

EXPLOITING EFFICIENT REPRESENTATIONS IN LARGE-SCALE TENSOR DECOMPOSITIONS*

NICO VERVLIE[†], OTTO DEBAL[‡], AND LIEVEN DE LATHAUWER[†]

Abstract. Decomposing tensors into simple terms is often an essential step toward discovering and understanding underlying processes or toward compressing data. However, storing the tensor and computing its decomposition is challenging in a large-scale setting. Though in many cases a tensor is structured, it can be represented using few parameters: a sparse tensor is determined by the positions and values of its nonzeros, a polyadic decomposition by its factor matrices, a Tensor Train by its core tensors, a Hankel tensor by its generating vector, etc. The complexity of tensor decomposition algorithms can be reduced significantly in terms of time and memory if these efficient representations are exploited directly. Only a few core operations such as norms and inner products need to be specialized to achieve this, thereby avoiding the explicit construction of multiway arrays. To improve the interpretability of tensor models, constraints are often imposed or multiple datasets are fused through joint factorizations. While imposing these constraints prohibits the use of traditional compression techniques, our framework allows constraints and compression, as well as other efficient representations, to be handled trivially as the underlying optimization variables do not change. To illustrate this, large-scale nonnegative tensor factorization is performed using MLSVD and Tensor Train compression. We also show how vector and matrix data can be analyzed using tensorization while keeping a vector or matrix complexity through the concept of implicit tensorization, as illustrated for Hankelization and Löwnerization. The concepts and numerical properties are extensively investigated through experiments.

Key words. tensor decompositions, canonical polyadic decomposition, tensorization, Tucker, large-scale, optimization

AMS subject classification. 15A69

DOI. 10.1137/17M1152371

1. Introduction. In signal processing and data analysis, tensors are often given as multiway arrays of numerical values. Tensor decompositions are then used to discover patterns, separate signals, model behavior, cluster similar phenomena, detect anomalies, predict missing data, etc. More concrete examples include intrusion detection in computer networks [67], analysis of food samples [6], crop classification using hyperspectral imaging [95], direction of arrival estimation using a grid of antennas [74], and modeling melting temperatures of alloys [93]. In these contexts, the canonical polyadic decomposition (CPD) and the block term decomposition (BTD) are common. The components, e.g., the rank-1 terms in a CPD, recovered by these decompositions can often be interpreted directly thanks to relatively mild uniqueness results compared to matrices. If prior knowledge is available, constraints such as nonnega-

*Submitted to the journal's Methods and Algorithms for Scientific Computing section October 16, 2017; accepted for publication (in revised form) December 3, 2018; published electronically March 19, 2019.

<http://www.siam.org/journals/sisc/41-2/M115237.html>

Funding: The first author's work was supported by an Aspirant Grant from the Research Foundation - Flanders (FWO). The second author's work was supported by a Ph.D. grant of the Agency for Innovation by Science and Technology (IWT). This work was supported by the Flemish Government, FWO projects G.0830.14N and G.0881.14N, by the European Research Council under the European Union's Seventh Framework Programme (FP7/2007-2013) / ERC Advanced grant BIOTENSORS (n° 339804), by KU Leuven Internal Funds C16/15/059. This paper reflects only the authors' views and the Union is not liable for any use that may be made of the contained information.

[†]Department of Electrical Engineering ESAT/STADIUS, KU Leuven, B-3001 Leuven, Belgium; and Group Science, Engineering and Technology, KU Leuven Kulak, 8500 Kortrijk, Belgium (nico.vervliet@esat.kuleuven.be, otto.debals@gmail.com, Lieven.DeLathauwer@kuleuven.be).

tivity, orthogonality, sparsity, smoothness, Vandermonde structure, (block) Toeplitz structure, and so on can be imposed on the components to further improve their interpretation in specific applications. When multiple datasets describing (partially) the same phenomenon are available, coupling or data fusion allows more meaningful patterns to be extracted. More examples can be found in several overview papers; see, e.g., [14, 52, 75]. While most publications tackle a single type of structure or constraint, we discuss a general framework in this paper.

Given as dense arrays of numerical values, the cost of storing and processing tensors tends to increase quickly, especially for high orders: for an N th-order tensor of size $I \times I \times \cdots \times I$ the number of entries is I^N . To deal with these large-scale tensors, various techniques have emerged to handle their decompositions, e.g., by deliberately sampling entries [92, 93], through randomization [89], through compression, and/or using parallelism, e.g., [13, 43, 55, 67, 76, 79], or by cross approximation [9, 56, 63, 66]. Such techniques rely on the fact that tensors are structured and can be represented efficiently using fewer parameters than the number of entries in the dense array. This is also our basic assumption in this paper. Moreover, we assume that structure is exploitable to reduce the computational and memory complexity from \mathcal{O} (tensor entries) to \mathcal{O} (parameters). Let us consider three important types of efficient representation that allow such exploitation: compact exact representation, compression-based representation, and (implicit) tensorization.

For the first type, the tensor has an *exact and known structure* that can be represented compactly. A sparse tensor, for instance, is represented compactly by the indices and values of the nonzeros, or more efficiently using tensor extensions of compressed sparse column or row formats [4, 77]. For sparse tensors, the efficient or compact representation has already been exploited extensively; see, e.g., [4, 13, 41, 42, 43, 44, 45, 53, 67, 77, 78, 79]. In some applications the tensor is given as a polyadic decomposition; the factor matrices then form the compact representation [57, 70].

Second, in the case of *compression-based representations*, the tensor is first approximated using a decomposition that is easier or cheaper to compute, and this approximation is then used to speed up subsequent computations. For example, in the experiment in subsection 4.2 a rank-10 CPD is computed directly from a $1024 \times 1024 \times 1024$ tensor which takes 39 s, while the computation time for the decomposition is reduced to 0.6 s after a multilinear singular value decomposition (MLSVD) compression step which takes 8 s. Common examples of this type are a truncated MLSVD [22], a Hierarchical Tucker approximation [31, 35], or a Tensor Train (TT) approximation [65]. MLSVD compression is used extensively as a preprocessing step to speed up the computation of an unconstrained CPD, which follows from the CANDELINC model [12, 50]: thanks to the orthogonal compression the rank-1 structure is preserved and the CPD can be computed from the compressed core tensor. To design effective updating algorithms, a previously computed compact representation of the old data combined with a newly arrived slice is used in [85]. In scientific computing, recompression or rounding is often used to lower the multilinear rank of a polyadic decomposition (PD) or a Tucker decomposition [48, 49, 64, 72, 73], or to lower the TT rank of a tensor given as a PD [65] or a TT approximation [37, 60]. Other examples include the approximation of a matrix or tensor by a greedy orthogonal rank-1 decomposition [51], a hierarchical Kronecker tensor product decomposition [34], or a two-level rank- (r_1, \dots, r_d) decomposition [47].

The third type is *implicit tensorization*. Tensorization involves mapping vectors or matrices to higher-order tensors [14, 25]. (Note that we use a much broader definition of tensorization than, e.g., the quantization concept in scientific computing; see, e.g., [32] and references therein.) By decomposing the obtained tensors, mild

uniqueness conditions of tensor decompositions can be leveraged, which are key in a variety of applications such as blind source separation (BSS). However, for some types of tensorization, the number of data points increases drastically. For example, third-order Hankelization of a vector of length M results in a tensor with $\mathcal{O}(M^3)$ entries [74], while Löwnerization [27] of K observations of mixed rational functions, sampled in M points, may result in a tensor of dimensions $M/2 \times M/2 \times K$. In both examples, the signal length M that can be handled is limited by the respective cubic and quadratic dependence of the number of entries on M . To avoid this data explosion, we propose for this type of problem the concept of implicit tensorization: by operating on the underlying data directly, no explicit construction of the tensor is required. This way, implicit tensorization alleviates computation and memory cost if the data underlying the tensor has few parameters, and the structure of this representation can be exploited easily, e.g., through fast Fourier transforms (FFT); see section 3. Examples are Hankelization, Löwnerization, Toeplitzization, segmentation with overlap, and outer product structures [5, 25]. Via implicit tensorization, deterministic BSS techniques become a viable alternative to classical methods such as ICA, even for datasets of realistic sizes. This is illustrated in subsection 4.3.

Optimization-based techniques are often used in, e.g., signal processing, data analysis, and machine learning to compute a low rank or a low multilinear rank approximation \mathcal{M} of a (structured) tensor \mathcal{T} . For example, for a CPD, $\mathcal{M} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$, and the following minimization problem is solved for the factor matrices \mathbf{A}, \mathbf{B} , and \mathbf{C} :

$$(1.1) \quad \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \frac{1}{2} \|\mathcal{M} - \mathcal{T}\|_{\text{F}}^2.$$

Often additional constraints are imposed to facilitate the interpretation of the resulting factor matrices, or multiple datasets are analyzed jointly. Many of these constraints can be implemented via projection, active set methods, penalties, or using parametric constraints; see [90] for an overview. As these constraints are usually imposed on factor matrices, classical complexity reduction techniques such as MLSVD compression [7] can often not be used as the constraint is not preserved. For example, if $\mathbf{A} \in \mathbb{R}^{I \times R}$ is constrained to be nonnegative and \mathbf{U} is an orthogonal basis for \mathbf{A} , then $\mathbf{A} = \mathbf{U}\hat{\mathbf{A}} \geq \mathbf{0}$ does not imply that $\hat{\mathbf{A}} \geq \mathbf{0}$. Hence, when using MLSVD compression, one cannot decompose the core tensor and impose nonnegativity. As the original factors remain the optimization variables in our approach, constrained decompositions can be computed using a compressed or structured tensor. An illustration in the case of compression is given in Figure 1.1.

1.1. Contributions. We explain how efficient representations originating from each of the three types of efficient representations (compact exact representation, compression-based representation, and (implicit) tensorization) can be exploited in optimization-based algorithms for the computation of a CPD, a decomposition into multilinear rank- $(L_r, L_r, 1)$ terms (LL1), a low-multilinear rank approximation (LMLRA), or a general BTD. To be able to exploit this efficient representation, the computation of the residual is avoided by expanding the least squares loss function into inner products and norms. For example, for the CPD we have

$$(1.2) \quad \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \frac{1}{2} \|\mathcal{M} - \mathcal{T}\|_{\text{F}}^2 = \min_{\mathbf{A}, \mathbf{B}, \mathbf{C}} \frac{1}{2} \|\mathcal{M}\|^2 - \langle \mathcal{M}, \mathcal{T} \rangle + \frac{1}{2} \|\mathcal{T}\|^2.$$

Each of these norms and the inner product, as well as the structured matrix-matrix products required for the gradient, can be implemented such that the structure un-

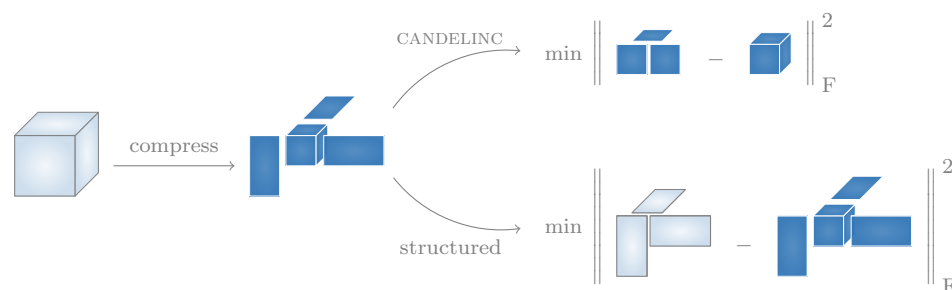


FIG. 1.1. In contrast to the classical CANDELINC model, the optimization variables remain the same in our framework when using compression. Instead of doing computations on a small core tensor, which is possible because the orthogonal transformation preserves low rank structure, we compute the decomposition using the compact Tucker model, factor matrices included. This requires that all operations are formulated in terms of the model and are implemented efficiently. (Figure adapted from [90].)

derlying the efficient representation is exploited. A number of papers have already proposed similar ideas for the computation of a CPD or LMLRA in the specific cases of sparse tensors [4, 13, 41, 42, 43, 44, 45, 53, 67, 77, 78, 79] or polyadic and Tucker formats [96, 97, 98]. We illustrate that these ideas can be extended to more types of structured tensors, including tensors resulting from tensorization techniques. Moreover, we show that, in contrast to previous results, both first- and second-order optimization algorithms can be used, and that constraints and joint factorizations are handled trivially, which we illustrate in the experiments by using Gauss–Newton type methods in combination with active set and parametric constraints. We also pay attention to often neglected numerical accuracy issues. Finally, we show that our approach enables enormous speedups for large-scale problems even when constraints are imposed or when multiple datasets are jointly factorized, in contrast to the traditional compression approach using the CANDELINC model [7], which cannot be used in these cases, as the constraints are not preserved by orthogonal compression. We demonstrate this experimentally for large-scale nonnegative tensor factorization using MLSVD and TT compression, and for implicit Hankelization of signals with up to 500 000 samples. MATLAB implementations for the framework are available in Tensorlab [94].

1.2. Outline. After discussing the notation in the remainder of this section, section 2 explains how efficient representations of structured tensors can be exploited by rewriting the Frobenius norm objective function and gradient for the CPD, LL1, LMLRA, and BTD. We give an overview of implementations of the four core operations—norm, inner product, matricized tensor times Khatri–Rao product, and matricized tensor times Kronecker product—in section 3 for structured tensors given in the polyadic, Tucker, or TT format, and derive implementations for Hankel and Löwner tensors. Finally, numerical experiments are given in section 4. We discuss the influence of the condition number and illustrate the performance for compression-based nonnegative CPD and blind separation of exponential polynomials through implicit Hankelization. In the last experiment, we show how ECG signals can be analyzed efficiently using Löwnerization.

1.3. Notation. Scalars, vectors, matrices, and tensors are denoted by lower case, e.g., a , bold lower case, e.g., \mathbf{a} , bold upper case, e.g., \mathbf{A} , and calligraphic, e.g.,

\mathcal{T} , letters, respectively. \mathbb{K} denotes either \mathbb{R} or \mathbb{C} . Sets are indexed by superscripts within parentheses, e.g., $\mathbf{A}^{(n)}$, $n = 1, \dots, N$. A mode- n vector is the generalization of a column (mode-1) and row (mode-2) vector and is defined by fixing all but the n th index of an N th-order tensor. The mode- n unfolding of a tensor \mathcal{T} is denoted by $\mathbf{T}_{(n)}$ and has the mode- n vectors as its columns. The mode- (m, n) unfolding is defined analogously and is denoted by $\mathbf{T}_{(m, n)}$. The vectorization operator $\text{vec}(\mathcal{T})$ stacks all mode-1 vectors into a column vector.

A number of products are needed. The mode- n tensor-matrix product is denoted by $\mathcal{T} \cdot_n \mathbf{A}$. The Kronecker and Hadamard, or elementwise, products are denoted by \otimes and $*$, respectively. The columnwise (rowwise) Khatri–Rao products between two matrices, denoted by \odot and \odot^T , respectively, are defined as the columnwise (rowwise) Kronecker product. As we assume that the first index runs faster than the second and so on, we define the following shorthand notations to simplify expressions:

$$(1.3) \quad \bigotimes_n \equiv \bigotimes_{n=N}^1, \quad \bigotimes_{k \neq n} \equiv \bigotimes_{\substack{k=N \\ k \neq n}}^1,$$

where N is the order of the tensor and indices are traversed in reverse order. Similar definitions are used for \odot , \odot^T , and $*$. To reduce the number of parentheses, we assume the matrix product takes precedence over \otimes , \odot , \odot^T , and $*$, e.g., $\mathbf{AB} \odot \mathbf{CD} \equiv (\mathbf{AB}) \odot (\mathbf{CD})$. The complex conjugate, transpose, and conjugated transpose are denoted by $\bar{\cdot}$, \cdot^T , and \cdot^H , respectively. The columnwise concatenation of vectors \mathbf{a} and \mathbf{b} is written as $\mathbf{x} = [\mathbf{a}; \mathbf{b}]$ and is shorthand for $\mathbf{x} = [\mathbf{a}^T \ \mathbf{b}^T]^T$. The inner product between \mathcal{A} and \mathcal{B} is denoted by $\langle \mathcal{A}, \mathcal{B} \rangle = \text{vec}(\mathcal{B})^H \text{vec}(\mathcal{A})$. The Frobenius norm is denoted by $\|\cdot\|$. $\text{Re}(\cdot)$ returns the real part of a complex number. \mathbf{I}_I is the $I \times I$ identity matrix and $\mathbf{1}_I$ is a length I column vector with ones. Finally, \mathbf{V}^f is the rowwise “flipped” version of $\mathbf{V} \in \mathbb{K}^{I \times J}$, i.e., $\mathbf{V}^f(i, :) = \mathbf{V}(I - i + 1, :)$ for $i = 1, \dots, I$.

2. Exploiting efficient representations. Many tensor decompositions used in signal processing and data analysis can be computed by minimizing the least squares (LS) error between the approximation and the given tensor \mathcal{T} :

$$(2.1) \quad \min_{\mathbf{z}} f(\mathbf{z}) \quad \text{with} \quad f(\mathbf{z}) = \frac{1}{2} \|\mathcal{M}(\mathbf{z}) - \mathcal{T}\|^2,$$

in which $\mathcal{M}(\mathbf{z})$ is a tensor decomposition with variables \mathbf{z} . Various algorithms have been proposed to solve (2.1), including alternating least squares (ALS) [8, 11, 23, 36, 54], nonlinear conjugate gradient (NCG) [1], quasi-Newton (qN) [82], and Gauss–Newton (GN) [82, 83, 84] algorithms. We first review the notation for common tensor decompositions in subsection 2.1. In subsection 2.2 we give an overview of LS optimization theory for algorithms using gradients (first-order information) and Hessian approximations (second-order information) and show that many common approximations to the Hessian do not depend on the data \mathcal{T} . Hence, second-order algorithms can be used without additional changes in implementation compared to first-order algorithms. Finally, the core operations required to reduce the computational complexity are introduced in subsection 2.3 by rewriting $f(\mathbf{z})$, similar to (1.2).

2.1. Overview of tensor decompositions. The techniques derived in this paper focus on four main decompositions, which can be grouped into two families from an optimization point-of-view: CPD-based and BTD-based algorithms. Note that we use different families compared to the usual division between decompositions computed using numerical linear algebra techniques such as singular value decompositions

(SVD), e.g., MLSVD, TT and Hierarchical Tucker, and decompositions usually computed through optimization (CPD, LL1, BTD). Here, we focus mainly on notation; for pointers on uniqueness results and applications, we refer the reader to the overview papers [14, 32, 33, 52, 75].

CPD-based algorithms. The CPD writes an N th-order tensor as a (minimal) sum of R rank-1 terms, each of which is an outer product, denoted by \otimes , of N nonzero factor vectors $\mathbf{a}_r^{(n)}$:

$$(2.2) \quad \mathcal{M}_{\text{CPD}}(\mathbf{z}) = \sum_{r=1}^R \mathbf{a}_r^{(1)} \otimes \cdots \otimes \mathbf{a}_r^{(N)} := \llbracket \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \rrbracket.$$

Each factor matrix $\mathbf{A}^{(n)}$ contains $\mathbf{a}_r^{(n)}$, $r = 1, \dots, R$, as its columns, and the variables \mathbf{z} are the vectorized factor matrices $\mathbf{A}^{(n)}$, i.e., $\mathbf{z} = [\text{vec}(\mathbf{A}^{(1)}) ; \dots ; \text{vec}(\mathbf{A}^{(N)})]$. (The underlying variables will always be denoted by \mathbf{z} .) Depending on the field, the CPD is also known as PARAFAC, CANDECOMP, or tensor/separation rank decomposition.

The decomposition into multilinear rank- $(L_r, L_r, 1)$ terms (LL1) is accommodated within this family as it is often computed as a constrained CPD [8, 23, 82]:

$$(2.3) \quad \mathcal{M}_{\text{LL1}}(\mathbf{z}) = \sum_{r=1}^R (\mathbf{A}_r \mathbf{B}_r^T) \otimes \mathbf{c}_r = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket,$$

in which $\mathbf{A}_r \in \mathbb{K}^{I_1 \times L_r}$, $\mathbf{B}_r \in \mathbb{K}^{I_2 \times L_r}$ for $r = 1, \dots, R$, $\mathbf{A} = [\mathbf{A}_1, \dots, \mathbf{A}_R]$, $\mathbf{B} = [\mathbf{B}_1, \dots, \mathbf{B}_R]$, $\mathbf{C} = [\mathbf{c}_1, \dots, \mathbf{c}_R]$, and \mathbf{P} is a binary matrix replicating column \mathbf{c}_r L_r times, i.e., $\mathbf{P} = \text{blockdiag}(\mathbf{1}_{L_1}^T, \dots, \mathbf{1}_{L_R}^T)$. Similar to the CPD, $\mathbf{z} = [\text{vec}(\mathbf{A}); \text{vec}(\mathbf{B}); \text{vec}(\mathbf{C})]$. The LL1 decomposition is similar to PARALIND [8]. CONFAC generalizes this decomposition by allowing an arbitrary full row rank matrix \mathbf{P} in every mode [17]. In [82], the decomposition into $(\text{rank-}L_r \otimes \text{rank-}1)$ terms is discussed as one generalization of \mathcal{M}_{LL1} to higher-order tensors.

BTD-based algorithms. The LMLRA or Tucker decomposition writes a tensor as a multilinear transformation of a core tensor $\mathcal{S} \in \mathbb{K}^{J_1 \times \cdots \times J_N}$ by the factor matrices $\mathbf{U}^{(n)} \in \mathbb{K}^{I_n \times J_n}$, $n = 1, \dots, N$:

$$(2.4) \quad \mathcal{M}_{\text{LMLRA}}(\mathbf{z}) = \mathcal{S} \cdot_1 \mathbf{U}^{(1)} \cdots_N \mathbf{U}^{(N)}.$$

An LMLRA can be computed using SVDs of tensors unfoldings, i.e., as a (truncated) MLSVD [21, 88], by higher-order orthogonal iteration [22] or through optimization [39, 40, 82, 83].

The more general block term decomposition writes a tensor as a sum of R multilinear rank- $(J_{r1}, J_{r2}, \dots, J_{rN})$ terms [18]:

$$(2.5) \quad \mathcal{M}_{\text{BTD}}(\mathbf{z}) = \sum_{r=1}^R \mathcal{S}^{(r)} \cdot_1 \mathbf{U}^{(r,1)} \cdots_N \mathbf{U}^{(r,N)}.$$

The difference between $\mathcal{M}_{\text{LMLRA}}$ and \mathcal{M}_{BTD} is the additional summation over $r = 1, \dots, R$. For simplicity of notation, we only use $\mathcal{M}_{\text{LMLRA}}$ in the derivations in this paper.

2.2. Optimization for least squares problems. In gradient-based optimization, an initial guess for the variables \mathbf{z}_0 is iteratively refined as

$$(2.6) \quad \mathbf{z}_k = \mathbf{z}_{k-1} + \alpha_k \mathbf{p}_k,$$

with the step direction \mathbf{p}_k found by solving

$$(2.7) \quad \mathbf{H}_k \mathbf{p}_k = -\bar{\mathbf{g}}_k,$$

in which \mathbf{g}_k and \mathbf{H}_k are the gradient and the Hessian (or an approximation) of $f(\mathbf{z})$ w.r.t. \mathbf{z} in (2.1), respectively, assuming $\mathcal{M}(\mathbf{z})$ is continuous in \mathbf{z} and independent of $\bar{\mathbf{z}}$ [81]. The step size α_k ensures sufficient decrease in the objective function value, e.g., through line search. Alternatively, \mathbf{z} can be updated using trust-region approaches or using plane search [58, 80, 82]. The linear system (2.7) can be solved using direct (pseudo)inversion or using (preconditioned) conjugate gradients [58].

For the LS objective function (2.1), various Hessian approximations are used depending on the algorithm, e.g.,

- $\mathbf{H}_k = \mathbf{I}$ for gradient descent,
- $\mathbf{H}_k = \mathbf{H}_{k-1} + \mathbf{U}_{k-1} + \mathbf{V}_{k-1}$ with \mathbf{U}_{k-1} and \mathbf{V}_{k-1} symmetric rank-1 matrices constructed using previous gradients and updates of \mathbf{z} for BFGS,
- $\mathbf{H}_k = \mathbf{J}_k^H \mathbf{J}_k$ with $\mathbf{J}_k = \frac{\partial}{\partial \mathbf{z}} \text{vec}(\mathcal{M}(\mathbf{z}) - \mathcal{T})|_{\mathbf{z}=\mathbf{z}_k} = \frac{\partial}{\partial \mathbf{z}} \text{vec}(\mathcal{M}(\mathbf{z}))|_{\mathbf{z}=\mathbf{z}_k}$ for Gauss-Newton,
- $\mathbf{H}_k = \mathbf{J}_k^H \mathbf{J}_k + \lambda_k \mathbf{I}$ for Levenberg-Marquardt.

None of these approximations of the Hessian depend on the tensor \mathcal{T} . Therefore, only the objective function evaluation and the gradient computation require adaptation to exploit structure in \mathcal{T} . Note that this also holds for ALS algorithms.

2.3. Exploiting efficient representations. Due to the explicit construction of the residual tensor $\mathcal{F} = \mathcal{M}(\mathbf{z}) - \mathcal{T}$ in (2.1) and the computation of the gradient, the per-iteration complexity is governed by the number of entries in the tensor, i.e., $\mathcal{O}(\prod_n I_n)$. To reduce the complexity, we assume \mathcal{T} is a function of a limited number of parameters. The parameters can, for instance, be the set of positions and values of nonzero entries in a sparse tensor, the factor matrices of a CPD, or the underlying signal in a Hankelized tensor. To exploit the structure in \mathcal{T} and to avoid the creation of the residual tensor, we rewrite the objective function and gradients, similar to (1.2), which reveals the core operations. More concretely, we have

$$(2.8) \quad f(\mathbf{z}) = \frac{1}{2} \|\mathcal{M}(\mathbf{z})\|^2 - \langle \mathcal{M}(\mathbf{z}), \mathcal{T} \rangle + \frac{1}{2} \|\mathcal{T}\|^2,$$

which requires the efficient computation of the norms of the model and the tensor, and the inner product between the model and the tensor. Similarly, the gradient, which is given by $\bar{\mathbf{g}} = \mathbf{J}^H \text{vec}(\mathcal{F})$, with $\mathbf{J} = \frac{\partial \text{vec}(\mathcal{F})}{\partial \mathbf{z}}$ the Jacobian matrix, is rewritten by grouping the model and data dependent terms as

$$(2.9) \quad \bar{\mathbf{g}} = \underbrace{\mathbf{J}^H \text{vec}(\mathcal{M}(\mathbf{z}))}_{\bar{\mathbf{g}}_{\mathcal{M}}} - \underbrace{\mathbf{J}^H \text{vec}(\mathcal{T})}_{\bar{\mathbf{g}}_{\mathcal{T}}}$$

$$(2.10) \quad = \bar{\mathbf{g}}_{\mathcal{M}} - \bar{\mathbf{g}}_{\mathcal{T}}.$$

(In the derivations, we work with \mathbf{g} rather than $\bar{\mathbf{g}}$.) The core operations are the MTKRPROD for the CPD-based algorithms and the MTKRONPROD for the BTD-based algorithms as shown below. Extensions to constrained and coupled constraints conclude this section.

CPD-based algorithms. Let us partition \mathbf{g} according to the variables \mathbf{z} , i.e., in the case when $\mathcal{M}(\mathbf{z})$ is a CPD, we have

$$(2.11) \quad \mathbf{g} = [\text{vec}(\mathbf{G}^{(1)}); \dots; \text{vec}(\mathbf{G}^{(N)})],$$

with

$$(2.12) \quad \text{vec}(\mathbf{G}^{(n)}) = \frac{\partial f}{\partial \text{vec}(\mathbf{A}^{(n)})}, \quad n = 1, \dots, N.$$

By grouping the model and data dependent terms as in (2.10), the gradients become

$$(2.13) \quad \mathbf{G}_{\mathcal{M}}^{(n)} = \bar{\mathbf{A}}^{(n)} \left(\underset{k \neq n}{*} \mathbf{A}^{(k)\text{H}} \mathbf{A}^{(k)} \right),$$

$$(2.14) \quad \mathbf{G}_{\mathcal{T}}^{(n)} = \bar{\mathbf{T}}_{(n)} \left(\underset{k \neq n}{\odot} \mathbf{A}^{(k)} \right).$$

It is clear that $\mathbf{G}_{\mathcal{M}}^{(n)}$ can be computed efficiently after forming the inner products $\mathbf{A}^{(k)\text{H}} \mathbf{A}^{(k)}$, $k = 1, \dots, N$. We therefore focus on the operation in (2.14), which involves the common matricized tensor times Khatri–Rao product (MTKRPROD):

$$(\text{MTKRPROD}) \quad \bar{\mathbf{T}}_{(n)} \left(\mathbf{A}^{(N)} \odot \mathbf{A}^{(N-1)} \odot \dots \odot \mathbf{A}^{(n+1)} \odot \mathbf{A}^{(n-1)} \odot \dots \odot \mathbf{A}^{(1)} \right).$$

In the case when $\mathcal{M}(\mathbf{z})$ is an LL1 decomposition, $\mathbf{A}^{(3)} = \mathbf{C}\mathbf{P}$ is used to form the inner products and to compute the MTKRPROD. From the chain rule it follows that for $n = 3$, $\mathbf{G}^{(n)}$, which is now the derivative w.r.t. \mathbf{C} , i.e., $\mathbf{G}^{(3)} = \frac{\partial f}{\partial \mathbf{C}}$, is given by

$$(2.15) \quad \mathbf{G}_{\mathcal{M}}^{(3)} = \bar{\mathbf{A}}^{(3)} \left(\underset{k \neq 3}{*} \mathbf{A}^{(k)\text{H}} \mathbf{A}^{(k)} \right) \mathbf{P}^{\text{T}},$$

$$(2.16) \quad \mathbf{G}_{\mathcal{T}}^{(3)} = \bar{\mathbf{T}}_{(3)} \left(\underset{k \neq 3}{\odot} \mathbf{A}^{(k)} \right) \mathbf{P}^{\text{T}}.$$

BTD-based algorithms. Similarly, in the case $\mathcal{M}(\mathbf{z})$ is an LMLRA, the gradient is partitioned as

$$(2.17) \quad \mathbf{g} = [\text{vec}(\mathcal{G}^{(0)}); \text{vec}(\mathbf{G}^{(1)}); \dots; \text{vec}(\mathbf{G}^{(N)})],$$

with

$$(2.18) \quad \text{vec}(\mathcal{G}^{(0)}) = \frac{\partial f}{\partial \text{vec}(\mathcal{S})},$$

$$(2.19) \quad \text{vec}(\mathbf{G}^{(n)}) = \frac{\partial f}{\partial \text{vec}(\mathbf{U}^{(n)})}, \quad n = 1, \dots, N.$$

More specifically, the gradient expressions for the core tensors are given by

$$(2.20) \quad \mathcal{G}_{\mathcal{M}}^{(0)} = \bar{\mathcal{S}} \cdot_1 \left(\mathbf{U}^{(1)\text{T}} \bar{\mathbf{U}}^{(1)} \right) \dots \cdot_N \left(\mathbf{U}^{(N)\text{T}} \bar{\mathbf{U}}^{(N)} \right),$$

$$(2.21) \quad \mathcal{G}_{\mathcal{T}}^{(0)} = \bar{\mathcal{T}} \cdot_1 \mathbf{U}^{(1)\text{T}} \dots \cdot_N \mathbf{U}^{(N)\text{T}},$$

and the expressions for the factors matrices by

$$(2.22) \quad \mathbf{G}_{\mathcal{M}}^{(n)} = \bar{\mathbf{U}}^{(n)} \bar{\mathbf{S}}_{(n)} \left(\bigotimes_{k \neq n} \mathbf{U}^{(k)\mathrm{H}} \mathbf{U}^{(k)} \right) \mathbf{S}_{(n)}^{\mathrm{T}},$$

$$(2.23) \quad \mathbf{G}_{\mathcal{T}}^{(n)} = \bar{\mathbf{T}}_{(n)} \left(\bigotimes_{k \neq n} \mathbf{U}^{(k)} \right) \mathbf{S}_{(n)}^{\mathrm{T}},$$

for $n = 1, \dots, N$. In the case of more general BTD with R terms, an additional summation is introduced; see, e.g., [83]. As (2.21) can be written as $\mathrm{vec} \left(\mathcal{G}_{\mathcal{T}}^{(0)} \right)^{\mathrm{T}} = \mathrm{vec}(\mathcal{T})^{\mathrm{H}} (\bigotimes_n \mathbf{U}^{(n)})$, both (2.21) and (2.23) require a vectorized or matricized tensor times Kronecker product (MTKRONPROD):

$$(\text{MTKRONPROD}) \quad \bar{\mathbf{T}}_{(n)} \left(\mathbf{U}^{(N)} \otimes \mathbf{U}^{(N-1)} \otimes \dots \otimes \mathbf{U}^{(n+1)} \otimes \mathbf{U}^{(n-1)} \otimes \dots \otimes \mathbf{U}^{(1)} \right).$$

Note that

$$(2.24) \quad \mathrm{vec} \left(\mathcal{G}_{\mathcal{T}}^{(0)} \right)^{\mathrm{T}} = \mathrm{vec} \left(\mathbf{U}^{(1)\mathrm{T}} \bar{\mathbf{T}}_{(1)} \left(\bigotimes_{k \neq 1} \mathbf{U}^{(k)} \right) \right);$$

hence the MTKRONPROD from (2.23) can be reused, although this may lead to a higher complexity if the construction of the $I \times R^{N-1}$ matrix $(\bigotimes_{k \neq n} \mathbf{U}^{(k)}) \mathbf{S}_{(n)}^{\mathrm{T}}$ can be avoided such as, e.g., for the polyadic format.

Core operations, constraints, and coupling. In order to reduce the complexity from \mathcal{O} (tensor entries) to \mathcal{O} (parameters), structure exploiting implementations of the Frobenius norm and inner products with CPD and BTD models are required to evaluate the objective function (2.8) as well as implementations for the MTKRPROD and the MTKRONPROD for CPD and BTD gradients, respectively. Examples of such implementations are discussed in section 3. In contrast to approaches taken in, e.g., [7, 15], the original factor matrices and core tensors remain the optimization variables. Therefore, algorithms for coupling datasets, e.g., [2, 83], only require new implementations of the core operations. Similarly, parametric constraints such as nonnegativity, orthogonality, or Vandermonde structure can be handled using the chain rule as in [83], as well as other types of constraints involving projecting onto a feasible set; see, e.g., [38, 46]. An overview of techniques for constrained optimization for tensor decompositions is given in [90]. Concrete examples are given in the experiments in section 4.

3. Operations on efficient representations. In the previous section, the computation of the norm, inner product, MTKRPROD, and MTKRONPROD were identified as the four key operations that allow tensor decomposition algorithms to work on structured tensors. In this section, we show that it is relatively straightforward to exploit efficient representations and to avoid the construction of the full tensor. We show the complexity is indeed governed by the number of parameters in the efficient representation, rather than the number of entries in the tensor. Unless stated otherwise, only third-order tensors $\mathcal{T} \in \mathbb{K}^{I \times I \times I}$ are considered to simplify notation and complexity expressions. In this section, we first give an overview of results for the polyadic (subsection 3.1) and Tucker (subsection 3.2) formats as efficient representations of \mathcal{T} , since expressions involving one of these formats occur in all of the discussed decompositions. Then, efficient implementations for the TT format (subsection 3.3) are discussed and two implicit tensorizations—Hankelization (subsection 3.4)

and Löwnerization (subsection 3.5)—are derived. These and other implementations are available in Tensorlab [94].

We derive expressions for the two computational families: the rank- R CPD $\mathcal{M}_{\text{CPD}} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{C} \rrbracket$ and the multilinear rank- (R, R, R) LMLRA $\mathcal{M}_{\text{LMLRA}} = \mathcal{S} \cdot_1 \mathbf{U} \cdot_2 \mathbf{V} \cdot_3 \mathbf{W}$. The extension to a BTD introduces additional summations and is omitted here; see [83]. For the MTKRPROD and MTKRONPROD operations, only the case $n = 1$ is shown, unless the expressions do not generalize trivially for $n = 2, \dots, N$. Table 3.1 gives a summary of the computational per-iteration complexity when computing a rank- R CPD and clearly illustrates independence on the total number of entries I^N .

TABLE 3.1

Computational per-iteration complexity when computing a rank- R CPD of an N th-order $I \times \dots \times I$ tensor given in its efficient representation. For Hankelization and Löwnerization, second-order tensorization of K signals is assumed, resulting in a tensor of size $I \times I \times K$. The number of samples is $M = 2I - 1$ for Hankel and $M = 2I$ for Löwner. The number of iterations depends on the optimization algorithm used: when using second-order information, e.g., in GN, one typically needs far fewer iterations than when using only first-order information or ALS.

Structure	Function	Complexity
Dense	Objective	$\mathcal{O}(RI^N)$
	Gradient	$\mathcal{O}(NRI^N)$
Polyadic	Parameters	NIF
	Objective	$\mathcal{O}(NIFR + NIR^2)$
	Gradient	$\mathcal{O}(NIFR + NIR^2)$
Tucker	Parameters	$\mathcal{O}(NIJ + J^N)$
	Objective	$\mathcal{O}(NIJR + J^N R + NIR^2)$
	Gradient	$\mathcal{O}(NIJR + J^N NR + NIR^2)$
TT	Parameters	$\mathcal{O}(NIr^2)$
	Objective	$\mathcal{O}(NIr^2 R + NIR^2)$
	Gradient	$\mathcal{O}(NIr^2 R + NIR^2)$
Hankel	Parameters	MK
	Objective	$\mathcal{O}(M(K + \log_2 M)R + (M + K)R^2)$
	Gradient	$\mathcal{O}(M(K + \log_2 M)R + (M + K)R^2)$
Löwner	Parameters	$M(K + 1)$
	Objective	$\mathcal{O}(M(K + \log_2 M)R + (M + K)R^2)$
	Gradient	$\mathcal{O}(M(K + \log_2 M)R + (M + K)R^2)$

3.1. Polyadic format. In the polyadic format, a third-order tensor \mathcal{T} is given by a set of factor matrices \mathbf{X} , \mathbf{Y} , and \mathbf{Z} with F columns, i.e., $\mathcal{T} = \llbracket \mathbf{X}, \mathbf{Y}, \mathbf{Z} \rrbracket$. The number of variables is therefore $F \sum_n I_n$ instead of $\prod_n I_n$. Efficient implementations in the case \mathcal{T} admits a PD, with F not necessarily equal to R , have been presented in [4] (see Kruskal tensor) and are here extended to complex data:

$$(3.1) \quad \|\mathcal{T}\|^2 = \mathbf{1}^T (\mathbf{X}^H \mathbf{X} * \mathbf{Y}^H \mathbf{Y} * \mathbf{Z}^H \mathbf{Z}) \mathbf{1},$$

$$(3.2) \quad \langle \mathcal{M}_{\text{CPD}}, \mathcal{T} \rangle = \mathbf{1}^T (\mathbf{X}^H \mathbf{A} * \mathbf{Y}^H \mathbf{B} * \mathbf{Z}^H \mathbf{C}) \mathbf{1},$$

$$(3.3) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{C} \odot \mathbf{B}) = \bar{\mathbf{X}} (\mathbf{Y}^H \mathbf{B} * \mathbf{Z}^H \mathbf{C}).$$

The complexity of these operations is governed by the construction of the inner products which require $\mathcal{O}(NIF^2)$ and $\mathcal{O}(NIFR)$ flop. Each MTKRPROD in (3.3) required for $\mathbf{G}_{\mathcal{T}}^{(n)}$ can then be computed using $\mathcal{O}(IFR)$ flop.

When computing an LMLRA, the following inner product is required:

$$(3.4) \quad \langle \mathcal{M}_{\text{LMLRA}}, \mathcal{T} \rangle = \text{trace}(\mathcal{S} \cdot_1 \mathbf{X}^H \mathbf{U} \cdot_2 \mathbf{Y}^H \mathbf{V} \cdot_3 \mathbf{Z}^H \mathbf{W}),$$

as well as the MTKRONPROD operation

$$(3.5) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{W} \otimes \mathbf{V}) = \bar{\mathbf{X}}(\mathbf{W}^H \mathbf{Z} \odot \mathbf{V}^H \mathbf{Y})^H,$$

$$(3.6) \quad \text{vec}(\mathcal{T})^H(\mathbf{W} \otimes \mathbf{V} \otimes \mathbf{U}) = \text{vec}(\llbracket \mathbf{U}^H \mathbf{X}, \mathbf{V}^H \mathbf{Y}, \mathbf{W}^H \mathbf{Z} \rrbracket)^H.$$

The computation of the trace in (3.4) requires only the diagonal entries of the LMLRA and takes $\mathcal{O}(NIFR + FR^N)$ flop, including the construction of the inner products. To compute the actual gradient $\mathbf{G}_{\mathcal{T}}^{(1)}$, the result of (3.5) is multiplied by $\mathbf{S}_{(1)}^T$, requiring the product $(\mathbf{W}^H \mathbf{Z} \odot \mathbf{V}^H \mathbf{Y})^H \mathbf{S}_{(1)}^T$, which is a transposed MTKRPROD involving usually small matrices. (In many practical applications the rank or multilinear rank is small compared to the tensor dimensions.) This product can be computed efficiently using, e.g., [69, 87] and requires $\mathcal{O}(NFR^N + NIFR)$ flop. The complexity of (3.6) is the same.

3.2. Tucker format. In the Tucker format a third-order tensor $\mathcal{T} = \mathcal{Q} \cdot_1 \mathbf{X} \cdot_2 \mathbf{Y} \cdot_3 \mathbf{Z}$ is defined by a third-order core tensor $\mathcal{Q} \in \mathbb{K}^{J_1 \times J_2 \times J_3}$ and the factor matrices $\mathbf{X} \in \mathbb{K}^{I \times J_1}$, $\mathbf{Y} \in \mathbb{K}^{I \times J_2}$, $\mathbf{Z} \in \mathbb{K}^{I \times J_3}$. The Tucker format often follows from the computation of an LMLRA or an MLSVD, e.g., as the result of a compression step [7]. For the complexity analysis, we assume $J_n = J$ for $n = 1, \dots, N$. We assume that the factor matrices have orthonormal columns, which is always possible through normalization.

The norm of \mathcal{T} is computed in $\mathcal{O}(J^N)$ operations as

$$(3.7) \quad \|\mathcal{T}\|^2 = \|\mathcal{Q}\|^2.$$

The inner product with \mathcal{M}_{CPD} and MTKRPROD is similar to (3.4) and (3.5):

$$(3.8) \quad \langle \mathcal{M}_{\text{CPD}}, \mathcal{T} \rangle = \text{trace}(\bar{\mathcal{Q}} \cdot_1 \mathbf{A}^T \bar{\mathbf{X}} \cdot_2 \mathbf{B}^T \bar{\mathbf{Y}} \cdot_3 \mathbf{C}^T \bar{\mathbf{Z}}),$$

$$(3.9) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{C} \odot \mathbf{B}) = \bar{\mathbf{X}} \bar{\mathbf{Q}}_{(1)}(\mathbf{Z}^H \mathbf{C} \odot \mathbf{Y}^H \mathbf{B}),$$

and both require $\mathcal{O}(NIJR + J^N R)$ flop.

To compute an LMLRA, the following expressions for the inner product and the MTKRONPROD in (2.23) can be used:

$$(3.10) \quad \langle \mathcal{M}_{\text{LMLRA}}, \mathcal{T} \rangle = \langle \mathcal{S}, \mathcal{Q} \cdot_1 \mathbf{U}^H \mathbf{X} \cdot_2 \mathbf{V}^H \mathbf{Y} \cdot_3 \mathbf{W}^H \mathbf{Z} \rangle,$$

$$(3.11) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{W} \otimes \mathbf{V}) = \bar{\mathbf{X}} \bar{\mathbf{Q}}_{(1)}(\mathbf{Z}^H \mathbf{W} \otimes \mathbf{Y}^H \mathbf{V}).$$

Similarly, the MTKRONPROD in (2.21) is computed as

$$(3.12) \quad \text{vec}(\mathcal{T})^H(\mathbf{W} \otimes \mathbf{V} \otimes \mathbf{U}) = \text{vec}(\mathcal{Q} \cdot_1 \mathbf{U}^H \mathbf{X} \cdot_2 \mathbf{V}^H \mathbf{Y} \cdot_3 \mathbf{W}^H \mathbf{Z})^H.$$

All computations involving the LMLRA require $\mathcal{O}(NIJR + J^N R + JR^N)$ flop.

3.3. Tensor Train format. An N th-order tensor $\mathcal{T} \in \mathbb{K}^{I_1 \times \dots \times I_N}$ can be represented by N core matrices or tensors $\mathcal{Q}^{(n)} \in \mathbb{K}^{r_{n-1} \times I_n \times r_n}$ serially linked by $N-1$ indices such that

$$(3.13) \quad t_{i_1, \dots, i_N} = \sum_{s_1=1}^{r_1} \cdots \sum_{s_{N-1}=1}^{r_{N-1}} q_{1, i_1, s_1}^{(1)} q_{s_1, i_2, s_2}^{(2)} \cdots q_{s_{N-1}, i_N, 1}^{(N)} \quad \text{or}$$

$$(3.14) \quad \mathcal{T} = \text{tt}(\mathcal{G}^{(1)}, \dots, \mathcal{G}^{(N)}).$$

The parameters r_n , $n = 0, \dots, N$, are the compression ranks with $r_0 = r_N = 1$. This formulation is known as matrix product states [59] in computational chemistry and as Tensor Trains (TT) [60] in scientific computing and can be computed using SVDs, cross approximation, or optimization [60, 61, 62, 66]. Assuming all dimensions and compression ranks are equal, i.e., $r_n = r$ for $n = 1, \dots, N-1$, the total number of parameters is $\mathcal{O}(NIr^2)$.

The squared norm of \mathcal{T} is given by $\|\mathcal{T}\|^2 = |\mathbf{F}^{(N)}|$, in which the scalar $\mathbf{F}^{(N)}$ is constructed using the recursive formula

$$(3.15) \quad \mathbf{F}^{(0)} = 1,$$

$$(3.16) \quad \mathbf{F}^{(n)} = \left(\mathbf{Q}^{(n)} \cdot_1 \mathbf{F}^{(n-1)} \right)_{(1,2)}^H \mathbf{Q}_{(1,2)}^{(n)}, \quad n = 1, \dots, N.$$

This way $\|\mathcal{T}\|^2$ can be computed in $\mathcal{O}(NIr^3)$ flop. Note that this recursive formula is essentially identical to the formula derived in [66].

To compute inner products with \mathcal{M}_{CPD} and $\mathcal{M}_{\text{LMLRA}}$, we first introduce the auxiliary core tensors $\hat{\mathcal{Q}}^{(n)}$. Let $\mathbf{A}^{(1)} = \mathbf{A}$, $\mathbf{A}^{(2)} = \mathbf{B}$, and $\mathbf{A}^{(3)} = \mathbf{C}$ for \mathcal{M}_{CPD} , and let $\mathbf{A}^{(1)} = \mathbf{U}$, $\mathbf{A}^{(2)} = \mathbf{V}$, and $\mathbf{A}^{(3)} = \mathbf{W}$ for $\mathcal{M}_{\text{LMLRA}}$; then $\hat{\mathcal{Q}}^{(n)}$ is computed using $\mathcal{O}(Ir^2R)$ flop as

$$(3.17) \quad \hat{\mathcal{Q}}^{(n)} = \bar{\mathcal{Q}}^{(n)} \cdot_2 \mathbf{A}^{(n)T}.$$

The inner products are then computed as

$$(3.18) \quad \langle \mathcal{M}_{\text{CPD}}, \mathcal{T} \rangle = \text{trace} \left(\text{tt}(\hat{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)}) \right),$$

$$(3.19) \quad \langle \mathcal{M}_{\text{LMLRA}}, \mathcal{T} \rangle = \left\langle \text{tt}(\hat{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)}), \bar{\mathcal{S}} \right\rangle.$$

As only the diagonal entries of $\text{tt}(\hat{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)})$ are required to compute the trace in (3.18), the construction of the $R \times \dots \times R$ tensor can be avoided. The construction of the diagonal entries in (3.18) and the inner product in (3.19) requires $\mathcal{O}(NRr^2)$ and $\mathcal{O}(R^N r^2)$ flop, respectively.

To compute an MTKRPROD, the auxiliary cores $\hat{\mathcal{Q}}^{(n)}$, $n = 1, \dots, N$, are used again. Define $\mathcal{P}^{(1)} = \text{tt}(\bar{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)})$ and $\mathbf{p}_r^{(1)}$ as the mode-1 vectors $\mathcal{P}^{(1)}(:, r, r)$, $r = 1, \dots, R$; then

$$(3.20) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{C} \odot \mathbf{B}) = \begin{bmatrix} \mathbf{p}_1^{(1)} & \mathbf{p}_2^{(1)} & \dots & \mathbf{p}_R^{(1)} \end{bmatrix}.$$

The complexity of computing one MTKRPROD is therefore $\mathcal{O}(Ir^2R)$, assuming auxiliary cores have been constructed.

Similarly, the MTKRONPROD can be formed efficiently as

$$(3.21) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{W} \otimes \mathbf{V}) = \left(\text{tt}(\bar{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)}) \right)_{(1)}.$$

To compute the gradient terms $\mathbf{G}_{\mathcal{T}}^{(n)}$ in (2.23) the unfolding $(\text{tt}(\bar{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)}))_{(1)}$ can be constructed and multiplied with $\mathbf{S}_{(1)}^T$, which requires $\mathcal{O}(IR^N + IR^{N-1}r^2)$ flop. However, by exploiting the TT structure, the product can be computed directly in $\mathcal{O}(rR^N + r^2R^{N-1} + Ir^2R)$ flop. The precise implementation is beyond the scope of this paper. Finally, to compute $\mathcal{G}_{\mathcal{T}}^{(0)}$ in $\mathcal{O}(r^2R^N)$ flop, we can use

$$(3.22) \quad \text{vec}(\mathcal{T})^H (\mathbf{W} \otimes \mathbf{V} \otimes \mathbf{U}) = \text{vec} \left(\text{tt}(\hat{\mathcal{Q}}^{(1)}, \hat{\mathcal{Q}}^{(2)}, \hat{\mathcal{Q}}^{(3)}) \right)^T.$$

3.4. Implicit Hankelization. Hankelization of a (possibly complex) signal vector $\mathbf{m} \in \mathbb{K}^M$ of length $M = \sum_{n=1}^N I_n - N + 1$ yields an $I_1 \times I_2$ matrix \mathbf{T} or $I_1 \times \cdots \times I_N$ tensor \mathcal{T} with constant antidiagonals ($t_{i_1 i_2} = m_{i_1 + i_2 - 1}$) or constant antidiagonal hyperplanes in the tensor ($t_{i_1 \dots i_N} = m_{i_1 + \dots + i_N - N + 1}$) for $N = 2$ and $N > 2$, respectively; see [25, 68] for some background on signal processing. In BSS, a set of K vectors $\mathbf{m}_k \in \mathbb{K}^M$ is often given [16, 25]. Each vector \mathbf{m}_k can be Hankelized separately, and the resulting matrices or tensors can be concatenated along the $(N + 1)$ th mode. To simplify the expressions,¹ we only consider second-order Hankelization and assume $I_1 = I_2 = I$. Let us stack the vectors \mathbf{m}_k as the columns of \mathbf{M} . Each entry of the $I \times I \times K$ tensor \mathcal{T} is given by $t_{ijk} = m_{i+j-1,k}$. As multiplying a Hankel matrix by a vector corresponds to a convolution, FFT can be used to speed up computations [3, 28]. Let \mathfrak{F} denote the M -point FFT, and let \mathfrak{F}^{-1} denote the inverse M -point FFT. Only the last I rows of the inverse M -point FFT are retained when the operator $\underline{\mathfrak{F}}^{-1}$ is used, i.e., $\underline{\mathfrak{F}}^{-1}\mathbf{x} = [\mathbf{0}_{I \times (M-I)} \quad \mathbf{I}_I] (\mathfrak{F}^{-1}\mathbf{x})$ for $\mathbf{x} \in \mathbb{C}^M$.

The squared norm can be computed in $\mathcal{O}(MK)$ flop as

$$(3.23) \quad \|\mathcal{T}\|^2 = \mathbf{w}^T (\mathbf{M} * \overline{\mathbf{M}}) \mathbf{1}_K,$$

in which \mathbf{w} is a vector that contains the number of occurrences of each entry of \mathbf{M} in \mathcal{T} and can be computed as $\mathbf{w} = \text{convolution}(\mathbf{1}_I, \mathbf{1}_I)$. The inner products can be computed as

$$(3.24) \quad \langle \mathcal{M}_{\text{CPD}}, \mathcal{T} \rangle = \langle \mathfrak{F}^{-1}(\mathfrak{F}\mathbf{A} * \mathfrak{F}\mathbf{B})\mathbf{C}^T, \mathbf{M} \rangle,$$

$$(3.25) \quad \langle \mathcal{M}_{\text{LMLRA}}, \mathcal{T} \rangle = \langle \mathfrak{F}^{-1}((\mathfrak{F}\mathbf{V} \odot^T \mathfrak{F}\mathbf{U})\mathbf{S}_{(1,2)}\mathbf{W}^T), \mathbf{M} \rangle.$$

The computational cost is $\mathcal{O}(RM(K + \log_2 M))$ and $\mathcal{O}((R + K)M \log_2 M + MR^3)$ flop for the inner product with \mathcal{M}_{CPD} and $\mathcal{M}_{\text{LMLRA}}$, respectively.

To compute the MTKRPROD, the Hankel structure can be exploited as

$$(3.26) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{C} \odot \mathbf{B}) = \underline{\mathfrak{F}}^{-1}(\mathfrak{F}\bar{\mathbf{M}}\mathbf{C} * \mathfrak{F}\mathbf{B}^f),$$

$$(3.27) \quad \bar{\mathbf{T}}_{(2)}(\mathbf{C} \odot \mathbf{A}) = \underline{\mathfrak{F}}^{-1}(\mathfrak{F}\bar{\mathbf{M}}\mathbf{C} * \mathfrak{F}\mathbf{A}^f),$$

$$(3.28) \quad \bar{\mathbf{T}}_{(3)}(\mathbf{B} \odot \mathbf{A}) = \mathbf{M}^H \underline{\mathfrak{F}}^{-1}(\mathfrak{F}\mathbf{A} * \mathfrak{F}\mathbf{B}),$$

requiring $\mathcal{O}(RM(K + \log_2 M))$ flop each. (\mathbf{A}^f and \mathbf{B}^f are the flipped versions of \mathbf{A} and \mathbf{B} ; see subsection 1.3.)

The MTKRONPROD in $\mathbf{G}_{\mathcal{T}}^{(n)}$ in (2.23) can be computed in $\mathcal{O}(R^2 M \log_2 M + RMK)$ flop as

$$(3.29) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{W} \otimes \mathbf{V}) = \underline{\mathfrak{F}}^{-1}(\mathfrak{F}\bar{\mathbf{M}}\mathbf{W} \odot^T \mathfrak{F}\mathbf{V}^f),$$

$$(3.30) \quad \bar{\mathbf{T}}_{(2)}(\mathbf{W} \otimes \mathbf{U}) = \underline{\mathfrak{F}}^{-1}(\mathfrak{F}\bar{\mathbf{M}}\mathbf{W} \odot^T \mathfrak{F}\mathbf{U}^f),$$

$$(3.31) \quad \bar{\mathbf{T}}_{(3)}(\mathbf{V} \otimes \mathbf{U}) = \mathbf{M}^H \underline{\mathfrak{F}}^{-1}(\mathfrak{F}\mathbf{V} \odot^T \mathfrak{F}\mathbf{U}).$$

The multiplication with $\mathbf{S}_{(n)}^T$ in (2.23) additionally takes $\mathcal{O}(IR^3)$ flop for $n = 1, 2$ and $\mathcal{O}(KR^3)$ flop for $n = 3$. However, by first computing the multiplication of the row-wise Khatri–Rao product with $\mathbf{S}_{(n)}^T$, i.e., before performing the inverse FFT, the total cost can be reduced to $\mathcal{O}(RM(K + \log_2 M) + MR^3)$. Finally, the MTKRONPROD in all modes can be computed in $\mathcal{O}(R^2 M \log_2 M + R^3 M + RMK)$ flop using

$$(3.32) \quad \text{vec}(\mathcal{T})^H (\mathbf{W} \otimes \mathbf{V} \otimes \mathbf{U}) = \text{vec}((\mathfrak{F}^{-1}(\mathfrak{F}\mathbf{V} \odot^T \mathfrak{F}\mathbf{U}))^T \bar{\mathbf{M}}\mathbf{W})^T.$$

¹More general implementations are available in Tensorlab 3.0 [94].

3.5. Implicit Löwnerization. Löwner matrices and tensors have attractive properties for applications involving rational functions [25, 27]. Given a function $h : \mathbb{K} \rightarrow \mathbb{K}$ evaluated at M points $t_i \in T = \{t_1, \dots, t_M\}$, we partition T into two disjoint point sets $X = \{x_1, \dots, x_I\}$ and $Y = \{y_1, \dots, y_J\}$ such that $T = X \cup Y$ and $M = I + J$. The entries in a Löwner matrix $\mathbf{L} \in \mathbb{K}^{I \times J}$ are then given by

$$(3.33) \quad l_{ij} = \frac{h(x_i) - h(y_j)}{x_i - y_j}, \quad i = 1, \dots, I, \quad \text{and} \quad j = 1, \dots, J.$$

While higher-order generalizations of the Löwner transformation exist [25, 26], we focus on third-order tensors in which each k th frontal slice is a Löwner matrix constructed from a function h_k evaluated at the points in X and Y . Let $\mathbf{P} \in \mathbb{K}^{I \times K}$ and $\mathbf{Q} \in \mathbb{K}^{J \times K}$ contain sampled function values at the points in X and Y , respectively, i.e.,

$$(3.34) \quad p_{ik} = h_k(x_i), \quad i = 1, \dots, I, \quad \text{and} \quad k = 1, \dots, K,$$

$$(3.35) \quad q_{jk} = h_k(y_j), \quad j = 1, \dots, J, \quad \text{and} \quad k = 1, \dots, K.$$

The tensor $\mathcal{T} \in \mathbb{K}^{I \times J \times K}$ is then fully determined by \mathbf{P} , \mathbf{Q} , X , and Y . To simplify complexity expressions, we take $I = J$. Each frontal slice \mathbf{T}_k , $k = 1, \dots, K$, can be written as

$$(3.36) \quad \mathbf{T}_k = \text{Diag}(\mathbf{p}_k) \mathbf{M} - \mathbf{M} \text{Diag}(\mathbf{q}_k) = \mathbf{p}_k \odot^T \mathbf{M} - \mathbf{M}(\mathbf{q}_k \odot^T \mathbf{I}_J),$$

in which \mathbf{M} is a Cauchy matrix with

$$(3.37) \quad m_{ij} = \frac{1}{x_i - y_j}, \quad 1 \leq i, j \leq I.$$

By assuming that all points in X and Y are equidistant, \mathbf{M} is also a Toeplitz matrix. As multiplying a Toeplitz matrix with a vector \mathbf{x} can be seen as a convolution of the generating vector $\mathbf{v} = [m_{1,I}; m_{1,I-1}; \dots; m_{1,1}; m_{2,1}; \dots; m_{I,1}] \in \mathbb{K}^{2I-1}$ with \mathbf{x} , it is not necessary to construct \mathbf{M} , and the multiplication can be performed using FFTs. Let \mathfrak{F} be the $(2I-1)$ -point FFT with zero padding for shorter vectors, let \mathfrak{F}^{-1} be the inverse $(2I-1)$ -point FFT, and define $\underline{\mathfrak{F}}^{-1}$ to be the last I rows of the result of the inverse FFT. For a matrix $\mathbf{X} \in \mathbb{K}^{I \times K}$, $\mathbf{M}\mathbf{X}$ and $\mathbf{M}^T \mathbf{X}$ can then be computed as

$$(3.38) \quad \mathbf{M}\mathbf{X} = \underline{\mathfrak{F}}^{-1} (\mathfrak{F} \mathbf{v} \mathbf{1}_K^T * \mathfrak{F} \mathbf{X}),$$

$$(3.39) \quad \mathbf{M}^T \mathbf{X} = \underline{\mathfrak{F}}^{-1} (\mathfrak{F} \mathbf{v}^f \mathbf{1}_K^T * \mathfrak{F} \mathbf{X}).$$

If $\mathfrak{F} \mathbf{v}$ is precomputed, the cost of this multiplication is $\mathcal{O}(2KI + 4KI \log_2 2I)$ flop.

We now define the different operations for Löwner tensors. As the unfoldings $\mathbf{T}_{(n)}$, $n = 1, 2, 3$, are given by

$$(3.40) \quad \mathbf{T}_{(1)} = \mathbf{P} \odot^T \mathbf{M} - \mathbf{M}(\mathbf{Q} \odot^T \mathbf{I}_J),$$

$$(3.41) \quad \mathbf{T}_{(2)} = \mathbf{M}^T (\mathbf{P} \odot^T \mathbf{I}_I) - (\mathbf{Q} \odot^T \mathbf{M}^T),$$

$$(3.42) \quad \mathbf{T}_{(3)} = \mathbf{P}^T (\mathbf{M} \odot^T \mathbf{I}_I) - \mathbf{Q}^T (\mathbf{I}_J \odot^T \mathbf{M}^T),$$

the expressions for the MTKRPROD and MTKRONPROD can be derived easily using multilinear algebra identities, e.g., $(\mathbf{A} \odot^T \mathbf{B})(\mathbf{C} \odot \mathbf{D}) = \mathbf{AC} * \mathbf{BD}$. The MTKRPROD

can be computed efficiently as

$$(3.43) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{C} \odot \mathbf{B}) = \bar{\mathbf{P}}\mathbf{C} * \bar{\mathbf{M}}\mathbf{B} - \bar{\mathbf{M}}(\bar{\mathbf{Q}}\mathbf{C} * \mathbf{B}),$$

$$(3.44) \quad \bar{\mathbf{T}}_{(2)}(\mathbf{C} \odot \mathbf{A}) = \mathbf{M}^H(\bar{\mathbf{P}}\mathbf{C} * \mathbf{A}) - \bar{\mathbf{Q}}\mathbf{C} * \mathbf{M}^H\mathbf{A},$$

$$(3.45) \quad \bar{\mathbf{T}}_{(3)}(\mathbf{B} \odot \mathbf{A}) = \mathbf{P}^H(\bar{\mathbf{M}}\mathbf{B} * \mathbf{A}) - \mathbf{Q}^H(\mathbf{B} * \mathbf{M}^H\mathbf{A}).$$

As each MTKRPROD requires two FFT-based multiplications with \mathbf{M} and two regular matrix-matrix multiplications with \mathbf{P} and \mathbf{Q} , the complexity is $\mathcal{O}(RI(K + \log_2 I))$. Similarly, the MTKRONPROD in (2.23) can be computed as

$$(3.46) \quad \bar{\mathbf{T}}_{(1)}(\mathbf{W} \otimes \mathbf{V}) = \bar{\mathbf{P}}\mathbf{W} \odot^T \bar{\mathbf{M}}\mathbf{V} - \bar{\mathbf{M}}(\bar{\mathbf{Q}}\mathbf{W} \odot^T \mathbf{V}),$$

$$(3.47) \quad \bar{\mathbf{T}}_{(2)}(\mathbf{W} \otimes \mathbf{U}) = \mathbf{M}^H(\bar{\mathbf{P}}\mathbf{W} \odot^T \mathbf{U}) - \bar{\mathbf{Q}}\mathbf{W} \odot^T \mathbf{M}^H\mathbf{U},$$

$$(3.48) \quad \bar{\mathbf{T}}_{(3)}(\mathbf{V} \otimes \mathbf{U}) = \mathbf{P}^H(\bar{\mathbf{M}}\mathbf{V} \odot^T \mathbf{U}) - \mathbf{Q}^H(\mathbf{V} \odot^T \mathbf{M}^H\mathbf{U}),$$

in $\mathcal{O}(R^2I + R^2I \log_2 I + RKI)$ flop for $n = 1, 2$ and $\mathcal{O}(RI \log_2 I + R^2IK)$ flop for $n = 3$. This can be reduced further to $\mathcal{O}(RI \log_2 I + RKI)$ flop for $n = 1, 2$, by exploiting the rowwise Khatri–Rao products when computing the subsequent multiplication with $\mathbf{S}_{(n)}^T$, which costs $\mathcal{O}(IR^3)$. A similar improvement can be made for $n = 3$. The MTKRONPROD in (2.21) is computed as

$$(3.49) \quad \text{vec}(\mathcal{T})^H(\mathbf{W} \otimes \mathbf{V} \otimes \mathbf{U}) = \text{vec}(\mathbf{U}^T(\bar{\mathbf{P}}\mathbf{W} \odot^T \bar{\mathbf{M}}\mathbf{V}) - (\mathbf{M}^H\mathbf{U})^T(\bar{\mathbf{Q}}\mathbf{W} \odot^T \mathbf{V}))^T,$$

and requires $\mathcal{O}(RI \log_2 I + R^3I)$ flop.

The squared Frobenius norm is computed as

$$(3.50) \quad \|\mathcal{T}\|^2 = \mathbf{1}_I^T(\mathbf{M} * \bar{\mathbf{M}})^T((\mathbf{P} * \bar{\mathbf{P}})\mathbf{1}_K)$$

$$(3.51) \quad + \mathbf{1}_I^T(\mathbf{M} * \bar{\mathbf{M}})((\mathbf{Q} * \bar{\mathbf{Q}})\mathbf{1}_K)$$

$$(3.52) \quad - 2\text{Re}(\mathbf{1}_K^T((\mathbf{M} * \bar{\mathbf{M}})^T\mathbf{P} * \bar{\mathbf{Q}})\mathbf{1}_K).$$

As only $K + 2$ FFTs are required, the total computational cost is $\mathcal{O}(4IK \log_2 2I)$ flop. To compute the inner products with \mathcal{M}_{CPD} and $\mathcal{M}_{\text{LMLRA}}$, the results from MTKRPROD and MTKRONPROD are reused:

$$(3.53) \quad \langle \mathcal{M}_{\text{CPD}}, \mathcal{T} \rangle = \langle \bar{\mathbf{T}}_{(3)}(\mathbf{B} \odot \mathbf{A}), \bar{\mathbf{C}} \rangle,$$

$$(3.54) \quad \langle \mathcal{M}_{\text{LMLRA}}, \mathcal{T} \rangle = \langle \mathbf{X}^T \bar{\mathbf{T}}_{(1)}(\mathbf{W} \otimes \mathbf{V}), \bar{\mathbf{S}}_{(1)} \rangle.$$

Apart from the cost of the MTKRPROD and the MTKRONPROD, an additional cost of $\mathcal{O}(KR)$ and $\mathcal{O}(R^N)$ flop is incurred for $\langle \mathcal{M}_{\text{CPD}}, \mathcal{T} \rangle$ and $\langle \mathcal{M}_{\text{LMLRA}}, \mathcal{T} \rangle$, respectively.

In the case when the point sets X or Y do not contain equidistant points, \mathbf{M} no longer admits a Toeplitz structure. By exploiting the low displacement rank of Cauchy matrices, it is again not necessary to explicitly construct the tensor. The computational complexity of the multiplications with \mathbf{M} in (3.38) and (3.39) is slightly increased to $\mathcal{O}(4KI \log_2^2 2I)$ flop [29, 30].

4. Experiments. The following experiments illustrate the scalability of constrained CPD, LL1, and BTM algorithms that exploit efficient representations of tensors and their effect on the accuracy of the results. The latter is discussed in subsection 4.1 for ill-conditioned problems and in subsection 4.3 for Hankelized signals. The former is illustrated by computing the nonnegative CPD of a compressed tensor

in subsection 4.2, and for the unconstrained LL1 decomposition and the constrained BTD in subsection 4.3. The CPD error E_{CPD} is defined as

$$(4.1) \quad E_{\text{CPD}} = \max_n \frac{\|\mathbf{A}^{(n)} - \hat{\mathbf{A}}^{(n)}\|}{\|\mathbf{A}^{(n)}\|},$$

in which $\mathbf{A}^{(n)}$ ($\hat{\mathbf{A}}^{(n)}$) are the exact (estimated) factor matrices, $n = 1, \dots, N$, and scaling and permutation indeterminacies are assumed to be resolved using the congruence maximization heuristic; see `cpderr` [94]. All timing results are total computation times for a complete run of the optimization algorithm, in contrast to the derived per-iteration complexities in section 3. Tensorlab 3.0 [94] is used for all experiments in combination with MATLAB 2016b running on a dual socket 20 core Intel Xeon E5-2660 v2 machine with 128GB of RAM running CentOS 7.

4.1. Accuracy and conditioning. In this first experiment, we study the accuracy of recovered factor matrices for the full tensor and its efficient representation. More specifically, the CPD of a tensor with an exact rank- R structure is computed while varying the relative condition number κ , which is defined as in [86]. Concretely, a rank-5 tensor of size $25 \times 25 \times 25$ is constructed using random factor matrices $\mathbf{A}^{(n)}$, $n = 1, 2, 3$. Each random factor vector $\mathbf{a}_r^{(n)}$ has norm one and a fixed angle α w.r.t. the other factor vectors in the same factor matrix, i.e., for $n = 1, 2, 3$, $\|\mathbf{a}_r^{(n)}\| = 1$, and $\cos \alpha = \mathbf{a}_r^{(n)\top} \mathbf{a}_s^{(n)}$ for $r, s = 1, \dots, R$, $r \neq s$. As α decreases, the rank-1 terms become more collinear and the condition number κ increases. A rank-5 CPD is computed from the full tensor and from the structured tensor in the polyadic format using inexact GN, implemented in `cpd_nls`, starting from a perturbed exact solution, i.e., the initial guesses for factor matrices $\mathbf{U}^{(n)}$ are constructed as $\mathbf{U}^{(n)} = \mathbf{A}^{(n)} + \sigma^{(n)} \mathbf{N}^{(n)}$, with $\mathbf{N}^{(n)}$ a random matrix with entries drawn from the normal distribution and $\sigma^{(n)} = 0.1 \cdot \|\mathbf{A}^{(n)}\| / \|\mathbf{N}^{(n)}\|$. Figure 4.1 shows the error E_{CPD} on the recovered factor matrices for $\alpha = \pi/2, \dots, \pi/180$ (using a logarithmic scale). As expected, the error increases when the condition number κ worsens. For well-conditioned problems, E_{CPD} is of the order of the machine precision $\epsilon \approx 10^{-16}$ if the full tensor is used, while E_{CPD} is higher for the structured tensor. This can be explained by the fact that changes smaller than $\sqrt{\epsilon} \approx 10^{-8}$ cannot be distinguished when using the structured tensor, as the computation of the objective function (2.8) requires the difference of squared, almost equal numbers. For very ill-conditioned problems, E_{CPD} may be undesirably high when using the structured tensor. If the tensor admits an exact decomposition *and* if this exact decomposition is required up to machine precision, using the full tensor may be necessary. However, the structured tensor can still be used as an initialization to reduce the overall computational cost. Note that for very large-scale tensors, the explicit construction of the tensor may not be feasible anyway, and that in many applications, the condition number is limited as indicated by the reference points in Figure 4.1.

4.2. Compression for nonnegative CPD. The following experiments illustrate how compression can be used to reduce the complexity of constrained tensor decompositions. Here, we focus on the widely used nonnegative CPD, which is sped up by first computing a truncated MLSVD or a TT approximation. Two approaches to enforce the constraints are used: a projected GN approach generating updates such that the variables are always positive using active sets [46], and using parameter-based constraints, i.e., each entry is the square of an underlying variable [71]. The former

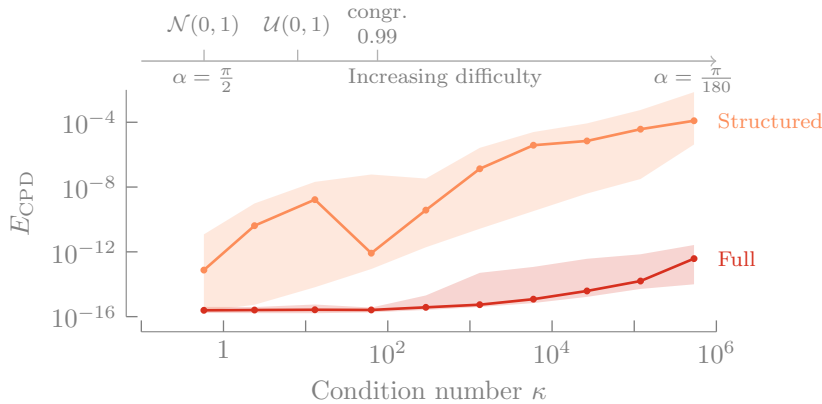


FIG. 4.1. By reducing the angle α between the factor vectors, the condition of the CPD worsens, resulting in a higher error E_{CPD} . For very ill-conditioned problems, using the structured tensor approach may result in an undesirably large error, and a few additional iterations using the full tensor may be required to improve the accuracy. Above the graph, three typical test setups are indicated as references: factor entries drawn from a normal distribution $\mathcal{N}(0,1)$, from a uniform distribution $\mathcal{U}(0,1)$, or a triple bottleneck setup such that the congruence between all factor vectors is 0.99 in each mode. The latter is typically considered very difficult. Results are medians over 100 experiments. The shaded areas give the minimum and maximum errors over all experiments.

approach is implemented using a boundary constrained GN algorithm `nlsb_gnd1`; the latter is implemented using the structured data fusion framework [83] (`sdf_nls` and `struct_nonneg`).

For the first experiment, let \mathcal{T} be a nonnegative rank-10 tensor of size $I \times I \times I$ constructed using random factor matrices $\mathbf{A}^{(n)}$ with entries drawn from a uniform distribution $\mathcal{U}(0,1)$. Gaussian i.i.d. noise is added such that the signal-to-noise ratio (SNR) is 20 dB. Using a randomized MLSVD (`mlsvd_rsi` [91]) the tensor is compressed such that the core \mathcal{G} has size $10 \times 10 \times 10$, i.e., $\mathcal{T} \approx \llbracket \mathcal{G}; \mathbf{U}^{(1)}, \mathbf{U}^{(2)}, \mathbf{U}^{(3)} \rrbracket$. The core tensor \mathcal{G} and the factors $\mathbf{U}^{(n)}$ are then used as structured format to compute a rank-10 CPD. Random initial factor matrices are drawn from the same distribution as $\mathbf{A}^{(n)}$. Figure 4.2 shows the time required to compute an unconstrained CPD (hence relying on CPD uniqueness), a constrained CPD using projected GN (`cpd_nls` with `nlsb_gnd1` solver), and a constrained CPD using parametrization (`sdf_nls`). In the three cases, the computation time for the original tensor \mathcal{T} rises cubically in the dimension I , for $I = 2^6, 2^7, \dots, 2^{10}$, while the time rises linearly using the structured approach. Hence, existing CPD algorithms can be scaled to handle large problems provided that the MLSVD approximation can be computed and the rank is modest. Note that the compression time, which also rises cubically in I from 20 ms for $I = 2^6$ to 8 s for $I = 2^{10}$, is not included. The total time for compression and decomposition is still (far) less than the time required for the decomposition of the full tensor, though. Moreover, in many signal processing and data analysis applications, the compression cost can be amortized by reusing the representation for multiple initializations, ranks, or experiments. The compression cost can be reduced further by resorting to parallelism and/or cross approximation techniques; see, e.g., [9, 56, 63].

In a second experiment, the complexity in function of the order N of the tensor is investigated. We compare the required time to compute a nonnegative CPD of an N th order rank-5 tensor with dimensions $10 \times \dots \times 10$ using the full tensor and its TT approximation. The data is generated by constructing N random factor matrices with

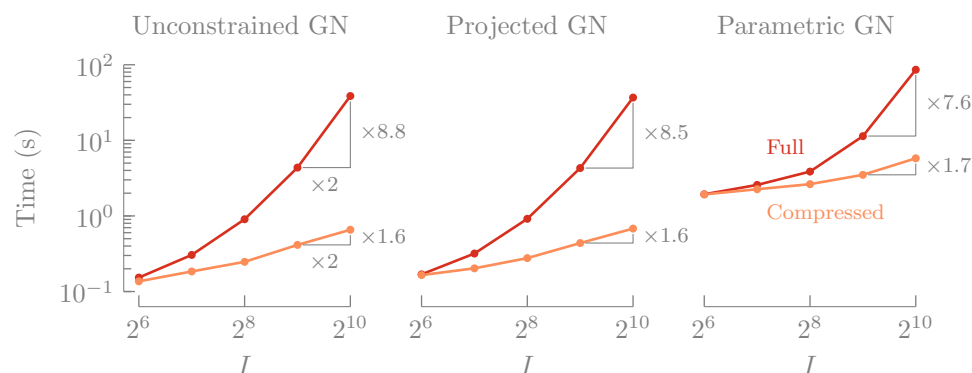


FIG. 4.2. When using a compressed tensor instead of the full tensor, the computational cost scales linearly in the tensor dimensions instead of cubically for a rank-10 nonnegative tensor of size $I \times I \times I$. MLSVD compression with a core size of $10 \times 10 \times 10$ is used. The increase for the structured tensor is actually less than linear ($\times 1.6 - \times 1.7$ instead of $\times 2$), which is caused by an improved multicore usage. The time required to compute the MLSVD increases cubically from 20 ms for $I = 2^6$ to 8 s for $I = 2^{10}$ and is not included. Medians over 100 experiments for each parameter are reported. All algorithms stop if the relative change in function value is smaller than 10^{-7} . The parametric GN approach has a larger overhead and is not preconditioned, explaining the higher computation time.

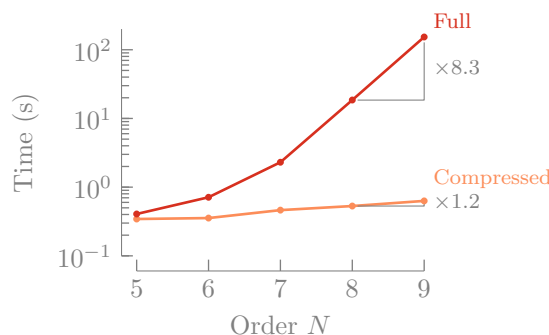


FIG. 4.3. Using a TT approximation instead of the full $10 \times 10 \times \dots \times 10$ tensor removes the exponential dependence on the order N when computing a rank-5 nonnegative CPD: from $N = 8$ to $N = 9$ the time increases with a factor $1.2 \approx 9/8$ using the TT approximation, while the time increase is $8.3 \approx 10$ using the full tensor. The median TT compression time is not included and increases exponentially from 20 ms for $N = 5$ to 114 s for $N = 9$. The results are medians over 100 experiments.

entries drawn from the uniform distribution $\mathcal{U}(0, 1)$. Random Gaussian i.i.d. noise is added to the generated tensor such that the SNR is 20 dB. The core tensors $\mathcal{G}^{(n)}$ corresponding to the TT approximation are computed using `tt_tensor` from the TT Toolbox [62], which we slightly adapted such that all compression ranks $r_n = 5$, $n = 1, \dots, N - 1$. Random factor matrices using the same distribution are used as initialization. Nonnegativity is enforced using the projected GN approach (`cpd_nls` with `nlsb_gnd1` solver). The timing results in Figure 4.3 clearly show that using the TT approximation avoids the curse of dimensionality as the time increases linearly in N in contrast to the time required for the decomposition of the full tensor. Note that we expect the time to increase by a factor 10 if the order increases by one, but the actual time increase is lower as seen in Figure 4.3. This is due to the fact

that the decomposition problem becomes easier as relatively more data is available per variable. This is also reflected in the accuracy: for $N = 5$ the median accuracy $E_{\text{CPD}} = 4.7 \cdot 10^{-3}$, while $E_{\text{CPD}} = 3.8 \cdot 10^{-5}$ for $N = 9$.

4.3. Signal separation through Hankelization. We use the BSS setup

$$(4.2) \quad \mathbf{X} = \mathbf{S}\mathbf{M}^T + \mathbf{N}$$

to recover the unknown sources \mathbf{S} and unknown mixing matrix \mathbf{M} from the given signals \mathbf{X} . The additive Gaussian i.i.d. noise $\mathbf{N} \in \mathbb{R}^{M \times K}$ in the experiments is scaled such that a given SNR is attained. In general, \mathbf{M} and \mathbf{S} cannot be recovered uniquely. To achieve uniqueness, the classical independent component analysis assumes statistically independent sources, while deterministic BSS techniques assume that the sources are, e.g., exponential polynomials or rational functions. In the case when source signals are sums of exponential polynomials, these signals can be recovered using Hankelization [19]. Hankelization along the columns of \mathbf{X} leads to a third-order tensor \mathcal{T} of size $\lfloor \frac{M+1}{2} \rfloor \times \lceil \frac{M+1}{2} \rceil \times K$ in which $\lfloor \cdot \rfloor$ and $\lceil \cdot \rceil$ are the floor and ceil operators, respectively. Hence, if the number of samples M doubles, the number of entries in \mathcal{T} quadruples. Given this fast increase, we show that exploiting the structure is crucial to apply deterministic BSS methods to problems of realistic sizes. If each of the R source signals is a sum of Q_r exponential polynomials of degree d_{qr} , $q = 1, \dots, Q_r$, \mathcal{T} can be decomposed into a sum of R multilinear rank- $(L_r, L_r, 1)$ terms with $L_r = \sum_{q=1}^{Q_r} (d_{qr} + 1)$ [19], i.e.,

$$(4.3) \quad \mathcal{T} = \sum_{r=1}^R (\mathbf{A}_r \mathbf{B}_r^T) \otimes \mathbf{c}_r,$$

in which \mathbf{c}_r estimates the r th mixing vector $\mathbf{M}(:, r)$. The source signals can be estimated up to scaling and permutation by dehankelizing $\mathbf{A}_r \mathbf{B}_r^T$, e.g., by averaging over the antidiagonals. The example from [19] is slightly adapted in this experiment: two sources are mixed using

$$(4.4) \quad \mathbf{M} = \begin{bmatrix} 2 & 1 \\ -1 & 1 \end{bmatrix},$$

i.e., $K = R = 2$. The sources are given by

$$(4.5) \quad s_1(t) = 0.5 \sin(6\pi t),$$

$$(4.6) \quad s_2(t) = (4t^2 - 2.8t) \exp(-t);$$

hence $L_1 = 2$ and $L_2 = 3$. Estimating R and L_r is beyond the scope of this paper; see, e.g., [19]. The true R and L_r are therefore used. M equidistant samples are taken between 0 s and 1 s.

In the first experiment, the SNR is varied for a fixed number of $M = 501$ samples, and the relative error E on the mixing matrix is compared when using the full Hankelized tensor or the implicitly Hankelized tensor, i.e., in the structured format. Hankelization is performed using `hankelize`, which returns both the explicit and implicit tensorization. The resulting tensors of size $251 \times 251 \times 2$ are decomposed using `ll1_nls` starting from a random initialization. From the factor matrices, the signals are recovered using `dehankelize`, without constructing the full tensor. As shown in Figure 4.4, when increasing the SNR from 0 dB to 300 dB, the error decreases at the

same rate for the explicit and the implicit tensorization until the SNR is 180 dB: while the error continues to decrease for the explicit tensorization, the error stagnates at $E \approx 10^{-10}$ for the implicit tensorization. This loss of accuracy is explained in subsection 4.1 and only occurs for a very high SNR and in the case when the sources are exact sums of exponential polynomials. As illustrated in the following experiments, a trade-off between computational cost and accuracy can be made. In the case of a high SNR, the solution using the implicit tensorization can be used to initialize the algorithm with the full tensor.

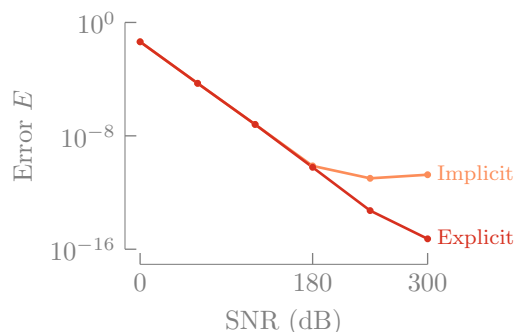


FIG. 4.4. For low and medium SNR, the errors on the mixing matrix for the explicit and implicit Hankelization approaches are equal. While the error for the implicit tensorization stagnates for SNR larger than 180 dB, the error continues to decrease when using the explicit tensorization. The errors are medians over 100 experiments, each using a best-out-of-five initializations strategy.

In the second experiment, the scaling behavior in the function of the number of samples is illustrated for $M = 501, 5001, 50001$, and 500001 points. The resulting Hankelized tensors require 953KB, 95.3MB, 9.53GB, and 953GB of memory, respectively, if formed explicitly. The SNR is fixed at 20 dB. We compute an unconstrained LL1 decomposition (`ll1_nls`) and a constrained BTM, which models the tensor as

$$(4.7) \quad \mathcal{T} = \sum_{r=1}^R \left(\mathbf{V}^{(r)} \mathbf{G}^{(r)} \mathbf{V}^{(r)\top} \right) \otimes \mathbf{c}_r,$$

in which $\mathbf{V}^{(r)}$ are confluent Vandermonde matrices and $\mathbf{G}^{(r)}$ are upper antitriangular matrices; see [19] for details. By imposing the structure in $\mathbf{V}^{(r)}$, the recovered signals are guaranteed to be exponential polynomials, and the underlying poles can be recovered easily. These parametric constraints are modeled in the SDF framework [83] using `struct_confvander` and `sdf_nls`. Figures 4.5 and 4.6 show the median time and error on the mixing matrix over 50 noise realizations. Each algorithm is initialized three and six times using random variables for the LL1 and BTM model, respectively. The time and error for the best initialization, i.e., the one resulting in the lowest error on the mixing matrix, is retained. The timing results in Figure 4.5 clearly show that implicit Hankelization drastically reduces the computation time, allowing longer signals to be analyzed. The computation of the constrained BTM is sensitive to the initialization and is difficult from an optimization point-of-view, as the factors are ill-conditioned due to the generalized Vandermonde structure. This results in higher computation times and lower accuracy for random initializations. Being able to analyze longer signal allows one to improve the accuracy for a fixed SNR, as is clear from Figure 4.6: by taking 100 times as many samples, E decreases by a factor 10.

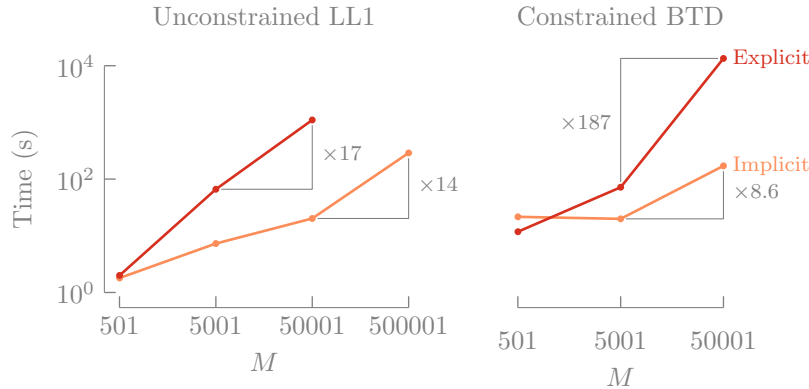


FIG. 4.5. For both the unconstrained LL1 decomposition and the constrained BTM, the computation time using implicit tensorization increases more slowly, enabling large-scale applications. For the LL1 decomposition, the time increases with a factor $14 \approx 13$ ($\mathcal{O}(M \log_2 M)$) using the efficient representation, and with a factor 17 when the full tensor is used. The latter factor is better than the expected factor 100 as fewer iterations are required. Computing the constrained BTM is more difficult due to ill-conditioned factors, and the variation in time is large causing larger deviations from the expected scaling factors. The results are medians over 50 experiments with multiple initializations.

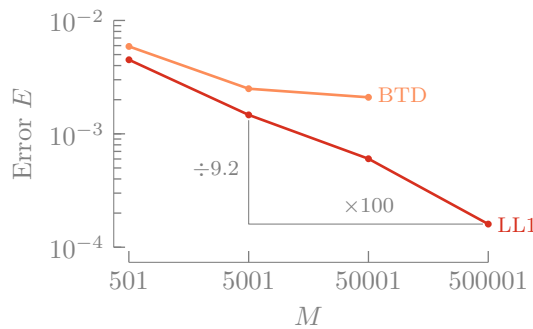


FIG. 4.6. The relative error E of the estimated mixing matrix decreases when the number of samples M increases. In the case of the unconstrained LL1 decomposition, E decreases with a factor $9.2 \approx 10$ when 100 times as many samples are used. The constrained BTM is more difficult from an optimization point-of-view and involves ill-conditioned factor matrices, resulting in a modest improvement in terms of error. The errors are computed using implicit tensorization and are medians over 50 noise realizations with multiple initializations. The results for full tensors are similar.

4.4. Fetal ECG extraction. If the underlying sources can be modeled by rational functions, BSS techniques based on Löwnerization can be used to recover the mixing matrix \mathbf{M} and the sources \mathbf{S} in (4.2). In this last experiment, the goal is to separate the heart beats of the mother from those of the fetus given five measured signals from five abdominal electrodes; see [20]. (The data is publicly available at [24].) As illustrated in [27], each signal can be modeled by a rational function, and the LL1 decomposition of the Löwnerized data can be used to recover the mixing matrix. Concretely, the signal matrix $\mathbf{X} \in \mathbb{R}^{2500 \times 5}$ is explicitly and implicitly mapped to a tensor $\mathcal{T} \in \mathbb{R}^{1250 \times 1250 \times 5}$ and an LL1 decomposition with $L_1 = L_2 = 56$, and $L_3 = 45$ is used to recover the 5×3 mixing matrix \mathbf{M} using the GN implementation in `ll1_nls`. (The values for L_r have been determined experimentally.) Note that we need more than

one source to model the mother's heartbeat, which is a three-dimensional signal [10]. The sources are recovered as $\mathbf{S} = \mathbf{X}\mathbf{M}^{\top\dagger}$ and are shown in Figure 4.7. The difference in computation time is relatively limited (81 s and 51 s for the explicit and implicit tensorization, respectively), as the multilinear ranks determined by L_r , $r = 1, 2, 3$, are quite high in this application. However, to store \mathcal{T} explicitly, 60 MB is required, while only 160 KB is needed to store the signal matrix and evaluation points in the efficient representation, allowing the signals to be processed on the limited resources of medical hardware.

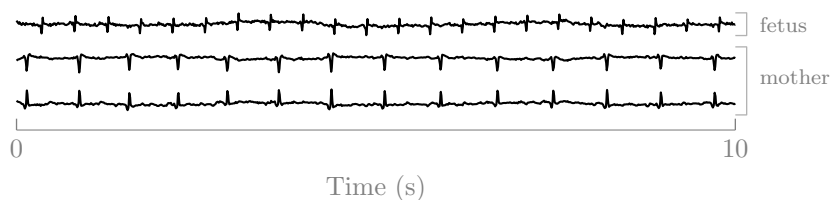


FIG. 4.7. The heart beats of the fetus and the mother can be separated using Löwnerization. An $LL1$ decomposition with $L_1 = L_2 = 56$, and $L_3 = 45$ is computed from the implicitly Löwnerized measured signals, after which the signals are reconstructed (shown) by pseudoinversion of mixing matrix.

5. Conclusion. By rewriting the objective function and gradient expressions commonly used to compute tensor decompositions, four core operations are exposed: squared norm, inner product, matricized tensor times Khatri–Rao product, and matricized tensor times Kronecker product. By specializing these operations for efficient representations of tensors, both the computational and memory complexity are reduced from \mathcal{O} (tensor entries) to \mathcal{O} (parameters) as illustrated for the polyadic, Tucker, and Tensor Train formats, as well as for Hankelization and Löwnerization. The operations can be used in many decomposition algorithms and frameworks, including ALS, second-order algorithms, and algorithms for constrained and/or coupled decompositions. The numerical consequences of exploiting efficient representations are studied in the case when a highly accurate solution is required. Finally, important concepts, such as Tucker or TT compression for constrained decompositions and implicit tensorization, allow large-scale datasets to be handled easily, as illustrated for nonnegative CPD and the Hankelization of mixtures of exponentials.

REFERENCES

- [1] E. ACAR, D. M. DUNLAVY, AND T. G. KOLDA, *A scalable optimization approach for fitting canonical tensor decompositions*, J. Chemometrics, 25 (2011), pp. 67–86, <https://doi.org/10.1002/cem.1335>.
- [2] E. ACAR, T. G. KOLDA, AND D. M. DUNLAVY, *All-at-once Optimization for Coupled Matrix and Tensor Factorizations*, preprint, <https://arxiv.org/abs/1105.3422>, 2011.
- [3] R. BADEAU AND R. BOYER, *Fast multilinear singular value decomposition for structured tensors*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 1008–1021, <https://doi.org/10.1137/060655936>.
- [4] B. BADER AND T. G. KOLDA, *Efficient MATLAB computations with sparse and factored tensors*, SIAM J. Sci. Comput., 30 (2007), pp. 205–231, <https://doi.org/10.1137/060676489>.
- [5] H. N. BHARATH, D. SIMA, N. SAUWEN, U. HIMMELREICH, L. DE LATHAUWER, AND S. VAN HUFEL, *Nonnegative canonical polyadic decomposition for tissue-type differentiation in gliomas*, IEEE J. Biomed. Health Inform., 21 (2017), pp. 1124–1132, <https://doi.org/10.1109/JBHI.2016.2583539>.

- [6] R. BRO, *Multi-way Analysis in the Food Industry: Models, Algorithms, and Applications*, Ph.D. thesis, University of Amsterdam, Amsterdam, The Netherlands, 1998.
- [7] R. BRO AND C. A. ANDERSSON, *Improving the speed of multiway algorithms: Part II: Compression*, Chemometr. Intell. Lab., 42 (1998), pp. 105–113, [https://doi.org/10.1016/S0169-7439\(98\)00011-2](https://doi.org/10.1016/S0169-7439(98)00011-2).
- [8] R. BRO, R. A. HARSHMAN, N. D. SIDIROPOULOS, AND M. E. LUNDY, *Modeling multi-way data with linearly dependent loadings*, J. Chemometrics, 23 (2009), pp. 324–340, <https://doi.org/10.1002/cem.1206>.
- [9] C. CAIAFA AND A. CICHOCKI, *Generalizing the column-row matrix decomposition to multi-way arrays*, Linear Algebra Appl., 433 (2010), pp. 557–573, <https://doi.org/10.1016/j.laa.2010.03.020>.
- [10] D. CALLAERTS, B. DE MOOR, J. VANDEWALLE, W. SANSSEN, G. VANTRAPPEN, AND J. JANSSENS, *Comparison of SVD methods to extract the foetal electrocardiogram from cutaneous electrode signals*, Med. Biol. Eng. Comput., 28 (1990), pp. 217–224, <https://doi.org/10.1007/BF02442670>.
- [11] J. D. CARROLL AND J.-J. CHANG, *Analysis of individual differences in multidimensional scaling via an n -way generalization of “Eckart–Young” decomposition*, Psychometrika, 35 (1970), pp. 283–319, <https://doi.org/10.1007/bf02310791>.
- [12] J. D. CARROLL, S. PRUZANSKY, AND J. B. KRUSKAL, *CANDELINC: A general approach to multidimensional analysis of many-way arrays with linear constraints on parameters*, Psychometrika, 45 (1980), pp. 3–24, <https://doi.org/10.1007/bf02293596>.
- [13] J. H. CHOI AND S. V. N. VISHWANATHAN, *DFacTo: Distributed factorization of tensors*, in Proceedings of the 27th International Conference on Neural Information Processing Systems, NIPS’14, Cambridge, MA, USA, 2014, MIT, Cambridge, MA, pp. 1296–1304.
- [14] A. CICHOCKI, D. MANDIC, A.-H. PHAN, C. CAIAFA, G. ZHOU, Q. ZHAO, AND L. DE LATHAUWER, *Tensor decompositions for signal processing applications: From two-way to multiway component analysis*, IEEE Signal Process. Mag., 32 (2015), pp. 145–163, <https://doi.org/10.1109/msp.2013.2297439>.
- [15] J. E. COHEN, R. CABRAL FARIAS, AND P. COMON, *Fast decomposition of large nonnegative tensors*, IEEE Signal Process. Lett., 22 (2015), pp. 862–866, <https://doi.org/10.1109/lsp.2014.2374838>.
- [16] P. COMON AND C. JUTTEN, *Handbook of Blind Source Separation: Independent Component Analysis and Applications*, Academic Press, Oxford, 2010.
- [17] A. DE ALMEIDA, G. FAVIER, AND J. C. M. MOTA, *A constrained factor decomposition with application to MIMO antenna systems*, IEEE Trans. Signal Process., 56 (2008), pp. 2429–2442, <https://doi.org/10.1109/TSP.2008.917026>.
- [18] L. DE LATHAUWER, *Decompositions of a higher-order tensor in block terms — Part II: Definitions and uniqueness*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 1033–1066, <https://doi.org/10.1137/070690729>.
- [19] L. DE LATHAUWER, *Blind separation of exponential polynomials and the decomposition of a tensor in rank- $(L_r, L_r, 1)$ terms*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1451–1474, <https://doi.org/10.1137/100805510>.
- [20] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *Fetal electrocardiogram extraction by blind source subspace separation*, IEEE Trans. Biomed. Eng., 47 (2000), pp. 567–572, <https://doi.org/10.1109/10.841326>.
- [21] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *A multilinear singular value decomposition*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1253–1278, <https://doi.org/10.1137/S0895479896305696>.
- [22] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *On the best rank-1 and rank- (R_1, R_2, \dots, R_N) approximation of higher-order tensors*, SIAM J. Matrix Anal. Appl., 21 (2000), pp. 1324–1342, <https://doi.org/10.1137/S0895479898346995>.
- [23] L. DE LATHAUWER AND D. NION, *Decompositions of a higher-order tensor in block terms — Part III: Alternating least squares algorithms*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 1067–1083, <https://doi.org/10.1137/070690730>.
- [24] B. DE MOOR, *DaISy: Database for the identification of systems*, Department of Electrical Engineering, ESAT/SISTA, KU Leuven, Belgium, <http://www.esat.kuleuven.ac.be/sista/daisy/>, July, 1. 2018. (Used dataset: cutaneous potential recordings of a pregnant woman, biomedical systems, 96-012).
- [25] O. DEBALS AND L. DE LATHAUWER, *The Concept of Tensorization*, Technical Report 17–99, ESAT-STADIUS, KU Leuven, Belgium, 2017, <ftp://ftp.esat.kuleuven.be/pub/stadius/odebals/debals2017concept.pdf>.

- [26] O. DEBALS, L. DE LATHAUWER, AND M. VAN BAREL, *About Higher-Order Löwner Tensors*, Technical Report 17–98, ESAT-STADIUS, KU Leuven, Leuven, Belgium, 2017, <ftp://ftp.esat.kuleuven.be/pub/stadius/odebals/debals2017higher.pdf>.
- [27] O. DEBALS, M. VAN BAREL, AND L. DE LATHAUWER, *Löwner-based blind signal separation of rational functions with applications*, IEEE Trans. Signal Process., 64 (2016), pp. 1909–1918, <https://doi.org/10.1109/tsp.2015.2500179>.
- [28] W. DING, L. QI, AND Y. WEI, *Fast Hankel tensor-vector product and its application to exponential data fitting*, Numer. Linear Algebra Appl., 22 (2015), pp. 814–832, <https://doi.org/10.1002/nla.1970>.
- [29] I. GOHBERG AND V. OLSHEVSKY, *Complexity of multiplication with vectors for structured matrices*, Linear Algebra Appl., 202 (1994), pp. 163–192, [https://doi.org/10.1016/0024-3795\(94\)90189-9](https://doi.org/10.1016/0024-3795(94)90189-9).
- [30] I. GOHBERG AND V. OLSHEVSKY, *Fast algorithms with preprocessing for matrix-vector multiplication problems*, J. Complexity, 10 (1994), pp. 411–427, <https://doi.org/10.1006/jcom.1994.1021>.
- [31] L. GRASEDYCK, *Hierarchical singular value decomposition of tensors*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 2029–2054, <https://doi.org/10.1137/090764189>.
- [32] L. GRASEDYCK, D. KRESSNER, AND C. TOBLER, *A literature survey of low-rank tensor approximation techniques*, GAMM-Mitt., 36 (2013), pp. 53–78, <https://doi.org/10.1002/gamm.201310004>.
- [33] W. HACKBUSCH, *Tensor Spaces and Numerical Tensor Calculus*, Springer Ser. Comput. Math. 42, Springer, Heidelberg, 2012, <https://doi.org/10.1007/978-3-642-28027-6>.
- [34] W. HACKBUSCH, B. N. KHOROMSKIJ, AND E. E. TYRTYSHNIKOV, *Hierarchical Kronecker tensor-product approximations*, J. Numer. Math., 13 (2005), pp. 119–156, <https://doi.org/10.1515/1569395054012767>.
- [35] W. HACKBUSCH AND S. KÜHN, *A new scheme for the tensor representation*, J. Fourier Anal. Appl., 15 (2009), pp. 706–722, <https://doi.org/10.1007/s00041-009-9094-9>.
- [36] F. L. HITCHCOCK, *The expression of a tensor or a polyadic as a sum of products*, Journal of Mathematics and Physics, 6 (1927), pp. 164–189, <https://doi.org/10.1002/sapm192761164>.
- [37] S. HOLTZ, T. ROHWEDDER, AND R. SCHNEIDER, *The alternating linear scheme for tensor optimization in the Tensor Train format*, SIAM J. Sci. Comput., 34 (2012), pp. A683–A713, <https://doi.org/10.1137/100818893>.
- [38] K. HUANG, N. D. SIDIROPOULOS, AND A. P. LIAVAS, *A flexible and efficient algorithmic framework for constrained matrix and tensor factorization*, IEEE Trans. Signal Process., 64 (2016), pp. 5052–5065, <https://doi.org/10.1109/TSP.2016.2576427>.
- [39] M. ISHTEVA, P.-A. ABSIL, S. VAN HUFFEL, AND L. DE LATHAUWER, *Best low multilinear rank approximation of higher-order tensors, based on the Riemannian trust-region scheme*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 115–135, <https://doi.org/10.1137/090764827>.
- [40] M. ISHTEVA, L. DE LATHAUWER, P.-A. ABSIL, AND S. VAN HUFFEL, *The best rank- (R_1, R_2, R_3) approximation of tensors by means of a geometric Newton method*, AIP Conference Proceedings, 1048 (2008), pp. 274–277, <https://doi.org/10.1063/1.2990911>.
- [41] B. JEON, I. JEON, L. SAEL, AND U. KANG, *SCouT: Scalable coupled matrix-tensor factorization — algorithm and discoveries*, in 2016 IEEE 32nd International Conference on Data Engineering (ICDE), IEEE, 2016, pp. 811–822, <https://doi.org/10.1109/icde.2016.7498292>.
- [42] I. JEON, E. E. PAPAEXAKIS, U. KANG, AND C. FALOUTSOS, *HaTen2: Billion-scale tensor decompositions*, in 2015 IEEE 31st International Conference on Data Engineering, IEEE, 2015, pp. 1047–1058, <https://doi.org/10.1109/icde.2015.7113355>.
- [43] U. KANG, E. E. PAPAEXAKIS, A. HARPALE, AND C. FALOUTSOS, *GigaTensor: Scaling tensor analysis up by 100 times - algorithms and discoveries*, Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '12, ACM, New York, 2012, pp. 316–324, <https://doi.org/10.1145/2339530.2339583>.
- [44] 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 on - SC '15, ACM, New York, 2015, 77, <https://doi.org/10.1145/2807591.2807624>.
- [45] O. KAYA AND B. UÇAR, *High performance parallel algorithms for the Tucker decomposition of sparse tensors*, in 2016 45th International Conference on Parallel Processing (ICPP), IEEE, 2016, pp. 103–112, <https://doi.org/10.1109/icpp.2016.19>.
- [46] C. KELLEY, *Iterative Methods for Optimization*, SIAM, Philadelphia, 1999, <https://doi.org/10.1137/1.9781611970920>.

- [47] B. N. KHOROMSKIJ, *Structured rank- (r_1, \dots, r_D) decomposition of function-related tensors in \mathbb{R}^D* , Comput. Methods Appl. Math., 6 (2006), pp. 194–220, <https://doi.org/10.2478/cmam-2006-0010>.
- [48] B. N. KHOROMSKIJ AND V. KHOROMSKAIA, *Low rank Tucker-type tensor approximation to classical potentials*, Cent. Eur. J. Math., 5 (2007), pp. 523–550, <https://doi.org/10.2478/s11533-007-0018-0>.
- [49] B. N. KHOROMSKIJ AND V. KHOROMSKAIA, *Multigrid accelerated tensor approximation of function related multidimensional arrays*, SIAM J. Sci. Comput., 31 (2009), pp. 3002–3026, <https://doi.org/10.1137/080730408>.
- [50] H. A. L. KIERS AND R. A. HARSHMAN, *Relating two proposed methods for speedup of algorithms for fitting two- and three-way principal component and related multilinear models*, Chemometr. Intell. Lab., 36 (1997), pp. 31–40, [https://doi.org/10.1016/s0169-7439\(96\)00074-3](https://doi.org/10.1016/s0169-7439(96)00074-3).
- [51] T. G. KOLDA, *Orthogonal tensor decompositions*, SIAM J. Matrix Anal. Appl., 23 (2001), pp. 243–255, <https://doi.org/10.1137/S0895479800368354>.
- [52] T. G. KOLDA AND B. BADER, *Tensor Decompositions and Applications*, SIAM Rev., 51 (2009), pp. 455–500, <https://doi.org/10.1137/07070111X>.
- [53] T. G. KOLDA AND J. SUN, *Scalable tensor decompositions for multi-aspect data mining*, 2008 Eighth IEEE International Conference on Data Mining, 2008, <https://doi.org/10.1109/icdm.2008.89>.
- [54] P. M. KROONENBERG AND J. DE LEEUW, *Principal component analysis of three-mode data by means of alternating least squares algorithms*, Psychometrika, 45 (1980), pp. 69–97, <https://doi.org/10.1007/bf02293599>.
- [55] A. P. LIAVAS AND N. D. SIDIROPOULOS, *Parallel algorithms for constrained tensor factorization via alternating direction method of multipliers*, IEEE Trans. Signal Process., 63 (2015), pp. 5450–5463, <https://doi.org/10.1109/tsp.2015.2454476>.
- [56] M. W. MAHONEY, M. MAGGIONI, AND P. DRINEAS, *Tensor-CUR decompositions for tensor-based data*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 957–987, <https://doi.org/10.1137/060665336>.
- [57] T. MÜLLER, K. KRUPPA, G. LICHTENBERG, AND N. RÉHAULT, *Fault detection with qualitative models reduced by tensor decomposition methods*, IFAC-PapersOnLine, 48 (2015), pp. 416–421, <https://doi.org/10.1016/j.ifacol.2015.09.562>.
- [58] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, 2nd ed., Springer, New York, 2006.
- [59] R. ORÚS, *A practical introduction to tensor networks: Matrix product states and projected entangled pair states*, Ann. Physics, 349 (2014), pp. 117–158, <https://doi.org/10.1016/j.aop.2014.06.013>.
- [60] I. V. OSELEDETS, *Tensor-train decomposition*, SIAM J. Sci. Comput., 33 (2011), pp. 2295–2317, <https://doi.org/10.1137/090752286>.
- [61] I. V. OSELEDETS AND S. DOLGOV, *Solution of linear systems and matrix inversion in the TT-format*, SIAM J. Sci. Comput., 34 (2012), pp. A2718–A2739, <https://doi.org/10.1137/110833142>.
- [62] I. V. OSELEDETS, S. DOLGOV, V. KAZEYEV, D. V. SAVOSTYANOV, O. LEBEDEVA, P. ZHLOBICH, T. MACH, AND L. SONG, *TT-Toolbox*, available online at <https://github.com/oseledets/TT-Toolbox>.
- [63] I. V. OSELEDETS, D. V. SAVOSTYANOV, AND E. E. TYRTYSHNIKOV, *Tucker dimensionality reduction of three-dimensional arrays in linear time*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 939–956, <https://doi.org/10.1137/060655894>.
- [64] I. V. OSELEDETS, D. V. SAVOSTYANOV, AND E. E. TYRTYSHNIKOV, *Linear algebra for tensor problems*, Computing, 85 (2009), pp. 169–188, <https://doi.org/10.1007/s00607-009-0047-6>.
- [65] I. V. OSELEDETS AND E. E. TYRTYSHNIKOV, *Breaking the curse of dimensionality, or how to use SVD in many dimensions*, SIAM J. Sci. Comput., 31 (2009), pp. 3744–3759, <https://doi.org/10.1137/090748330>.
- [66] I. V. OSELEDETS AND E. E. TYRTYSHNIKOV, *TT-cross approximation for multidimensional arrays*, Linear Algebra Appl., 432 (2010), pp. 70–88, <https://doi.org/10.1016/j.laa.2009.07.024>.
- [67] E. E. PAPALEXAKIS, C. FALOUTSOS, AND N. D. SIDIROPOULOS, *ParCube: Sparse parallelizable CANDECOMP-PARAFAC tensor decomposition*, ACM Trans. Knowl. Discov. Data, 10 (2015), pp. 1–25, <https://doi.org/10.1145/2729980>.
- [68] J. M. PAPPY, L. DE LATHAUWER, AND S. VAN HUFFEL, *Exponential data fitting using multilinear algebra: The single-channel and multi-channel case*, Numer. Linear Algebra Appl., 12 (2005), pp. 809–826, <https://doi.org/10.1002/nla.453>.

- [69] A.-H. PHAN, P. TICHAVSKÝ, AND A. CICHOCKI, *Fast alternating LS algorithms for high order CANDECOMP/PARAFAC tensor factorizations*, IEEE Trans. Signal Process., 61 (2013), pp. 4834–4846, <https://doi.org/10.1109/tsp.2013.2269903>.
- [70] M. J. REYNOLDS, G. BEYLKIN, AND A. DOOSTAN, *Optimization via separated representations and the canonical tensor decomposition*, J. Comput. Phys., 348 (2017), pp. 220–230, <https://doi.org/10.1016/j.jcp.2017.07.012>.
- [71] J.-P. ROYER, N. THIRION-MOREAU, AND P. COMON, *Computing the polyadic decomposition of nonnegative third order tensors*, Signal Processing, 91 (2011), pp. 2159–2171, <https://doi.org/10.1016/j.sigpro.2011.03.006>.
- [72] D. V. SAVOSTYANOV, *Fast revealing of mode ranks of tensor in canonical form*, Numer. Math. Theor. Meth. Appl., 2 (2009), pp. 439–444, <https://doi.org/10.4208/nmtma.2009.m9006s>.
- [73] D. V. SAVOSTYANOV, E. E. TYRTYSHNIKOV, AND N. L. ZAMARASHKIN, *Fast truncation of mode ranks for bilinear tensor operations*, Numer. Linear Algebra Appl., 19 (2011), pp. 103–111, <https://doi.org/10.1002/nla.765>.
- [74] N. D. SIDIROPOULOS, R. BRO, AND G. B. GIANNAKIS, *Parallel factor analysis in sensor array processing*, IEEE Trans. Signal Process., 48 (2000), pp. 2377–2388, <https://doi.org/10.1109/78.852018>.
- [75] N. D. SIDIROPOULOS, L. DE LATHAUWER, X. FU, K. HUANG, E. E. PAPALEXAKIS, AND C. FALOUTSOS, *Tensor decomposition for signal processing and machine learning*, IEEE Trans. Signal Process., 65 (2017), pp. 3551–3582, <https://doi.org/10.1109/TSP.2017.2690524>.
- [76] N. D. SIDIROPOULOS, E. E. PAPALEXAKIS, AND C. FALOUTSOS, *Parallel randomly compressed cubes: A scalable distributed architecture for big tensor decomposition*, IEEE Signal Process. Mag., 31 (2014), pp. 57–70, <https://doi.org/10.1109/MSP.2014.2329196>.
- [77] S. SMITH AND G. KARYPIS, *Tensor-matrix products with a compressed sparse tensor*, in Proceedings of the 5th Workshop on Irregular Applications Architectures and Algorithms - IA3 '15, ACM, New York, 2015, 5, <https://doi.org/10.1145/2833179.2833183>.
- [78] S. SMITH AND G. KARYPIS, *A medium-grained algorithm for sparse tensor factorization*, in 2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS), IEEE, 2016, pp. 902–911, <https://doi.org/10.1109/ipdps.2016.113>.
- [79] S. SMITH, N. RAVINDRAN, N. D. SIDIROPOULOS, AND G. KARYPIS, *SPLATT: Efficient and parallel sparse tensor-matrix multiplication*, in 2015 IEEE International Parallel and Distributed Processing Symposium, IEEE, 2015, pp. 61–70, <https://doi.org/10.1109/ipdps.2015.27>.
- [80] L. SORBER, I. DOMANOV, M. VAN BAREL, AND L. DE LATHAUWER, *Exact line and plane search for tensor optimization*, Comput. Optim. Appl., 63 (2015), pp. 121–142, <https://doi.org/10.1007/s10589-015-9761-5>.
- [81] L. SORBER, M. VAN BAREL, AND L. DE LATHAUWER, *Unconstrained optimization of real functions in complex variables*, SIAM J. Optim., 22 (2012), pp. 879–898, <https://doi.org/10.1137/110832124>.
- [82] L. SORBER, M. VAN BAREL, AND L. DE LATHAUWER, *Optimization-based algorithms for tensor decompositions: Canonical polyadic decomposition, decomposition in rank- $(L_r, L_r, 1)$ terms, and a new generalization*, SIAM J. Optim., 23 (2013), pp. 695–720, <https://doi.org/10.1137/120868323>.
- [83] L. SORBER, M. VAN BAREL, AND L. DE LATHAUWER, *Structured data fusion*, IEEE J. Sel. Topics Signal Process., 9 (2015), pp. 586–600, <https://doi.org/10.1109/JSTSP.2015.2400415>.
- [84] G. TOMASI AND R. BRO, *A comparison of algorithms for fitting the PARAFAC model*, Comput. Stat. Data Anal., 50 (2006), pp. 1700 – 1734, <https://doi.org/10.1016/j.csda.2004.11.013>.
- [85] M. VANDECAPPELLE, N. VERVLIT, AND L. DE LATHAUWER, *Nonlinear least squares updating of the canonical polyadic decomposition*, in 2017 25th European Signal Processing Conference (EUSIPCO17), 2017, pp. 693–697, <https://doi.org/10.23919/EUSIPCO.2017.8081290>.
- [86] N. VANNIEUWENHOVEN, *Condition numbers for the tensor rank decomposition*, Linear Algebra Appl., 535 (2017), pp. 35–86, <https://doi.org/10.1016/j.laa.2017.08.014>.
- [87] N. VANNIEUWENHOVEN, K. MEERBERGEN, AND R. VANDEBRIL, *Computing the gradient in optimization algorithms for the CP decomposition in constant memory through tensor blocking*, SIAM J. Sci. Comput., 37 (2015), pp. C415–C438, <https://doi.org/10.1137/14097968x>.
- [88] N. VANNIEUWENHOVEN, R. VANDEBRIL, AND K. MEERBERGEN, *A new truncation strategy for the higher-order singular value decomposition*, SIAM J. Sci. Comput., 34 (2012), pp. A1027–A1052, <https://doi.org/10.1137/110836067>.
- [89] N. VERVLIT AND L. DE LATHAUWER, *A randomized block sampling approach to canonical polyadic decomposition of large-scale tensors*, IEEE J. Sel. Topics Signal Process., 10 (2016), pp. 284–295, <https://doi.org/10.1109/JSTSP.2015.2503260>.

- [90] N. VERVLIET AND L. DE LATHAUWER, *Numerical Optimization Based Algorithms for Data Fusion*, in Data Fusion Methodology and Applications, M. Cocchi, ed., Elsevier, 2018, accepted for publication.
- [91] N. VERVLIET, O. DEBALS, AND L. DE LATHAUWER, *Tensorlab 3.0—Numerical optimization strategies for large-scale constrained and coupled matrix/tensor factorization*, in 2016 50th Asilomar Conference on Signals, Systems and Computers, 2016, pp. 1733–1738, <https://doi.org/10.1109/ACSSC.2016.7869679>.
- [92] N. VERVLIET, O. DEBALS, AND L. DE LATHAUWER, *Canonical Polyadic Decomposition of Incomplete Tensors with Linearly Constrained Factors*, Technical Report 16–172, ESAT-STADIUS, KU Leuven, Belgium, 2017, <ftp://ftp.esat.kuleuven.be/pub/stadius/nvervliet/vervliet2017cpdli.pdf>.
- [93] N. VERVLIET, O. DEBALS, L. SORBER, AND L. DE LATHAUWER, *Breaking the curse of dimensionality using decompositions of incomplete tensors: Tensor-based scientific computing in big data analysis*, IEEE Signal Process. Mag., 31 (2014), pp. 71–79, <https://doi.org/10.1109/MSP.2014.2329429>.
- [94] N. VERVLIET, O. DEBALS, L. SORBER, M. VAN BAREL, AND L. DE LATHAUWER, *Tensorlab 3.0*, 2016, available online at <https://www.tensorlab.net>.
- [95] Z. YAWEN, D. GUANGJUN, AND X. ZHIXIANG, *Hyperspectral image tensor feature extraction based on fusion of multiple spectral-spatial features*, in Proceedings of the 2016 International Conference on Intelligent Information Processing - ICIIP '16, ACM, New York, 2016, 43, <https://doi.org/10.1145/3028842.3028885>.
- [96] G. ZHOU, A. CICHOCKI, AND S. XIE, *Decomposition of Big Tensors with Low Multilinear Rank*, preprint, 2014, <https://arxiv.org/abs/1412.1885>.
- [97] G. ZHOU, A. CICHOCKI, Q. ZHAO, AND S. XIE, *Nonnegative matrix and tensor factorizations: An algorithmic perspective*, IEEE Signal Process. Mag., 31 (2014), pp. 54–65, <https://doi.org/10.1109/msp.2014.2298891>.
- [98] G. ZHOU, A. CICHOCKI, Q. ZHAO, AND S. XIE, *Efficient nonnegative Tucker decompositions: Algorithms and uniqueness*, IEEE Trans. Image Process., 24 (2015), pp. 4990–5003, <https://doi.org/10.1109/tip.2015.2478396>.