

# Tensor Decompositions for Signal Processing Applications

## From Two-way to Multiway Component Analysis

A. Cichocki, D. Mandic, A-H. Phan, C. Caiafa, G. Zhou, Q. Zhao, and  
L. De Lathauwer

### Summary

The widespread use of multi-sensor technology and the emergence of big datasets has highlighted the limitations of standard flat-view matrix models and the necessity to move towards more versatile data analysis tools. We show that higher-order tensors (i.e., multiway arrays) enable such a fundamental paradigm shift towards models that are essentially polynomial and whose uniqueness, unlike the matrix methods, is guaranteed under very mild and natural conditions. Benefiting from the power of multilinear algebra as their mathematical backbone, data analysis techniques using tensor decompositions are shown to have great flexibility in the choice of constraints that match data properties, and to find more general latent components in the data than matrix-based methods. A comprehensive introduction to tensor decompositions is provided from a signal processing perspective, starting from the algebraic foundations, via basic Canonical Polyadic and Tucker models, through to advanced cause-effect and multi-view data analysis schemes. We show that tensor decompositions enable natural generalizations of some commonly used signal processing paradigms, such as canonical correlation and subspace techniques, signal separation, linear regression, feature extraction and classification. We also cover computational aspects, and point out how ideas from compressed sensing and scientific computing may be used for addressing the otherwise unmanageable storage and manipulation problems associated with big datasets. The concepts are supported by illustrative real world case studies illuminating the benefits of the tensor framework, as efficient and promising tools for modern signal processing, data analysis and machine learning applications; these benefits also extend to vector/matrix data through tensorization.

### INTRODUCTION

**Historical notes.** The roots of multiway analysis can be traced back to studies of homogeneous polynomials in the 19th century, contributors include Gauss, Kronecker, Cayley, Weyl and Hilbert — in modern day interpretation these are fully symmetric tensors. Decompositions of non-symmetric tensors have been studied since the early 20th century [1], whereas the benefits of using more than two matrices in factor analysis [2] became apparent in several communities since the 1960s. The Tucker decomposition for tensors was introduced in psychometrics [3], [4], while the Canonical Polyadic Decomposition (CPD) was independently rediscovered and put into an application context under the names of Canonical Decomposition (CANDECOMP) in psychometrics [5] and Parallel Factor Model (PARAFAC) in linguistics [6]. Tensors were subsequently adopted in diverse branches of data analysis such as chemometrics, food industry and social sciences [7], [8]. When it comes to Signal Processing, the early 1990s saw a considerable interest in Higher-Order Statistics (HOS) [9] and it was soon realized that for the multivariate case HOS are effectively higher-order tensors; indeed, algebraic approaches to Independent Component Analysis (ICA) using HOS [10]–[12] were inherently tensor-based. Around 2000, it was realized that the Tucker decomposition represents a MultiLinear Singular Value Decomposition (MLSVD) [13]. Generalizing the matrix SVD, the workhorse of numerical linear algebra, the MLSVD spurred the interest in tensors in applied mathematics and scientific computing in very high dimensions [14]–[16]. In parallel, CPD was successfully adopted as a tool for sensor array processing and deterministic signal separation in wireless communication [17], [18]. Subsequently, tensors have been used in audio, image and video processing, machine

learning and biomedical applications, to name but a few. The significant interest in tensors and their fast emerging applications are reflected in books [7], [8], [12], [19]–[21] and tutorial papers [22]–[29] covering various aspects of multiway analysis.

**From a matrix to a tensor.** Approaches to two-way (matrix) component analysis are well established, and include Principal Component Analysis (PCA), Independent Component Analysis (ICA), Nonnegative Matrix Factorization (NMF) and Sparse Component Analysis (SCA) [12], [19], [30]. These techniques have become standard tools for e.g., blind source separation (BSS), feature extraction, or classification. On the other hand, large classes of data arising from modern heterogeneous sensor modalities have a multiway character and are therefore naturally represented by multiway arrays or tensors (see Section **Tensorization**).

Early multiway data analysis approaches reformatted the data tensor as a matrix and resorted to methods developed for classical two-way analysis. However, such a “flattened” view of the world and the rigid assumptions inherent in two-way analysis are not always a good match for multiway data. It is only through higher-order tensor decomposition that we have the opportunity to develop sophisticated models capturing multiple interactions and couplings, instead of standard pairwise interactions. In other words, we can only discover hidden components within multiway data if the analysis tools account for intrinsic multi-dimensional patterns present — motivating the development of multilinear techniques.

In this article, we emphasize that tensor decompositions are not just matrix factorizations with additional subscripts — multilinear algebra is much structurally richer than linear algebra. For example, even basic notions such as rank have a more subtle meaning, uniqueness conditions of higher-order tensor decompositions are more relaxed and accommodating than those for matrices [31], [32], while matrices and tensors also have completely different geometric properties [20]. This boils down to matrices representing linear transformations and quadratic forms, while tensors are connected with multilinear mappings and multivariate polynomials [29].

## NOTATIONS AND CONVENTIONS

A tensor can be thought of as a multi-index numerical array, whereby the order of a tensor is the number of its “modes” or “dimensions”, these may include space, time, frequency, trials, classes, and dictionaries. A real-valued tensor of order  $N$  is denoted by  $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  and its entries by  $a_{i_1, i_2, \dots, i_N}$ . Then, an  $N \times 1$  vector  $\mathbf{a}$  is considered a tensor of order one, and an  $N \times M$  matrix  $\mathbf{A}$  a tensor of order two. Subtensors are parts of the original data tensor, created when only a fixed subset of indices is used. Vector-valued subtensors are called *fibers*, defined by fixing every index but one, and matrix-valued subtensors are called *slices*, obtained by fixing all but two indices (see Table I). Manipulation of tensors often requires their reformatting (*reshaping*); a particular case of reshaping tensors to matrices is termed matrix unfolding or matricization (see Figure 4 (left)). Note that a mode- $n$  multiplication of a tensor  $\mathcal{A}$  with a matrix  $\mathbf{B}$  amounts to the multiplication of all mode- $n$  vector fibers with  $\mathbf{B}$ , and that in linear algebra the tensor (or outer) product appears in the expression for a rank-1 matrix:  $\mathbf{a}\mathbf{b}^T = \mathbf{a} \circ \mathbf{b}$ . Basic tensor notations are summarized in Table I, while Table II outlines several types of products used in this paper.

### INTERPRETABLE COMPONENTS IN TWO-WAY DATA ANALYSIS

The aim of blind source separation (BSS), factor analysis (FA) and latent variable analysis (LVA) is to decompose a data matrix  $\mathbf{X} \in \mathbb{R}^{I \times J}$  into the factor matrices  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_R] \in \mathbb{R}^{I \times R}$  and  $\mathbf{B} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_R] \in \mathbb{R}^{J \times R}$  as:

$$\begin{aligned} \mathbf{X} &= \mathbf{A} \mathbf{D} \mathbf{B}^T + \mathbf{E} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \mathbf{b}_r^T + \mathbf{E} \\ &= \sum_{r=1}^R \lambda_r \mathbf{a}_r \circ \mathbf{b}_r + \mathbf{E}, \end{aligned} \quad (1)$$

where  $\mathbf{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_R)$  is a scaling (normalizing) matrix, the columns of  $\mathbf{B}$  represent the unknown source signals (factors or latent variables depending on the tasks in hand), the columns of  $\mathbf{A}$  represent the associated mixing vectors (or factor loadings), while  $\mathbf{E}$  is noise due to an unmodelled data part or model error. In other words, model (1) assumes that the data matrix  $\mathbf{X}$  comprises hidden components  $\mathbf{b}_r$  ( $r = 1, 2, \dots, R$ ) that are mixed together in an unknown manner through coefficients  $\mathbf{A}$ , or, equivalently, that data contain *factors* that have

TABLE I: Basic notation.

$\mathcal{A}, \mathbf{A}, \mathbf{a}, a$	tensor, matrix, vector, scalar
$\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_R]$	matrix $\mathbf{A}$ with column vectors $\mathbf{a}_r$
$\mathbf{a}(:, i_2, i_3, \dots, i_N)$	fiber of tensor $\mathcal{A}$ obtained by fixing all but one index
$\mathbf{A}(:, :, i_3, \dots, i_N)$	matrix slice of tensor $\mathcal{A}$ obtained by fixing all but two indices
$\mathcal{A}(:, :, :, i_4, \dots, i_N)$	tensor slice of $\mathcal{A}$ obtained by fixing some indices
$\mathcal{A}(\mathcal{I}_1, \mathcal{I}_2, \dots, \mathcal{I}_N)$	subtensor of $\mathcal{A}$ obtained by restricting indices to belong to subsets $\mathcal{I}_n \subseteq \{1, 2, \dots, I_n\}$
$\mathbf{A}_{(n)} \in \mathbb{R}^{I_n \times I_1 I_2 \dots I_{n-1} I_{n+1} \dots I_N}$	mode- $n$ matricization of tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ whose entry at row $i_n$ and column $(i_1 - 1)I_2 \dots I_{n-1} I_{n+1} \dots I_N + \dots + (i_{N-1} - 1)I_N + i_N$ is equal to $a_{i_1 i_2 \dots i_N}$
$\text{vec}(\mathcal{A}) \in \mathbb{R}^{I_N I_{N-1} \dots I_1}$	vectorization of tensor $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ with the entry at position $i_1 + \sum_{k=2}^N [(i_k - 1)I_1 I_2 \dots I_{k-1}]$ equal to $a_{i_1 i_2 \dots i_N}$
$\mathbf{D} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_R)$	diagonal matrix with $d_{rr} = \lambda_r$
$\mathcal{D} = \text{diag}_N(\lambda_1, \lambda_2, \dots, \lambda_R)$	diagonal tensor of order $N$ with $d_{rr \dots r} = \lambda_r$
$\mathbf{A}^T, \mathbf{A}^{-1}, \mathbf{A}^\dagger$	transpose, inverse, and Moore-Penrose pseudo-inverse

an associated *loading* for every data channel. Figure 2 (top) depicts the model (1) as a dyadic decomposition, whereby the terms  $\mathbf{a}_r \circ \mathbf{b}_r = \mathbf{a}_r \mathbf{b}_r^T$  are rank-1 matrices.

The well-known indeterminacies intrinsic to this model are: (i) arbitrary scaling of components, and (ii) permutation of the rank-1 terms. Another indeterminacy is related to the physical meaning of the factors: if the model in (1) is unconstrained, it admits infinitely many combinations of  $\mathbf{A}$  and  $\mathbf{B}$ . Standard matrix factorizations in linear algebra, such as the QR-factorization, Eigenvalue Decomposition (EVD), and Singular Value Decomposition (SVD), are only special cases of (1), and owe their uniqueness to hard and restrictive constraints such as triangularity and orthogonality. On the other hand, certain properties of the factors in (1) can be represented by appropriate constraints, making possible unique estimation or extraction of such factors. These constraints include statistical independence, sparsity, nonnegativity, exponential structure, uncorrelatedness, constant modulus, finite alphabet, smoothness and unimodality. Indeed, the first four properties form the basis of Independent Component Analysis (ICA) [12], [33], [34], Sparse Component Analysis (SCA)

[30], Nonnegative Matrix Factorization (NMF) [19], and harmonic retrieval [35].

#### TENSORIZATION — BLESSING OF DIMENSIONALITY

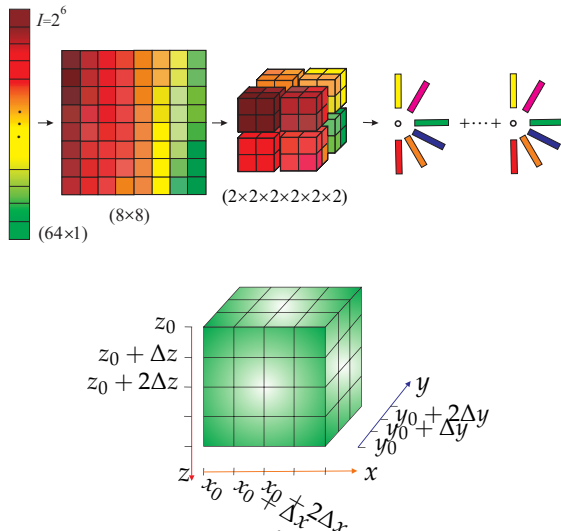
While one-way (vectors) and two-way (matrices) algebraic structures were respectively introduced as natural representations for segments of scalar measurements and measurements on a grid, tensors were initially used purely for the mathematical benefits they provide in data analysis; for instance, it seemed natural to stack together excitation-emission spectroscopy matrices in chemometrics into a third-order tensor [7].

The procedure of creating a data tensor from lower-dimensional original data is referred to as *tensorization*, and we propose the following taxonomy for tensor generation:

- 1) *Rearrangement of lower dimensional data structures.* Large-scale vectors or matrices are readily tensorized to higher-order tensors, and can be compressed through tensor decompositions if they admit a low-rank tensor approximation; this principle facilitates big data analysis [21], [27], [28]

TABLE II: Definition of products.

$\mathcal{C} = \mathcal{A} \times_n \mathbf{B}$	mode- $n$ product of $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ and $\mathbf{B} \in \mathbb{R}^{J_n \times I_n}$ yields $\mathcal{C} \in \mathbb{R}^{I_1 \times \dots \times I_{n-1} \times J_n \times I_{n+1} \times \dots \times I_N}$ with entries $c_{i_1 \dots i_{n-1} j_n i_{n+1} \dots i_N} = \sum_{i_n=1}^{I_n} a_{i_1 \dots i_{n-1} i_n i_{n+1} \dots i_N} b_{j_n i_n}$ and matrix representation $\mathbf{C}_{(n)} = \mathbf{B} \mathbf{A}_{(n)}$
$\mathcal{C} = [\mathcal{A}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \dots, \mathbf{B}^{(N)}]$	full multilinear product, $\mathcal{C} = \mathcal{A} \times_1 \mathbf{B}^{(1)} \times_2 \mathbf{B}^{(2)} \dots \times_N \mathbf{B}^{(N)}$
$\mathcal{C} = \mathcal{A} \circ \mathcal{B}$	tensor or outer product of $\mathcal{A} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ and $\mathcal{B} \in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_M}$ yields $\mathcal{C} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N \times J_1 \times J_2 \times \dots \times J_M}$ with entries $c_{i_1 i_2 \dots i_N j_1 j_2 \dots j_M} = a_{i_1 i_2 \dots i_N} b_{j_1 j_2 \dots j_M}$
$\mathcal{X} = \mathbf{a}^{(1)} \circ \mathbf{a}^{(2)} \circ \dots \circ \mathbf{a}^{(N)}$	tensor or outer product of vectors $\mathbf{a}^{(n)} \in \mathbb{R}^{I_n}$ ( $n = 1, \dots, N$ ) yields a rank-1 tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ with entries $x_{i_1 i_2 \dots i_N} = a_{i_1}^{(1)} a_{i_2}^{(2)} \dots a_{i_N}^{(N)}$
$\mathbf{C} = \mathbf{A} \otimes \mathbf{B}$	Kronecker product of $\mathbf{A} \in \mathbb{R}^{I_1 \times I_2}$ and $\mathbf{B} \in \mathbb{R}^{J_1 \times J_2}$ yields $\mathbf{C} \in \mathbb{R}^{I_1 \times I_2 \times J_1 \times J_2}$ with entries $c_{(i_1-1)J_1+j_1, (i_2-1)J_2+j_2} = a_{i_1 i_2} b_{j_1 j_2}$
$\mathbf{C} = \mathbf{A} \odot \mathbf{B}$	Khatri-Rao product of $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_R] \in \mathbb{R}^{I \times R}$ and $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_R] \in \mathbb{R}^{J \times R}$ yields $\mathbf{C} \in \mathbb{R}^{I \times J \times R}$ with columns $\mathbf{c}_r = \mathbf{a}_r \otimes \mathbf{b}_r$



**Figure 1:** Construction of tensors. *Top:* Tensorization of a vector or matrix into the so-called *quantized format*; in scientific computing this facilitates super-compression of large-scale vectors or matrices. *Bottom:* Tensor formed through the discretization of a trivariate function  $f(x, y, z)$ .

(see Figure 1 (top)). For instance, a one-way exponential signal  $x(k) = az^k$  can be rearranged into a rank-1 Hankel matrix or a Hankel tensor [36]:

$$\mathbf{H} = \begin{pmatrix} x(0) & x(1) & x(2) & \dots \\ x(1) & x(2) & x(3) & \dots \\ x(2) & x(3) & x(4) & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix} = \mathbf{a} \mathbf{b} \circ \mathbf{b}, \quad (2)$$

where  $\mathbf{b} = [1, z, z^2, \dots]^T$ .

Also, in sensor array processing, tensor structures naturally emerge when combining snapshots from identical subarrays [17].

- 2) *Mathematical construction.* Among many such examples, the  $N$ th-order moments (cumulants) of a vector-valued random variable form an  $N$ th-order tensor [9], while in second-order ICA snapshots of data statistics (covariance matrices) are effectively slices of a third-order tensor [12], [37]. Also, a (*channel*  $\times$  *time*) data matrix can be transformed into a (*channel*  $\times$  *time*  $\times$  *frequency*) or (*channel*  $\times$  *time*  $\times$  *scale*) tensor via time-frequency or wavelet representations, a powerful procedure in multichannel EEG analysis in brain science [19], [38].
- 3) *Experiment design.* Multi-faceted data can be naturally stacked into a tensor; for instance, in wireless communications the so-called signal diversity (temporal, spatial, spectral, ...) corresponds to the order of the tensor [18]. In the same spirit, the standard EigenFaces can be generalized to TensorFaces by combining images with different illuminations, poses, and expressions [39], while the common modes in EEG recordings across subjects, trials, and conditions are best analyzed when combined together into a tensor [26].
- 4) *Naturally tensor data.* Some data sources are readily generated as tensors (e.g., RGB color images, videos, 3D light field dis-

plays) [40]. Also in scientific computing we often need to evaluate a discretized multivariate function; this is a natural tensor, as illustrated in Figure 1 (bottom) for a trivariate function  $f(x, y, z)$  [21], [27], [28].

The high dimensionality of the tensor format is associated with blessings — these include possibilities to obtain compact representations, uniqueness of decompositions, flexibility in the choice of constraints, and generality of components that can be identified.

#### CANONICAL POLYADIC DECOMPOSITION

**Definition.** A Polyadic Decomposition (PD) represents an  $N$ th-order tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  as a linear combination of rank-1 tensors in the form

$$\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{b}_r^{(1)} \circ \mathbf{b}_r^{(2)} \circ \dots \circ \mathbf{b}_r^{(N)}. \quad (3)$$

Equivalently,  $\mathcal{X}$  is expressed as a multilinear product with a diagonal core:

$$\begin{aligned} \mathcal{X} &= \mathcal{D} \times_1 \mathbf{B}^{(1)} \times_2 \mathbf{B}^{(2)} \dots \times_N \mathbf{B}^{(N)} \\ &= \llbracket \mathcal{D}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \dots, \mathbf{B}^{(N)} \rrbracket, \end{aligned} \quad (4)$$

where  $\mathcal{D} = \text{diag}_N(\lambda_1, \lambda_2, \dots, \lambda_R)$  (*cf.* the matrix case in (1)). Figure 2 (bottom) illustrates these two interpretations for a third-order tensor. The tensor rank is defined as the smallest value of  $R$  for which (3) holds exactly; the minimum rank PD is called *canonical* (CPD) and is desired in signal separation. The term CPD may also be considered as an abbreviation of CANDECOMP/PARAFAC decomposition, see **Historical notes**. The matrix/vector form of CPD can be obtained via the Khatri-Rao products as:

$$\begin{aligned} \mathbf{X}_{(n)} &= \mathbf{B}^{(n)} \mathbf{D} \left( \mathbf{B}^{(N)} \odot \dots \odot \mathbf{B}^{(n+1)} \right. \\ &\quad \left. \odot \mathbf{B}^{(n-1)} \odot \dots \odot \mathbf{B}^{(1)} \right)^T \\ \text{vec}(\mathcal{X}) &= [\mathbf{B}^{(N)} \odot \mathbf{B}^{(N-1)} \odot \dots \odot \mathbf{B}^{(1)}] \mathbf{d}. \end{aligned} \quad (5)$$

where  $\mathbf{d} = (\lambda_1, \lambda_2, \dots, \lambda_R)^T$ .

**Rank.** As mentioned earlier, rank-related properties are very different for matrices and tensors. For instance, the number of complex-valued rank-1 terms needed to represent a higher-order tensor can be strictly less than the number of real-valued rank-1 terms [20], while the determination of tensor rank is in general NP-hard [41]. Fortunately, in signal processing applications, rank estimation most often corresponds to determining the number of

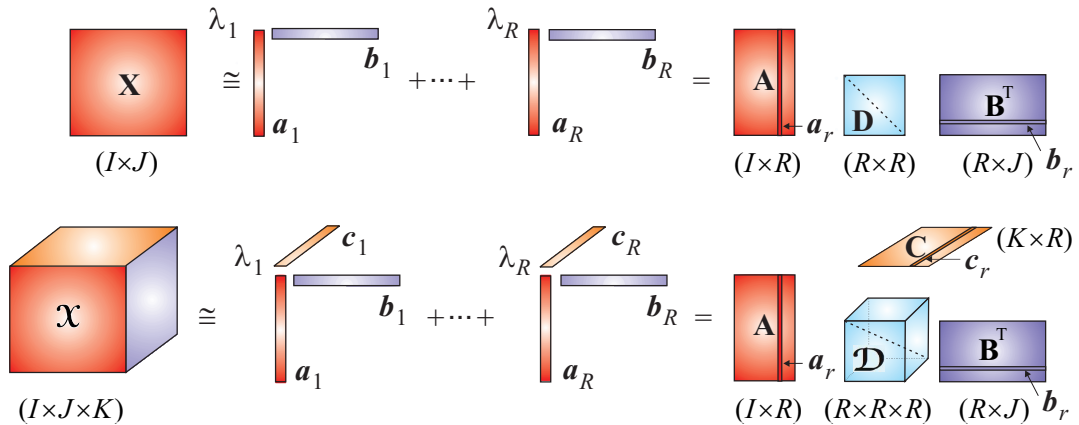
tensor components that can be retrieved with sufficient accuracy, and often there are only a few data components present. A pragmatic first assessment of the number of components may be through the inspection of the multilinear singular value spectrum (see Section **Tucker Decomposition**), which indicates the size of the core tensor in Figure 2 (bottom-right). Existing techniques for rank estimation include the CORCONDIA algorithm (core consistency diagnostic) which checks whether the core tensor is (approximately) diagonalizable [7], while a number of techniques operate by balancing the approximation error versus the number of degrees of freedom for a varying number of rank-1 terms [42]–[44].

**Uniqueness.** Uniqueness conditions give theoretical bounds for exact tensor decompositions. A classical uniqueness condition is due to Kruskal [31], which states that for third-order tensors the CPD is unique up to unavoidable scaling and permutation ambiguities, provided that  $k_{\mathbf{B}^{(1)}} + k_{\mathbf{B}^{(2)}} + k_{\mathbf{B}^{(3)}} \geq 2R + 2$ , where the Kruskal rank  $k_{\mathbf{B}}$  of a matrix  $\mathbf{B}$  is the maximum value ensuring that any subset of  $k_{\mathbf{B}}$  columns is linearly independent. In sparse modeling, the term  $(k_{\mathbf{B}} + 1)$  is also known as the spark [30]. A generalization to  $N$ th-order tensors is due to Sidiropoulos and Bro [45], and is given by:

$$\sum_{n=1}^N k_{\mathbf{B}^{(n)}} \geq 2R + N - 1. \quad (6)$$

More relaxed uniqueness conditions can be obtained when one factor matrix has full column rank [46]–[48]; for a thorough study of the third-order case, we refer to [32]. This all shows that, compared to matrix decompositions, CPD is unique under more natural and relaxed conditions, that only require the components to be “sufficiently different” and their number not unreasonably large. These conditions do not have a matrix counterpart, and are at the heart of tensor based signal separation.

**Computation.** Certain conditions, including Kruskal’s, enable explicit computation of the factor matrices in (3) using linear algebra (essentially, by solving sets of linear equations and by computing (generalized) Eigenvalue Decomposition) [6], [47], [49], [50]. The presence of noise in data means that CPD is rarely exact, and we need to fit a CPD model to the data by minimizing a suitable cost function. This is typically achieved by minimizing the Frobenius norm of the difference between the



**Figure 2:** Analogy between dyadic (top) and polyadic (bottom) decompositions; the Tucker format has a diagonal core. The uniqueness of these decompositions is a prerequisite for blind source separation and latent variable analysis.

given data tensor and its CP approximation, or alternatively by least absolute error fitting when the noise is Laplacian [51]. Theoretical Cramér-Rao Lower Bound (CRLB) and Cramér-Rao Induced Bound (CRIB) for the assessment of CPD performance were derived in [52] and [53].

Since the computation of CPD is intrinsically multilinear, we can arrive at the solution through a sequence of linear sub-problems as in the Alternating Least Squares (ALS) framework, whereby the LS cost function is optimized for one component matrix at a time, while keeping the other component matrices fixed [6]. As seen from (5), such a conditional update scheme boils down to solving overdetermined sets of linear equations.

While the ALS is attractive for its simplicity and satisfactory performance for a few well separated components and at sufficiently high SNR, it also inherits the problems of alternating algorithms and is not guaranteed to converge to a stationary point. This can be rectified by only updating the factor matrix for which the cost function has most decreased at a given step [54], but this results in an  $N$ -times increase in computational cost per iteration. The convergence of ALS is not yet completely understood — it is quasi-linear close to the stationary point [55], while it becomes rather slow for ill-conditioned cases; for more detail we refer to [56], [57].

Conventional all-at-once algorithms for numerical optimization such as nonlinear conjugate gradients, quasi-Newton or nonlinear least squares [58], [59] have been shown to often outperform ALS for ill-conditioned cases and

to be typically more robust to overfactoring, but come at a cost of a much higher computational load per iteration. More sophisticated versions use the rank-1 structure of the terms within CPD to perform efficient computation and storage of the Jacobian and (approximate) Hessian; their complexity is on par with ALS while for ill-conditioned cases the performance is often superior [60], [61].

An important difference between matrices and tensors is that the existence of a best rank- $R$  approximation of a tensor of rank greater than  $R$  is not guaranteed [20], [62] since the set of tensors whose rank is at most  $R$  is not closed. As a result, cost functions for computing factor matrices may only have an infimum (instead of a minimum) so that their minimization will approach the boundary of that set without ever reaching the boundary point. This will cause two or more rank-1 terms go to infinity upon convergence of an algorithm, however, numerically the diverging terms will almost completely cancel one another while the overall cost function will still decrease along the iterations [63]. These diverging terms indicate an inappropriate data model: the mismatch between the CPD and the original data tensor may arise due to an underestimated number of components, not all tensor components having a rank-1 structure, or data being too noisy.

**Constraints.** As mentioned earlier, under quite mild conditions the CPD is unique by itself, without requiring additional constraints. However, in order to enhance the accuracy and robustness with respect to noise, prior knowledge of data properties (e.g., statistical inde-

pendence, sparsity) may be incorporated into the constraints on factors so as to facilitate their physical interpretation, relax the uniqueness conditions, and even simplify computation [64]–[66]. Moreover, the orthogonality and non-negativity constraints ensure the existence of the minimum of the optimization criterion used [63], [64], [67].

**Applications.** The CPD has already been established as an advanced tool for signal separation in vastly diverse branches of signal processing and data analysis, such as in audio and speech processing, biomedical engineering, chemometrics, and machine learning [7], [22], [23], [26]. Note that algebraic ICA algorithms are effectively based on the CPD of a tensor of the statistics of recordings; the statistical independence of the sources is reflected in the diagonality of the core tensor in Figure 2, that is, in vanishing cross-statistics [11], [12]. The CPD is also heavily used in exploratory data analysis, where the rank-1 terms capture essential properties of dynamically complex signals [8]. Another example is in wireless communication, where the signals transmitted by different users correspond to rank-1 terms in the case of line-of-sight propagation [17]. Also, in harmonic retrieval and direction of arrival type applications, real or complex exponentials have a rank-1 structure, for which the use of CPD is natural [36], [65].

**Example 1.** Consider a sensor array consisting of  $K$  displaced but otherwise identical subarrays of  $I$  sensors, with  $\tilde{I} = KI$  sensors in total. For  $R$  narrowband sources in the far field, the baseband equivalent model of the array output becomes  $\mathbf{X} = \mathbf{A}\mathbf{S}^T + \mathbf{E}$ , where  $\mathbf{A} \in \mathbb{C}^{\tilde{I} \times R}$  is the global array response,  $\mathbf{S} \in \mathbb{C}^{J \times R}$  contains  $J$  snapshots of the sources, and  $\mathbf{E}$  is noise. A single source ( $R = 1$ ) can be obtained from the best rank-1 approximation of the matrix  $\mathbf{X}$ , however, for  $R > 1$  the decomposition of  $\mathbf{X}$  is not unique, and hence the separation of sources is not possible without incorporating additional information. Constraints on the sources that may yield a unique solution are, for instance, constant modulus or statistical independence [12], [68].

Consider a row-selection matrix  $\mathbf{J}_k \in \mathbb{C}^{I \times \tilde{I}}$  that extracts the rows of  $\mathbf{X}$  corresponding to the  $k$ -th subarray,  $k = 1, \dots, K$ . For two identical subarrays, the generalized EVD of the matrices  $\mathbf{J}_1\mathbf{X}$  and  $\mathbf{J}_2\mathbf{X}$  corresponds to the well-known ESPRIT [69]. For the case  $K > 2$ , we

shall consider  $\mathbf{J}_k\mathbf{X}$  as slices of the tensor  $\mathcal{X} \in \mathbb{C}^{I \times J \times K}$  (see Section **Tensorization**). It can be shown that the signal part of  $\mathcal{X}$  admits a CPD as in (3)–(4), with  $\lambda_1 = \dots = \lambda_R = 1$ ,  $\mathbf{J}_k\mathbf{A} = \mathbf{B}^{(1)}\text{diag}(\mathbf{b}_{k1}^{(3)}, \dots, \mathbf{b}_{kR}^{(3)})$  and  $\mathbf{B}^{(2)} = \mathbf{S}$  [17], and the consequent source separation under rather mild conditions — its uniqueness does not require constraints such as statistical independence or constant modulus. Moreover, the decomposition is unique even in cases when the number of sources  $R$  exceeds the number of subarray sensors  $I$ , or even the total number of sensors  $\tilde{I}$ . Notice that particular array geometries, such as linearly and uniformly displaced subarrays, can be converted into a constraint on CPD, yielding a further relaxation of the uniqueness conditions, reduced sensitivity to noise, and often faster computation [65].

#### TUCKER DECOMPOSITION

Figure 3 illustrates the principle of Tucker decomposition which treats a tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  as a multilinear transformation of a (typically dense but small) core tensor  $\mathcal{G} \in \mathbb{R}^{R_1 \times R_2 \times \dots \times R_N}$  by the factor matrices  $\mathbf{B}^{(n)} = [\mathbf{b}_1^{(n)}, \mathbf{b}_2^{(n)}, \dots, \mathbf{b}_{R_n}^{(n)}] \in \mathbb{R}^{I_n \times R_n}$ ,  $n = 1, 2, \dots, N$  [3], [4], given by

$$\mathcal{X} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \dots \sum_{r_N=1}^{R_N} g_{r_1 r_2 \dots r_N} (\mathbf{b}_{r_1}^{(1)} \circ \mathbf{b}_{r_2}^{(2)} \circ \dots \circ \mathbf{b}_{r_N}^{(N)}) \quad (7)$$

or equivalently

$$\begin{aligned} \mathcal{X} &= \mathcal{G} \times_1 \mathbf{B}^{(1)} \times_2 \mathbf{B}^{(2)} \dots \times_N \mathbf{B}^{(N)} \\ &= \llbracket \mathcal{G}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \dots, \mathbf{B}^{(N)} \rrbracket. \end{aligned} \quad (8)$$

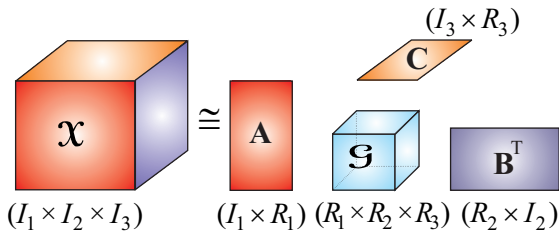
Via the Kronecker products (see Table II) Tucker decomposition can be expressed in a matrix/vector form as:

$$\begin{aligned} \mathbf{X}_{(n)} &= \mathbf{B}^{(n)} \mathbf{G}_{(n)} (\mathbf{B}^{(N)} \otimes \dots \otimes \mathbf{B}^{(n+1)} \otimes \mathbf{B}^{(n-1)} \otimes \dots \otimes \mathbf{B}^{(1)})^T \\ \text{vec}(\mathcal{X}) &= [\mathbf{B}^{(N)} \otimes \mathbf{B}^{(N-1)} \otimes \dots \otimes \mathbf{B}^{(1)}] \text{vec}(\mathcal{G}). \end{aligned}$$

Although Tucker initially used the orthogonality and ordering constraints on the core tensor and factor matrices [3], [4], we can also employ other meaningful constraints (see below).

**Multilinear rank.** For a core tensor of minimal size,  $R_1$  is the column rank (the dimension of the subspace spanned by mode-1 fibers),  $R_2$  is the row rank (the dimension of the subspace spanned by mode-2 fibers), and so on. A remarkable difference from matrices is that the values of  $R_1, R_2, \dots, R_N$  can be different for  $N \geq$





**Figure 3:** Tucker decomposition of a third-order tensor. The column spaces of  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$  represent the signal subspaces for the three modes. The core tensor  $\mathcal{G}$  is nondiagonal, accounting for possibly complex interactions among tensor components.

3. The  $N$ -tuple  $(R_1, R_2, \dots, R_N)$  is consequently called the *multilinear rank* of the tensor  $\mathcal{X}$ .

**Links between CPD and Tucker decomposition.** Eq. (7) shows that Tucker decomposition can be considered as an expansion in rank-1 terms (polyadic but not necessary canonical), while (4) represents CPD as a multilinear product of a core tensor and factor matrices (but the core is not necessary minimal); Table III shows various other connections. However, despite the obvious interchangeability of notation, the CP and Tucker decompositions serve different purposes. In general, the Tucker core cannot be diagonalized, while the number of CPD terms may not be bounded by the multilinear rank. Consequently, in signal processing and data analysis, CPD is typically used for factorizing data into easy to interpret components (i.e., the rank-1 terms), while the goal of unconstrained Tucker decompositions is most often to compress data into a tensor of smaller size (i.e., the core tensor) or to find the subspaces spanned by the fibers (i.e., the column spaces of the factor matrices).

**Uniqueness.** The unconstrained Tucker decomposition is in general not unique, that is, factor matrices  $\mathbf{B}^{(n)}$  are rotation invariant. However, physically, the subspaces defined by the factor matrices in Tucker decomposition are unique, while the bases in these subspaces may be chosen arbitrarily — their choice is compensated for within the core tensor. This becomes clear upon realizing that any factor matrix in (8) can be post-multiplied by any nonsingular (rotation) matrix; in turn, this multiplies the core

**TABLE III:** Different forms of CPD and Tucker representations of a third-order tensor  $\mathcal{X} \in \mathbb{R}^{I \times J \times K}$ .

CPD	Tucker Decomposition
Tensor representation, outer products	
$\mathcal{X} = \sum_{r=1}^R \lambda_r \mathbf{a}_r \circ \mathbf{b}_r \circ \mathbf{c}_r$	$\mathcal{X} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} g_{r_1 r_2 r_3} \mathbf{a}_{r_1} \circ \mathbf{b}_{r_2} \circ \mathbf{c}_{r_3}$
Tensor representation, multilinear products	
$\mathcal{X} = \mathcal{D} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}$	$\mathcal{X} = \mathcal{G} \times_1 \mathbf{A} \times_2 \mathbf{B} \times_3 \mathbf{C}$
Matrix representations	
$\mathbf{X}_{(1)} = \mathbf{A} \mathbf{D} (\mathbf{C} \odot \mathbf{B})^T$	$\mathbf{X}_{(1)} = \mathbf{A} \mathbf{G}_{(1)} (\mathbf{C} \otimes \mathbf{B})^T$
$\mathbf{X}_{(2)} = \mathbf{B} \mathbf{D} (\mathbf{C} \odot \mathbf{A})^T$	$\mathbf{X}_{(2)} = \mathbf{B} \mathbf{G}_{(2)} (\mathbf{C} \otimes \mathbf{A})^T$
$\mathbf{X}_{(3)} = \mathbf{C} \mathbf{D} (\mathbf{B} \odot \mathbf{A})^T$	$\mathbf{X}_{(3)} = \mathbf{C} \mathbf{G}_{(3)} (\mathbf{B} \otimes \mathbf{A})^T$
Vector representation	
$\text{vec}(\mathcal{X}) = (\mathbf{C} \odot \mathbf{B} \odot \mathbf{A}) \mathbf{d}$	$\text{vec}(\mathcal{X}) = (\mathbf{C} \otimes \mathbf{B} \otimes \mathbf{A}) \text{vec}(\mathcal{G})$
Scalar representation	
$x_{ijk} = \sum_{r=1}^R \lambda_r a_{ir} b_{jr} c_{kr}$	$x_{ijk} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \sum_{r_3=1}^{R_3} g_{r_1 r_2 r_3} a_{i r_1} b_{j r_2} c_{k r_3}$
Matrix slices $\mathbf{X}_k = \mathbf{X}(:, :, k)$	
$\mathbf{X}_k = \mathbf{A} \text{diag}(c_{k1}, c_{k2}, \dots, c_{kR}) \mathbf{B}^T$	$\mathbf{X}_k = \mathbf{A} \sum_{r_3=1}^{R_3} c_{k r_3} \mathbf{G}(:, :, r_3) \mathbf{B}^T$

tensor by its inverse, that is

$$\begin{aligned}
 \mathcal{X} &= \llbracket \mathcal{G}; \mathbf{B}^{(1)}, \mathbf{B}^{(2)}, \dots, \mathbf{B}^{(N)} \rrbracket \\
 &= \llbracket \mathcal{H}; \mathbf{B}^{(1)} \mathbf{R}^{(1)}, \mathbf{B}^{(2)} \mathbf{R}^{(2)}, \dots, \mathbf{B}^{(N)} \mathbf{R}^{(N)} \rrbracket, \\
 \mathcal{H} &= \llbracket \mathcal{G}; \mathbf{R}^{(1)-1}, \mathbf{R}^{(2)-1}, \dots, \mathbf{R}^{(N)-1} \rrbracket, \quad (9)
 \end{aligned}$$

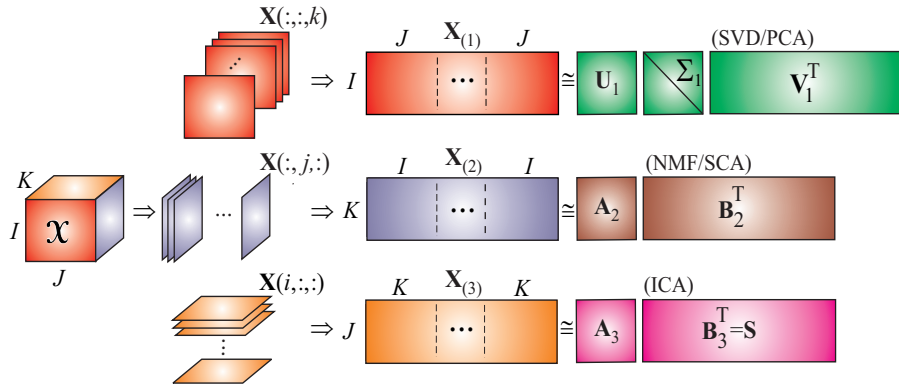
where  $\mathbf{R}^{(n)}$  are invertible.

**Multilinear SVD (MLSVD).** Orthonormal bases in a constrained Tucker representation can be obtained via the SVD of the mode- $n$  matrixed tensor  $\mathbf{X}_{(n)} = \mathbf{U}_n \Sigma_n \mathbf{V}_n^T$  (i.e.,  $\mathbf{B}^{(n)} = \mathbf{U}_n$ ,  $n = 1, 2, \dots, N$ ). Due to the orthonormality, the corresponding core tensor becomes

$$\mathcal{S} = \mathcal{X} \times_1 \mathbf{U}_1^T \times_2 \mathbf{U}_2^T \cdots \times_N \mathbf{U}_N^T. \quad (10)$$

Then, the singular values of  $\mathbf{X}_{(n)}$  are the Frobenius norms of the corresponding slices of the core tensor  $\mathcal{S}$ :  $(\Sigma_n)_{r_n, r_n} = \|\mathcal{S}(:, :, \dots, r_n, :, \dots, :)\|$ , with slices in the same mode being mutually orthogonal, i.e., their inner products





**Figure 4:** Multiway Component Analysis (MWCA) for a third-order tensor, assuming that the components are: principal and orthogonal in the first mode, nonnegative and sparse in the second mode and statistically independent in the third mode.

are zero. The columns of  $\mathbf{U}_n$  may thus be seen as multilinear singular vectors, while the norms of the slices of the core are multilinear singular values [13]. As in the matrix case, the multilinear singular values govern the multilinear rank, while the multilinear singular vectors allow, for each mode separately, an interpretation as in PCA [8].

**Low multilinear rank approximation.** Analogous to PCA, a large-scale data tensor  $\mathcal{X}$  can be approximated by discarding the multilinear singular vectors and slices of the core tensor that correspond to small multilinear singular values, that is, through truncated matrix SVDs. Low multilinear rank approximation is always well-posed, however, the truncation is not necessarily optimal in the LS sense, although a good estimate can often be made as the approximation error corresponds to the degree of truncation. When it comes to finding the best approximation, the ALS type algorithms exhibit similar advantages and drawbacks to those used for CPD [8], [70]. Optimization-based algorithms exploiting second-order information have also been proposed [71], [72].

**Constraints and Tucker-based multiway component analysis (MWCA).** Besides orthogonality, constraints that may help to find unique basis vectors in a Tucker representation include statistical independence, sparsity, smoothness and nonnegativity [19], [73], [74]. Components of a data tensor seldom have the same properties in its modes, and for physically meaningful representation different constraints may be required in different modes, so as to match the properties of the data at hand. Figure 4 illus-

trates the concept of MWCA and its flexibility in choosing the mode-wise constraints; a Tucker representation of MWCA naturally accommodates such diversities in different modes.

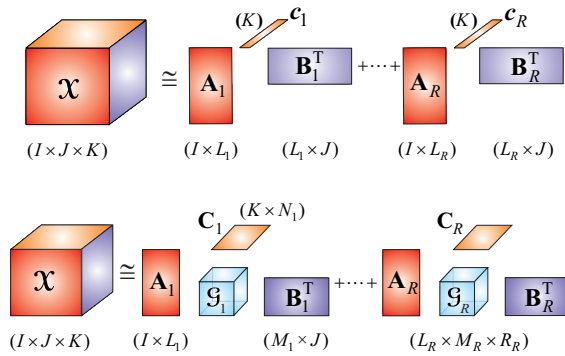
**Other applications.** We have shown that Tucker decomposition may be considered as a multilinear extension of PCA [8]; it therefore generalizes signal subspace techniques, with applications including classification, feature extraction, and subspace-based harmonic retrieval [25], [39], [75], [76]. For instance, a low multilinear rank approximation achieved through Tucker decomposition may yield a higher Signal-to-Noise Ratio (SNR) than the SNR in the original raw data tensor, making Tucker decomposition a very natural tool for compression and signal enhancement [7], [8], [24].

## BLOCK TERM DECOMPOSITIONS

We have already shown that CPD is unique under quite mild conditions, a further advantage of tensors over matrices is that it is even possible to relax the rank-1 constraint on the terms, thus opening completely new possibilities in e.g. BSS. For clarity, we shall consider the third-order case, whereby, by replacing the rank-1 matrices  $\mathbf{b}_r^{(1)} \circ \mathbf{b}_r^{(2)} = \mathbf{b}_r^{(1)} \mathbf{b}_r^{(2)T}$  in (3) by low-rank matrices  $\mathbf{A}_r \mathbf{B}_r^T$ , the tensor  $\mathcal{X}$  can be represented as (Figure 5, top):

$$\mathcal{X} = \sum_{r=1}^R (\mathbf{A}_r \mathbf{B}_r^T) \circ \mathbf{c}_r. \quad (11)$$

Figure 5 (bottom) shows that we can even use terms that are only required to have a low



**Figure 5:** Block Term Decompositions (BTDs) find data components that are structurally more complex than the rank-1 terms in CPD. *Top:* Decomposition into terms with multilinear rank  $(L_r, L_r, 1)$ . *Bottom:* Decomposition into terms with multilinear rank  $(L_r, M_r, N_r)$ .

multilinear rank (see also Section **Tucker Decomposition**), to give:

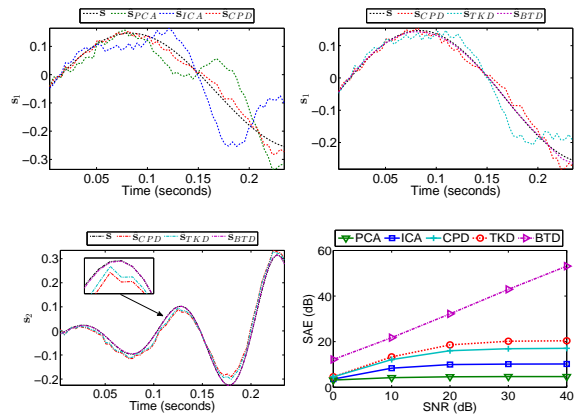
$$\mathcal{X} = \sum_{r=1}^R \mathcal{G}_r \times_1 \mathbf{A}_r \times_2 \mathbf{B}_r \times_3 \mathbf{C}_r. \quad (12)$$

These so-called Block Term Decompositions (BTD) admit the modelling of more complex signal components than CPD, and are unique under more restrictive but still fairly natural conditions [77]–[79].

**Example 3.** To compare some standard and tensor approaches for the separation of short duration correlated sources, BSS was performed on five linear mixtures of the sources  $s_1(t) = \sin(6\pi t)$  and  $s_2(t) = \exp(10t) \sin(20\pi t)$ , which were contaminated by white Gaussian noise, to give the mixtures  $\mathbf{X} = \mathbf{A}\mathbf{S} + \mathbf{E} \in \mathbb{R}^{5 \times 60}$ , where  $\mathbf{S}(t) = [s_1(t), s_2(t)]^T$  and  $\mathbf{A} \in \mathbb{R}^{5 \times 2}$  was a random matrix whose columns (mixing vectors) satisfy  $\mathbf{a}_1^T \mathbf{a}_2 = 0.1$ ,  $\|\mathbf{a}_1\| = \|\mathbf{a}_2\| = 1$ . The 3Hz sine wave did not complete a full period over the 60 samples, so that the two sources had a correlation degree of  $\frac{|\mathbf{s}_1^T \mathbf{s}_2|}{\|\mathbf{s}_1\|_2 \|\mathbf{s}_2\|_2} = 0.35$ . The tensor approaches, CPD, Tucker decomposition and BTD employed a third-order tensor  $\mathcal{X}$  of size  $24 \times 37 \times 5$  generated from five Hankel matrices whose elements obey  $\mathcal{X}(i, j, k) = \mathcal{X}(k, i + j - 1)$  (see Section **Tensorization**). The average squared angular error (SAE) was used as the performance measure. Figure 6 shows the simulation results, illustrating that:

- **PCA** failed since the mixing vectors were not orthogonal and the source signals were correlated, both violating the assumptions for PCA.

- **ICA** (using the JADE algorithm [10]) failed because the signals were not statistically independent, as assumed in ICA.
- **Low rank tensor approximation:** a rank-2 CPD was used to estimate  $\mathbf{A}$  as the third factor matrix, which was then inverted to yield the sources. The accuracy of CPD was compromised as the components of tensor  $\mathcal{X}$  cannot be represented by rank-1 terms.
- **Low multilinear rank approximation:** Tucker decomposition (TKD) for the multilinear rank  $(4, 4, 2)$  was able to retrieve the column space of the mixing matrix but could not find the individual mixing vectors due to the non-uniqueness of TKD.
- **BTD in multilinear rank- $(2, 2, 1)$  terms** matched the data structure [78], and it is remarkable that the sources were recovered using as few as 6 samples in the noise-free case.



**Figure 6:** Blind separation of the mixture of a pure sine wave and an exponentially modulated sine wave using PCA, ICA, CPD, Tucker decomposition (TKD) and BTD. The sources  $s_1$  and  $s_2$  are correlated and of short duration; the symbols  $\hat{s}_1$  and  $\hat{s}_2$  denote the estimated sources.

## HIGHER-ORDER COMPRESSED SENSING

The aim of *Compressed Sensing* (CS) is to provide faithful reconstruction of a signal of interest when the set of available measurements is (much) smaller than the size of the original signal [80]–[83]. Formally, we have available  $M$  (compressive) data samples  $\mathbf{y} \in \mathbb{R}^M$ , which are assumed to be linear transformations of the original signal  $\mathbf{x} \in \mathbb{R}^I$  ( $M < I$ ). In other words,  $\mathbf{y} = \mathbf{\Phi}\mathbf{x}$ , where the *sensing matrix*  $\mathbf{\Phi} \in \mathbb{R}^{M \times I}$  is usually random. Since the projections are of a lower dimension than the original data, the reconstruction is an ill-posed inverse

problem, whose solution requires knowledge of the physics of the problem converted into constraints. For example, a 2D image  $\mathbf{X} \in \mathbb{R}^{I_1 \times I_2}$  can be vectorized as a long vector  $\mathbf{x} = \text{vec}(\mathbf{X}) \in \mathbb{R}^I$  ( $I = I_1 I_2$ ) that admits sparse representation in a known *dictionary*  $\mathbf{B} \in \mathbb{R}^{I \times I}$ , so that  $\mathbf{x} = \mathbf{B}\mathbf{g}$ , where the matrix  $\mathbf{B}$  may be a wavelet or discrete cosine transform (DCT) dictionary. Then, faithful recovery of the original signal  $\mathbf{x}$  requires finding the sparsest vector  $\mathbf{g}$  such that:

$$\mathbf{y} = \mathbf{W}\mathbf{g}, \text{ with } \|\mathbf{g}\|_0 \leq K, \quad \mathbf{W} = \Phi\mathbf{B}, \quad (13)$$

where  $\|\cdot\|_0$  is the  $\ell_0$ -norm (number of non-zero entries) and  $K \ll I$ .

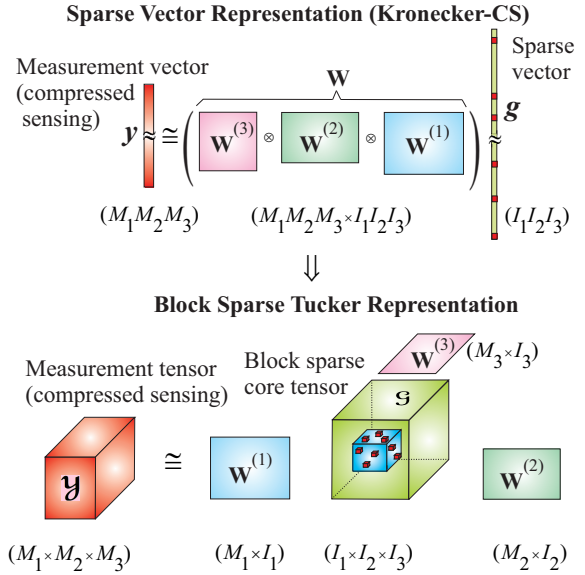
Since the  $\ell_0$ -norm minimization is not practical, alternative solutions involve iterative refinements of the estimates of vector  $\mathbf{g}$  using greedy algorithms such as the Orthogonal Matching Pursuit (OMP) algorithm, or the  $\ell_1$ -norm minimization algorithms ( $\|\mathbf{g}\|_1 = \sum_{i=1}^I |g_i|$ ) [83]. Low coherence of the composite dictionary matrix  $\mathbf{W}$  is a prerequisite for a satisfactory recovery of  $\mathbf{g}$  (and hence  $\mathbf{x}$ ) — we need to choose  $\Phi$  and  $\mathbf{B}$  so that the correlation between the columns of  $\mathbf{W}$  is minimum [83].

When extending the CS framework to tensor data, we face two obstacles:

- *Loss of information*, such as spatial and contextual relationships in data, when a tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  is vectorized.
- *Data handling*, since the size of vectorized data and the associated dictionary  $\mathbf{B} \in \mathbb{R}^{I \times I}$  easily becomes prohibitively large (see Section **Curse of Dimensionality**), especially for tensors of high order.

Fortunately, tensor data are typically highly structured – a perfect match for compressive sampling – so that the CS framework relaxes data acquisition requirements, enables compact storage, and facilitates data completion (inpainting of missing samples due to a broken sensor or unreliable measurement).

**Kronecker-CS for fixed dictionaries.** In many applications, the dictionary and the sensing matrix admit a Kronecker structure (Kronecker-CS model), as illustrated in Figure 7 (top) [84]. In this way, the global *composite dictionary matrix* becomes  $\mathbf{W} = \mathbf{W}^{(N)} \otimes \mathbf{W}^{(N-1)} \otimes \dots \otimes \mathbf{W}^{(1)}$ , where each term  $\mathbf{W}^{(n)} = \Phi^{(n)}\mathbf{B}^{(n)}$  has a reduced dimensionality since  $\mathbf{B}^{(n)} \in \mathbb{R}^{I_n \times I_n}$  and  $\Phi^{(n)} \in \mathbb{R}^{M_n \times I_n}$ . Denote  $M = M_1 M_2 \dots M_N$  and  $I = I_1 I_2 \dots I_N$ , and since  $M_n \leq I_n$ ,  $n = 1, 2, \dots, N$ , this reduces storage requirements by a factor  $\frac{\sum_n I_n M_n}{MI}$ . The computation of  $\mathbf{W}\mathbf{g}$  is affordable



**Figure 7:** Compressed sensing with a Kronecker-structured dictionary. *Top:* Vector representation. *Bottom:* Tensor representation; Orthogonal Matching Pursuit (OMP) can perform faster if the sparse entries belong to a small subtensor, up to permutation of the columns of  $\mathbf{W}^{(1)}$ ,  $\mathbf{W}^{(2)}$ ,  $\mathbf{W}^{(3)}$ .

since  $\mathbf{g}$  is sparse, however, computing  $\mathbf{W}^T \mathbf{y}$  is expensive but can be efficiently implemented through a sequence of products involving much smaller matrices  $\mathbf{W}^{(n)}$  [85]. We refer to [84] for links between the coherence of factors  $\mathbf{W}^{(n)}$  and the coherence of the global composite dictionary matrix  $\mathbf{W}$ .

Figure 7 and Table III illustrate that the Kronecker-CS model is effectively a vectorized Tucker decomposition with a sparse core. The tensor equivalent of the CS paradigm in (13) is therefore to find the sparsest core tensor  $\mathcal{G}$  such that:

$$\mathcal{Y} \cong \mathcal{G} \times_1 \mathbf{W}^{(1)} \times_2 \mathbf{W}^{(2)} \dots \times_N \mathbf{W}^{(N)}, \quad (14)$$

with  $\|\mathcal{G}\|_0 \leq K$ , for a given set of mode-wise dictionaries  $\mathbf{B}^{(n)}$  and sensing matrices  $\Phi^{(n)}$  ( $n = 1, 2, \dots, N$ ). Working with several small dictionary matrices, appearing in a Tucker representation, instead of a large global dictionary matrix, is an example of the use of tensor structure for efficient representation, see also Section **Curse of Dimensionality**.

A higher-order extension of the OMP algorithm, referred to as the Kronecker-OMP algorithm [85], requires  $K$  iterations to find the  $K$  non-zero entries of the core tensor  $\mathcal{G}$ . Additional

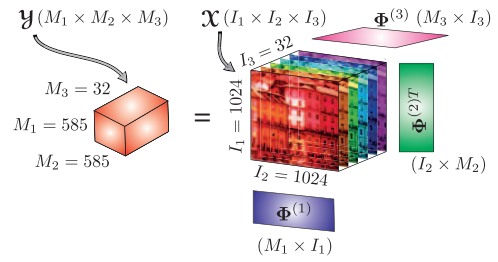
computational advantages can be gained if it can be assumed that the  $K$  non-zero entries belong to a small subtensor of  $\mathcal{G}$ , as shown in Figure 7 (bottom); such a structure is inherent to e.g., hyperspectral imaging [85], [86] and 3D astrophysical signals. More precisely, if the  $K = L^N$  non-zero entries are located within a subtensor of size  $(L \times L \times \dots \times L)$ , where  $L \ll I_n$ , then the so-called  $N$ -way Block OMP algorithm (N-BOMP) requires at most  $NL$  iterations, which is linear in  $N$  [85]. The Kronecker-CS model has been applied in Magnetic Resonance Imaging (MRI), hyper-spectral imaging, and in the inpainting of multiway data [84], [86].

#### Approaches without fixed dictionaries.

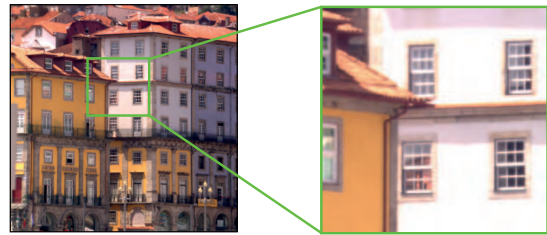
In Kronecker-CS the mode-wise dictionaries  $\mathbf{B}^{(n)} \in \mathbb{R}^{I_n \times I_n}$  can be chosen so as best to represent physical properties or prior knowledge about the data. They can also be learned from a large ensemble of data tensors, for instance in an ALS type fashion [86]. Instead of the total number of sparse entries in the core tensor, the size of the core (i.e., the multilinear rank) may be used as a measure for sparsity so as to obtain a low-complexity representation from compressively sampled data [87], [88]. Alternatively, a PD representation can be used instead of a Tucker representation. Indeed, early work in chemometrics involved excitation-emission data for which part of the entries was unreliable because of scattering; the CPD of the data tensor is then computed by treating such entries as missing [7]. While CS variants of several CPD algorithms exist [59], [89], the ‘‘oracle’’ properties of tensor-based models are still not as well understood as for their standard models; a notable exception is CPD with sparse factors [90].

**Example 2.** Figure 8 shows an original 3D  $(1024 \times 1024 \times 32)$  hyperspectral image  $\mathcal{X}$  which contains scene reflectance measured at 32 different frequency channels, acquired by a low-noise Peltier-cooled digital camera in the wavelength range of 400–720 nm [91]. Within the Kronecker-CS setting, the tensor of compressive measurements  $\mathcal{Y}$  was obtained by multiplying the frontal slices by random Gaussian sensing matrices  $\Phi^{(1)} \in \mathbb{R}^{M_1 \times 1024}$  and  $\Phi^{(2)} \in \mathbb{R}^{M_2 \times 1024}$  ( $M_1, M_2 < 1024$ ) in the first and second mode, respectively, while  $\Phi^{(3)} \in \mathbb{R}^{32 \times 32}$  was the identity matrix (see Figure 8 (top)). We used Daubechies wavelet factor matrices  $\mathbf{B}^{(1)} = \mathbf{B}^{(2)} \in \mathbb{R}^{1024 \times 1024}$  and  $\mathbf{B}^{(3)} \in \mathbb{R}^{32 \times 32}$ , and employed N-BOMP to recover the small

#### Kronecker-CS of a 32-channel hyperspectral image $\mathcal{X}$



Original hyperspectral image - RGB display  
(1024 x 1024 x 32) (256 x 256 x 32)



Reconstruction (SP=33%, PSNR = 35.51dB) - RGB display  
(1024 x 1024 x 32) (256 x 256 x 32)



**Figure 8:** Multidimensional compressed sensing of a 3D hyperspectral image using Tucker representation with a small sparse core in wavelet bases.

sparse core tensor and, subsequently, reconstruct the original 3D image as shown in Figure 8 (bottom). For the Sampling Ratio SP=33% ( $M_1 = M_2 = 585$ ) this gave the Peak Signal to Noise Ratio (PSNR) of 35.51dB, while taking 71 minutes to compute the required  $N_{iter} = 841$  sparse entries. For the same quality of reconstruction (PSNR=35.51dB), the more conventional Kronecker-OMP algorithm found 0.1% of the wavelet coefficients as significant, thus requiring  $N_{iter} = K = 0.001 \times (1024 \times 1024 \times 32) = 33,555$  iterations and days of computation time.

#### LARGE-SCALE DATA AND CURSE OF DIMENSIONALITY

The sheer size of tensor data easily exceeds the memory or saturates the processing capability of standard computers, it is therefore natural to ask ourselves how tensor decompositions can



be computed if the tensor dimensions in all or some modes are large or, worse still, if the tensor order is high. The term *curse of dimensionality*, in a general sense, was introduced by Bellman to refer to various computational bottlenecks when dealing with high-dimensional settings. In the context of tensors, the *curse of dimensionality* refers to the fact that the number of elements of an  $N$ th-order ( $I \times I \times \dots \times I$ ) tensor,  $I^N$ , scales *exponentially* with the tensor order  $N$ . For example, the number of values of a discretized function in Figure 1 (bottom), quickly becomes unmanageable in terms of both computations and storing as  $N$  increases. In addition to their standard use (signal separation, enhancement, etc.), tensor decompositions may be elegantly employed in this context as efficient representation tools. The first question is then which type of tensor decomposition is appropriate.

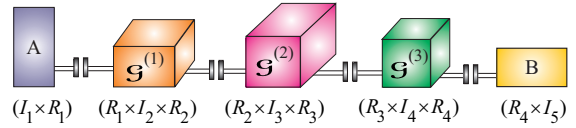
**Efficient data handling.** If all computations are performed on a *CP representation* and not on the raw data tensor itself, then instead of the original  $I^N$  raw data entries, the number of parameters in a CP representation reduces to *NIR*, which scales *linearly* in  $N$  (see Table IV). This effectively bypasses the *curse of dimensionality*, while giving us the freedom to choose the rank  $R$  as a function of the desired accuracy [14]; on the other hand the CP approximation may involve numerical problems (see Section **Canonical Polyadic Decomposition**).

Compression is also inherent to *Tucker decomposition*, as it reduces the size of a given data tensor from the original  $I^N$  to  $(NIR + R^N)$ , thus exhibiting an approximate compression ratio of  $(\frac{I}{R})^N$ . We can then benefit from the well understood and reliable approximation by means of matrix SVD, however, this is only meaningful for low  $N$ .

**TABLE IV:** Storage complexities of tensor models for an  $N$ th-order tensor  $\mathcal{X} \in \mathbb{R}^{I \times I \times \dots \times I}$ , whose original storage complexity is  $\mathcal{O}(I^N)$ .

1. CPD	$\mathcal{O}(NIR)$
2. Tucker	$\mathcal{O}(NIR + R^N)$
3. TT	$\mathcal{O}(NIR^2)$
4. QTT	$\mathcal{O}(NR^2 \log_2(I))$

**Tensor networks.** A numerically reliable way to tackle curse of dimensionality is through a concept from scientific computing and quan-

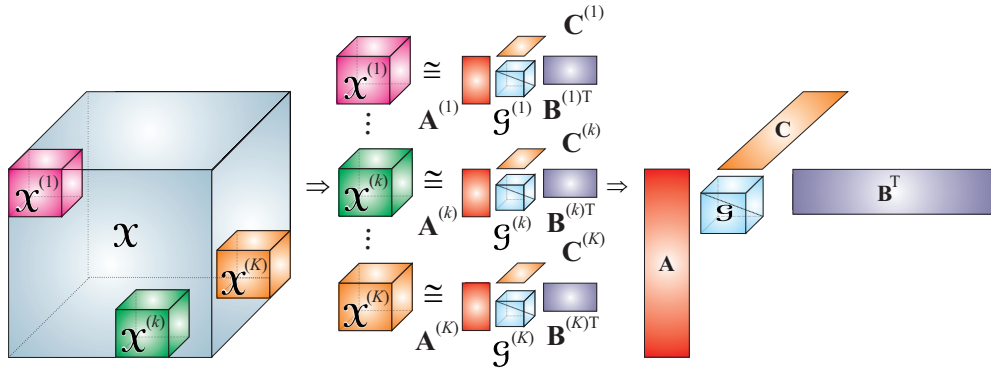


**Figure 9:** Tensor Train (TT) decomposition of a fifth-order tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_5}$ , consisting of two matrix carriages and three third-order tensor carriages. The five carriages are connected through tensor contractions, which can be expressed in a scalar form as  $x_{i_1, i_2, i_3, i_4, i_5} = \sum_{r_1=1}^{R_1} \sum_{r_2=1}^{R_2} \dots \sum_{r_4=1}^{R_4} a_{i_1, r_1} g_{r_1, i_2, r_2}^{(1)} g_{r_2, i_3, r_3}^{(2)} g_{r_3, i_4, r_4}^{(3)} b_{r_4, i_5}$ .

tum information theory, termed *tensor networks*, which represents a tensor of a possibly very high order as a set of sparsely interconnected matrices and core tensors of low order (typically, order 3). These low-dimensional cores are interconnected via tensor contractions to provide a highly compressed representation of a data tensor. In addition, existing algorithms for the approximation of a given tensor by a tensor network have good numerical properties, making it possible to control the error and achieve any desired accuracy of approximation. For example, tensor networks allow for the representation of a wide class of discretized multivariate functions even in cases where the number of function values is larger than the number of atoms in the universe [21], [27], [28].

Examples of tensor networks are the hierarchical Tucker (HT) decompositions and Tensor Trains (TT) (see Figure 9) [15], [16]. The TTs are also known as Matrix Product States (MPS) and have been used by physicists more than two decades (see [92], [93] and references therein). The PARATREE algorithm was developed in signal processing and follows a similar idea, it uses a polyadic representation of a data tensor (in a possibly nonminimal number of terms), whose computation then requires only the matrix SVD [94].

For very large-scale data that exhibit a well-defined structure, an even more radical approach can be employed to achieve a parsimonious representation — through the concept of *quantized or quantic tensor networks* (QTN) [27], [28]. For example, a huge vector  $x \in \mathbb{R}^I$  with  $I = 2^L$  elements can be quantized and tensorized through reshaping into a  $(2 \times 2 \times \dots \times 2)$  tensor  $\mathcal{X}$  of order  $L$ , as illustrated in Figure 1 (top). If  $x$  is an exponential signal,  $x(k) = az^k$ , then  $\mathcal{X}$  is a symmetric rank-1 tensor



**Figure 10:** Efficient computation of the CP and Tucker decompositions, whereby tensor decompositions are computed in parallel for sampled blocks, these are then merged to obtain the global components  $\mathbf{A}$ ,  $\mathbf{B}$ ,  $\mathbf{C}$  and a core tensor  $\mathcal{G}$ .

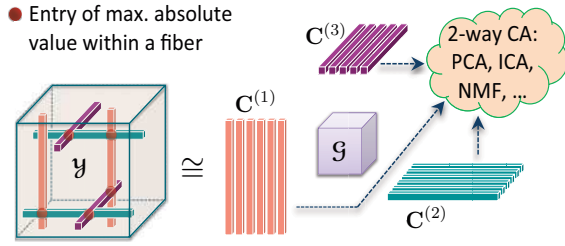
that can be represented by two parameters: the scaling factor  $a$  and the generator  $z$  (cf. (2) in Section **Tensorization**). Non-symmetric terms provide further opportunities, beyond the sum-of-exponential representation by symmetric low-rank tensors. Huge matrices and tensors may be dealt with in the same manner. For instance, an  $N$ th-order tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ , with  $I_n = q^{L_n}$ , can be quantized in all modes simultaneously to yield a  $(q \times q \times \dots \times q)$  quantized tensor of higher order. In QTN,  $q$  is small, typically  $q = 2, 3, 4$ , for example, the binary encoding ( $q = 2$ ) reshapes an  $N$ th-order tensor with  $(2^{L_1} \times 2^{L_2} \times \dots \times 2^{L_N})$  elements into a tensor of order  $(L_1 + L_2 + \dots + L_N)$  with the same number of elements. The tensor train decomposition applied to quantized tensors is referred to as the quantized TT (QTT); variants for other tensor representations have also been derived [27], [28]. In scientific computing, such formats provide the so-called *super-compression* — a logarithmic reduction of storage requirements:  $\mathcal{O}(I^N) \rightarrow \mathcal{O}(N \log_q(I))$ .

**Computation of the decomposition/representation.** Now that we have addressed the possibilities for efficient tensor representation, the question that needs to be answered is how these representations can be computed from the data in an efficient manner. The first approach is to process the data in smaller blocks rather than in a batch manner [95]. In such a “divide-and-conquer” approach, different blocks may be processed in parallel and their decompositions carefully recombined (see Figure 10) [95], [96]. In fact, we may even compute the decomposition through recursive updating, as new data arrive [97]. Such recursive techniques may be used

for efficient computation and for tracking decompositions in the case of nonstationary data.

The second approach would be to employ compressed sensing ideas (see Section **Higher-Order Compressed Sensing**) to fit an algebraic model with a limited number of parameters to possibly large data. In addition to completion, the goal here is a significant reduction of the cost of data acquisition, manipulation and storage — breaking the Curse of Dimensionality being an extreme case.

While algorithms for this purpose are available both for low rank and low multilinear rank representation [59], [87], an even more drastic approach would be to directly adopt sampled fibers as the bases in a tensor representation. In the Tucker decomposition setting we would choose the columns of the factor matrices  $\mathbf{B}^{(n)}$  as mode- $n$  fibers of the tensor, which requires addressing the following two problems: (i) how to find fibers that allow us to best represent the tensor, and (ii) how to compute the corresponding core tensor at a low cost (i.e., with minimal access to the data). The matrix counterpart of this problem (i.e., representation of a large matrix on the basis of a few columns and rows) is referred to as the *pseudoskeleton approximation* [98], where the optimal representation corresponds to the columns and rows that intersect in the submatrix of maximal volume (maximal absolute value of the determinant). Finding the optimal submatrix is computationally hard, but quasi-optimal submatrices may be found by heuristic so-called “cross-approximation” methods that only require a limited, partial exploration of the data matrix. Tucker variants of this approach have



**Figure 11:** Tucker representation through fiber sampling and cross-approximation: the columns of factor matrices are sampled from the fibers of the original data tensor  $\mathcal{X}$ . Within MWCA the selected fibers may be further processed using BSS algorithms.

been derived in [99]–[101] and are illustrated in Figure 11, while cross-approximation for the TT format has been derived in [102]. Following a somewhat different idea, a tensor generalization of the CUR decomposition of matrices samples fibers on the basis of statistics derived from the data [103].

#### MULTIWAY REGRESSION — HIGHER ORDER PLS (HOPLS)

**Multivariate regression.** Regression refers to the modelling of one or more *dependent variables* (responses),  $Y$ , by a set of *independent data* (predictors),  $X$ . In the simplest case of conditional MSE estimation,  $\hat{y} = E(y|x)$ , the response  $y$  is a linear combination of the elements of the vector of predictors  $x$ ; for multivariate data the Multivariate Linear Regression (MLR) uses a matrix model,  $Y = XP + E$ , where  $P$  is the matrix of coefficients (loadings) and  $E$  the residual matrix. The MLR solution gives  $P = (X^T X)^{-1} X^T Y$ , and involves inversion of the moment matrix  $X^T X$ . A common technique to stabilize the inverse of the moment matrix  $X^T X$  is *principal component regression* (PCR), which employs low rank approximation of  $X$ .

#### Modelling structure in data — the PLS.

Notice that in stabilizing multivariate regression PCR uses only information in the  $X$ -variables, with no feedback from the  $Y$ -variables. The idea behind the Partial Least Squares (PLS) method is to account for structure in data by assuming that the underlying system is governed by a small number,  $R$ , of specifically constructed latent variables, called *scores*, that are shared between the  $X$ - and  $Y$ -variables; in estimating the number  $R$ , PLS compromises between fitting  $X$  and predicting  $Y$ . Figure 12 illustrates that the PLS procedure: (i) uses eigenanalysis to perform

$$\begin{aligned} \mathbf{X} &\cong \mathbf{T} \mathbf{P}^T = \sum_{r=1}^R \mathbf{t}_r \mathbf{p}_r^T \\ &(I \times N) \quad (I \times R) \quad (R \times N) \\ \mathbf{Y} &\cong \mathbf{U} \mathbf{Q}^T = \sum_{r=1}^R \mathbf{u}_r \mathbf{q}_r^T \\ &(I \times M) \quad (I \times R) \quad (R \times M) \end{aligned}$$

**Figure 12:** The basic PLS model performs joint sequential low-rank approximation of the matrix of predictors  $\mathbf{X}$  and the matrix of responses  $\mathbf{Y}$ , so as to share (up to the scaling ambiguity) the latent components — columns of the *score matrices*  $\mathbf{T}$  and  $\mathbf{U}$ . The matrices  $\mathbf{P}$  and  $\mathbf{Q}$  are the *loading matrices* for predictors and responses, and  $\mathbf{E}$  and  $\mathbf{F}$  are the corresponding *residual matrices*.

*contraction* of the data matrix  $\mathbf{X}$  to the principal eigenvector *score matrix*  $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_R]$  of rank  $R$ ; (ii) ensures that the  $\mathbf{t}_r$  components are maximally correlated with the  $\mathbf{u}_r$  components in the approximation of the responses  $\mathbf{Y}$ , this is achieved when the  $\mathbf{u}_r$ 's are scaled versions of the  $\mathbf{t}_r$ 's. The  $Y$ -variables are then regressed on the matrix  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_R]$ . Therefore, PLS is a multivariate model with inferential ability that aims to find a representation of  $\mathbf{X}$  (or a part of  $\mathbf{X}$ ) that is relevant for predicting  $\mathbf{Y}$ , using the model

$$\mathbf{X} = \mathbf{T} \mathbf{P}^T + \mathbf{E} = \sum_{r=1}^R \mathbf{t}_r \mathbf{p}_r^T + \mathbf{E}, \quad (15)$$

$$\mathbf{Y} = \mathbf{U} \mathbf{Q}^T + \mathbf{F} = \sum_{r=1}^R \mathbf{u}_r \mathbf{q}_r^T + \mathbf{F}. \quad (16)$$

The score vectors  $\mathbf{t}_r$  provide an LS fit of  $\mathbf{X}$ -data, while at the same time the maximum correlation between  $\mathbf{t}$ - and  $\mathbf{u}$ -scores ensures a good predictive model for  $Y$ -variables. The predicted responses  $\mathbf{Y}_{new}$  are then obtained from new data  $\mathbf{X}_{new}$  and the loadings  $\mathbf{P}$  and  $\mathbf{Q}$ .

In practice, the score vectors  $\mathbf{t}_r$  are extracted sequentially, by a series of orthogonal projections followed by the deflation of  $\mathbf{X}$ . Since the rank of  $\mathbf{Y}$  is not necessarily decreased with each new  $\mathbf{t}_r$ , we may continue deflating until the rank of the  $\mathbf{X}$ -block is exhausted, so as to balance between prediction accuracy and model order.

The PLS concept can be generalized to tensors in the following ways:

- 1) *By unfolding multiway data.* For example



$\mathcal{X}(I \times J \times K)$  and  $\mathcal{Y}(I \times M \times N)$  can be flattened into long matrices  $\mathbf{X}(I \times JK)$  and  $\mathbf{Y}(I \times MN)$ , so as to admit matrix-PLS (see Figure 12). However, the flattening prior to standard *bilinear* PLS obscures structure in multiway data and compromises the interpretation of latent components.

- 2) *By low rank tensor approximation.* The so-called N-PLS attempts to find score vectors having maximal covariance with response variables, under the constraints that tensors  $\mathcal{X}$  and  $\mathcal{Y}$  are decomposed as a sum of rank-one tensors [104].
- 3) *By a BTD-type approximation,* as in the Higher Order PLS (HOPLS) model shown in Figure 13 [105]. The use of block terms within HOPLS equips it with additional flexibility, together with a more realistic analysis than unfolding-PLS and N-PLS.

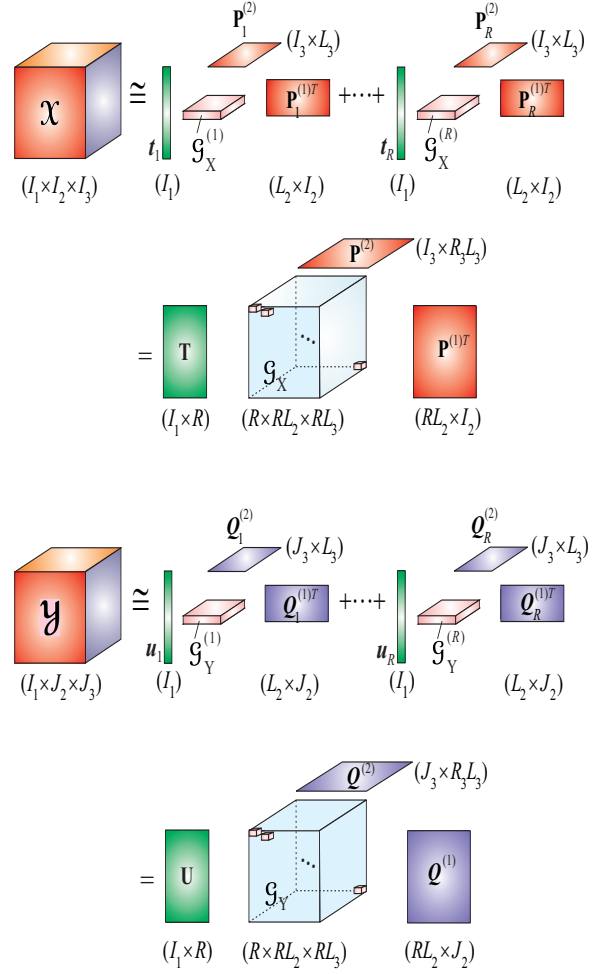
The principle of HOPLS can be formalized as a set of sequential approximate decompositions of the independent tensor  $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$  and the dependent tensor  $\mathcal{Y} \in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_M}$  (with  $I_1 = J_1$ ), so as to ensure maximum similarity (correlation) between the scores  $\mathbf{t}_r$  and  $\mathbf{u}_r$  within the loadings matrices  $\mathbf{T}$  and  $\mathbf{U}$ , based on

$$\mathcal{X} \cong \sum_{r=1}^R \mathcal{G}_X^{(r)} \times_1 \mathbf{t}_r \times_2 \mathbf{P}_r^{(1)} \dots \times_N \mathbf{P}_r^{(N-1)} \quad (17)$$

$$\mathcal{Y} \cong \sum_{r=1}^R \mathcal{G}_Y^{(r)} \times_1 \mathbf{u}_r \times_2 \mathbf{Q}_r^{(1)} \dots \times_N \mathbf{Q}_r^{(M-1)}. \quad (18)$$

A number of data-analytic problems can be reformulated as either regression or “similarity analysis” (ANOVA, ARMA, LDA, CCA), so that both the matrix and tensor PLS solutions can be generalized across exploratory data analysis.

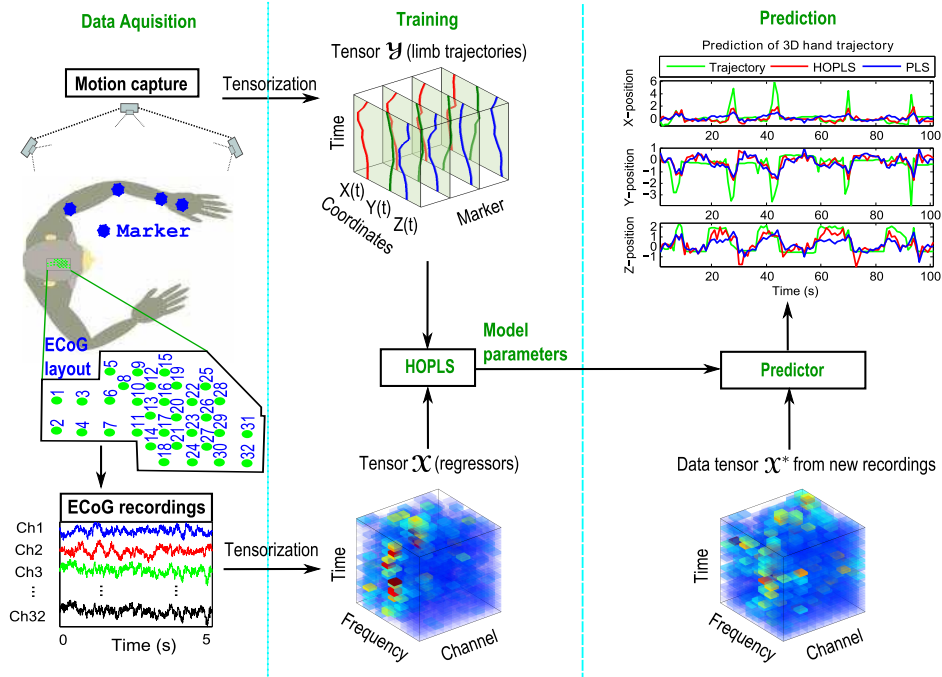
**Example 4: Decoding of a 3D hand movement trajectory from the electrocorticogram (ECoG).** The predictive power of tensor-based PLS is illustrated on a real-world example of the prediction of arm movement trajectory from ECoG. Fig. 14(left) illustrates the experimental setup, whereby 3D arm movement of a monkey was captured by an optical motion capture system with reflective markers affixed to the left shoulder, elbow, wrist, and hand; for full detail see (<http://neurotycho.org>). The predictors (32 ECoG channels) naturally build a fourth-order tensor  $\mathcal{X}$  (time  $\times$  channel\_no  $\times$  epoch\_length  $\times$  frequency) while the movement trajectories for the four markers (response) can be



**Figure 13:** The principle of Higher Order PLS (HOPLS) for third-order tensors. The core tensors  $\mathcal{G}_X$  and  $\mathcal{G}_Y$  are block-diagonal. The BTD-type structure allows for the modelling of general components that are highly correlated in the first mode.

represented as a third-order tensor  $\mathcal{Y}$  (time  $\times$  3D\_marker\_position  $\times$  marker\_no). The goal of the training stage is to identify the HOPLS parameters:  $\mathcal{G}_X^{(r)}, \mathcal{G}_Y^{(r)}, \mathbf{P}_r^{(n)}, \mathbf{Q}_r^{(n)}$ , see also Figure 13. In the test stage, the movement trajectories,  $\mathcal{Y}^*$ , for the new ECoG data,  $\mathcal{X}^*$ , are predicted through multilinear projections: (i) the new scores,  $\mathbf{t}_r^*$ , are found from new data,  $\mathcal{X}^*$ , and the existing model parameters:  $\mathcal{G}_X^{(r)}, \mathbf{P}_r^{(1)}, \mathbf{P}_r^{(2)}, \mathbf{P}_r^{(3)}$ , (ii) the predicted trajectory is calculated as  $\mathcal{Y}^* \approx \sum_{r=1}^R \mathcal{G}_Y^{(r)} \times_1 \mathbf{t}_r^* \times_2 \mathbf{Q}_r^{(1)} \times_3 \mathbf{Q}_r^{(2)} \times_4 \mathbf{Q}_r^{(3)}$ . In the simulations, standard PLS was applied in the same way to the unfolded tensors.

Figure 14(right) shows that although the standard PLS was able to predict the movement cor-



**Figure 14:** Prediction of arm movement from brain electrical responses. *Left:* Experiment setup. *Middle:* Construction of the data and response tensors and training. *Right:* The new data tensor (bottom) and the predicted 3D arm movement trajectories ( $X$ ,  $Y$ ,  $Z$  coordinates) obtained by tensor-based HOPLS and standard matrix-based PLS (top).

responding to each marker individually, such prediction is quite crude as the two-way PLS does not adequately account for mutual information among the four markers. The enhanced predictive performance of the BTD-based HOPLS (red line in Fig.14(right)) is therefore attributed to its ability to model interactions between complex latent components of both predictors and responses.

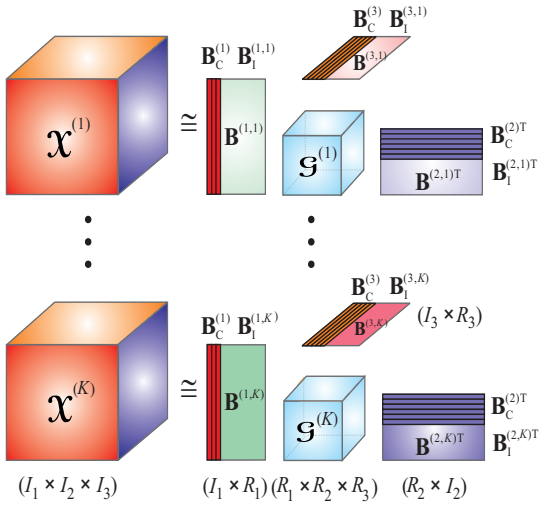
#### LINKED MULTIWAY COMPONENT ANALYSIS AND TENSOR DATA FUSION

Data fusion concerns joint analysis of an ensemble of data sets, such as multiple “views” of a particular phenomenon, where some parts of the “scene” may be visible in only one or a few data sets. Examples include fusion of visual and thermal images in low visibility conditions, or the analysis of human electrophysiological signals in response to a certain stimulus but from different subjects and trials; these are naturally analyzed together by means of matrix/tensor factorizations. The “coupled” nature of the analysis of multiple datasets ensures that there may be common factors across the datasets, and that some components are not shared (e.g., processes that are independent of excitations or stimuli/tasks).

The linked multiway component analysis (LMWCA) [106], shown in Figure 15, performs such decomposition into shared and individual factors, and is formulated as a set of approximate joint Tucker decompositions of a set of data tensors  $\mathcal{X}^{(k)} \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ , ( $k = 1, 2, \dots, K$ ):

$$\mathcal{X}^{(k)} \cong \mathcal{G}^{(k)} \times_1 \mathbf{B}^{(1,k)} \times_2 \mathbf{B}^{(2,k)} \dots \times_N \mathbf{B}^{(N,k)}, \quad (19)$$

where each factor matrix  $\mathbf{B}^{(n,k)} = [\mathbf{B}_C^{(n)}, \mathbf{B}_I^{(n,k)}] \in \mathbb{R}^{I_n \times R_n}$  has: (i) components  $\mathbf{B}_C^{(n)} \in \mathbb{R}^{I_n \times C_n}$  (with  $0 \leq C_n \leq R_n$ ) that are common (i.e., maximally correlated) to all tensors, and (ii) components  $\mathbf{B}_I^{(n,k)} \in \mathbb{R}^{I_n \times (R_n - C_n)}$  that are tensor-specific. The objective is to estimate the common components  $\mathbf{B}_C^{(n)}$ , the individual components  $\mathbf{B}_I^{(n,k)}$ , and, via the core tensors  $\mathcal{G}^{(k)}$ , their mutual interactions. As in MWCA (see Section **Tucker Decomposition**), constraints may be imposed to match data properties [73], [76]. This enables a more general and flexible framework than group ICA and Independent Vector Analysis, which also perform linked analysis of multiple data sets but assume that: (i) there exist only common components



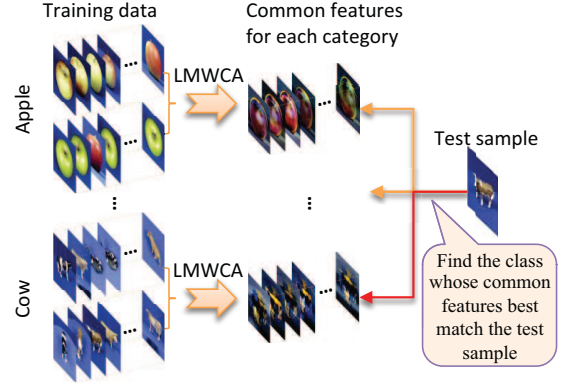
**Figure 15:** Coupled Tucker decomposition for linked multiway component analysis (LMWCA). The data tensors have both shared and individual components. Constraints such as orthogonality, statistical independence, sparsity and non-negativity may be imposed where appropriate.

and (ii) the corresponding latent variables are statistically independent [107], [108], both quite stringent and limiting assumptions. As an alternative to Tucker decompositions, coupled tensor decompositions may be of a polyadic or even block term type [89], [109].

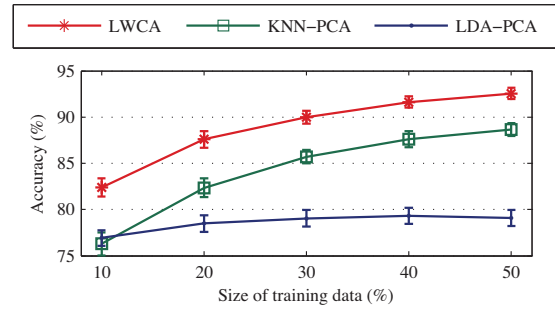
**Example 5: Feature extraction and classification of objects using LMWCA.** Classification based on common and distinct features of natural objects from the ETH-80 database (<http://www.d2.mpi-inf.mpg.de/Datasets>) was performed using LMWCA, whereby the discrimination among objects was performed using only the common features. This dataset consists of 3280 images in 8 categories, each containing 10 objects with 41 views per object. For each category, the training data were organized in two distinct fourth-order ( $128 \times 128 \times 3 \times I_4$ ) tensors, where  $I_4 = 10 \times 41 \times 0.5p$ , with  $p$  the fraction of training data. LMWCA was applied to these two tensors to find the common and individual features, with the number of common features set to 80% of  $I_4$ . In this way, eight sets of common features were obtained for each category. The test sample label was assigned to the category whose common features matched the new sample best (evaluated by canonical correlations) [110]. Figure 16 shows the results over 50 Monte Carlo runs and compares LMWCA with the standard



Sample images from different and same categories



Classification based on LMWCA



Performance comparison

**Figure 16:** Classification of color objects belonging to different categories. Due to using only common features, LMWCA achieves a high classification rate, even when the training set is small.

K-NN and LDA classifiers, the latter using 50 principal components as features. The enhanced classification results for LMWCA are attributed to the fact that the classification only makes use of the common components and is not hindered by components that are not shared across objects or views.

## SOFTWARE

The currently available software resources for tensor decompositions include:

- The Tensor Toolbox, a versatile framework for basic operations on sparse and dense tensors, including CPD and Tucker formats [111].
- The TDALAB and TENSORBOX, which provide a user-friendly interface and ad-

vanced algorithms for CPD, nonnegative Tucker decomposition and MWCA [112], [113].

- The Tensorlab toolbox builds upon the complex optimization framework and offers numerical algorithms for computing the CPD, BTD and Tucker decompositions. The toolbox includes a library of constraints (e.g. nonnegativity, orthogonality) and the possibility to combine and jointly factorize dense, sparse and incomplete tensors [89].
- The *N*-Way Toolbox, which includes (constrained) CPD, Tucker decomposition and PLS in the context of chemometrics applications [114]. Many of these methods can handle constraints (e.g., nonnegativity, orthogonality) and missing elements.
- The TT Toolbox, the Hierarchical Tucker Toolbox and the Tensor Calculus library provide tensor tools for scientific computing [115]–[117].
- Code developed for multiway analysis is also available from the Three-Mode Company [118].

#### CONCLUSIONS AND FUTURE DIRECTIONS

We live in a world overwhelmed by data, from multiple pictures of Big Ben on various social web links to terabytes of data in multiview medical imaging, while we may need to repeat the scientific experiments many times to obtain ground truth. Each snapshot gives us a somewhat incomplete view of the same object, and involves different angles, illumination, lighting conditions, facial expressions, and noise.

We have cast a light on tensor decompositions as a perfect match for exploratory analysis of such multifaceted data sets, and have illustrated their applications in multi-sensor and multi-modal signal processing. Our emphasis has been to show that tensor decompositions and multilinear algebra open completely new possibilities for component analysis, as compared with the “flat view” of standard two-way methods.

Unlike matrices, tensors are multiway arrays of data samples whose representations are typically overdetermined (fewer parameters in the decomposition than the number of data entries). This gives us an enormous flexibility in finding hidden components in data and the ability to enhance both robustness to noise and tolerance to missing data samples and faulty sensors.

We have also discussed multilinear variants of several standard signal processing tools such as multilinear SVD, ICA, NMF and PLS, and have shown that tensor methods can operate in a deterministic way on signals of very short duration.

At present the uniqueness conditions of standard tensor models are relatively well understood and efficient computation algorithms do exist, however, for future applications several challenging problems remain to be addressed in more depth:

- A whole new area emerges when several decompositions which operate on different datasets are coupled, as in multiview data where some details of interest are visible in only one mode. Such techniques need theoretical support in terms of existence, uniqueness, and numerical properties.
- As the complexity of advanced models increases, their computation requires efficient iterative algorithms, extending beyond the ALS class.
- Estimation of the number of components in data, and the assessment of their dimensionality would benefit from automation, especially in the presence of noise and outliers.
- Both new theory and algorithms are needed to further extend the flexibility of tensor models, e.g., for the constraints to be combined in many ways, and tailored to the particular signal properties in different modes.
- Work on efficient techniques for saving and/or fast processing of ultra large-scale tensors is urgent, these now routinely occupy tera-bytes, and will soon require peta-bytes of memory.
- Tools for rigorous performance analysis and rule of thumb performance bounds need to be further developed across tensor decomposition models.
- Our discussion has been limited to tensor models in which all entries take values independently of one another. Probabilistic versions of tensor decompositions incorporate prior knowledge about complex variable interaction, various data alphabets, or noise distributions, and so promise to model data more accurately and efficiently [119], [120].

It is fitting to conclude with a quote from Marcel Proust “*The voyage of discovery is not in*

seeking new landscapes but in having new eyes". We hope to have helped to bring to the eyes of the Signal Processing Community the multi-disciplinary developments in tensor decompositions, and to have shared our enthusiasm about tensors as powerful tools to discover new landscapes. The future computational, visualization and interpretation tools will be important next steps in supporting the different communities working on large-scale and big data analysis problems.

#### BIOGRAPHICAL NOTES

**Andrzej Cichocki** received the Ph.D. and Dr.Sc. (habilitation) degrees, all in electrical engineering, from the Warsaw University of Technology (Poland). He is currently a Senior Team Leader of the laboratory for Advanced Brain Signal Processing, at RIKEN Brain Science Institute (JAPAN) and Professor at Systems Research Institute, Polish Academy of Science (POLAND). He has authored of more than 400 publications and 4 monographs in the areas of signal processing and computational neuroscience. He serves as Associate Editor for the IEEE Transactions on Signal Processing and Journal Neuroscience Methods.

**Danilo P. Mandic** is a Professor of signal processing at Imperial College London, London, U.K. and has been working in the area of nonlinear and multidimensional adaptive signal processing and time-frequency analysis. His publication record includes two research monographs titled Recurrent Neural Networks for Prediction (West Sussex, U.K.: Wiley, August 2001) and Complex Valued Nonlinear Adaptive Filters: Noncircularity, Widely Linear and Neural Models, an edited book titled Signal Processing for Information Fusion, and more than 200 publications on signal and image processing.

**Anh Huy Phan** received the Ph.D. degree from the Kita Kyushu Institute of Technology, Japan in 2011. He worked as Deputy Head of the Research and Development Department, Broadcast Research and Application Center, Vietnam Television, and is currently a Research Scientist at the Laboratory for Advanced Brain Signal Processing, and a Visiting Research Scientist with the Toyota Collaboration Center, Brain Science Institute, RIKEN. He has served on the Editorial Board of International Journal of Computational Mathematics. His research interests include multilinear algebra, tensor computation, blind source separation, and brain computer interface.

**Cesar F. Caiafa** received the Ph.D. degree in engineering from the Faculty of Engineering, University of Buenos Aires, in 2007. He is currently Adjunct Researcher with the Argentinean Radioastronomy Institute (IAR) - CONICET and Assistant Professor with Faculty of Engineering, University of Buenos Aires. He is also Visiting Scientist at Lab. for Advanced Brain Signal Processing, BSI - RIKEN, Japan.

**Guoxu Zhou** received his Ph.D. degree in intelligent signal and information processing from South China University of Technology, Guangzhou, China, in 2010. He is currently a Research Scientist of the Laboratory for Advanced Brain Signal Processing, at RIKEN Brain Science Institute, Japan. His research interests include statistical signal processing, tensor analysis, intelligent information processing, and machine learning.

**Qibin Zhao** received his Ph.D. degree from the Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China, in 2009. He is currently a research scientist at the Laboratory for Advanced Brain Signal Processing in RIKEN Brain Science Institute, Japan and a visiting research scientist in BSI TOYOTA Collaboration

Center, RIKEN-BSI. His research interests include multiway data analysis, brain computer interface and machine learning.

**Lieven De Lathauwer** received the Ph.D. degree from the Faculty of Engineering, KU Leuven, Belgium, in 1997. From 2000 to 2007 he was Research Associate with the Centre National de la Recherche Scientifique, France. He is currently Professor with KU Leuven. He is affiliated with both the Group Science, Engineering and Technology of Kulak, with the Stadius Center for Dynamical Systems, Signal Processing and Data Analytics of the Electrical Engineering Department (ESAT) and with iMinds Future Health Department. He is Associate Editor of the SIAM Journal on Matrix Analysis and Applications and has served as Associate Editor for the IEEE Transactions on Signal Processing. His research concerns the development of tensor tools for engineering applications.

#### REFERENCES

- [1] F. L. Hitchcock, "Multiple invariants and generalized rank of a p-way matrix or tensor," *Journal of Mathematics and Physics*, vol. 7, pp. 39–79, 1927.
- [2] R. Cattell, "Parallel proportional profiles and other principles for determining the choice of factors by rotation," *Psychometrika*, vol. 9, pp. 267–283, 1944.
- [3] L. R. Tucker, "The extension of factor analysis to three-dimensional matrices," in *Contributions to Mathematical Psychology*, H. Gulliksen and N. Frederiksen, Eds. New York: Holt, Rinehart and Winston, 1964, pp. 110–127.
- [4] —, "Some mathematical notes on three-mode factor analysis," *Psychometrika*, vol. 31, no. 3, pp. 279–311, September 1966.
- [5] J. Carroll and J.-J. Chang, "Analysis of individual differences in multidimensional scaling via an  $n$ -way generalization of 'Eckart-Young' decomposition," *Psychometrika*, vol. 35, no. 3, pp. 283–319, September 1970.
- [6] R. A. Harshman, "Foundations of the PARAFAC procedure: Models and conditions for an explanatory multimodal factor analysis," *UCLA Working Papers in Phonetics*, vol. 16, pp. 1–84, 1970.
- [7] A. Smilde, R. Bro, and P. Geladi, *Multi-way Analysis: Applications in the Chemical Sciences*. New York: John Wiley & Sons Ltd, 2004.
- [8] P. Kroonenberg, *Applied Multiway Data Analysis*. New York: John Wiley & Sons Ltd, 2008.
- [9] C. Nikias and A. Petropulu, *Higher-Order Spectra Analysis: A Nonlinear Signal Processing Framework*. Prentice Hall, 1993.
- [10] J.-F. Cardoso and A. Souloumiac, "Blind beamforming for non-Gaussian signals," in *IEE Proceedings F (Radar and Signal Processing)*, vol. 140, no. 6. IET, 1993, pp. 362–370.
- [11] P. Comon, "Independent component analysis, a new concept?" *Signal Processing*, vol. 36, no. 3, pp. 287–314, 1994.
- [12] P. Comon and C. Jutten, Eds., *Handbook of Blind Source Separation: Independent Component Analysis and Applications*. Academic Press, 2010.
- [13] L. De Lathauwer, B. De Moor, and J. Vandewalle, "A multilinear singular value decomposition," *SIAM Journal of Matrix Analysis and Applications*, vol. 24, pp. 1253–1278, 2000.
- [14] G. Beylkin and M. Mohlenkamp, "Algorithms for numerical analysis in high dimensions," *SIAM J. Scientific Computing*, vol. 26, no. 6, pp. 2133–2159, 2005.
- [15] J. Ballani, L. Grasedyck, and M. Kluge, "Black box approximation of tensors in hierarchical Tucker format," *Linear Algebra and its Applications*, vol. 438, no. 2, pp. 639–657, 2013.

- [16] I. V. Oseledets, "Tensor-train decomposition," *SIAM J. Scientific Computing*, vol. 33, no. 5, pp. 2295–2317, 2011.
- [17] N. Sidiropoulos, R. Bro, and G. Giannakis, "Parallel factor analysis in sensor array processing," *IEEE Transactions on Signal Processing*, vol. 48, no. 8, pp. 2377–2388, 2000.
- [18] N. Sidiropoulos, G. Giannakis, and R. Bro, "Blind PARAFAC receivers for DS-CDMA systems," *IEEE Transactions on Signal Processing*, vol. 48, no. 3, pp. 810–823, 2000.
- [19] A. Cichocki, R. Zdunek, A.-H. Phan, and S. Amari, *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-way Data Analysis and Blind Source Separation*. Chichester: Wiley, 2009.
- [20] J. Landsberg, *Tensors: Geometry and Applications*. AMS, 2012.
- [21] W. Hackbusch, *Tensor Spaces and Numerical Tensor Calculus*, ser. Springer series in computational mathematics. Heidelberg: Springer, 2012, vol. 42.
- [22] E. Acar and B. Yener, "Unsupervised multiway data analysis: A literature survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 21, pp. 6–20, 2009.
- [23] T. Kolda and B. Bader, "Tensor decompositions and applications," *SIAM Review*, vol. 51, no. 3, pp. 455–500, September 2009.
- [24] P. Comon, X. Luciani, and A. L. F. de Almeida, "Tensor decompositions, Alternating Least Squares and other Tales," *Jour. Chemometrics*, vol. 23, pp. 393–405, 2009.
- [25] H. Lu, K. Plataniotis, and A. Venetsanopoulos, "A survey of multilinear subspace learning for tensor data," *Pattern Recognition*, vol. 44, no. 7, pp. 1540–1551, 2011.
- [26] M. Mørup, "Applications of tensor (multiway array) factorizations and decompositions in data mining," *Wiley Interdisc. Rev.: Data Mining and Knowledge Discovery*, vol. 1, no. 1, pp. 24–40, 2011.
- [27] B. Khoromskij, "Tensors-structured numerical methods in scientific computing: Survey on recent advances," *Chemometrics and Intelligent Laboratory Systems*, vol. 110, no. 1, pp. 1–19, 2011.
- [28] L. Grasedyck, D. Kessner, and C. Tobler, "A literature survey of low-rank tensor approximation techniques," *CGAMM-Mitteilungen*, vol. 36, pp. 53–78, 2013.
- [29] P. Comon, "Tensors: A brief survey," *IEEE Signal Processing Magazine*, p. (accepted), 2014.
- [30] A. Bruckstein, D. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *SIAM Review*, vol. 51, no. 1, pp. 34–81, 2009.
- [31] J. Kruskal, "Three-way arrays: rank and uniqueness of trilinear decompositions, with application to arithmetic complexity and statistics," *Linear Algebra and its Applications*, vol. 18, no. 2, pp. 95 – 138, 1977.
- [32] I. Domanov and L. De Lathauwer, "On the uniqueness of the canonical polyadic decomposition of third-order tensors — part i: Basic results and uniqueness of one factor matrix and part ii: Uniqueness of the overall decomposition," *SIAM J. Matrix Anal. Appl.*, vol. 34, no. 3, pp. 855–903, 2013.
- [33] A. Cichocki and S. Amari, *Adaptive Blind Signal and Image Processing*. John Wiley, Chichester, 2003.
- [34] A. Hyvärinen, "Independent component analysis: recent advances," *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, vol. 371, no. 1984, 2013.
- [35] M. Elad, P. Milanfar, and G. H. Golub, "Shape from moments – an estimation theory perspective," *Signal Processing, IEEE Transactions on*, vol. 52, no. 7, pp. 1814–1829, 2004.
- [36] N. Sidiropoulos, "Generalizing Caratheodory's uniqueness of harmonic parameterization to N dimensions," *IEEE Trans. Information Theory*, vol. 47, no. 4, pp. 1687–1690, 2001.
- [37] A. Belouchrani, K. Abed-Meraim, J.-F. Cardoso, and É. Moulines, "A blind source separation technique using second-order statistics," *IEEE Trans. Signal Processing*, vol. 45, no. 2, pp. 434–444, 1997.
- [38] F. Miwakeichi, E. Martinez-Montes, P. Valds-Sosa, N. Nishiyama, H. Mizuhara, and Y. Yamaguchi, "Decomposing EEG data into space–time–frequency components using parallel factor analysis," *NeuroImage*, vol. 22, no. 3, pp. 1035–1045, 2004.
- [39] M. Vasilescu and D. Terzopoulos, "Multilinear analysis of image ensembles: Tensorfaces," in *Proc. European Conf. on Computer Vision (ECCV)*, vol. 2350, Copenhagen, Denmark, May 2002, pp. 447–460.
- [40] M. Hirsch, D. Lanman, G. Wetzstein, and R. Raskar, "Tensor displays," in *Int. Conf. on Computer Graphics and Interactive Techniques, SIGGRAPH 2012, Los Angeles, CA, USA, Aug. 5-9, 2012, Emerging Technologies Proceedings*, 2012, pp. 24–42.
- [41] J. Hastad, "Tensor rank is NP-complete," *Journal of Algorithms*, vol. 11, no. 4, pp. 644–654, 1990.
- [42] M. Timmerman and H. Kiers, "Three mode principal components analysis: Choosing the numbers of components and sensitivity to local optima," *British Journal of Mathematical and Statistical Psychology*, vol. 53, no. 1, pp. 1–16, 2000.
- [43] E. Ceulemans and H. Kiers, "Selecting among three-mode principal component models of different types and complexities: A numerical convex-hull based method," *British Journal of Mathematical and Statistical Psychology*, vol. 59, no. 1, pp. 133–150, May 2006.
- [44] M. Mørup and L. K. Hansen, "Automatic relevance determination for multiway models," *Journal of Chemometrics, Special Issue: In Honor of Professor Richard A. Harshman*, vol. 23, no. 7-8, pp. 352 – 363, 2009. [Online]. Available: <http://www2.imm.dtu.dk/pubdb/p.php?25806>
- [45] N. Sidiropoulos and R. Bro, "On the uniqueness of multilinear decomposition of N-way arrays," *J. Chemometrics*, vol. 14, no. 3, pp. 229–239, 2000.
- [46] T. Jiang and N. D. Sidiropoulos, "Kruskal's permutation lemma and the identification of CANDECOMP/PARAFAC and bilinear models," *IEEE Trans. Signal Processing*, vol. 52, no. 9, pp. 2625–2636, 2004.
- [47] L. De Lathauwer, "A link between the canonical decomposition in multilinear algebra and simultaneous matrix diagonalization," *SIAM J. Matrix Analysis Applications*, vol. 28, no. 3, pp. 642–666, 2006.
- [48] A. Stegeman, "On uniqueness conditions for Candecom/Parafac and Indscal with full column rank in one mode," *Linear Algebra and its Applications*, vol. 431, no. 1–2, pp. 211–227, 2009.
- [49] E. Sanchez and B. Kowalski, "Tensorial resolution: a direct trilinear decomposition," *J. Chemometrics*, vol. 4, pp. 29–45, 1990.
- [50] I. Domanov and L. De Lathauwer, "Canonical polyadic decomposition of third-order tensors: Reduction to generalized eigenvalue decomposition," ESAT, KU Leuven, ESAT-SISTA Internal Report 13-36, 2013.
- [51] S. Vorobyov, Y. Rong, N. Sidiropoulos, and A. Gershman, "Robust iterative fitting of multilinear models," *IEEE Transactions Signal Processing*, vol. 53, no. 8, pp. 2678–2689, 2005.
- [52] X. Liu and N. Sidiropoulos, "Cramer-Rao lower bounds for low-rank decomposition of multidimen-



- sional arrays," *IEEE Trans. on Signal Processing*, vol. 49, no. 9, pp. 2074–2086, Sep. 2001.
- [53] P. Tichavsky, A. Phan, and Z. Koldovsky, "Cramér-rao-induced bounds for candecomp/parafac tensor decomposition," *IEEE Transactions on Signal Processing*, vol. 61, no. 8, pp. 1986–1997, 2013.
- [54] B. Chen, S. He, Z. Li, and S. Zhang, "Maximum block improvement and polynomial optimization," *SIAM Journal on Optimization*, vol. 22, no. 1, pp. 87–107, 2012.
- [55] A. Uschmajew, "Local convergence of the alternating least squares algorithm for canonical tensor approximation," *SIAM J. Matrix Anal. Appl.*, vol. 33, no. 2, pp. 639–652, 2012.
- [56] M. J. Mohlenkamp, "Musings on multilinear fitting," *Linear Algebra and its Applications*, vol. 438, no. 2, pp. 834–852, 2013.
- [57] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 1126–1153, 2013.
- [58] P. Paatero, "The multilinear engine: A table-driven least squares program for solving multilinear problems, including the n-way parallel factor analysis model," *Journal of Computational and Graphical Statistics*, vol. 8, no. 4, pp. 854–888, Dec. 1999.
- [59] E. Acar, D. Dunlavy, T. Kolda, and M. Mørup, "Scalable tensor factorizations for incomplete data," *Chemometrics and Intelligent Laboratory Systems*, vol. 106 (1), pp. 41–56, 2011. [Online]. Available: <http://www2.imm.dtu.dk/pubdb/p.php?5923>
- [60] A.-H. Phan, P. Tichavsky, and A. Cichocki, "Low complexity Damped Gauss-Newton algorithms for CANDECOMP/PARAFAC," *SIAM Journal on Matrix Analysis and Applications (SIMAX)*, vol. 34, no. 1, pp. 126–147, 2013.
- [61] 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. Optimization*, vol. 23, no. 2, 2013.
- [62] V. de Silva and L.-H. Lim, "Tensor rank and the ill-posedness of the best low-rank approximation problem," *SIAM J. Matrix Anal. Appl.*, vol. 30, pp. 1084–1127, September 2008.
- [63] W. Krijnen, T. Dijkstra, and A. Stegeman, "On the non-existence of optimal solutions and the occurrence of "degeneracy" in the Candecomp/Parafac model," *Psychometrika*, vol. 73, pp. 431–439, 2008.
- [64] M. Sørensen, L. De Lathauwer, P. Comon, S. Icart, and L. Deneire, "Canonical Polyadic Decomposition with orthogonality constraints," *SIAM J. Matrix Anal. Appl.*, vol. 33, no. 4, pp. 1190–1213, 2012.
- [65] M. Sørensen and L. De Lathauwer, "Blind signal separation via tensor decomposition with Vandermonde factor: Canonical polyadic decomposition," *IEEE Trans. Signal Processing*, vol. 61, no. 22, pp. 5507–5519, Nov. 2013.
- [66] G. Zhou and A. Cichocki, "Canonical Polyadic Decomposition based on a single mode blind source separation," *IEEE Signal Processing Letters*, vol. 19, no. 8, pp. 523–526, 2012.
- [67] L.-H. Lim and P. Comon, "Nonnegative approximations of nonnegative tensors," *Journal of Chemometrics*, vol. 23, no. 7–8, pp. 432–441, 2009.
- [68] A. van der Veen and A. Paulraj, "An analytical constant modulus algorithm," *IEEE Transactions Signal Processing*, vol. 44, pp. 1136–1155, 1996.
- [69] R. Roy and T. Kailath, "Esprit-estimation of signal parameters via rotational invariance techniques," *Acoustics, Speech and Signal Processing, IEEE Transactions on*, vol. 37, no. 7, pp. 984–995, 1989.
- [70] 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 Journal of Matrix Analysis and Applications*, vol. 21, no. 4, pp. 1324–1342, 2000.
- [71] B. Savas and L.-H. Lim, "Quasi-Newton methods on Grassmannians and multilinear approximations of tensors," *SIAM J. Scientific Computing*, vol. 32, no. 6, pp. 3352–3393, 2010.
- [72] 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 Analysis Applications*, vol. 32, no. 1, pp. 115–135, 2011.
- [73] G. Zhou and A. Cichocki, "Fast and unique Tucker decompositions via multiway blind source separation," *Bulletin of Polish Academy of Science*, vol. 60, no. 3, pp. 389–407, 2012.
- [74] A. Cichocki, "Generalized Component Analysis and Blind Source Separation Methods for Analyzing Multichannel Brain Signals," in *Statistical and Process Models for Cognitive Neuroscience and Aging*. Lawrence Erlbaum Associates, 2007, pp. 201–272.
- [75] M. Haardt, F. Roemer, and G. D. Galdo, "Higher-order SVD based subspace estimation to improve the parameter estimation accuracy in multi-dimensional harmonic retrieval problems," *IEEE Trans. Signal Processing*, vol. 56, pp. 3198 – 3213, Jul. 2008.
- [76] A. Phan and A. Cichocki, "Tensor decompositions for feature extraction and classification of high dimensional datasets," *Nonlinear Theory and Its Applications, IEICE*, vol. 1, no. 1, pp. 37–68, 2010.
- [77] L. De Lathauwer, "Decompositions of a higher-order tensor in block terms – Part I and II," *SIAM Journal on Matrix Analysis and Applications (SIMAX)*, vol. 30, no. 3, pp. 1022–1066, 2008, special Issue on Tensor Decompositions and Applications. [Online]. Available: <http://publi-etis.ensea.fr/2008/De08e>
- [78] 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 Analysis Applications*, vol. 32, no. 4, pp. 1451–1474, 2011.
- [79] L. De Lathauwer, "Block component analysis, a new concept for blind source separation," in *Proc. 10th International Conf. LVA/ICA, Tel Aviv, March 12-15,, 2012*, pp. 1–8.
- [80] E. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: exact signal reconstruction from highly incomplete frequency information," *IEEE Transactions on Information Theory*, vol. 52, no. 2, pp. 489–509, 2006.
- [81] E. J. Candes and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?" *Information Theory, IEEE Transactions on*, vol. 52, no. 12, pp. 5406–5425, 2006.
- [82] D. L. Donoho, "Compressed sensing," *Information Theory, IEEE Transactions on*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [83] Y. Eldar and G. Kutyniok, "Compressed Sensing: Theory and Applications," *New York: Cambridge Univ. Press*, vol. 20, p. 12, 2012.
- [84] M. F. Duarte and R. G. Baraniuk, "Kronecker compressive sensing," *IEEE Transactions on Image Processing*, vol. 21, no. 2, pp. 494–504, 2012.
- [85] C. Caiafa and A. Cichocki, "Computing sparse representations of multidimensional signals using Kronecker bases," *Neural Computaion*, vol. 25, no. 1, pp. 186–220, 2013.
- [86] —, "Multidimensional compressed sensing and their applications," *WIRES Data Mining and Knowledge Discovery*, 2013 (accepted).



- [87] S. Gandy, B. Recht, and I. Yamada, "Tensor completion and low-n-rank tensor recovery via convex optimization," *Inverse Problems*, vol. 27, no. 2, 2011.
- [88] M. Signoretto, Q. T. Dinh, L. De Lathauwer, and J. A. Suykens, "Learning with tensors: a framework based on convex optimization and spectral regularization," *Machine Learning*, pp. 1–49, 2013.
- [89] L. Sorber, M. Van Barel, and L. De Lathauwer, "Tensorlab v1.0," Feb. 2013. [Online]. Available: <http://esat.kuleuven.be/sista/tensorlab/>
- [90] N. Sidiropoulos and A. Kyrillidis, "Multi-way compressed sensing for sparse low-rank tensors," *IEEE Signal Processing Letters*, vol. 19, no. 11, pp. 757–760, 2012.
- [91] D. Foster, K. Amano, S. Nascimento, and M. Foster, "Frequency of metamerism in natural scenes," *Journal of the Optical Society of America A*, vol. 23, no. 10, pp. 2359–2372, 2006.
- [92] A. Cichocki, "Era of big data processing: A new approach via tensor networks and tensor decompositions, (invited talk)," in *Proc. of Int. Workshop on Smart Info-Media Systems in Asia (SISA2013), Nagoya, Japan, Sept.30–Oct.2*, 2013.
- [93] R. Orus, "A Practical Introduction to Tensor Networks: Matrix Product States and Projected Entangled Pair States," *The Journal of Chemical Physics*, 2013.
- [94] J. Salmi, A. Richter, and V. Koivunen, "Sequential unfolding SVD for tensors with applications in array signal processing," *IEEE Transactions on Signal Processing*, vol. 57, pp. 4719–4733, 2009.
- [95] A.-H. Phan and A. Cichocki, "PARAFAC algorithms for large-scale problems," *Neurocomputing*, vol. 74, no. 11, pp. 1970–1984, 2011.
- [96] S. K. Suter, M. Makhynia, and R. Pajarola, "Tamresh - tensor approximation multiresolution hierarchy for interactive volume visualization," *Comput. Graph. Forum*, vol. 32, no. 3, pp. 151–160, 2013.
- [97] D. Nion and N. Sidiropoulos, "Adaptive algorithms to track the PARAFAC decomposition of a third-order tensor," *IEEE Trans. on Signal Processing*, vol. 57, no. 6, pp. 2299–2310, Jun. 2009.
- [98] S. A. Goreinov, N. L. Zamarashkin, and E. E. Tyrtyshnikov, "Pseudo-skeleton approximations by matrices of maximum volume," *Mathematical Notes*, vol. 62, no. 4, pp. 515–519, 1997.
- [99] C. Caiafa and A. Cichocki, "Generalizing the column-row matrix decomposition to multi-way arrays," *Linear Algebra and its Applications*, vol. 433, no. 3, pp. 557–573, 2010.
- [100] S. A. Goreinov, "On cross approximation of multi-index array," *Doklady Math.*, vol. 420, no. 4, pp. 404–406, 2008.
- [101] I. Oseledets, D. V. Savostyanov, and E. Tyrtyshnikov, "Tucker dimensionality reduction of three-dimensional arrays in linear time," *SIAM J. Matrix Analysis Applications*, vol. 30, no. 3, pp. 939–956, 2008.
- [102] I. Oseledets and E. Tyrtyshnikov, "TT-cross approximation for multidimensional arrays," *Linear Algebra and its Applications*, vol. 432, no. 1, pp. 70–88, 2010.
- [103] M. W. Mahoney, M. Maggioni, and P. Drineas, "Tensor-CUR decompositions for tensor-based data," *SIAM Journal on Matrix Analysis and Applications*, vol. 30, no. 3, pp. 957–987, 2008.
- [104] R. Bro, "Multiway calibration. Multilinear PLS," *Journal of Chemometrics*, vol. 10, pp. 47–61, 1996.
- [105] Q. Zhao, C. Caiafa, D. Mandic, Z. Chao, Y. Nagasaka, N. Fujii, L. Zhang, and A. Cichocki, "Higher-order partial least squares (HOPLS): A generalized multilinear regression method," *IEEE Trans on Pattern Analysis and machine Intelligence (PAMI)*, vol. 35, no. 7, pp. 1660–1673, 2013.
- [106] A. Cichocki, "Tensors decompositions: New concepts for brain data analysis?" *Journal of Control, Measurement, and System Integration (SICE)*, vol. 47, no. 7, pp. 507–517, 2011.
- [107] V. Calhoun, J. Liu, and T. Adali, "A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data," *Neuroimage*, vol. 45, pp. 163–172, 2009.
- [108] Y.-O. Li, T. Adali, W. Wang, and V. Calhoun, "Joint blind source separation by multiset canonical correlation analysis," *IEEE Transactions on Signal Processing*, vol. 57, no. 10, pp. 3918–3929, oct. 2009.
- [109] E. Acar, T. Kolda, and D. Dunlavy, "All-at-once optimization for coupled matrix and tensor factorizations," *CoRR*, vol. abs/1105.3422, 2011.
- [110] G. Zhou, A. Cichocki, S. Xie, and D. Mandic, "Beyond Canonical Correlation Analysis: Common and individual features analysis," *IEEE Transactions on PAMI*, 2013. [Online]. Available: <http://arxiv.org/abs/1212.3913>, 2012
- [111] B. Bader, T. G. Kolda *et al.*, "MATLAB tensor toolbox version 2.5," Available online, Feb. 2012. [Online]. Available: <http://www.sandia.gov/~tgkolda/TensorToolbox>
- [112] G. Zhou and A. Cichocki, "TDALAB: Tensor Decomposition Laboratory," LABSP, Wako-shi, Japan, 2013. [Online]. Available: <http://bsp.brain.riken.jp/TDALAB/>
- [113] A.-H. Phan, P. Tichavský, and A. Cichocki, "Tensorbox: a matlab package for tensor decomposition," LABSP, RIKEN, Japan, 2012. [Online]. Available: <http://www.bsp.brain.riken.jp/~phan/tensorbox.php>
- [114] C. Andersson and R. Bro, "The N-way toolbox for MATLAB," *Chemometrics Intell. Lab. Systems*, vol. 52, no. 1, pp. 1–4, 2000. [Online]. Available: <http://www.models.life.ku.dk/nwaytoolbox>
- [115] I. Oseledets, "TT-toolbox 2.2," 2012. [Online]. Available: <https://github.com/oseledets/TT-Toolbox>
- [116] D. Kressner and C. Tobler, "htucker—A MATLAB toolbox for tensors in hierarchical Tucker format," *MATHICSE, EPF Lausanne*, 2012. [Online]. Available: <http://anchp.epfl.ch/htucker>
- [117] M. Espig, M. Schuster, A. Killaitis, N. Waldren, P. Wähnert, S. Handschuh, and H. Auer, "Tensor Calculus library," 2012. [Online]. Available: <http://gitorious.org/tensorcalculus>
- [118] P. Kroonenberg, "The three-mode company. A company devoted to creating three-mode software and promoting three-mode data analysis." [Online]. Available: <http://three-mode.leidenuniv.nl/>
- [119] S. Zhe, Y. Qi, Y. Park, I. Molloy, and S. Chari, "DinTucker: Scaling up Gaussian process models on multidimensional arrays with billions of elements," *PAMI (in print) arXiv preprint arXiv:1311.2663*, 2014.
- [120] K. Yilmaz and A. T. Cemgil, "Probabilistic latent tensor factorisation," in *Proc. of International Conference on Latent Variable analysis and Signal Separation*, vol. 6365, 2010, pp. 346–353, cPCI-S.