node t. Similar to Algorithm 2, it starts by checking if \mathbf{V}_t is already computed, 496497 and returns immediately in that case. Otherwise, it calls SMP-DTREE-TTV on the parent node $\mathcal{P}(t)$ to make sure that parent's tensor values $\mathbf{V}_{\mathcal{P}(t)}$ are available. Once 498 $\mathbf{V}_{\mathcal{P}(t)}$ is ready, the algorithm proceeds with computing \mathbf{V}_t for each nonzero index 499 $\mathcal{I}_t(i) = (i_1, \ldots, i_N)$. As for each such index the reduction set is defined during the 500symbolic TTV, $\mathbf{V}_t(i, :)$ can be independently updated in parallel. In performing this 501

⁵⁰² update, for each element $\mathbf{V}_{\mathcal{P}(t)}(j,:)$ of the parent, the algorithm multiplies this vector ⁵⁰³ with the rows of the corresponding matrices of the TTMV in $\delta(t)$ of t, then adds it ⁵⁰⁴ to $\mathbf{V}_t(i,:)$.

For shared-memory parallel CP-ALS, we replace Line 9 of Algorithm 3 with a call to SMP-DTREE-TTV (l_n) . The parallelization of the rest of the computations is trivial. In computing the matrices $\mathbf{W}^{(n)}$ and $\mathbf{U}^{(n)}$ at Lines 3, 13 and 15, we use parallel dense BLAS kernels. Computing the matrix $\mathbf{H}^{(n)}$ at Line 12 and normalizing the columns of $\mathbf{U}^{(n)}$ are embarrassingly parallel element-wise matrix operations. We skip the details of the parallel convergence check whose cost is negligible.

4.2. Distributed memory parallelism. Parallelizing CP-ALS in a distributed memory setting involves defining unit parallel tasks, data elements, and their interdependencies. Following to this definition, we partition and distribute tensor elements and factor matrices to all available processes. We discuss a *fine-grain* and a *mediumgrain* parallel task model together with the associated distributed memory parallel algorithms.

517 We start the discussion with the following straightforward lemma that enables us 518 to distribute tensor nonzeros for parallelization.

519 LEMMA 5 (Distributive property of TTVs). Let \mathcal{X}, \mathcal{Y} , and \mathcal{Z} be tensors in 520 $\mathbb{R}^{I_1 \times \cdots \times I_N}$ with $\mathcal{X} = \mathcal{Y} + \mathcal{Z}$. Then, for any $n \in \mathbb{N}_N$ and $u \in \mathbb{R}^{I_n} \mathcal{X} \times_n u =$ 521 $\mathcal{Y} \times_n u + \mathcal{Z} \times_n u$ holds.

522 Proof. Using (1) we express the element-wise result of $\mathcal{X} \times_n \mathbf{u}$ as

523
$$(\mathcal{X} \times_{n} \mathbf{u})_{i_{1},...,i_{N},...,i_{N}} = \sum_{j=1}^{I_{n}} u_{j}(x_{i_{1},...,j_{,...,i_{N}}} - z_{i_{1},...,j_{,...,i_{N}}} + z_{i_{1},...,j_{,...,i_{N}}})$$
524
$$= \sum_{j=1}^{I_{n}} u_{j}(x_{i_{1},...,j_{,...,i_{N}}} - z_{i_{1},...,j_{,...,i_{N}}}) + \sum_{j=1}^{I_{n}} u_{j}z_{i_{1},...,j_{,...,i_{N}}}$$

$$= (\boldsymbol{\mathcal{Y}} \times_n \mathbf{u})_{i_1,\dots,i_N} + (\boldsymbol{\mathcal{Z}} \times_n \mathbf{u})_{i_1,\dots,i_N}.$$

527 By extending the previous lemma to P summands and all but one mode TTV, 528 we obtain the next corollary.

529 COROLLARY 6. Let \mathcal{X} and $\mathcal{X}_1, \ldots, \mathcal{X}_P$ be tensors in $\mathbb{R}^{I_1 \times \cdots \times I_N}$ with $\sum_{i=1}^{P} \mathcal{X}_i =$ 530 \mathcal{X} . Then, for any $n \in \mathbb{N}_N$ and $u^{(i)} \in \mathbb{R}^{I_i}$ for $i \in \mathbb{N}_N \setminus \{n\}$ we obtain $\mathcal{X}_{\times_{i\neq n}} =$ 531 $\mathcal{X}_1 \times_{i\neq n} u^{(i)} + \cdots + \mathcal{X}_P \times_{i\neq n} u^{(i)}$.

Froof. Multiplying the tensors \mathcal{X} and $\mathcal{X}_1, \ldots, \mathcal{X}_P$ in any mode $n' \neq n$ in the equation gives tensors $\mathcal{X}' = \mathcal{X} \times_{n'} \mathbf{u}^{(n')}$ and $\mathcal{X}'_i = \mathcal{X}_i \times_{n'} \mathbf{u}^{(n')}$, and $\mathcal{X}' = \sum_{i=1}^P \mathcal{X}'_i$ holds by the distributive property. The same process is repeated in the remaining modes to obtain the desired result.

536 4.2.1. Fine-grain parallelism. Corollary 6 allows us to partition the tensor \mathcal{X} in the sum form $\mathcal{X}_1 + \cdots + \mathcal{X}_P$ for any P > 1, then perform TTMVs on each tensor part \mathcal{X}_p independently, finally sum up these results to obtain the TTMV result 538 for \mathcal{X} multiplied in N-1 modes. As \mathcal{X} is sparse, an intuitive way to achieve this 539decomposition is by partitioning its nonzeros to P tensors where P is the number of 540available distributed processes. This way, for any dimension n, we can perform the 541TTMV of \mathcal{X}_p with the columns of the set of factor matrices $\{\mathbf{U}^{(1)},\ldots,\mathbf{U}^{(N)}\}$ in all 542modes except n. This yields a "local" matrix $\mathbf{M}_{p}^{(n)}$ at each process p, and all these local 543matrices must subsequently be "assembled" by summing up their rows corresponding 544

to the same row indices. In order to perform this assembly of rows, we also partition 545 the rows of matrices $\mathbf{M}^{(n)}$ so that each row is "owned" by a process that is responsible 546for holding the final row sum. We represent this partition with a vector $\boldsymbol{\sigma}^{(n)} \in \mathbb{R}^{I_n}$ where $\sigma_i^{(n)} = p$ implies that the final value of $\mathbf{M}^{(n)}(i, :)$ resides at the process p. We assume the same partition given by $\sigma^{(n)}$ on the corresponding factor matrices $\mathbf{U}^{(n)}$, as this enables each process to compute the rows of $\mathbf{U}^{(n)}$ that it owns using the rows

of $\mathbf{M}^{(n)}$ belonging to that process without incurring any communication. 551

547

548

549

550

Algorithm 5 DMP-DTREE-CP-ALS: Dimension tree-based CP-ALS algorithm

Input: \mathcal{X}_p : A part of an *N*-mode tensor \mathcal{X} R: The rank of CP decomposition For all $n \in \mathbb{N}_N$: $\sigma^{(n)}$: The vector indicating the owner process of each row of $\mathbf{U}^{(n)}$ $\mathcal{F}_p^{(n)}$: The index set with elements $i \in \mathcal{F}_p^{(n)}$ where $\sigma_i^{(n)} = p$ $\mathcal{G}_p^{(n)}$: The set of all unique indices of the nonzeros of \mathcal{X}_p in mode n $\mathbf{U}^{(n)}(\mathcal{G}_p^{(n)},:)$: Distributed initial matrix (rows needed by process p) **Output:** $[\lambda; \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}]$: Rank-*R* CP decomposition of $\boldsymbol{\mathcal{X}}$ with distributed $\mathbf{U}^{(n)}$ ▶ The tree has leaves $\{l_1, \ldots, l_N\}$. 1: $\mathcal{T} \leftarrow \text{CONSTRUCT-DIMENSION-TREE}(\boldsymbol{\mathcal{X}}_p)$ 2: for n = 2...N do $\mathbf{W}^{(n)} \leftarrow \text{All-Reduce}(\mathbf{U}^{(n)}(\mathcal{F}_p^{(n)},:)^T \mathbf{U}^{(n)}(\mathcal{F}_p^{(n)},:))$ 3: 4: repeat for $n = 1 \dots N$ do 5: for all $t \in \mathcal{T}$ do \blacktriangleright Invalidate tree tensors that are multiplied in mode n. 6: if $n \in \mu'(t)$ then 7: \blacktriangleright Destroy all tensors that are multiplied by $\mathbf{U}^{(n)}$. $DESTROY(\mathbf{V}_t)$ 8: SMP-DTREE-TTV (l_n) ▶ Perform the TTMVs for the leaf node tensors. 9: $\mathbf{M}^{(n)}(\mathcal{G}_p^{(n)},:) \leftarrow \mathbf{V}_{l_n}$ \blacktriangleright Form $\mathbf{M}^{(n)}$ using leaf tensors (done implicitly). 10: COMM-FACTOR-MATRIX($\mathbf{M}^{(n)}$, "fold", $\mathcal{G}_p^{(d)}, \mathcal{F}_p^{(d)}, \boldsymbol{\sigma}^{(d)}$) \blacktriangleright Assemble $\mathbf{M}^{(n)}(\mathcal{F}_p^{(n)}, :)$. 11: $\mathbf{H}^{(n)} \leftarrow *_{i \neq n} \mathbf{W}^{(i)}$ 12: $\mathbf{U}^{(n)}(\mathcal{F}_{p}^{(n)},:) \leftarrow \mathbf{M}^{(n)}(\mathcal{F}_{p}^{(n)},:)\mathbf{H}^{(n)^{\dagger}} \\ \boldsymbol{\lambda} \leftarrow \text{COLUMN-NORMALIZE}(\mathbf{U}^{(n)})$ 13:14:COMM-FACTOR-MATRIX $(\mathbf{U}^{(n)}, \text{"expand"}, \mathcal{G}_p^{(d)}, \mathcal{F}_p^{(d)}, \boldsymbol{\sigma}^{(d)}) \rightarrow \text{Send/Receive } \mathbf{U}^{(n)}.$ 15: $\mathbf{W}^{(n)} \leftarrow \text{All-Reduce}([\mathbf{U}^{(n)}(\mathcal{F}_p^{(n)},:)]^T \mathbf{U}^{(n)}(\mathcal{F}_p^{(n)},:))$ 16: 17: until converge is achieved or the maximum number of iterations is reached. 18: return $[\![\boldsymbol{\lambda}; \mathbf{U}^{(1)}, \dots, \mathbf{U}^{(N)}]\!]$

This approach amounts to a fine-grain parallelism where each fine-grain computational task corresponds to performing TTMV operations due to a nonzero element $x_{i_1,\ldots,i_N} \in \mathcal{X}$. Specifically, according to (10), the process p needs the matrix rows $\mathbf{U}^{(1)}(i_1,:),\ldots,\mathbf{U}^{(N)}(i_N,:)$ for each nonzero x_{i_1,\ldots,i_N} in its local tensor \mathcal{X}_p in order 554555to perform its local TTMVs. For each dimension n, we represent the union of all 556these "required" row indices for the process p by $\mathcal{G}_p^{(n)}$. Similarly, we represent the set of "owned" rows by the process p by $\mathcal{F}_p^{(n)}$. In this situation, the set $\mathcal{G}_p^{(n)} \setminus \mathcal{F}_p^{(n)}$ correspond to the rows of $\mathbf{M}^{(n)}$ for which the process p generates a partial TTMV 558 559result, which need to be sent to their owner processes. Equally, it represents the set of 560rows of $\mathbf{U}^{(n)}$ that are not owned by the process p and are needed in its local TTMVs 561according to (10). These rows of $\mathbf{U}^{(n)}$ are similarly to be received from their owners 562in order to carry out the TTMVs at process p. Hence, a "good" partition in general 563 involves a significant overlap of $\mathcal{G}_p^{(n)}$ and $\mathcal{F}_p^{(n)}$ to minimize the cost of communication. 564In Algorithm 5, we describe the fine-grain parallel algorithm that operates in 565

O. KAYA AND B. UCAR

this manner at process p. The elements \mathcal{X}_p , $\mathcal{F}_p^{(n)}$, and $\mathcal{G}_p^{(n)}$ are determined in the partitioning phase, and are provided as input to the algorithm. Each process starts 566 567 with the subset $\mathcal{G}_p^{(n)}$ of rows of each factor matrix $\mathbf{U}^{(n)}$ that it needs for its local computations. Similar to Algorithm 3, at Line 1 we start by forming the dimension 568 569 tree for the local tensor \mathcal{X}_p . We then compute the matrices $\mathbf{W}^{(n)}$ corresponding to 570 $\mathbf{U}^{(n)}{}^{T}\mathbf{U}^{(n)}$ using the initial factor matrices. We do this step in parallel in which each process computes the local contribution $[\mathbf{U}^{(n)}(\mathcal{F}_{p}^{(n)},:)]^{T}\mathbf{U}^{(n)}(\mathcal{F}_{p}^{(n)},:)$ due to its owned rows. Afterwards, we perform an ALL-REDUCE communication to sum up these local 571572573results to obtain a copy of $\mathbf{W}^{(n)}$ at each process. The cost of this communication 574is typically negligible as $\mathbf{W}^{(n)}$ is a small matrix of size $R \times R$. The main CP-ALS 575subiteration for mode n begins with destroying tensors in the tree that will become 576 invalid after updating $\mathbf{U}^{(n)}$. Next, we perform SMP-DTREE-TTV on the leaf node l_n , Algorithm 6 COMM-FACTOR-MATRIX: Communication routine for factor matrices

Input: M: Distributed factor matrix to be communicated comm: The type of communication. "fold" sums up the partial results in owner processes, whereas "expand" communicates the final results from owners to all others. σ : The ownership of each row of **M** \mathcal{G}_p : The rows used by process p \mathcal{F}_p : The rows owned by process pOutput: Rows of M are properly communicated. 1: if comm = "fold" then 2: for all $i \in \mathcal{G}_p \setminus \mathcal{F}_p$ do ▶ Send non-owned rows to their owners. 3: Send $\mathbf{M}(i, :)$ to the process σ_i . for all $i \in \mathcal{F}_p$ do 4: ▶ Gather all partial results of owned rows together. Receive and sum up all partial results for $\mathbf{M}(i, :)$. 5:6: else if *comm* = "*expand*" then 7: for all $i \in \mathcal{F}_p$ do ▶ Send owned rows to all processes in need. Send $\mathbf{M}(i,:)$ to the all processes p' with $i \in F_{p'}$. 8: for all $i \in \mathcal{G}_p \setminus \mathcal{F}_p$ do 9: ▶ Receive rows that are needed locally. Receive $\mathbf{M}(i,:)$ from the process σ_i . 10:

577

and obtain the "local" matrix $\mathbf{M}^{(n)}$. Then, the partial results for the rows of $\mathbf{M}^{(n)}$ 578are communicated to be assembled at their owner processes. We name this as the 579 fold communication step following the convention from the fine-grain parallel sparse 580 matrix computations. Afterwards, we form the matrix $\mathbf{H}^{(n)}$ locally at each process 581p in order to compute the owned part $\mathbf{U}^{(n)}(\mathcal{F}_p^{(n)},:)$ using the recently assembled 582 $\mathbf{M}^{(n)}(\mathcal{F}_p^{(n)},:)$. Once the new distributed $\mathbf{U}^{(n)}$ is computed, we normalize it column-583wise and obtain the vector $\boldsymbol{\lambda}$ of norms. The computational and the communication 584costs of this step are negligible. The new $\mathbf{U}^{(n)}$ is finalized after the normalization, 585 and we then perform an *expand* communication step in which we send the rows of 586 $\mathbf{U}^{(n)}$ from the owner processes to all others in need. This is essentially the inverse 587 of the *fold* communication step in the sense that each process p that sends a partial 588 row result of $\mathbf{M}^{(n)}(i, :)$ to another process q in the fold step receives the final result 589 for the corresponding row $\mathbf{U}^{(n)}(i:)$ from the process q in the expand communication. 590Finally, we update the matrix $\mathbf{W}^{(n)}$ using the new $\mathbf{U}^{(n)}$ in parallel. 591

The expand and the fold communications at Lines 11 and 15 constitute the most 592 expensive communication steps. We outline these communications in Algorithm 6. In 593 the expand communication, the process p sends the partial results for the set $\mathcal{G}_p \setminus \mathcal{F}_p$ of 594rows to their owner processes, while similarly receiving all partial results for its set \mathcal{F}_p

of owned rows and summing them up. Symmetrically, in the fold communication, the 596 process p sends the rows with indices \mathcal{F}_p , and receives the rows with indices $\mathcal{G}_p \setminus \mathcal{F}_p$. 597 The exact set of row indices that needs to be communicated in fold and expand steps 598depends on the partitioning of \mathcal{X} and the factor matrices. As this partition does not 599 change once determined, the communicated rows between p and q stays the same in 600 CP-ALS iterations. Therefore, in our implementation we determine this row set once 601 outside the main CP-ALS iteration, and reuse it at each iteration. Another advantage 602 of determining this set of rows to be communicated in advance is that it eliminates 603 the need to store the vector $\boldsymbol{\sigma}^{(n)} \in \mathbb{R}^{I_n}$ which could otherwise be costly for a large 604 tensor dimension. 605

4.2.2. Medium-grain parallelism. For an N-mode tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ 606 and using $P = \prod_{i=1}^{N} P_i$ processes, the medium-grain decomposition imposes a par-607 tition with $P_1 \times \cdots \times P_N$ Cartesian topology on the dimensions of \mathcal{X} . Specifically, 608 for each dimension n, the index set \mathbb{N}_{I_n} is partitioned into P_n sets $\mathcal{S}_1^{(n)}, \ldots, \mathcal{S}_{P_n}^{(n)}$. With this partition, the process with the index $(p_1, \ldots, p_N) \in P_1 \times \cdots \times P_N$ gets 609 610 $\mathcal{X}(\mathcal{S}_{p_1}^{(1)},\ldots,\mathcal{S}_{p_N}^{(n)})$ as its local tensor. Each factor matrix $\mathbf{U}^{(n)}$ is also partitioned fol-611 lowing this topology where the set of rows $\mathbf{U}^{(n)}(\mathcal{S}_j^{(n)},:)$ is owned by the processes with index (p_1,\ldots,p_N) where $p_n = j, j \in \mathbb{N}_{P_n}$, even though these rows are to be 612 613 further partitioned among the processes having $p_n = j$. As a result, one advantage 614 615 of the medium-grain partition is that only the processes with $p_n = j$ need to communicate with each other in mode n. This does not necessarily reduce the volume of 616 communication, but it can reduce the number of messages by a factor of P_n in the 617 nth dimension. 618

619 One can design an algorithm specifically for the medium-grain decomposition [40]. 620 However, using the fine-grain algorithm on a medium-grain partition effectively pro-621 vides a medium-grain algorithm. For this reason, we do not need nor provide a sepa-622 rate algorithm for the medium-grain task model, and use the fine-grain algorithm with 623 a proper medium-grain partition instead, which equally benefits from the topology.

4.2.3. Partitioning. The distributed memory algorithms that we described require partitioning the data and the computations, as in any distributed memory algorithm. In order to reason about their computational load balance and communication cost, we use hypergraph models. Once the models are built, different hypergraph partitioning methods can be used to partition the data and the computations. We discuss a few partitioning alternatives.

4.2.4. Partitioning for the fine-grain parallelism. We propose a hypergraph model to capture the computational load and the communication volume of the fine-grain parallelization given in Algorithm 5. For the simplicity of the discussion, we present the model for a 3rd order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ and factor matrices $\mathbf{U}^{(1)} \in \mathbb{R}^{I_1 \times R}, \mathbf{U}^{(2)} \in \mathbb{R}^{I_2 \times R}$, and $\mathbf{U}^{(3)} \in \mathbb{R}^{I_3 \times R}$. For these inputs, we construct a hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ with the vertex set \mathcal{V} and the hyperedge set \mathcal{E} . The generalization of the model to higher order tensors should be clear from this construction.

The vertex set $\mathcal{V} = \mathcal{V}^{(1)} \cup \mathcal{V}^{(2)} \cup \mathcal{V}^{(3)} \cup \mathcal{V}^{(\mathcal{X})}$ of the hypergraph involves four types of vertices. The first three types correspond to the rows of the matrices $\mathbf{U}^{(1)}, \mathbf{U}^{(2)},$ and $\mathbf{U}^{(3)}$. In particular, we have vertices $v_i^{(1)} \in \mathcal{V}^{(1)}$ for $i \in \mathbb{N}_{I_1}, v_j^{(2)} \in \mathcal{V}^{(2)}$ for $j \in \mathbb{N}_{I_2}$, and $v_k^{(3)} \in \mathcal{V}^{(3)}$ for $k \in \mathbb{N}_{I_3}$. These vertices represent the "ownership" of the corresponding matrix rows, and we assign unit weight to each such vertex. The fourth type of vertices are denoted by $v_{i,j,k}^{(\mathcal{X})}$, which we define for each nonzero $x_{i,j,k} \in \mathcal{X}$.



Fig. 2: Fine-grain hypergraph model for the $3 \times 3 \times 3$ tensor $\mathcal{X} = \{(1,2,3), (2,3,1), (3,1,2)\}$ and a 3-way partition of the hypergraph. The objective is to minimize the cutsize of the partition while maintaining a balance on the total part weights corresponding to each vertex type (shown with different colors).

This vertex type relates to the number of operations performed in TTMV due to the 643 nonzero element $x_{i,j,k} \in \mathcal{X}$ in using (10) in all modes. In the N-dimensional case, 644this includes up to N vector Hadamard products involving the value of the nonzero 645 $x_{i,j,k}$, and the corresponding matrix rows. The exact number of performed Hadamard 646 products depends on how nonzero indices coincide as TTVs are carried out, and cannot 647 be determined before a partitioning takes place. In our earlier work [26], this cost 648 was exactly N Hadamard products per nonzero, as the MTTKRPs were computed 649 without reusing partial results and without index compression after each TTMV. In 650 the current case, we assign a cost of N to each vertex $v_{i,j,k}^{(\boldsymbol{\chi})}$ to represent an upper bound on the computational cost, and expect this to lead to a good load balance in 651 652 practice. With these vertex definitions, one can use multi-constraint partitioning (6) 653 with one constraint per vertex type. In this case, the first, the second, and the third 654 655 types have unit weights in the first, second, and third constraints, respectively, and zero weight in all other constraints. The fourth vertex type also gets a unit weight (N, N)656 or equivalently, 1) in the fourth constraint, and zero weight for others. Here, balancing 657 the first three constraints corresponds to balancing the number of matrix rows at each 658 process (which provides the memory balance as well as the computational balance in 659 dense matrix operations), whereas balancing the fourth type corresponds to balancing 660 the computational load due to TTMVs. 661

As TTMVs are carried out using (10), data dependencies to the rows of $\mathbf{U}^{(1)}(i,:)$. 662 $\mathbf{U}^{(2)}(j,:)$, and $\mathbf{U}^{(3)}(k,:)$ take place when performing Hadamard products due to each nonzero $x_{i,j,k}$. We introduce three types of hyperedges in $\mathcal{E} = \mathcal{E}^{(1)} \cup \mathcal{E}^{(2)} \cup \mathcal{E}^{(3)}$ 663 664 to represent these dependencies as follows: $\mathcal{E}^{(1)}$ contains a hyperedge $n_i^{(1)}$ for each 665 matrix row $\mathbf{U}^{(1)}(i,:), \mathcal{E}^{(2)}$ contains a hyperedge $n_j^{(2)}$ for each row $\mathbf{U}^{(2)}(j,:)$, and $\mathcal{E}^{(3)}$ 666 contains a hyperedge $n_k^{(3)}$ for each row $\mathbf{U}^{(3)}(k,:)$. Initially, $n_i^{(1)}$, $n_i^{(2)}$ and $n_k^{(3)}$ contain 667 the corresponding vertices $v_i^{(1)}$, $v_j^{(2)}$, and $v_k^{(3)}$, as the owner of a matrix row has a dependency to it by default. In computing the MTTKRP using (10), each nonzero 668 669 $x_{i,j,k}$ requires access to $\mathbf{U}^{(1)}(i,:)$, $\mathbf{U}^{(2)}(j,:)$, and $\mathbf{U}^{(3)}(k,:)$. Therefore, we add the vertex $v_{i,j,k}^{(\boldsymbol{\chi})}$ to the hyperedges $n_i^{(1)}$, $n_j^{(2)}$ and $n_k^{(3)}$ to model this dependency. In Figure 2, we demonstrate this fine-grain hypergraph model on a sample tensor $\boldsymbol{\chi} =$ 670 671 672

 $\{(1,2,3),(2,3,1),(3,1,2)\}$, yet we exclude the vertex weights for simplicity. Each vertex type and hyperedge type is shown using a different color in the figure.

Consider now a P-way partition of the vertices of $\mathcal{H} = (\mathcal{V}, \mathcal{E})$ where each part is 675 associated with a unique process to obtain a P-way parallel execution of Algorithm 5. 676 We consider the first subiteration of Algorithm 5 that updates $\mathbf{U}^{(1)}$, and assume that 677 each process already has all data elements to carry out the local TTMVs at Line 9. 678 Now suppose that the nonzero $x_{i,j,k}$ is owned by the process p and the matrix row 679 $\mathbf{U}^{(1)}(i,:)$ is owned by the process q. Then, the process p computes a partial result 680 for $\mathbf{M}^{(1)}(i, :)$ which needs to be sent to the process q at Line 3 of Algorithm 6. By 681 construction of the hypergraph, we have $v_i^{(1)} \in n_i^{(1)}$ which resides at the process q, and due to the nonzero $x_{i,j,k}$ we have $v_{i,j,k}^{(1)} \in n_i^{(1)}$ which resides at the process p; therefore, this communication is accurately represented in the *connectivity* $\kappa_{n_i^{(1)}}$ of the 682 683 684 hyperedge $n_i^{(1)}.$ In general, the hyperedge $n_i^{(1)}$ incurs $\kappa_{n_i^{(1)}}-1$ messages to transfer the 685 partial results for the matrix row $\mathbf{M}^{(1)}(i,:)$ to the process q at Line 3. Therefore, the 686 connectivity-1 cutsize metric (5) over the hyperedges exactly encodes the total volume 687 of messages sent at Line 3, if we set $c[\cdot] = R$. Since the send operations at Line 8 are 688 duals of the send operations at Line 3, the total volume of messages sent at Line 8 for 689 the first mode is also equal to this number. By extending this reasoning to all other 690 691 modes, we obtain that the cumulative (over all modes) volume of communication in one iteration of Algorithm 5 equals to the connectivity-1 cut-size metric. As the 692 communication due to each mode take place in different stages, one might alternatively 693 use a multi-objective hypergraph model to minimize the communication volume due 694 to each mode (or equivalently, hyperedge type) independently. 695

As discussed above, the proper model for partitioning the data and the compu-696 tations for the fine-grain parallelism calls for a multi-constraint and a multi-objective 697 partitioning formulation to achieve the load balance and minimize the communication 698 cost with a single call to a hypergraph partitioning routine. Since these formulations 699 are expensive, we follow a two-step approach. In the first step, we partition only the 700 nonzeros of the tensor on the hypergraph $\mathcal{H} = (\mathcal{V}^{(\mathcal{X})}, \mathcal{E})$ using just one load constraint 701 due to the vertices in $\mathcal{V}^{(\boldsymbol{\mathcal{X}})}$, and we thereby avoid multi-constraint partitioning. We 702 also avoid multi-objective partitioning by treating all hyperedge types as the same, 703 and thereby aim to minimize the total communication volume across all dimensions, 704 which works well in practice. Once the nonzero partitioning is settled, we partition 705 the rows of the factor matrices in a way to *balance* the communication, which is not 706 achievable using standard partitioning tools. 707

708 We now discuss three methods for partitioning the described hypergraph.

Random: This approach visits the vertices of the hypergraph and assigns each 709 visited vertex to a part chosen uniformly at random. It is expected to balance the 710 TTMV work assigned to each process while ignoring the cost of communication. We 711use random partitioning only as a "worst case" point of reference for other methods. 712 **Standard:** In this standard approach, we feed the hypergraph to a standard 713 hypergraph partitioning tool to obtain balance on the number of tensor nonzeros and 714 715 the amount of TTMV work assigned to a process, while minimizing the communication volume. This approach promises significant reductions in communication cost with 716 respect to the others, yet imposes high computational and memory requirements. 717

T18 Label propagation-like: Given that the standard partitioning approach is too costly in practice, we developed a fast hypergraph partitioning heuristic which has reasonable memory and computational costs. The method is based on the balanced

O. KAYA AND B. UÇAR

1 label-propagation algorithm [38, 45], and includes some additional adaptations to handle hypergraphs [9, 20]. The heuristic starts with an initial assignment of vertices to parts, and then proceeds with multiple passes over the hypergraph. At each pass, the vertices are visited in an order, and are possibly moved to other parts in order to reduce the cutsize while respecting the balance constraints.

For the heuristic to be efficient on hypergraphs, some adaptations are needed. 726 Each pass involves two types of updates. In the first step, each hyperedge chooses 727 a "preferred part" by considering the current part of its vertices. Next, each vertex 728 updates its part according to the preferred parts of the hyperedges that include the 729 vertex. In both steps, the most dominant part index is chosen for the update. The 730 heuristic runs in linear time on the size of the hypergraph per iteration, and requires 731 732 a memory of $2|\mathcal{V}| + |\mathcal{E}| + 4P$. Running the algorithm for a few iterations provides reasonably good partitions. This basic algorithm can have many variants. In one 733 variant, we visit the vertices in an order imposed by an increasing ordering by size of 734 the hyperedges. This variant has an overhead of sorting the hyperedges. In another 735 variant, we reweigh the preference of a hyperedge of size s by the multiplier $(1-\frac{1}{D})^{s-1}$. 736 This last variant has a memory overhead for storing the weights for efficiency purposes; 737 for each size s, the value $(1 - \frac{1}{P})^{s-1}$ is needed. 738

4.2.5. Partitioning for the medium-grain parallelism. Similar to the fine-739 grain model, one can use a hypergraph model for the medium-grain parallel compu-740 741 tations to reduce the communication volume using hypergraph partitioners. However, medium-grain variant is analogous to checkerboard partitioning designed for 742 matrices [13, 14], and calls for a multi-constrained partitioning. Specifically, for a 743 3-dimensional tensor with a process topology $P_1 \times P_2 \times P_3$ where $P = P_1 P_2 P_3$, the 744 hypergraph is to be partitioned in three phases; using one load constraint in the first 745 phase, P_1 constraints (where each constraint is obtained from the first phase) in the 746 second phase, and P_1P_2 constraints (obtained from the second phase partitioning) 747 in the third phase. As P can be large, the number of constraints P_1 and P_1P_2 can 748 similarly get large, and in this case the state of the art partitioners do not perform 749 well both in terms of partition quality and speed. For higher dimensional tensors, this 750 situation only gets worse. That is why explicit communication reduction using hyper-751 graph partitioning for the medium-grain algorithm is not feasible in practice. Hence, 752we use the partitioning heuristic by Smith and Karypis [40] to partition medium-grain 753 hypergraphs for load balance, and to expect a communication reduction due to par-754tition topology indirectly. We also determine the partition topology by choosing P_1 , 755 P_2 , and P_3 proportional to the tensor dimensions I_1 , I_2 , and I_3 . 756

4.2.6. Mode partitioning. Once the nonzero partitioning is obtained for the 757 given fine- or medium-grain parallelism, we proceed with partitioning the mode in-758 dices (or, equivalently, the rows of the factor matrices) using a similar heuristic com-759 mon in similar work [26, 40]. For each matrix row i in dimension n, we identify the 760 processes that have a data dependency to that row. These are exactly the processes 761 which have at least one nonzero with index i in the *n*th dimension. Next, all row 762 indices are sorted in increasing order of the number of dependent processes. Finally, 763 each row is greedily assigned to the process having the minimum total communication 764 volume among all processes dependent to that row. 765

5. Related work. There has been many recent advances in the efficient computations of tensor factorizations in general, and CP decomposition in particular. We briefly mention these here and refer the reader to the original sources for details. In [4],

Bader and Kolda show how to efficiently carry out MTTKRP as well as other funda-769 770mental tensor operations on sparse tensors in MATLAB. GigaTensor [22] is a parallel implementation of CP-ALS using the Map-Reduce framework. DFacTo [16] is a C++ 771 implementation with distributed memory parallelism using MPI, and it uses a partic-772 773 ular formulation of MTTKRP using sparse matrix-vector multiplication. SPLATT [40] is an efficient parallelization of MTTKRP and CP-ALS both in shared [42] and dis-774 tributed memory [40] environments using OpenMP and MPI, and is implemented in 775 C. It uses a medium-grain distributed parallelism with a Cartesian partitioning of the 776 tensor, and generalizes this technique to the tensor completion problem [41]. It is the 777 fastest publicly available CP-ALS implementation in the existing literature, and their 778 approach translates to performing N(N-1) TTMVs in performing the MTTKRP in 779 780 the main CP-ALS iteration. Karlsson et al. similarly discuss a parallel computation of the tensor completion problem using CP formulation [23] in which they replicate 781 the entire factor matrix $\mathbf{U}^{(n)}$ among MPI processes unlike our approach, and report 782 scalability results only up to 100 cores. For computing the CP decomposition of dense 783 tensors. Phan et al. [35] proposes a scheme that divides the tensor modes into two 784 sets, pre-computes the TTMVs for each mode set, and finally reuses these partial 785 786 results to obtain the final MTTKRP result in each mode. This provides a factor of 2 improvement in the number of TTMVs over the traditional approach, and our dimen-787 sion tree-based framework can be considered as the generalization of this approach 788 that provides a factor of $N/\log N$ improvement. 789

While this paper was under evaluation, another paper appeared [33]. Li et al. use the same idea of storing intermediate tensors but use a different formulation based on tensor times tensor multiplication and a tensor times matrix through Hadamard products for shared memory systems. The overall approach is similar to that by Phan et al. [35], where the difference lies in the application of the method to sparse tensors and auto-tuning to better control memory use and gains in the operation counts.

Aside from CP decomposition, Baskaran et al. [6] provide a shared-memory parallel implementation for the Tucker decomposition of sparse tensors. We [27] provide efficient shared and distributed memory parallelization of the Tucker decomposition for sparse tensors using OpenMP and MPI. Austin et al. [3] discuss a high performance distributed memory parallelization of dense Tucker factorization in the context of data compression. Finally, Perros et al. [34] investigate an efficient computation of hierarchical Tucker decomposition for sparse tensors.

6. Experiments. We first investigate how CP-ALS implementations compare using a single thread to assess the algorithmic impact of using a BDT in the same implementation. Then, we compare these implementations using multiple threads to evaluate their shared memory parallel performance. Finally, we compare the mediumand the fine-grain distributed memory parallel algorithms.

808 **6.1.** Dataset and environment. We experimented with five real-world tensors whose sizes are shown in Table 1. Netflix tensor has user \times movie \times time dimensions, 809 which we formed from the data of the Netflix Prize competition [7]. In this tensor, 810 811 nonzeros correspond to the user reviews for movies, and the review date extends the data to the third dimension. The values of the nonzeros are determined by the cor-812 813 responding review scores given by the users. We obtained the NELL tensor from the Never Ending Language Learning (NELL) knowledge database of the "Read the 814 Web" project [10], which consists of tuples of the form (entity, relation, entity) such 815 as ("Chopin", "plays musical instrument", "piano"). The nonzeros of this tensor cor-816respond to these entries discovered by NELL from the web, and the values are set to 817

Tensor	I_1	I_2	I_3	I_4	#nonzeros
Delicious	1.4K	532K	17M	2.4M	140M
Flickr	731	319K	28M	1.6M	112M
Netflix	480K	17K	2K	-	100M
NELL	3.2M	301	638K	-	78M
Amazon	6.6M	2.4M	23K	-	1.3B

Table 1: Real-world tensors used in the experiments.

be the "belief" scores given by the learning algorithms used in NELL. Delicious and 818 Flickr are the datasets for the web-crawl of Delicious.com and Flickr.com during 2006 819 and 2007, which are formed by Görlitz et al. [17]. These datasets consist of tuples of 820 form (time \times users \times resources \times tags); hence, we form 4-mode tensors out of these 821 tuples. We obtained the Amazon review dataset from SNAP [32], which contains 822 product review texts by users. We first processed this dataset with the standard text 823 processing routines. We used the nltk package [8] in Python to tokenize the review 824 text, to discard the stop words, to apply Porter stemmer, and to keep the words that 825 are in the US, GB, or CA dictionaries. Afterwards, we retained only the words with 826 at least five occurrences in the whole review set. Then, we created a three dimen-827 sional tensor whose dimensions correspond to the users, products, and retained words. 828 Numerical values are set to the frequency of a word in a review. 829

We conducted experiments on a shared memory and a separate distributed mem-830 ory system. The shared memory system has two CPU sockets (Intel(R) Xeon(R))831 E5-2695 v3) each having 14 cores at a clock speed of 2.30GHz with Turbo Boost dis-832 abled. The system has a total memory of size 768GB, and each socket has L1, L2, 833 L3 caches of sizes of 32KB, 256KB, and 35MB, respectively. All codes are compiled 834 with gcc/g++5.3.0 using OpenMP directives and compiler options -O3, -ffast-math, 835 -funroll-loops, -ftree-vectorize, -fstrict-aliasing on this shared memory system. The 836 837 distributed memory system is an IBM Blue Gene/Q cluster. This system consists of 6 racks of 1024 nodes with each node having 16GB of memory and a 16-core IBM 838 PowerPC A2 processor running at 1.6GHz. We ran our experiments up to 256 nodes 839 (4096 cores). Each core of PowerPC A2 can handle one arithmetic and memory oper-840 ation simultaneously; therefore, we assigned 32 threads per node (2 threads per core) 841 for better performance. On this system, all codes were compiled using the Clang C++ 842 compiler (version 3.5.2) with IBM MPI wrapper using the same optimization flags, 843 and linked against IBM ESSL library for LAPACK and BLAS routines. 844

We also used synthetic tensors created randomly having 4, 8, 16, and 32 dimensions. In these random tensors, each dimension is of size 10M, and there are 100M nonzeros with a uniform random distribution of indices. Using these tensors, we measure the effect of tensor dimensionality on the performance.

We provide the dimension-tree based CP-ALS implementation in our tensor factorization library called HYPERTENSOR. It is a C++11 implementation providing shared and distributed memory parallelism through OpenMP and MPI libraries. We compared our code against SPLATT v1.1.1 [40], a C code with OpenMP and MPI parallelizations. The "winner" for each test case in the results are highlighted with bold font.

6.2. Shared memory experiments. We compare the shared memory performance of the dimension tree-based CP-ALS algorithm with the state of the art. We
experimented with four methods called ht-tree2, ht-tree3, ht-tree, and splatt.
The ht-tree method, implemented in HYPERTENSOR, uses a full BDT to carry out

TTMVs. The ht-tree 2 method is the same implementation as ht-tree except that it 859 uses a 2-level dimension tree. In this tree, N leaf nodes are directly connected to the 860 root, hence no intermediate results are generated. However, TTMVs are performed 861 one mode at a time to benefit from the index compression to reduce the operation 862 count. As a result, this method performs N-1 TTMVs for each mode in an iteration 863 just as SPLATT; we thus expect comparable performance. The ht-tree3 method uses 864 a three level BDT whose second level has two nodes holding the partial TTMV results 865 corresponding to the first and the second half of the set of dimensions. This is anal-866 ogous to the approach by Phan et al. [35] for computing dense CP decompositions, 867 but it also employs our data structure and shared memory parallelization for sparse 868 tensors. Storing partial results in the second level reduces the number of TTMVs by 869 870 a factor of 2, but for an N-mode tensor, this method still performs $O(N^2)$ TTMVs per iteration. Note that for 3- and 4-dimensional tensors, ht-tree3 and ht-tree use 871 identical trees, thus give the same results. These results for ht-tree3 are indicated 872 with an asterisk in the tables. Finally, splatt corresponds to the parallel CP-ALS 873 implementation in SPLATT. We ran all algorithms for 20 iterations with the rank of 874 approximation R = 20 (except for the sequential execution of 16- and 32-dimensional 875 random tensors, which are run for 2 iterations due to their cost), and recorded the 876 average time spent per CP-ALS iteration. Test instances in which a method gets out 877 of memory are indicated with a dash symbol. 878

6.2.1. Sequential execution. In Table 2, we give the sequential per-iteration 879 880 run time of all methods. We report the run time in seconds for **splatt**, and the relative speedup with respect to **splatt** for the other three methods. We first note 881 that ht-tree2 runs slightly slower than splatt on three dimensional Amazon and 882 NELL tensors (0.99x and 0.87x), and notably slower on Netflix tensor (0.60x). This 883 is so because SPLATT has a specially tuned implementation for 3-dimensional tensors, 884 whereas we use a single code for all dimensions. On all higher dimensional tensors, 885 886 ht-tree2 performs significantly better than splatt, up to 2.08x on Random8D, which shows the efficiency of our implementation for N-dimensional tensors even before 887 using a BDT. The gap between **splatt** and **ht-tree2** narrows using Random16D, 888 as ht-tree2 depletes the memory in one NUMA node, which is discussed more in 889 subsection 6.2.4, and starts accessing the distant memory in the other NUMA node. 890 **ht-tree2** gets out of memory using Random32D as it stores $O(N^2)$ index arrays. 891

892 We now measure the effect of dimension trees by comparing **ht-tree** with **ht**tree2 in Table 2. These two methods use the same TTMV implementation, whereas 893 ht-tree uses a full BDT. On Delicious, Flickr, Netflix, and NELL, ht-tree obtains 894 1.78x, 1.61x, 1.63x, and 1.47x speedup over ht-tree2 thanks to the BDT. Likewise, 895 896 on random tensors, we observed 1.43x, 1.91x, 3.01x speedup on tensors Random4D, Random8D, and Random16D, respectively, using ht-tree. This validates our perfor-897 mance expectation (Theorem 2) that as the dimensionality of the tensor increases, a 898 BDT results in significantly fewer TTMVs hence better performance. 899

Comparing ht-tree with splatt similarly yields a speedup of 1.98x, 1.98x, 1.28x, 2.05x, 3.97x, 3.94x, and 5.96x on tensors Delicious, Flickr, NELL, Random4D, Random8D, Random16D, and Random32D, respectively, which similarly meets our expectation of performance gain from Theorem 2. On Amazon, ht-tree was only 2% faster, whereas on Netflix, splatt was only 2% faster, which was the only instance in which splatt had a slight edge over ht-tree.

Finally, we note that the performance gap between **ht-tree** and **ht-tree3** widens significantly as the tensor gets higher dimensional. Using Random8D, Random16D,

Table 2:	Sequential	l CP-ALS ru	1 time pe	er iteration.	Timings a	re in secor	nds for \mathbf{sp}	latt,
whereas	we report	the relative	speedup	with respect	t to splat	t for othe	r method	s.

	splatt	ht-tree 2	ht-tree3	ht-tree
Delicious	66.6	1.11	*	1.98
Flickr	43.6	1.23	*	1.98
Netflix	8.2	0.60	*	0.98
NELL	8.3	0.87	*	1.28
Amazon	214.6	0.99	*	1.02
Random4D	224.7	1.43	*	2.05
Random8D	1527.1	2.08	2.70	3.97
Random16D	4401.6	1.31	2.02	3.94
Random32D	19919.9	-	2.38	5.96

Table 3: Shared memory parallel CP-ALS run time per iteration (in seconds). Timings are in seconds for **splatt**, whereas we report the relative speedup with respect to **splatt** for other methods.

	\mathbf{splatt}	ht-tree2	ht-tree3	ht-tree
Delicious	8.3	0.93	*	2.00
Flickr	5.8	1.09	*	1.81
Netflix	0.7	0.55	*	0.87
NELL	1.3	1.11	*	1.46
Amazon	24.4	0.86	*	0.95
Random4D	20.1	0.91	*	1.47
Random8D	86.9	0.81	1.69	2.18
Random16D	349.2	0.82	1.73	3.47
Random32D	1601.8	-	2.18	5.65

and Random32D, ht-tree is 1.47x, 1.95x, and 2.50x faster than ht-tree3 as it incurs
 significantly fewer TTMVs. These results suggest that using a full BDT is indeed the

910 ideal choice for performance.

6.2.2. Shared memory parallel execution. In Table 3, we give the run time 911 results of all methods with shared memory parallelism using 14 threads. We first note 912 that in HYPERTENSOR, using dimension trees consistently yields better execution 913 914 times. Using ht-tree, we obtain 2.15x, 1.66x, 1.58x, 1.32x, and 1.10x speedup over ht-tree2 on Delicious, Flickr, Netflix, NELL, and Amazon tensors, respectively. For 915random tensors, we get 1.62x, 2.69x, and 4.23x speedup on Random4D, Random8D, 916 and Random16D. Comparing ht-tree with splatt, we observe a speedup of 2.00x, 917 918 1.81x, 1.46x, 1.47x, 2.18x, 3.47x, and 5.65x on tensors Delicious, Flickr, NELL, Random4D, Random8D, Random16D, and Random32D, respectively. This demonstrates 919 that the use of a BDT in CP-ALS computations can be effectively parallelized in a 920 shared memory setting, on top of significantly reducing the amount of TTMV work. 921 On three dimensional Amazon and Netflix tensors, splatt has a slight edge over ht-922 923 tree by 5% and 14% faster executions, respectively. Another point to note is that splatt has somewhat better parallel speedup in general (over its own sequential run 924 925 time) than **ht-tree2** and **ht-tree**. This is mostly due to the fact that TTMV is a memory-bound computation; hence, once the memory bandwidth is fully utilized, one 926 cannot expect further speedup through multi-threading. When performing TTMVs, 927 our implementation makes slightly more memory accesses due to extra pointer ar-928 929 rays involved in the dimension tree nodes, which saturates the bandwidth earlier and

which	splatt is faste	er than th	splattht-tree2ht-tree3ht-tree87.32.41*1.39				
		splatt	ht-tree2	ht-tree3	ht-tree		
	Delicious	87.3	2.41	*	1.39		
	Flickr	57.1	2.70	*	1.32		

1.96

1.63

1.85

2.21

6.36

15.71

1.48

1.47

1.58

2.40

3.97

3.67

3.07

*

4.35

4.32

4.85

51.6

37.7

720.3

62.9

86.2

 $\mathbf{233}$

638.7

Netflix

NELL

Amazon

Random4D

Random8D

Random16D

Random32D

Table 4: Symbolic precomputation timings. We report the exact timing for **splatt** in seconds, and relative timing with respect to **splatt** for other methods, indicating the ratio at which **splatt** is faster than these methods this precomputation.

930	thereby affects	s the parall	el speedup t	o a certain	extent.	Neverthe	eless, usir	ıg ht-tre €
931	we achieve up	to $5.65x$ fas	ster runs ove	r splatt in	a shared	l memory	r parallel	execution
	~ ^ 1			-				

Conformally with the sequential case, **ht-tree** gets significantly faster than **ht**tree3 as the tensor dimensionality increases, up to 2.59x using Random32D, which demonstrates the effectiveness of using a full BDT.

935 **6.2.3.** Preprocessing cost. In Table 4, we provide symbolic TTV costs of our 936 methods as well as the precomputation cost of **splatt** for setting up its data structures. All runtimes are for a sequential execution; neither SPLATT nor HYPERTENSOR 937 parallelizes this step in their current version. We first note that **ht-tree** incurs sig-938 nificantly less cost than ht-tree2 in all instances. This is expected as ht-tree2 has 939 940 $O(N^2)$ index arrays to be sorted, whereas **ht-tree** has only $O(N \log N)$ of them. For the same reason, ht-tree gets notably faster than ht-tree3 as the tensor dimensional-941 ity increases to 32. The cost is comparable between **splatt** and **ht-tree** for Delicious, 942 Flickr, Netflix, NELL, and Amazon tensors, but splatt takes significantly less time 943 for higher dimensional random tensors, up to 3.97x on Random8D, as it sorts only 944 945O(N) arrays.

Comparing these timings with the iteration times in Table 2, we see that this precomputation is amortized in a few iterations, except for Netflix and NELL. In practice, CP-ALS is typically executed multiple times with different initial matrices and ranks of approximation using the same symbolic dimension tree construct, which should render this preprocessing cost less important even for these two tensors.

6.2.4. Memory usage. We provide the memory consumption of all methods in 951 952 Table 5. The first column corresponds to the amount of memory used to store factor matrices, which is common to all methods. We give the memory usage for storing 953 954 index arrays in GBs for **splatt**, and as the ratio to **splatt** for all other methods. We also give the memory consumption for storing the value matrices of intermedi-955 ate tensors for ht-tree3 and ht-tree in GBs. We first note that ht-tree2 uses the 956957 highest amount of memory to store index arrays as expected, up to 8.32 times more than splatt on Random16D. In all tensors, the amount of index memory used by ht-958 959 tree is only slightly higher than **splatt**, the worst case being Flickr tensor for which ht-tree consumes 1.35 times more memory in comparison. There is a multitude of 960 reasons for this observation. First, even though splatt uses only O(N) index arrays, 961 we realized upon inspecting the implementation that it uses two different representa-962963 tions of a tensor for faster execution, effectively doubling its memory requirements.

Table 5: Memory usage of different methods. Index usages of ht-tree2, ht-tree3, and ht-tree are reported with respect to that of **splatt**, whereas the memory usage of factors, value matrices, and **splatt** index are in GBs.

	factors	splatt	ht-tree2	ht-t	ree3	ht-tree	
	lactors	index	index	index	value	index	value
Delicious	3	7	2.44	*	*	1.21	5.7
Flickr	4.5	4.8	2.29	*	*	1.35	4.3
Netflix	0.1	3.4	1.85	*	*	1.29	1.4
NELL	0.6	2.5	1.88	*	*	1.28	0.5
Amazon	1.4	49.6	1.84	*	*	1.32	65.4
Random4D	6	9.1	2.21	*	*	1.12	10.2
Random8D	11.9	21	4.24	1.08	15.3	1.08	30.5
Random16D	23.8	44.9	8.32	1.47	15.3	1.23	45.7
Random32D	47.7	94.3	-	2.40	15.3	1.31	61

It also employs two arrays per dimension in its compressed sparse fiber (CSF) tensor storage [39], doubling the memory consumption. Finally, its memory efficiency depends heavily on the index overlaps after TTMVs, which happens rarely for high dimensional random tensors.

Aside from this, we note that the amount of memory used to store the values of intermediate tensors is reasonable, Random8D being the only exception in which the value matrix size exceeds the index size using **ht-tree**. Finally, **ht-tree3** uses more memory for index storage than **ht-tree** for high dimensional tensors, as it uses more index arrays similar to **ht-tree2**. It uses less space for value matrices, however, as it has only one level for intermediate tensors.

All in all, we conclude upon considering these results that **ht-tree** provides remarkable performance improvements with a reasonable increase in the memory usage.

976 **6.3.** Distributed memory experiments. We compare the performance and the scalability of the fine- and the medium-grain parallel CP-ALS algorithms. In 977 these experiments, we do not use SPLATT software to benchmark medium-grain par-978 allelization for two reasons. First, we would like to compare the effect of load balance 979 and communication cost in different algorithms using different partitionings, while 980 isolating the effects of the efficiency of local CP-ALS computations. Since SPLATT's 981 982 medium-grain implementation does not use BDTs for local TTMVs, and is slower, comparing it against HYPERTENSOR's fine-grain implementation which has faster 983 local TTMVs would not be fair, nor would correctly reveal the effect of different 984 partitioning strategies. Second, we were not able to get SPLATT to work on our 985 986 distributed system despite our full efforts. Therefore, we instead performed mediumgrain partitioning of tensors following the description of SPLATT's heuristic [40], and 987 ran HYPERTENSOR on these partitions which incurs the same cost in terms of the 988 communication volume and the number of messages as SPLATT, while using more effi-989 cient TTMV kernels. For local CP-ALS computations, we use the BDT-based method 990 991 **ht-tree** for shared memory parallelism, as it gives the best performance. This way, the experiments become more precise in terms of measuring the influence of medium-992 993 and fine-grain algorithms and associated partitionings on parallel scalability.

We investigate the performance in two tables. In Table 6, we give the strong scalability results of the medium- and the fine-grain algorithms up to 256 MPI ranks using 4096 cores. Since we achieved the maximum scalability in most tensors with 256 MPI ranks and 4096 cores, the discussion is mostly confined to this case. Especially Table 6: Per iteration speedup results for distributed memory parallel CP-ALS using different partitions. The best single threaded execution is given in seconds (shaded cells), and the relative speedups are reported for all other cases. # nodes and # cores correspond to the number of nodes (and equivalently, MPI ranks) and cores per node used in each instance, respectively.

#modeax#aamaa		Delic	ious			Flie	ckr	
$\#noues \times \#cores$	med-gd	fine-rd	fine-lb	fine-hp	med-gd	fine-rd	fine-lb	fine-hp
8×1	45.078	0.65	1.00	1.00	0.86	-	25.674	0.99
8×16	14.12	4.64	12.51	18.13	10.72	-	12.77	16.98
16×16	17.41	6.30	19.62	31.00	16.02	2.71	23.53	29.78
32×16	27.96	9.02	30.92	51.64	23.30	4.03	36.99	53.60
64×16	34.28	15.51	50.14	83.63	25.34	7.08	72.73	89.77
128×16	54.84	25.83	84.57	128.43	39.38	13.43	115.65	143.43
256×16	74.76	40.43	141.75	183.24	46.34	23.49	148.40	178.29
#nodes×#cores	Netflix					NE	LL	
# 1100C3 × # COTC3	med-gd	fine-rd	fine-lb	fine-hp	med-gd	fine-rd	fine-lb	fine-hp
4×1	27.656	0.86	0.98	0.94	19.540	0.87	0.97	0.98
4×16	25.26	16.41	23.22	24.18	16.41	10.14	16.22	17.62
8×16	44.82	22.18	39.68	40.20	26.09	13.43	27.91	28.69
16×16								
10 / 10	76.82	27.99	65.69	65.85	43.52	19.58	45.13	44.01
$\frac{10 \times 10}{32 \times 16}$	$\begin{array}{r} 76.82 \\ 124.58 \end{array}$	27.99 34.40	$65.69 \\ 103.58$	$65.85 \\ 105.96$	43.52 68.08	$19.58 \\ 26.66$	45.13 69.29	44.01 61.45
	$\frac{76.82}{124.58}\\200.41$	$ \begin{array}{r} 27.99 \\ 34.40 \\ 41.90 \end{array} $	65.69 103.58 155.37	65.85 105.96 159.86	43.52 68.08 109.16	$ \begin{array}{r} 19.58 \\ 26.66 \\ 36.39 \end{array} $	45.13 69.29 97.21	44.01 61.45 81.76
	76.82 124.58 200.41 264.49	$ \begin{array}{r} 27.99 \\ 34.40 \\ 41.90 \\ 52.68 \end{array} $	65.69 103.58 155.37 236.37	65.85 105.96 159.86 271.14	43.52 68.08 109.16 153.86	$ \begin{array}{r} 19.58 \\ 26.66 \\ 36.39 \\ 49.10 \end{array} $	45.13 69.29 97.21 127.71	$ \begin{array}{r} 44.01 \\ 61.45 \\ 81.76 \\ 126.06 \end{array} $

#nodes × # cores	Amazon					
# 1100C3 × # COTC3	med-gd	fine-rd	fine-lb			
64×1	-	0.55	36.303			
64×16	-	4.12	19.37			
128×16	31.82	5.86	36.05			
256×16	34.38	8.58	63.69			

998 in Table 7, we give the detailed load balance and communication cost metrics just for 999 this case.

In Table 6, we compare the execution of Algorithm 5 with three partitioning 1000 methods. The **fine-hp** and **fine-lb** correspond, respectively, to the standard hyper-1001 1002 graph partitioning and the label-propagation-like heuristic of subsection 4.2.4. For fine-hp, we used PaToH [12] with the default settings. On Amazon tensor, we could 1003 1004 not obtain results for **fine-hp** as the tensor was too big for PaToH to handle. For fine-lb, we ran the three alternatives, each for three passes, and chose the partitioning 1005with the smallest cut. The **med-gd** method corresponds to the medium-grain parti-1006 1007 tioning heuristic [40]. The **fine-rd** method refers to the random partitioning of the 1008 fine-grain hypergraph, which is given as a reference to illustrate the impact of a good partitioning. Due to memory constraints, we were not able to execute Algorithm 5 1009 on a single node, as the original tensors are large. Therefore, for each tensor, we give 1010the results starting from the minimum number of nodes needed, and for the same 1011 instance we also give the single threaded results. We use the run time of the single 1012 threaded execution of the fastest method as our baseline (highlighted with shaded 1013 cells) in computing the parallel speedup. In passing from one core per node to 16 1014 cores per node, we see some speedup over 16 in Table 6; this is because we use two 1015 threads per core (see subsection 6.1). 1016

In order to be able to analyze the speedup results of Table 6, we give the number of tensor nonzeros per part, computational load (the number of Hadamard products), communication volume, and the number of messages incurred by these three partitionings, using 256 MPI ranks in Table 7. For the four performance metrics, we give the maximum and the average value observed across all processes. We see in all instances except **fine-hp** on NELL that balancing the number of nonzeros per part

Dentitioning	Nr	ız	Comp.	Load	Comn	n. Vol.	Num	. Msg
Farmoning	Max.	Avg.	Max.	Avg.	Max.	Avg.	Max.	Avg.
			Deli	cious				
fine-hp	547K	547K	1947K	1807 K	199K	137K	2039	2018
fine-lb	564K	547K	1849K	1737K	309K	265K	2040	2040
fine-rd	550K	547K	2747K	2737K	1083K	1080K	2040	2040
medium-gd	598K	547K	2353K	2214K	624K	571K	1096	1096
			Fli	ckr				
fine-hp	441K	441K	1334K	1221K	54K	38K	1827	1588
fine-lb	454K	441K	1335K	1257K	107K	87K	2040	2030
fine-rd	443K	441K	2318K	2308K	1042K	1038K	2040	2040
medium-gd	443K	441K	1826K	1806K	576K	558K	1152	1152
			Net	tflix				
fine-hp	392K	392K	1439K	1211K	49K	19K	1380	1158
fine-lb	404K	392K	1252K	1208K	48K	39K	1530	1528
fine-rd	394K	392K	1681K	1674K	412K	411K	1530	1530
medium-gd	394K	393K	1184K	1177K	26K	24K	642	642
			NE	ELL				
fine-hp	307K	307K	1219K	787K	46K	23K	1513	1402
fine-lb	316K	307K	920K	888K	89K	84K	1522	1502
fine-rd	309K	307K	1211K	1205K	271K	269K	1530	1520
medium-gd	310K	307K	871K	855K	58K	51K	583	570
	•	•	Am	azon	•			•
fine-lb	5104K	4955K	16871K	16241K	211K	203K	1503	1503
fine-rd	4962K	4955K	20964K	20940 K	4656K	4651K	1530	1530
medium-gd	19984K	4955K	50023K	16255K	230K	170K	550	514

Table 7: Load balance and communication statistics for 256-way partitioning.

gracefully translates into balancing the actual computational load.

1024 Using 256 MPI ranks on Delicious, fine-hp and fine-lb are 2.5x and 1.9x faster over med-gd. We observe in Table 7 that this is due to better minimization of the 1025total and the maximum communication volume. On Flickr, fine-hp is 3.9x faster 1026 than **med-gd** at 256 MPI ranks with 14.7x and 10.6x less total and maximum com-1027 munication volume, while **fine-lb** shows a speedup of 3.2x over **med-gd** with 6.4x 1028 and 5.4x less total and maximum communication cost. In both tensors, med-gd re-1029 sults in about the half the communication volume of **fine-rd**. In overall, on Delicious 1030 fine-hp and fine-lb obtain 183x and 142x speedup using 4096 cores, whereas medgd and fine-rd give 75x and 40x speedup for the same tensor. On Flickr, fine-hp 1032and fine-lb similarly yield 178x and 148x speedup using 4096 cores, while med-gd 1033and fine-rd could achieve 46x and 23x speedup. For the Delicious and Flickr tensors, 1034 while passing from 8 nodes (with 8×16 cores) to 256 nodes (with 256×16 cores). med-gd results in 5.29x and 4.32x speedup. The fine-hp and fine-lb result in 11.33x 1036and 10.11x speedup for Delicious, and 11.62x and 10.50x speedup for Flickr in the 1037 same scenario. fine-rd is significantly slower than the other methods, incurring the 1038 1039 highest communication as shown in Table 7.

On Netflix and NELL, **med-gd** yields 321x and 197x speedup using 4096 cores. fine-hp shows a comparable performance with 261x and 164x speedup, whereas finelb is slightly slower than fine-hp with 236x and 158x speedup. In passing from 4 nodes (with 4×16 cores) to 256 nodes (with 256×16 cores), **med-gd** results in 12.73x and 12.03x speedup for Netflix and NELL, respectively. The fine-hp and finelb partitioning result in 10.18x and 11.18x speedup for Netflix, and 9.72x and 9.32x speedup for NELL in the same scenario. fine-rd is similarly the slowest of all methods giving 67x and 62x speedup for these two tensors. Using Netflix and NELL, **med-**



Fig. 3: Running time dissection of a parallel CP-ALS iteration using **fine-hp** and **ht-tree** scheme. The legends **ttv**, **dmat**, **comm**, and **fit** correspond to the time spent for TTMVs, dense matrix operations following the TTMVs, communication, and fit computation. The time for fit computation not discernible in the plot.

gd gets 23% and 20% faster than fine-hp, and 36% and 25% faster than fine-lb. In 1048 Table 7, we see that this is due to **med-gd** incurring smaller maximum communication 1049 1050 volume and having fewer messages. We investigated this outcome and observed that when a tensor is long in one mode and short in all others, the communication due 1051 to the long mode dominates the overall cost, and the communication for the small modes remains negligible in comparison. On Netflix with P = 256 MPI ranks, using 1053 **med-gd** with $P = p \times q \times r$ topology, the worst case communication volume for the 1054first mode is upper-bounded by 480K(qr-1), as $I_1 = 480K$ indices are distributed 1056 to p process "slices" each with qr processes. Similarly, for the second and the third modes, the worst case communication volumes are 17K(pr-1) and 2K(pq-1). In 1057 such cases, choosing a large p and smaller q and r significantly reduces the worst-1058 case communication cost in the first mode, while the cost in other modes stays low. 1059The medium-grain heuristic achieves this. Specifically, on Netflix, the medium-grain 1060 heuristic chooses a grid size of $64 \times 4 \times 1$. This advantage is lost when there are at 1061 1062 least two long dimensions (see Delicious and Flickr), as using more processes in one long mode can increase the communication significantly in the other long modes. 1063

1064 On Amazon tensor, med-gd starts to lose scalability at 256 MPI ranks. We observe in Table 7 that this is due to load imbalance. Amazon tensor has some relatively 1065"dense" slices that make load balancing difficult for the medium-grain heuristic. This 1066 problem never arises in the fine-grain partitioning due to finer granularity of tasks; as 1067 a result, fine-lb runs 1.85x faster than med-gd using 256 MPI ranks. In this tensor, 1068 1069 from 128 nodes to 256 nodes, **med-gd** displays a speedup of 1.08. With **fine-lb**, the parallel algorithm enjoys 1.86x speedup in passing from 64 nodes to 128 nodes, and 10703.2x speedup in passing from 64 to 256 nodes. In overall, med-gd gives 34x speedup 1071 over the baseline whereas **fine-lb** gets significantly faster with 63x speedup. 1072

In Figure 3, we present the dissection of the parallel run time for a CP-ALS iteration on Flickr tensor using 256 MPI ranks. We choose Flickr as representative, as it includes the highest proportion of dense matrix operations in comparison to all other tensors. Despite this fact and using a BDT for faster TTMVs, the TTMV O. KAYA AND B. UÇAR

step still remains to be the dominant computational cost. In this figure, we first 1077 1078observe that the workload due to TTMV and dense matrix computations decrease with the increasing number of processes. Second, we expect in general that having 1079 more processes increases the total communication volume; yet we observe in the plot 1080 that the communication cost declines until 128 MPI ranks. This is because a good 1081 partitioning can reduce the communication volume per process (while increasing the 1082 total communication volume). At 256 MPI ranks, however, communication cost starts 1083 to increase and become the bottleneck. The fit computation takes a negligible amount 1084 of time hence is not discernible in the plot. 1085

On Flickr tensor, the three variants of the fine-lb took 58.38, 89.66, and 65.04 1086 seconds to partition the hypergraph, med-gd took 190 seconds, and fine-hp took 1087 1088 207 minutes. In all data instances, fine-lb gives good results while being a fast partitioning heuristic. fine-hp consistently provides better partitions than fine-lb in all 1089 instances, yet the partitioning cost might render it impractical to use in real-world 1090 scenarios. med-gd heuristic is only effective when the tensor nonzeros are homoge-1091 neously distributed, and the tensor has only one large dimension. One might consider 1092 1093 reducing the communication volume on a medium-grain topology using hypergraph 1094 partitioning, yet the high number of constraints prevents this approach from being amenable. Therefore, we believe that fine-lb serves well in most practical situations. 1095

7. Conclusion. We investigated an efficient computation of successive tensor-1096 times-vector multiplication in the context of the well-known CP-ALS algorithm for 1097 sparse tensor factorization. We introduced a computational scheme using dimension 1098 trees that asymptotically reduces the computational cost of the TTMV operations 1099 1100 for higher order tensors while using a reasonable amount of memory. Our technique provides performance benefits for lower order tensors, and gets progressively better 1101 as the dimensionality of the tensor increases in comparison to the state of the art. 1102 We proposed an effective shared memory parallelization of this method with a pre-1103 computation step in order to efficiently carry out numerical computations within 1104the CP-ALS iterations. We introduced a fine-grain parallelization approach in the 1105 distributed memory setting, compared it against a recently proposed medium-grain 1106 variant, discussed good partitionings for both approaches, and validated these findings 1107 with experiments on real-world tensors. The proposed computational scheme can 1108 be applied to both dense and sparse tensors as well as other tensor decomposition 1109 algorithms involving successive tensor-times-vector and -matrix multiplications. We 1110are planning to investigate this potential in our future work. 1111

Acknowledgments. Some preliminary experiments were carried out using the workstations and the PSMN cluster at ENS Lyon. This work was performed using HPC resources from GENCI-[TGCC/CINES/IDRIS] (Grant 2016 - i2016067501).

1115

REFERENCES

- [1] E. ACAR, D. M. DUNLAVY, AND T. G. KOLDA, A scalable optimization approach for fitting canonical tensor decompositions, Journal of Chemometrics, 25 (2011), pp. 67–86.
 [2] G. A. A. M. DUNLAVY, AND T. G. KOLDA, A scalable optimization approach for fitting canonical tensor decompositions, Journal of Chemometrics, 25 (2011), pp. 67–86.
- [2] C. A. ANDERSSON AND R. BRO, *The N-way toolbox for MATLAB*, Chemometrics and Intelligent Laboratory Systems, 52 (2000), pp. 1–4.
- [3] W. AUSTIN, G. BALLARD, AND T. G. KOLDA, Parallel tensor compression for large-scale scientific data, in IEEE International Parallel and Distributed Processing Symposium (IPDPS), Chicago, IL, USA, May 23–27, 2016, pp. 912–922.
- [4] B. W. BADER AND T. G. KOLDA, Efficient MATLAB computations with sparse and factored tensors, SIAM Journal on Scientific Computing, 30 (2007), pp. 205–231.

- [5] B. W. BADER, T. G. KOLDA, ET AL., Matlab tensor toolbox version 2.6. Available online http://www.sandia.gov/~tgkolda/TensorToolbox/, February 2015.
- [6] M. BASKARAN, B. MEISTER, N. VASILACHE, AND R. LETHIN, Efficient and scalable computations with sparse tensors, in IEEE Conference on High Performance Extreme Computing (HPEC), Sept 2012, pp. 1–6.
- [7] J. BENNETT AND S. LANNING, *The Netflix Prize*, in Proceedings of KDD cup and workshop,
 vol. 2007, 2007, p. 35.
- [8] S. BIRD, E. LOPER, AND E. KLEIN, Natural Language Processing with Python, O'Reilly Media Inc., 2009.
- [9] J. BUURLAGE, Self-improving sparse matrix partitioning and bulk-synchronous pseudo streaming, master's thesis, Utrecht University, 2016.
- [10] A. CARLSON, J. BETTERIDGE, B. KISIEL, B. SETTLES, E. R. H. JR., AND T. M. MITCHELL,
 Toward an architecture for never-ending language learning, in AAAI, vol. 5, 2010, p. 3.
- [11] D. J. CARROLL AND J. CHANG, Analysis of individual differences in multidimensional scaling
 via an N-way generalization of "Eckart-Young" decomposition, Psychometrika, 35 (1970),
 pp. 283–319.
- [12] Ü. V. ÇATALYÜREK AND C. AYKANAT, PaToH: A Multilevel Hypergraph Partitioning Tool, Version 3.0, Bilkent University, Department of Computer Engineering, Ankara, 06533
 Turkey. http://bmi.osu.edu/~umit/software.htm, 1999.
- [13] Ü. V. ÇATALYÜREK AND C. AYKANAT, A hypergraph-partitioning approach for coarse-grain decomposition, in Supercomputing, ACM/IEEE 2001 Conference, Denver, Colorado, 2001, p. 42.
- [14] Ü. V. ÇATALYÜREK, C. AYKANAT, AND B. UÇAR, On two-dimensional sparse matrix partitioning: Models, methods, and a recipe, SIAM Journal on Scientific Computing, 32 (2010), pp. 656–683.
- [15] Ü. V. ÇATALYÜREK, Hypergraph Models for Sparse Matrix Partitioning and Reordering, PhD
 thesis, Bilkent University, Computer Engineering and Information Science, Nov 1999.
- [16] J. H. CHOI AND S. V. N. VISHWANATHAN, DFacTo: Distributed factorization of tensors, in
 27th Advances in Neural Information Processing Systems, Montreal, Quebec, Canada,
 2014, pp. 1296–1304.
- [17] O. GÖRLITZ, S. SIZOV, AND S. STAAB, *PINTS: Peer-to-peer infrastructure for tagging systems*, in Proceedings of the 7th International Conference on Peer-to-Peer Systems, Berkeley, CA, USA, 2008, USENIX Association, p. 19.
- [18] L. GRASEDYCK, *Hierarchical singular value decomposition of tensors*, SIAM Journal on Matrix
 Analysis and Applications, 31 (2010), pp. 2029–2054.
- [19] R. A. HARSHMAN, Foundations of the PARAFAC procedure: Models and conditions for an "explanatory" multi-modal factor analysis, UCLA Working Papers in Phonetics, 16 (1970), pp. 1–84.
- 1163 [20] V. HENNE, Label propagation for hypergraph partitioning, master's thesis, Karsruhe Institute 1164 of Technology, Germany, 2015.
- 1165 [21] J. HÅSTAD, Tensor rank is NP-complete, Journal of Algorithms, 11 (1990), pp. 644–654.
- [22] U. KANG, E. PAPALEXAKIS, A. HARPALE, AND C. FALOUTSOS, GigaTensor: Scaling tensor analysis up by 100 times - Algorithms and discoveries, in Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, New York, NY, USA, 2012, ACM, pp. 316–324.
- [23] L. KARLSSON, D. KRESSNER, AND A. USCHMAJEW, Parallel algorithms for tensor completion
 in the CP format, Parallel Computing, 57 (2016), pp. 222–234.
- [24] G. KARYPIS AND V. KUMAR, Multilevel algorithms for multi-constraint hypergraph partitioning,
 Tech. Report 99-034, University of Minnesota, Department of Computer Science/Army
 HPC Research Center, Minneapolis, MN 55455, November 1998.
- [25] O. KAYA AND B. UÇAR, High-performance parallel algorithms for the Tucker decomposition of higher order sparse tensors, Tech. Report RR-8801, Inria, Oct 2015.
- [26] O. KAYA AND B. UÇAR, Scalable sparse tensor decompositions in distributed memory systems, in Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, New York, NY, USA, 2015, ACM, pp. 77:1–77:11.
- [27] O. KAYA AND B. UÇAR, *High performance parallel algorithms for the Tucker decomposition of sparse tensors*, in 45th International Conference on Parallel Processing (ICPP '16), Aug
 2016, pp. 103–112.
- 1183 [28] T. G. KOLDA AND B. BADER, *The TOPHITS model for higher-order web link analysis*, in 1184 Proceedings of Link Analysis, Counterterrorism and Security, 2006.
- [29] T. G. KOLDA AND B. BADER, Tensor decompositions and applications, SIAM Review, 51 (2009),
 pp. 455–500.

O. KAYA AND B. UÇAR

- [30] L. D. LATHAUWER AND B. D. MOOR, From matrix to tensor: Multilinear algebra and signal processing, in Institute of Mathematics and Its Applications Conference Series, vol. 67, 1998, pp. 1–16.
- [31] T. LENGAUER, Combinatorial Algorithms for Integrated Circuit Layout, Wiley–Teubner, Chich ester, U.K., 1990.
- [32] J. LESKOVEC AND A. KREVL, SNAP Datasets: Stanford large network dataset collection. http: //snap.stanford.edu/data, June 2014.
- [33] J. LI, J. CHOI, I. PERROS, J. SUN, AND R. VUDUC, Model-driven sparse CP decomposition for higher-order tensors, in IPDPS 2017, 31th IEEE International Symposium on Parallel and Distributed Processing, Orlando, FL, USA, May 2017, pp. 1048–1057.
- [34] I. PERROS, R. CHEN, R. VUDUC, AND J. SUN, Sparse hierarchical Tucker factorization and its application to healthcare, in Data Mining (ICDM), 2015 IEEE International Conference on, Nov 2015, pp. 943–948.
- [35] A. H. PHAN, P. TICHAVSKÝ, AND A. CICHOCKI, Fast alternating LS algorithms for high order
 CANDECOMP/PARAFAC tensor factorizations, IEEE Transactions on Signal Processing,
 61 (2013), pp. 4834–4846.
- [36] S. RENDLE AND T. S. LARS, Pairwise interaction tensor factorization for personalized tag rec *ommendation*, in Proceedings of the Third ACM International Conference on Web Search
 and Data Mining, WSDM '10, New York, NY, USA, 2010, ACM, pp. 81–90.
- [37] S. RENDLE, B. M. LEANDRO, A. NANOPOULOS, AND L. SCHMIDT-THIEME, Learning optimal ranking with tensor factorization for tag recommendation, in Proceedings of the 15th ACM
 SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '09, New York, NY, USA, 2009, ACM, pp. 727–736.
- [38] G. M. SLOTA, K. MADDURI, AND S. RAJAMANICKAM, PuLP: Scalable multi-objective multiconstraint partitioning for small-world networks, in Proc. 2nd IEEE Int'l. Conf. on Big
 Data (BigData), IEEE, Oct. 2014, pp. 481–490.
- [39] S. SMITH AND G. KARYPIS, Tensor-matrix products with a compressed sparse tensor, in Pro ceedings of the 5th Workshop on Irregular Applications: Architectures and Algorithms,
 ACM, 2015, p. 7.
- [40] S. SMITH AND G. KARYPIS, A medium-grained algorithm for sparse tensor factorization, in
 2016 IEEE International Parallel and Distributed Processing Symposium, IPDPS 2016,
 Chicago, IL, USA, May 23-27, 2016, 2016, pp. 902–911.
- 1219[41] S. SMITH, J. PARK, AND G. KARYPIS, An exploration of optimization algorithms for high per-1220formance tensor completion, Proceedings of the 2016 ACM/IEEE conference on Super-1221computing, (2016).
- [42] S. SMITH, N. RAVINDRAN, N. D. SIDIROPOULOS, AND G. KARYPIS, SPLATT: Efficient and parallel sparse tensor-matrix multiplication, in 29th IEEE International Parallel & Distributed Processing Symposium, Hyderabad, India, May 2015, IEEE Computer Society, pp. 61–70.
- [43] P. SYMEONIDIS, A. NANOPOULOS, AND Y. MANOLOPOULOS, Tag recommendations based on tensor dimensionality reduction, in Proceedings of the 2008 ACM Conference on Recommender Systems, New York, NY, USA, 2008, ACM, pp. 43–50.
- 1228 [44] G. TOMASI AND R. BRO, A comparison of algorithms for fitting the parafac model, Computa-1229 tional Statistics & Data Analysis, 50 (2006), pp. 1700 – 1734.
- [45] J. UGANDER AND L. BACKSTROM, Balanced label propagation for partitioning massive graphs, in
 Proceedings of the Sixth ACM International Conference on Web Search and Data Mining,
 WSDM '13, New York, NY, USA, 2013, ACM, pp. 507–516.
- [46] M. A. O. VASILESCU AND D. TERZOPOULOS, Multilinear analysis of image ensembles: Tensor Faces, in Computer Vision—ECCV 2002, Springer, 2002, pp. 447–460.
- [47] N. ZHENG, Q. LI, S. LIAO, AND L. ZHANG, *Flickr group recommendation based on tensor de- composition*, in Proceedings of the 33rd International ACM SIGIR Conference on Research
 and Development in Information Retrieval, SIGIR '10, NY, USA, 2010, ACM, pp. 737–738.