

Tracking Online Low-Rank Approximations of Higher-Order Incomplete Streaming Tensors

Le Trung Thanh^{1,2}, Karim Abed-Meraim¹, Nguyen Linh
Trung^{2*} and Adel Hafiane¹

¹University of Orléans, INSA CVL, PRISME, France.

²VNU University of Engineering and Technology, Vietnam.

*Corresponding author(s). E-mail(s): linhtrung@vnu.edu.vn;

Contributing authors: trung-thanh.le@univ-orleans.fr;
karim-abed.meraim@univ-orleans.fr; adel.hafiane@insa-cvl.fr;

Abstract

In this paper, we propose two new provable algorithms for tracking online low-rank approximations of high-order streaming tensors with missing data. The first algorithm, dubbed adaptive Tucker decomposition (ATD), minimizes a weighted recursive least-squares cost function to obtain the tensor factors and the core tensor in an efficient way, thanks to the alternating minimization framework and the randomized sketching technique. Under the Canonical Polyadic (CP) model, the second algorithm called ACP is developed as a variant of ATD when the core tensor is imposed to be identity. Both algorithms are low-complexity tensor trackers that have fast convergence and low memory storage requirements. A unified convergence analysis is presented for ATD and ACP to justify their performance. Experiments indicate that the two proposed algorithms are capable of streaming tensor decomposition with competitive performance with respect to estimation accuracy and runtime on both synthetic and real data.

Keywords: Tensor decomposition; low-rank approximation; adaptive algorithm; streaming tensor; missing data; randomized method.

1 Introduction

The era of “Big Data”, which deals with massive datasets, has brought new analysis techniques for discovering new valuable information hidden in the data [1]. Among these techniques is multilinear low-rank approximation (LRA) of matrices and tensors, which has recently attracted considerable attention from engineers and researchers in the signal processing and machine learning communities [2]. A tensor is a multidimensional array and provides a natural representation of high-dimensional data. Low-rank approximation of tensors (t-LRA) can be considered as a multiway extension of LRA of matrices (which are two-way) to higher dimensions [3]. Generally, t-LRA is referred to as tensor decomposition which factorizes a tensor into a sequence of basic components [3]. As a result, t-LRA provides a useful tool for dealing with several large-scale multidimensional problems in modern data analysis that would otherwise be intractable by classical methods.

Two widely-used approaches for t-LRA are CP/PARAFAC decomposition¹ [4] and Tucker decomposition [5]. Under the CP format, a tensor can be represented as a sum of rank-1 tensors; each rank-1 tensor is formulated as the outer product of vectors. Under the Tucker format, a tensor is factorized into a sequence of factor matrices acting on a reduced-size core tensor. “Workhorse” algorithms are based on the method of alternating least-squares (ALS). Readers are referred to the work of Kolda and Bader [3] for a good review.

The characteristics of “Big Data” are often associated with the following three “V”s: volume, velocity and veracity [1]. Velocity and veracity are the focus of this paper. Velocity requires (near) real-time processing of data streams, while veracity demands robust algorithms to better deal with missing, noisy and inconsistent data. In online applications, data acquisition is often a time-varying process in which data are serially collected or changing with time. Besides, missing data are ubiquitous and more and more common in high-dimensional problems in which collecting all attributes of data is either too expensive or even impossible. However, well-known t-LRA algorithms either face high complexity or operate in batch mode, and thus, may not be suitable for such problems. This has led to defining a variant of t-LRA, namely online (adaptive) t-LRA.

In the literature, there are several studies related to the problem of tracking online t-LRA in the missing data context; the tensors are said to be both *streaming* and *incomplete*. Under the CP format, the very first adaptive tensor algorithms were proposed by Nion *et al.* in [6] more than 10 years ago. Since then, several adaptive CP decomposition methods have been introduced. We refer the reader to ref. [7] for a good survey. Among them, Mardani *et al.* proposed TeCPSGD [8], which is a first-order algorithm and uses the method of stochastic gradient descent (SGD) to track the CP decomposition of 3rd-order

¹In the literature, there exist some other names for the CP decomposition, such as PARAFAC (Parallel Factors), CPD (Canonical Polyadic Decomposition), and CANDECOMP or CAND (Canonical Decomposition).

streaming tensors with missing data. Leveraging the framework of alternating minimization, TeCPSGD can estimate directly all factors except the one corresponding to the dimension growing over time in an efficient way. Because of SGD, TeCPSGD is a low-complexity tensor tracker but with a slow convergence rate. Therefore, it is not really suitable for fast time-varying scenarios in which a class of methods with a fast rate of convergence is preferable. In [9], Kasai developed OLSTEC which is an efficient second-order algorithm and exploits the recursive least-squares technique. OLSTEC provides competitive performance in terms of estimation accuracy, but its computational complexity is much higher than that of TeCPSGD. In parallel, Chinh *et al.* proposed to first track the low-dimensional tensor subspace and then derive the loading factors from its Khatri-Rao structure [10]. Its performance, however, is sensitive to initialization, see Fig. 6 for an illustration.

All the adaptive CP decomposition algorithms above are specifically designed for factorizing 3rd-order streaming tensors (i.e., the temporal tensor slices aka data observations are matrices) which could limit their applications in practice where the underlying tensor is of higher order (i.e., greater than or equal to 4). One possible way is to reshape the underlying higher-order streaming tensor into a 3rd-order one and then apply the abovementioned algorithms for tracking. However, the high-dimensional structure of the original tensor might not be fully preserved in its reshaped variant, resulting in low estimation accuracy. Dealing with N th-order streaming tensors ($N > 3$) is nontrivial due to several issues. Some mathematical tools, transformations, and operations applied for 3rd-order streaming tensors are not straightforward for higher-order ones, such as the low-rank regularization, the tensor subspace/dictionary (used in [8, 9]) and the bi-iteration SVD procedure (used in [6, 10]). Particularly in [8, 9], the nuclear norm of a matrix \mathbf{Z} can be derived from $\min_{\mathbf{Z}=\mathbf{U}\mathbf{V}^\top} \frac{1}{2}(\|\mathbf{U}\|_F^2 + \|\mathbf{V}\|_F^2)$, thanks to Lemma 5.1 in [11]. This property is widely used by several matrix factorization methods for low-rank regularization [12]. Accordingly, the sum of squared Frobenius norms of two non-temporal factors of 3rd-order streaming tensors can be used as a regularization promoting the low-rank approximation of data streams. However, this property does not generally hold for higher-order tensors due to the presence of more than two factors and their multilinear connection. In [6, 10], the subspace-based algorithms track the underlying low-dimensional tensor subspace matrix $\mathbf{H}_t = \mathbf{U}_t^{(1)} \odot \mathbf{U}_t^{(2)} \odot \dots \odot \mathbf{U}_t^{(N-1)}$ where N is the tensor order, and then, estimate the tensor factors $\{\mathbf{U}_t^{(n)}\}_{n=1}^{N-1}$ by exploiting its Khatri-Rao structure. When $N = 3$, subspace tracking algorithms and bi-iteration SVD are specifically applied to estimate \mathbf{H}_t and $\{\mathbf{U}_t^{(1)}, \mathbf{U}_t^{(2)}\}$. However when $N > 3$, it becomes more complicated due to two main issues: (i) tracking the matrix \mathbf{H}_t having a “hierarchical” Khatri-Rao structure over time is nontrivial especially in noisy and time-varying environments and (ii) even when \mathbf{H}_t is assumed to be estimated correctly at each time t , the estimation of $\{\mathbf{U}_t^{(n)}\}_{n=1}^{N-1}$ might cost a high computational complexity, e.g., if bi-iteration SVD is used, we have to repeat

bi-iteration SVD recursively $N - 1$ times which is very expensive for streaming processing. These characteristics prevent us extending their methods for tracking higher-order streaming tensors efficiently and effectively. Accordingly, designing adaptive methods capable of directly tracking tensors of higher order is of great importance and it is our main concern in this study. Some adaptive methods have been developed for handling higher-order streaming tensors in the literature. For example, Dawon *et al.* in [13] introduced another online CP algorithm called STF which is capable of dealing with higher-order streaming tensors. The authors imposed a temporal regularization on the loading factors and used the SGD method to update them over time. In [14], Zhang *et al.* developed a Bayesian-based streaming method called BRST robust to outliers. To track and separate the low-rank and sparsity components of the underlying tensor, a Bayesian statistical model was applied. The computational complexity of BRST is however very high and, thus, the method becomes inefficient when handling high-dimensional and fast-arriving data streams. In [15], Lee *et al.* proposed another robust streaming CP algorithm called SOFIA which has the potential to handle real-world data streams with missing values and sparse outliers. Specifically, SOFIA exploits the well-known time-series forecasting model namely Holt-Winters for detecting outliers, temporal patterns, and hence factorizing the underlying tensor. In our past work [16], we developed a robust adaptive CP decomposition with missing data and outliers. Thanks to the recursive least-squares technique in adaptive filtering and the alternating direction method of multipliers (ADMM) method, RACP is capable of online detecting sparse outliers and tracking successfully the underlying CP model of data streams over time. In parallel, some adaptive CP decomposition algorithms, such as [17–21], are capable of handling higher-order tensors. However, they do not handle incomplete datasets.

Under the Tucker format, there are many adaptive methods capable of factorizing streaming tensors in online settings [7]. Particularly, several Tucker trackers were proposed to decompose streaming tensors having one mode/dimension evolving with time. Some of them work under the assumption that temporal slices of the streaming tensor interact with the same core tensor of fixed size, such as RPTucker [22], BASS-Tucker [23], RT-NTD [24], BK-NTD [24], and D-TuckerO [25]. Some others, on the other hand, assume that the core tensor has one temporal mode and that its temporal slices associate with data streams, such as STA [26], OTL [27], ORLTM [28], D-L1-Tucker [29], and ROLTD [30], to name a few. Among them, a few Tucker trackers can deal with data corruption. RPTucker [22] is specifically designed for dynamic tensor completion but its ability is limited to 3rd-order tensors. ORLTM [28], D-L1-Tucker [29], and ROLTD [30] are robust to sparse outliers. However, their design is not suitable for handling incomplete observations.

Some studies have been conducted to design efficient t-SVD algorithms for higher-order tensors, for example [31–34]. Most of them were designed for batch computation, and thus, are not suitable for dynamic models. Only TOUCAN [31] has the ability to track t-SVD over time. However, it is only useful

for 3rd-order streaming tensors. In parallel, it is well known that block-term decomposition (BTD) can be considered as a combination of CP and Tucker decompositions [35]. In the tensor literature, there are two adaptive BTD algorithms that are able to factorize streaming tensors, namely OnlineBTD [36] and O-BTD-RLS [37]. They are, however, sensitive to data corruption. With respect to tensor-train decomposition, we proposed TT-FOA [38], which is an adaptive tensor-train (TT) model for streaming tensors. Although TT-FOA and its stochastic version are capable of tracking the online low-rank tensor-train representation of large-scale and higher-order tensors, they are not designed to handle missing data. ROBOT was recently proposed in [39] to overcome this drawback. Very recently, Jinshi *et al.* in [40] proposed, for the first time, an online tensor completion based on a tensor-ring format. Its convergence, however, has not yet been mathematically proven.

In the multi-aspect streaming perspective of tensor analysis, Song *et al.* proposed an effective multi-aspect streaming tensor framework (MAST) [41], used for dynamic tensor completion. MAST can successfully track the multilinear LRA of incomplete tensors with dynamic growth in more than one tensor mode. A robust version of MAST for handling outliers, called outlier-robust multi-aspect streaming tensor completion and factorization (OR-MSTC), was proposed in [42]. Thanks to ADMM, OR-MSTC can estimate the low-rank component from measurements corrupted by outliers. A new inductive framework, called SIITA, has been proposed to incorporate side information into incremental tensor analysis [43]. SIITA can be seen as a counterpart of MAST for multi-aspect streaming Tucker decomposition. Although all these approaches provide good frameworks for the problem of dynamic tensor completion, they are either useful for third-order tensors only or are of high complexity, and hence, relatively inefficient in online applications with data streams. In addition, convergence analysis of these algorithms is not available.

This study considers the problem of tracking t-LRA of higher-order incomplete tensors using randomized sketching techniques. It is mainly motivated by the fact that randomized algorithms reduce the computational complexity and memory storage of their conventional counterparts [44]. As a result, they have recently attracted a great deal of attention and achieved success in large-scale data analysis in general, and in tensor decomposition in particular. For example, Wang *et al.* applied a sketching technique to develop a fast algorithm for orthogonal tensor decomposition [45]. Under mild conditions, the tensor sketch can be obtained without accessing all the data [46]. Battaglino *et al.* proposed a practical randomized CP decomposition [47]. Their work aimed to speed up the traditional ALS algorithm via randomized least-squares regressions. With respect to Tucker decomposition, Malik *et al.* proposed two randomized algorithms using TensorSketch for low-rank tensor decomposition [48]. Che *et al.* designed an effective randomized algorithm for computing the low-rank approximation of tensors under the sequentially truncated HOSVD (ST-HOSVD) model [49]. Recently, they provided an improved version of ST-HOSVD with a lower computational complexity and analyzed its probabilistic

error bound [50]. In parallel, two other randomized versions for HOSVD and ST-HOSVD were introduced by Rachel *et al.* in [51]. However, these algorithms only perform batch computation, so they are not appropriate for online processing. We refer the reader to ref. [52] for a good survey on randomized algorithms for tensor decomposition. This shortcoming motivates us to develop a new efficient randomized algorithm for the problem of tracking t-LRA.

The main contributions of this paper are two-fold. Firstly, under the Tucker format, we propose a novel *adaptive Tucker decomposition (ATD)* algorithm for tracking the online t-LRA of higher-order incomplete streaming tensors. ATD is a low-complexity tensor tracker and its convergence is fast, thanks to the alternating minimization and randomized sketching. It can handle incomplete tensors derived from infinite data streams because it performs Tucker decomposition with constant time and space complexity that are independent of time index t . A convergence analysis is then provided to establish performance guarantees. Secondly, under the CP format, we derive a second algorithm, namely *adaptive CP decomposition (ACP)* for the problem of online t-LRA. ACP is faster than ATD as the cost of both computation and memory storage is lower. ACP exhibits a competitive performance in terms of estimation accuracy and running time. To the best of our knowledge, ATD and ACP are the first of their kind capable of dealing with streaming tensors of higher orders with a “provable” convergence guarantee.

The rest of this paper is structured as follows. Section 2 presents a brief review of tensor operators and the t-LRA problem. Section 3 formulates the problem of tracking t-LRA for incomplete and streaming tensors. Sections 4 describes in detail the proposed method for tracking t-LRA and its convergence analysis. Section 5 conducts extensive experiments to demonstrate the effectiveness and efficiency of our algorithms, in comparison with state-of-the-art algorithms. Section 6 concludes the paper.

2 Background

2.1 Notations and Definitions

In this paper, we use the following notational conventions. Scalars and vectors are denoted by lowercase letters (e.g., x) and boldface lowercase letters (e.g., \mathbf{x}), respectively. Boldface capital and bold calligraphic letters denote matrices (e.g., \mathbf{X}) and tensors (e.g., \mathcal{X}). For index notation, the (i_1, i_2, \dots, i_N) -th entry of \mathcal{X} is indicated by x_{i_1, i_2, \dots, i_N} . The symbols \circ , \odot , \otimes and \circledast are used to denote the outer, Khatri-Rao, Kronecker, and Hadamard products respectively. The symbol $\lfloor \cdot \rfloor$ denotes the operator for rounding to the nearest integer. We use \mathbf{X}^\top and $\mathbf{X}^\#$ to represent the transpose and pseudo-inverse of \mathbf{X} . Also, $\|\cdot\|_F$ denotes the Frobenius norm of a vector, matrix, and tensor. The operator `randsample(n, k)` returns k integers sampled uniformly at random from the range $[1, n]$. In addition, we outline here some algebraic operators on matrices and tensors that are frequently used throughout this paper.

Considering a N -order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, the mode- n fibers of \mathcal{X} are I_n -dimensional vectors derived from fixing all but the i_n -th index. The mode- n unfolding of \mathcal{X} , written as $\underline{\mathbf{X}}^{(n)}$, is a matrix whose columns are the mode- n fibers of \mathcal{X} . We also use $\text{unfold}_n(\mathcal{X})$ to denote the unfolding operation along the n -th mode.

The n -mode product of \mathcal{X} with a matrix $\mathbf{U} \in \mathbb{R}^{J \times I_n}$, written as $\mathcal{X} \times_n \mathbf{U}$, yields a new tensor $\mathcal{Y} \in \mathbb{R}^{I_1 \times \cdots \times I_{n-1} \times J \times I_{(n+1)} \times \cdots \times I_N}$ such that $\underline{\mathbf{Y}}^{(n)} = \mathbf{U} \underline{\mathbf{X}}^{(n)}$. The product of \mathcal{X} with N matrices $\{\mathbf{U}^{(n)}\}_{n=1}^N$ along all N modes is denoted by

$$\mathcal{X} \prod_{n=1}^N \times_n \mathbf{U}^{(n)} = \mathcal{X} \times_1 \mathbf{U}^{(1)} \times_2 \cdots \times_N \mathbf{U}^{(N)}. \quad (1)$$

The concatenation of \mathcal{X} with a tensor $\mathcal{Y} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_{N-1}}$, written as $\mathcal{X} \boxplus \mathcal{Y}$, yields a new tensor $\mathcal{Z} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_{N+1}}$ such that

$$z_{i_1, i_2, \dots, i_N} = \begin{cases} x_{i_1, i_2, \dots, i_N}, & \text{if } i_N \leq I_N, \\ y_{i_1, i_2, \dots, i_{N-1}}, & \text{if } i_N = I_N + 1. \end{cases} \quad (2)$$

The Khatri-Rao and Kronecker products of a sequence of matrices in a reverse order are denoted by

$$\bigodot_{n=1}^N \mathbf{U}^{(n)} = \mathbf{U}^{(N)} \odot \mathbf{U}^{(N-1)} \odot \cdots \odot \mathbf{U}^{(1)}, \quad (3)$$

$$\bigotimes_{n=1}^N \mathbf{U}^{(n)} = \mathbf{U}^{(N)} \otimes \mathbf{U}^{(N-1)} \otimes \cdots \otimes \mathbf{U}^{(1)}. \quad (4)$$

For clarity, the frequently used acronyms and notational conventions are summarized in Tab. 1.

2.2 Low-Rank Approximations of Tensors (t-LRA)

Consider a N -order tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times \cdots \times I_N}$, the t-LRA of \mathcal{X} can be achieved by solving the following minimization:

$$\underset{\mathcal{Y}}{\text{argmin}} \quad \|\mathcal{X} - \mathcal{Y}\|_F^2 \quad \text{subject to} \quad \mathcal{Y} = \mathcal{G} \prod_{n=1}^N \times_n \mathbf{U}^{(n)}, \quad (5)$$

where $\mathbf{r} = [r_1, r_2, \dots, r_N]$ is the desired low multilinear rank, \mathcal{G} is the core tensor of size $r_1 \times r_2 \times \cdots \times r_N$, and $\{\mathbf{U}^{(n)}\}_{n=1}^N$ with $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times r_n}$ are called loading factors [53]. The two most well-known and widely-used approaches for the t-LRA are based on CP decomposition [4] and Tucker decomposition [5].

Tucker decomposition can be considered as a generalization of SVD for tensors, where the loading factors are orthogonal. Generally, this decomposition is not unique in the sense that we can rotate the columns of $\mathbf{U}^{(n)}$ by an

Table 1: Notational conventions.

$x, \mathbf{x}, \mathbf{X}, \mathcal{X}, \mathbb{X}$	scalar, vector, matrix, tensor, and set/subset/support
x_{i_1, i_2, \dots, i_N}	(i_1, i_2, \dots, i_N) -th entry of \mathcal{X}
$\mathbf{x} = \text{vec}(\mathbf{X})$	vectorization of \mathbf{X}
$\mathbf{X} = \text{diag}(\mathbf{x})$	diagonal matrix \mathbf{X} with \mathbf{x} on the main diagonal
$\text{tr}(\mathbf{X})$	trace of \mathbf{X}
$\mathbf{X}(i, :), \mathbf{X}(:, j)$	i -th row and j -th column of \mathbf{X}
$\mathbf{X}^\top, \mathbf{X}^{-1}, \mathbf{X}^\#$	transpose, inverse, and pseudo-inverse of \mathbf{X}
$\underline{\mathbf{X}}^{(n)}$	mode- n unfolding of \mathcal{X}
$\circ, \odot, \otimes, \oplus$	outer, Khatri-Rao, Kronecker, Hadamard product
$\mathcal{X} \boxplus \mathcal{Y}$	concatenation of \mathcal{X} with \mathcal{Y}
$\mathcal{X} \times_n \mathbf{U}$	n -mode product of \mathcal{X} with \mathbf{U}
$\mathcal{X} \prod_{n=1}^N \times_n \mathbf{U}^{(n)}$	$\mathcal{X} \times_1 \mathbf{U}^{(1)} \times_2 \dots \times_N \mathbf{U}^{(N)}$
$\odot_{n=1}^N \mathbf{U}^{(n)}$	$\mathbf{U}^{(N)} \odot \mathbf{U}^{(N-1)} \odot \dots \odot \mathbf{U}^{(1)}$
$\otimes_{n=1}^N \mathbf{U}^{(n)}$	$\mathbf{U}^{(N)} \otimes \mathbf{U}^{(N-1)} \otimes \dots \otimes \mathbf{U}^{(1)}$
$\ \cdot\ _F$	Frobenius norm
$\lfloor \cdot \rfloor$	operator for finding the nearest integer
$\text{randsample}(n, k)$	operator for selecting k integers randomly from $[1, n]$
r_{CP}	CP rank
r_{TD}	Tucker rank

orthogonal matrix $\mathbf{Q}^{(n)} \in \mathbb{R}^{r_n \times r_n}$ while still retaining the Tucker representation. Fortunately, the column space covering the factor $\mathbf{U}^{(n)}$ is unique, thus we can estimate subspaces of the loading factors instead [2, 3, 54].

CP decomposition allows us to represent the tensor \mathcal{X} by a sequence of factors having the same number of columns, i.e.,

$$\mathcal{X} = \sum_{i=1}^r \alpha_i \mathbf{u}_i^{(1)} \circ \mathbf{u}_i^{(2)} \circ \dots \circ \mathbf{u}_i^{(N)}, \quad (6)$$

where r is the tensor rank, $\mathbf{u}_i^{(n)}$ is the i -th column of the n -th factor $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times r}$. This decomposition is complex than Tucker decomposition in terms of representation and computation, but essentially unique under mild conditions [2, 3]. Note that parameters $\{\alpha_i\}_{i=1}^r$ can be absorbed in the loading factors, and hence the main interest of the CP decomposition is in finding $\{\mathbf{U}^{(n)}\}_{n=1}^N$.

3 Problem Statement

In this study, we investigate the problem of tracking t-LRA of an incomplete streaming tensor $\mathcal{X}[t] \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_{(N+1)}[t]}$ fixing all but the last dimension $I_{(N+1)}[t]$ (see illustration in Fig. 1 where the gray boxes represent missing data). Specifically, the t -th tensor slice $\mathcal{X}_t \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$ of $\mathcal{X}[t]$ is assumed

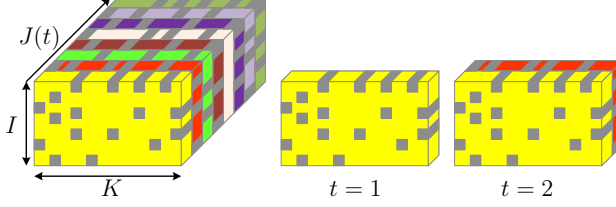


Fig. 1: Incomplete streaming tensors. The gray boxes represent missing data. At each time t , the underlying streaming tensor is obtained by appending the new data (i.e., temporal tensor slice) to the old observations along the time dimension. Particularly, the dimension $J(t)$ is increasing with time while the two dimensions I and K are fixed.

to be generated under the following model:

$$\mathcal{P}_t \circledast \mathcal{X}_t = \mathcal{P}_t \circledast (\mathcal{Y}_t + \mathcal{N}_t), \quad (7)$$

where \mathcal{P}_t is a binary observation mask, \mathcal{N}_t is a Gaussian noise tensor of the same size with \mathcal{X}_t , and \mathcal{Y}_t is the multilinear low-rank component. The mask \mathcal{P}_t shows whether the (i_1, i_2, \dots, i_N) -th entry of \mathcal{X}_t is missing or not, i.e., $p_{i_1, i_2, \dots, i_N} = 1$ if x_{i_1, i_2, \dots, i_N} is observed and $p_{i_1, i_2, \dots, i_N} = 0$ otherwise. The low-rank component \mathcal{Y}_t is given by²

$$\mathcal{Y}_t = \left(\mathcal{G} \prod_{n=1}^N \times_n \mathbf{U}^{(n)} \right) \times_{N+1} \mathbf{u}_t^\top, \quad (8)$$

where $\mathbf{r} = [r_1, r_2, \dots, r_{(N+1)}]$ is the desired low multilinear rank, $\mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times \dots \times r_{(N+1)}}$ is the core tensor, $\mathbf{U} = \{\mathbf{U}^{(n)}\}_{n=1}^N$ with $\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times r_n}$ contains the first N loading factors, and $\mathbf{u}_t \in \mathbb{R}^{r_{(N+1)}}$ is the weight vector. Specifically, the weight vector \mathbf{u}_t in (8) is indeed the t -th row of the last loading factor $\mathbf{U}^{(N+1)} \in \mathbb{R}^{I_{(N+1)}[t] \times r_{(N+1)}}$ of $\mathcal{X}[t]$. The tensor $\mathcal{X}[t]$ is derived by appending the new slice \mathcal{X}_t to the previous $\mathcal{X}[t-1]$ along the time dimension $\mathcal{X}[t] = \mathcal{X}[t-1] \boxplus \mathcal{X}_t$, where $I_{(N+1)}[t] = I_{(N+1)}[t-1] + 1$, as shown in Fig 1.

The problem of tracking t-LRA of the incomplete streaming tensor $\mathcal{X}[t]$ can be stated as follows:

Tracking t-LRA *At each time t , we observe a streaming tensor slice \mathcal{X}_t under the model of (7). We aim to estimate \mathcal{G}_t and \mathbf{u}_t that will provide a good multilinear low-rank approximation for $\mathcal{X}[t]$ in time.*

²In online settings, the tensor core \mathcal{G} and loading factors $\{\mathbf{U}^{(n)}\}$ may be slowly time-varying, i.e., $\mathcal{G} = \mathcal{G}_t$ and $\mathbf{U}^{(n)} = \mathbf{U}_t^{(n)}$, $n = 1, 2, \dots, N$. Our algorithms are capable of estimating \mathcal{G} and \mathbf{U} accurately, but also successfully tracking their variation along the time.

Applying batch methods to $\mathcal{X}[t]$ is possible, but these turns out to be inefficient for online (adaptive) settings. Our goal is to develop efficient one-pass algorithms, both in computational complexity and memory storage, for tracking the t-LRA of $\mathcal{X}[t]$ from past estimations at each time t . In particular, in an adaptive scheme, we propose to minimize the following exponentially weighted cost function:

$$\{\mathcal{G}_t, \mathbf{U}_t\} = \underset{\mathcal{G}, \mathbf{U}}{\operatorname{argmin}} \left[f_t(\mathcal{G}, \mathbf{U}) = \frac{1}{t} \sum_{k=1}^t \lambda^{t-k} \ell(\mathcal{G}, \mathbf{U}, \mathcal{P}_k, \mathcal{X}_k) \right], \quad (9)$$

where the loss function $\ell(\cdot)$ with respect to the k -th slice \mathcal{X}_k is given by

$$\ell(\mathcal{G}, \mathbf{U}, \mathcal{P}_k, \mathcal{X}_k) \triangleq \min_{\mathbf{u}_k \in \mathbb{R}^{r(N+1)}} \left\| \mathcal{P}_k \circledast \left(\mathcal{X}_k - \mathcal{G} \prod_{n=1}^N \times_n \mathbf{U}^{(n)} \times_{N+1} \mathbf{u}_k^\top \right) \right\|_F^2, \quad (10)$$

and $\lambda \in (0, 1]$ is the forgetting parameter. Here, all observations (i.e. tensor slices) in the time interval $[1, t]$ are taken into consideration in the estimation of the underlying low-rank component at each time t . The loss $\ell(\cdot)$ presents the residual for each observation which measures the difference between the observed value and the estimated value of the tensor slice.³ λ is used to discount the effect of past observations exponentially, and to ensure that observations in the distant past are substantially down-weighted in the cost function relative to the latest ones. Accordingly, when $\lambda < 1$, this can facilitate the tracking ability of estimators, especially in time-varying and non-stationary environments. The effective window length for $\lambda < 1$ is $(1 - \lambda)^{-1}$ when t is large. When $\lambda = 1$, (9) boils down to its counterpart of (5) in batch setting.

In the next two sections, we describe the two proposed algorithms for solving (9) under CP and Tucker decompositions. We make the following four assumptions for the convenience of deploying our algorithms as well as analyzing their performance.

- (A1) Observed tensor slices $\{\mathcal{X}_t\}_{t \geq 1}$ are independent and identically distributed from a data-generating distribution, which is the underlying distribution of the dataset, having a compact set \mathcal{V} . This assumption is very common for convergence analysis in online settings in general and adaptive tensor decomposition in particular, e.g., [8, 9, 55–57].⁴

³As there are many choices of low-rank estimations for a given data stream \mathcal{X}_k , $\ell(\cdot)$ in (10) is defined as the minimum loss over all possible choices, and hence, it results in the best low-rank tensor approximation to \mathcal{X}_k . On the arrival of a new data \mathcal{X}_t at each time t , thanks to (10), we can obtain a good estimation of the coefficient vector \mathbf{u}_t which is necessary to establish the first-order surrogate of the cost function $f_t(\cdot)$, to be detailed later in Section 4.

⁴(A1) is a strong assumption in our analysis, but it can be relaxed as follows: Observed tensor slices $\{\mathcal{X}_t\}_{t \geq 1}$ are Frobenius-norm bounded, i.e., $\|\mathcal{X}_t\|_F < M < \infty$. Low-rank components $\{\mathcal{Y}_t\}_{t \geq 1}$ of the observed tensor slices $\{\mathcal{X}_t\}_{t \geq 1}$ are assumed to be deterministic and bounded. Noise tensors $\{\mathcal{N}_t\}_{t \geq 1}$ are i.i.d. from a distribution having a compact support.

- (A2) Tensor slices $\{\mathbf{X}_t\}_{t \geq 1}$ follow the data model (7) where the true underlying loading factors $\{\mathbf{U}_t^{(n)}\}_{t \geq 1}$ are bounded, i.e., $\|\mathbf{U}_t^{(n)}\|_F \leq \kappa < \infty$. It prevents arbitrarily large values in $\mathbf{U}_t^{(n)}$ and ill-conditioned computation. Furthermore, we assume that $\mathbf{U}_t^{(n)}$ is full-column rank for every n and t . This constraint is useful to establish nice propositions in convergence analysis (e.g., boundedness of solutions and Lipschitz continuity of the objective function) as well as to improve the well-posedness of the tensor tracking problem. When (A1) holds, (A2) naturally holds.
- (A3) Observation mask tensors $\{\mathbf{P}_t\}_{t \geq 1}$ are independent of $\{\mathbf{X}_t\}_{t \geq 1}$ and their entries obey the uniform distribution. With respect to the imputation of missing values and recovery of low-rank components, the uniform randomness allows the sequence of binary masks $\{\mathbf{P}_t\}_{t \geq 1}$ to admit stable recovery which is defined as follows:

Definition 1 (Stable recovery [58]). *We say that the sequence of binary masks $\{\mathbf{P}_t\}_{t \geq 1}$ admits the stable recovery if it satisfies the following property: Assume two sequences $\{\mathbf{A}_t\}_{t \geq 1}$, $\{\mathbf{B}_t\}_{t \geq 1}$ where \mathbf{A}_t and \mathbf{B}_t share the same size as \mathbf{P}_t , the rank and the maximum value of \mathbf{A}_t and \mathbf{B}_t are bounded for every t . For any $\varepsilon > 0$, there exists $\delta > 0$ depending only on ε such that if $\limsup_{t \rightarrow \infty} \|\mathbf{P}_t \circledast (\mathbf{A}_t - \mathbf{B}_t)\|_{\bar{F}} \leq \delta$, then $\limsup_{t \rightarrow \infty} \|\mathbf{A}_t - \mathbf{B}_t\|_{\bar{F}} \leq \varepsilon$ where $\|\cdot\|_{\bar{F}}$ is the averaged Frobenius norm defined as $\|\mathbf{A}\|_{\bar{F}} = \|\mathbf{A}\|_F / \sqrt{I_1 I_2 \dots I_N}$.*

Moreover, the number of observed entries in \mathbf{X}_t is assumed to be larger than the lower bound $\mathcal{O}(rL \log(L))$, where $L = \sqrt{I_1 I_2 \dots I_N}$ and $r = \max(r_1, r_2, \dots, r_N)$,⁵ and every row of $\mathbf{X}_t^{(n)}$ is observed at least r entries for all n . The constraints are fundamental conditions to prevent the problem of completion/imputation from being underdetermined where available observations may be insufficient to cover missing entries.

- (A4) The low multilinear-rank model is either static or slowly time-varying, i.e., the core tensor and loading factors may vary slowly between two consecutive times $t-1$ and t : $\mathbf{G}_t \simeq \mathbf{G}_{t-1}$ and $\mathbf{U}_t^{(n)} \simeq \mathbf{U}_{t-1}^{(n)}$. The tensor rank is assumed to be known.

4 Proposed Methods

In this section, we first propose a fast adaptive Tucker algorithm called ATD for tracking the online t-LRA of incomplete streaming tensors. Then, a novel variant of ATD is presented based on the CP format, namely ACP. Next, we provide a performance analysis in terms of complexity and convergence to demonstrate their effectiveness and efficiency.

⁵It is indicated in [59] that no completion method can recover missing data of $\mathbf{M} \in \mathbb{R}^{n \times n}$ with rank $r = \mathcal{O}(1)$ unless the number of observed entries in \mathbf{M} satisfies $m \geq cn \log n$ for some positive constant $c > 0$ and this lower bound is the information theoretical limit.

4.1 Proposed ATD Algorithm

Algorithm 1 Adaptive Tucker Decomposition (ATD)

Input:

- Incomplete slices $\{\mathcal{P}_t \circledast \mathcal{X}_t\}_{t=1}^{\infty}$, $\mathcal{X}_t \in \mathbb{R}^{I_1 \times I_2 \times \dots \times I_N}$,
- Tucker rank $\mathbf{r}_{\text{TD}} = [r_1, r_2, \dots, r_{(N+1)}]$, forgetting factor $\lambda \in (0, 1]$,
- Parameters: $\alpha > 0$, $\delta > 0$, and $m > 0$.

Output: Loading factors $\{\mathbf{U}_t^{(n)}\}_{n=1}^N$ and the core tensor \mathcal{G}_t .

Initialization: $\{\mathbf{U}_0^{(n)}\}_{n=1}^N$ and \mathcal{G}_0 are initialized randomly and $\{\mathbf{S}_0^{(n)}\}_{n=1}^N = \delta \mathbf{I}_{r_n}$.

Main Program:

```

for  $t = 1, 2, \dots$  do
   $\mathcal{X}_{\Omega_t} = \mathcal{P}_t \circledast \mathcal{X}_t$ 
  // Stage 1: Estimation of  $\mathbf{u}_t$ 
   $\mathcal{S} = \text{randsample}(|\Omega_t|, \lfloor mr_{(N+1)} \log r_{(N+1)} \rfloor)$ 

   $\mathcal{H}_t = \mathcal{G}_{t-1} \prod_{n=1}^N \times_n \mathbf{U}_{t-1}^{(n)}$ 

   $\mathbf{u}_t = (\mathbf{H}_{\mathcal{S}_t}^\top \mathbf{H}_{\mathcal{S}_t} + \alpha \mathbf{I})^{-1} \mathbf{H}_{\mathcal{S}_t}^\top \mathbf{x}_{\mathcal{S}_t}$ 
   $\Delta \mathcal{X}_t = \mathcal{P}_t \circledast (\mathcal{X}_t - \mathcal{H}_t \times_{N+1} \mathbf{u}_t^\top)$ 

  // Stage 2: Estimation of  $\{\mathbf{U}_t^{(n)}\}_{n=1}^N$ 
  for  $n = 1, 2, \dots, N$  do
     $\mathbf{W}_t^{(n)} = (\mathbf{U}_{t-1}^{(n)})^\# \underline{\mathbf{X}}_{\Omega_t}^{(n)}$ 
     $\mathbf{S}_t^{(n)} = \lambda \mathbf{S}_{t-1}^{(n)} + \mathbf{W}_t^{(n)} (\mathbf{W}_t^{(n)})^\top$ 
     $\mathbf{V}_t^{(n)} = (\mathbf{S}_t^{(n)})^{-1} \mathbf{W}_t^{(n)}$ 
     $\mathbf{U}_t^{(n)} = \mathbf{U}_{t-1}^{(n)} + \Delta \underline{\mathbf{X}}_t^{(n)} (\mathbf{V}_t^{(n)})^\top$ 
  end for
  // Stage 3: Estimation of  $\mathcal{G}_t$ 
   $\mathbf{Z}_t = \mathbf{u}_t \otimes \left( \bigotimes_{n=2}^{N-1} \mathbf{U}_t^{(n)} \right)$ 
   $\Delta \mathbf{G}_t = (\mathbf{U}_t^{(1)})^\# \Delta \underline{\mathbf{X}}_t^{(1)} \mathbf{Z}_t^\#$ 
   $\Delta \mathcal{G}_t = \text{reshape}(\Delta \mathbf{G}_t, \mathbf{r}_{\text{TD}})$ 
   $\mathcal{G}_t = \mathcal{G}_{t-1} + \Delta \mathcal{G}_t$ 
end for

```

Leveraging past estimations of the loading factors and the core tensor, we propose to minimize the surrogate $g_t(\mathcal{G}, \mathcal{U})$ of $f_t(\mathcal{G}, \mathcal{U})$ instead, which is

defined, for a given value of $\{\mathbf{u}_k\}_{1 \leq k \leq t}$, by

$$g_t(\mathcal{G}, \mathbf{U}) = \frac{1}{t} \sum_{k=1}^t \lambda^{t-k} \left\| \mathcal{P}_k \circledast \left(\mathcal{X}_k - \mathcal{G} \prod_{n=1}^N \times_n \mathbf{U}^{(n)} \times_{N+1} \mathbf{u}_k^\top \right) \right\|_F^2, \quad (11)$$

The main motivation here stems from the following observations: First, it is easy to verify that $g_t(\mathcal{G}, \mathbf{U})$ provides an upper bound on $f_t(\mathcal{G}, \mathbf{U})$ (i.e., $f_t(\mathcal{G}, \mathbf{U}) \leq g_t(\mathcal{G}, \mathbf{U})$ for all \mathcal{G} and \mathbf{U} , and a fixed set of $\{\mathbf{u}_k\}_{1 \leq k \leq t}$). Also, the error function $e_t(\mathcal{G}, \mathbf{U}) = g_t(\mathcal{G}, \mathbf{U}) - f_t(\mathcal{G}, \mathbf{U})$ is L -smooth for some constant $L > 0$, i.e., it is differentiable and $\nabla e_t(\mathcal{G}, \mathbf{U})$ is L -Lipschitz continuous. As a result, $g_t(\mathcal{G}, \mathbf{U})$ is a *first-order surrogate* function of $f_t(\mathcal{G}, \mathbf{U})$ [60] and hence its theoretical convergence results can be achieved without making any strong assumptions on $f_t(\mathcal{G}, \mathbf{U})$. In particular, the sequence of surrogate values $\{g_t(\mathcal{G}_t, \mathbf{U}_t)\}_{t=1}^\infty$ is quasi-martingale and converges almost surely. Accordingly, under a simple assumption that the directional derivative of f_t exists in any direction at any \mathcal{G} and \mathbf{U} , $\{g_t(\mathcal{G}_t, \mathbf{U}_t)\}_{t=1}^\infty$ and $\{f_t(\mathcal{G}_t, \mathbf{U}_t)\}_{t=1}^\infty$ converge to the same limit. Indeed, the solution $\{\mathcal{G}_t, \mathbf{U}_t\}$ derived from minimizing $g_t(\mathcal{G}, \mathbf{U})$ converges to a stationary point of $f_t(\mathcal{G}, \mathbf{U})$ when t approaches infinity. Furthermore, $g_t(\mathcal{G}, \mathbf{U})$ can be effectively minimized with a convergence rate of $\mathcal{O}(1/t)$ and it is much simpler than minimizing $f_t(\mathcal{G}, \mathbf{U})$.

In order to obtain a low-complexity estimator, we exploit the fact that (11) can be efficiently solved using the alternating minimization framework whose iteration step coincides with the tensor slice's acquisition in time. In particular, it can be divided into three main stages: (i) estimate \mathbf{u}_t first, given the old estimation of \mathcal{G}_{t-1} and \mathbf{U}_{t-1} ; (ii) update the loading factor $\mathbf{U}_t^{(n)}$, given \mathbf{u}_t , \mathcal{G}_{t-1} , and the remaining factors and (iii) estimate the core tensor \mathcal{G}_t . The proposed ATD algorithm is summarized in Algorithm 1. In the following, we will describe the key steps of our algorithm for minimizing (11).

4.1.1 Estimation of \mathbf{u}_t

Under the assumption that the loading factors and the core tensor might be static or slowly time-varying (i.e., $\mathbf{U}_t \simeq \mathbf{U}_{t-1}$ and $\mathcal{G}_t \simeq \mathcal{G}_{t-1}$), the weight vector \mathbf{u}_t can be derived from the loss function $\ell(\cdot)$ in (10) at time t by

$$\mathbf{u}_t = \underset{\mathbf{u} \in \mathbb{R}^{r(N+1)}}{\operatorname{argmin}} \left\| \mathcal{P}_t \circledast (\mathcal{X}_t - \mathcal{H}_t \times_{N+1} \mathbf{u}^\top) \right\|_2^2, \quad (12)$$

where $\mathcal{H}_t = \mathcal{G}_{t-1} \prod_{n=1}^N \times_n \mathbf{U}_{t-1}^{(n)}$. Problem (12) can be readily converted into the standard form of

$$\mathbf{u}_t = \underset{\mathbf{u} \in \mathbb{R}^{r(N+1)}}{\operatorname{argmin}} \left\| \mathbf{P}_t (\mathbf{x}_t - \mathbf{H}_t \mathbf{u}) \right\|_2^2, \quad (13)$$

where $\mathbf{P}_t = \operatorname{diag}(\operatorname{vec}(\mathcal{P}_t))$, $\mathbf{x}_t = \operatorname{vec}(\mathcal{X}_t)$, and \mathbf{H}_t is the unfolding matrix of the tensor \mathcal{H}_t . For the sake of convenience, let Ω_t and \mathbf{x}_{Ω_t} be the set and

vector containing the observed entries of \mathcal{X}_t , while \mathbf{H}_{Ω_t} is the sub-matrix of \mathbf{H}_t obtained by selecting the rows corresponding to \mathbf{x}_{Ω_t} .

Generally, problem (13) is an overdetermined least-squares (LS) regression and requires $\mathcal{O}(|\Omega_t|r^2)$ with respect to (w.r.t.) computational complexity to compute the exact LS solution [61]. Thus, it costs time and effort when handling high-dimensional and high-order tensors. We propose to solve a regularized least-squares sketch of (13) instead, i.e.,

$$\mathbf{u}_t = \underset{\mathbf{u} \in \mathbb{R}^{r(N+1)}}{\operatorname{argmin}} \left\| \mathcal{L}(\mathbf{x}_{\Omega_t} - \mathbf{H}_{\Omega_t} \mathbf{u}) \right\|_2^2 + \alpha \|\mathbf{u}\|_2^2, \quad (14)$$

where α is a small positive parameter for regularization, $\mathcal{L}(\cdot)$ is a sketching map that helps reduce the sample size, and hence speed up the calculations. Accordingly, the updated rule for \mathbf{u}_t is given by

$$\mathbf{u}_t = \left(\mathbf{H}_{\mathcal{S}_t}^\top \mathbf{H}_{\mathcal{S}_t} + \alpha \mathbf{I} \right)^{-1} \mathbf{H}_{\mathcal{S}_t}^\top \mathbf{x}_{\mathcal{S}_t}, \quad (15)$$

where $\mathbf{H}_{\mathcal{S}_t}$ and $\mathbf{x}_{\mathcal{S}_t}$ are transformed versions of \mathbf{H}_{Ω_t} and \mathbf{x}_{Ω_t} under the sketching $\mathcal{L}(\cdot)$, respectively. Here, the introduction of $\alpha \|\mathbf{u}\|_2^2$ is in order to avoid the singular/ill-posed computation, multicollinearity, and other pathological phenomena in practice. For example, in some cases, the matrix $\mathbf{H}_{\mathcal{S}_t}$ in (18) may be rank-deficient or near singular. The computation of its inverse is prone to large numerical errors, and hence, the ordinary least-squares solution is no longer well defined. The presence of $\alpha \|\mathbf{u}\|_2^2$ with $\alpha > 0$ results in an additional positive term $\alpha \mathbf{I}$ on the diagonal of the moment matrix $\mathbf{H}_{\mathcal{S}_t}^\top \mathbf{H}_{\mathcal{S}_t}$ while it does not change the eigenvectors of $\mathbf{H}_{\mathcal{S}_t}^\top \mathbf{H}_{\mathcal{S}_t}$. Accordingly, it can ensure that all of the eigenvalues are strictly greater than 0. In other words, the introduction of the regularization term can prevent singularity and ill-posed problems. With respect to the interpretability aspect, the inclusion of $\alpha \|\mathbf{u}\|_2^2$ effectively eliminates multicollinearity – a phenomenon where two or more variables are highly correlated – that affects the interpretability of regression models, including the least-squares estimation [62, 63]. It stems from the fact that adding $\alpha \|\mathbf{u}\|_2^2$ introduces bias to unbiased least-squares estimation but it reduces the variance. Thereby, the regularized least-squares solver can result in a more precise and hence more interpretable estimation in the case when the multicollinearity problem exists in data.⁶

Thanks to the tensor structure of \mathcal{H}_t , the uniform row-sampling is effective in many cases especially when we deal with a high-order streaming tensor (N is large) and/or with some incoherent tensor factors, see Appendix A for details. In the presence of highly coherent factors, a preconditioning (mixing) step is

⁶Typically, the value of α can be chosen by cross validation in batch setting. However, it turns out to be inefficient for stream processing due to its high complexity. In this work, the optimal value of the coefficient vector \mathbf{u}_t is perfect but a good estimation of \mathbf{u}_t is reasonable for tensor tracking. A small α close to the noise level is enough to avoid singular/ill-posed problems during the tracking process. Therefore, we can choose its value in the range $[10^{-3}, 1]$ for reasonable performance in practice.

necessary to guarantee the incoherence. For instance, the subsampled randomized Hadamard transform (SRHT) is a good candidate which can produce a transformed matrix whose rows have (almost) uniform leverage scores [64]. In this context, we here emphasize that well-known randomized LS algorithms can help save much computational complexity while obtaining reasonable estimations of \mathbf{u}_t , especially for large-scale low-rank tensors. It is also worth noting that the update of \mathbf{u}_t costs the most computation time of every tensor tracker as it requires all loading factors $\{\mathbf{U}_{t-1}^{(n)}\}_{n=1}^N$ to form \mathcal{H}_t , and hence, solves the least-squares problem. Therefore, the proposed randomized technique here plays an important role in reducing the overall complexity.

4.1.2 Estimation of $\mathbf{U}_t^{(n)}$

The loading factor $\mathbf{U}_t^{(n)}$ can be updated by minimizing $g_t(\cdot)$ w.r.t. $\mathbf{U}^{(n)}$, as

$$\mathbf{U}_t^{(n)} = \underset{\mathbf{U}^{(n)} \in \mathbb{R}^{I_n \times r_n}}{\operatorname{argmin}} \left[\frac{1}{t} \sum_{k=1}^t \lambda^{t-k} \left\| \mathbf{P}_k^{(n)} \circledast \left(\mathbf{X}_k^{(n)} - \mathbf{U}^{(n)} \mathbf{W}_k^{(n)} \right) \right\|_F^2 \right], \quad (16)$$

where $\mathbf{X}_k^{(n)}$ (resp. $\mathbf{P}_k^{(n)}$) is the mode- n unfolding of \mathcal{X}_k (resp. \mathcal{P}_k) and the coefficient matrix $\mathbf{W}_k^{(n)}$ is the mode- n unfolding of the tensor \mathcal{W}_k which is defined by $\mathcal{W}_k = (\mathcal{G}_{t-1} \prod_{i=1, i \neq n}^N \times_i \mathbf{U}_{t-1}^{(i)}) \times_{N+1} \mathbf{u}_k^\top$.

Interestingly, we exploit the fact that minimization (16) can boil down to the problem of subspace tracking in the presence of missing data [65]. Particularly, the solution of (16) can be obtained by minimizing subproblems for each row $\mathbf{u}_m^{(n)}$ of $\mathbf{U}^{(n)}$, $m = 1, 2, \dots, I_n$ as

$$\mathbf{u}_{t,m}^{(n)} = \underset{\mathbf{u}_m^{(n)} \in \mathbb{R}^r}{\operatorname{argmin}} \left[\frac{1}{t} \sum_{k=1}^t \lambda^{t-k} \left\| \mathbf{P}_{k,m}^{(n)} \left((\mathbf{x}_{k,m}^{(n)})^\top - \mathbf{W}_k^{(n)} (\mathbf{u}_m^{(n)})^\top \right) \right\|_F^2 \right], \quad (17)$$

where $\mathbf{x}_{k,m}^{(n)}$ is the m -th row of $\mathbf{X}_k^{(n)}$ and $\mathbf{P}_{k,m}^{(n)} = \operatorname{diag}(\mathbf{P}_k^{(n)}(m, :))$. Thanks to the parallel scheme of the well-known PETRELS algorithm for subspace tracking [66], we derive an efficient estimator for minimizing the exponentially weighted LS cost function (16). Particularly, we first define two auxiliary matrices $\mathbf{S}_t^{(n)}$ and $\mathbf{V}_t^{(n)}$ as follows⁷

$$\mathbf{S}_t^{(n)} = \lambda \mathbf{S}_{t-1}^{(n)} + (\mathbf{W}_t^{(n)})^\top \mathbf{W}_t^{(n)} \text{ and } \mathbf{V}_t^{(n)} = (\mathbf{S}_t^{(n)})^{-1} (\mathbf{W}_t^{(n)})^\top. \quad (18)$$

The loading factor $\mathbf{U}_t^{(n)}$ is then updated recursively by

$$\mathbf{U}_t^{(n)} = \mathbf{U}_{t-1}^{(n)} + \Delta \mathbf{X}_t^{(n)} (\mathbf{V}_t^{(n)})^\top, \quad (19)$$

⁷To enable the recursive updating rule, the matrix $\mathbf{S}_0^{(n)}$ is initialized by a scaled identity matrix $\mathbf{S}_0^{(n)} = \delta_n \mathbf{I}_{r_n}$ with $\delta_n > 0$.

where the matrix $\Delta \underline{\mathbf{X}}_t^{(n)}$ is derived from the mode- n unfolding of the residual error tensor $\Delta \mathbf{X}_t = \underline{\mathbf{P}}_t \circledast (\mathbf{X}_t - \mathbf{H}_t \times_{N+1} \mathbf{u}_t^\top)$. This is not PETRELS, but a modified version. Here, we can utilize the already updated $\mathbf{U}_t^{(n)}$ to track the remaining factors, which can improve the rate of convergence. Also, we can estimate all the N factors in a parallel scheme which further reduces the cost when several computational units are available.

4.1.3 Estimation of \mathcal{G}_t

For the estimation of \mathcal{G}_t given the latest updated loading factors, (11) is reformulated as

$$\mathcal{G}_t = \underset{\mathcal{G}}{\operatorname{argmin}} \left[\frac{1}{t} \sum_{k=1}^t \lambda^{t-k} \left\| \underline{\mathbf{P}}_k^{(1)} \circledast \left(\underline{\mathbf{X}}_k^{(1)} - \mathbf{U}_t^{(1)} \underline{\mathbf{G}}^{(1)} \mathbf{Z}_k \right) \right\|_F^2 \right], \quad (20)$$

where the variable $\underline{\mathbf{G}}^{(1)}$ is the mode-1 unfolding of \mathcal{G} and the matrix \mathbf{Z}_k is given by $\mathbf{Z}_k = \mathbf{u}_k \otimes \left(\bigotimes_{n=2}^N \mathbf{U}_t^{(n)} \right)$.

When handling a streaming tensor with a huge number of slices (i.e., t is large) and a large number of unknown parameters in \mathcal{G} (i.e., $\prod_{n=1}^{N+1} r_n$ is large), applying batch gradient methods for (20) may be time-consuming despite the effect of the forgetting factor λ . Stochastic approximation is introduced as a good alternative [67]. In particular, we minimize the following function:

$$\mathcal{G}_t = \underset{\mathcal{G}}{\operatorname{argmin}} \left\| \underline{\mathbf{P}}_t^{(1)} \circledast \left(\underline{\mathbf{X}}_t^{(1)} - \mathbf{U}_t^{(1)} \underline{\mathbf{G}}^{(1)} \mathbf{Z}_t \right) \right\|_F^2. \quad (21)$$

Given the estimation of \mathbf{U}_t , the residual error between the newcoming tensor slice and the recovered one is given by $\Delta \underline{\mathbf{X}}_t^{(1)} = \underline{\mathbf{P}}_t^{(1)} \circledast (\underline{\mathbf{X}}_t^{(1)} - \mathbf{U}_t^{(1)} \underline{\mathbf{G}}_{t-1}^{(1)} \mathbf{Z}_t)$. Accordingly, we can derive the variation of \mathcal{G} at time t from

$$\Delta \underline{\mathbf{X}}_t^{(1)} = \underline{\mathbf{P}}_t^{(1)} \circledast (\mathbf{U}_t^{(1)} \Delta \underline{\mathbf{G}}_t^{(1)} \mathbf{Z}_t), \quad (22)$$

where $\Delta \underline{\mathbf{G}}_t^{(1)} = \underline{\mathbf{G}}_t^{(1)} - \underline{\mathbf{G}}_{t-1}^{(1)}$. In particular, $\Delta \underline{\mathbf{G}}_t$ is computed as

$$\Delta \underline{\mathbf{G}}_t^{(1)} = (\mathbf{U}_t^{(1)})^\# \Delta \underline{\mathbf{X}}_t^{(1)} \mathbf{Z}_t^\#. \quad (23)$$

As \mathbf{Z}_t is of the Kronecker structure, we can obtain the pseudo-inverse of \mathbf{Z}_t efficiently by using the following nice property [68]: $(\mathbf{A}_1 \otimes \mathbf{A}_2 \otimes \cdots \otimes \mathbf{A}_n)^\# = \mathbf{A}_1^\# \otimes \mathbf{A}_2^\# \otimes \cdots \otimes \mathbf{A}_n^\#$. After that, $\Delta \underline{\mathbf{G}}_t^{(1)}$ will be reshaped into a tensor $\Delta \mathcal{G}_t$ of size $r_1 \times r_2 \times \cdots \times r_{(N+1)}$. To sum up, we obtain the simple rule for updating \mathcal{G}_t as follows:

$$\mathcal{G}_t = \mathcal{G}_{t-1} + \Delta \mathcal{G}_t. \quad (24)$$

We note that for overdetermined cases, the rule for updating \mathbf{G}_t can be sped up by using the following “vector trick”: $\text{vec}(\mathbf{ABC}^\top) = (\mathbf{C} \otimes \mathbf{A}) \text{vec}(\mathbf{B})$, $\forall \mathbf{A}, \mathbf{B}, \mathbf{C}$. In particular, expression (22) can be cast into the standard least-squares format as follows:

$$\delta \mathbf{x}_t = \mathbf{P}_t \left(\mathbf{u}_t \otimes \left(\bigotimes_{n=1}^N \mathbf{U}_t^{(n)} \right) \right) \delta \mathbf{g}_t, \quad (25)$$

where $\delta \mathbf{x}_t = \text{vec}(\Delta \mathbf{X}_t^{(1)})$, $\delta \mathbf{g}_t = \text{vec}(\Delta \mathbf{G}_t^{(1)})$ and $\mathbf{P}_t = \text{diag}(\text{vec}(\mathbf{P}_t^{(1)}))$. Interestingly, (25) has a Kronecker structure, thus $\delta \mathbf{g}_t$ can be efficiently computed by applying randomized sketching techniques with a much lower complexity, e.g., the uniform sampling or the Kronecker product regression in [69].

4.2 Variants of ATD

4.2.1 Orthogonal ATD

In the cases where the orthogonality constraints are imposed on the loading factors, we add an orthogonalization step of $\mathbf{U}^{(n)}$ at each time t as follows

$$\mathbf{U}_t^{(n)} = \mathbf{U}_t^{(n)} [(\mathbf{U}_t^{(n)})^\top \mathbf{U}_t^{(n)}]^{-1/2}, \quad (26)$$

where $(\cdot)^{-1/2}$ represents the inverse square root or simply take the QR decomposition of $\mathbf{U}_t^{(n)}$. Accordingly, the update of $\Delta \mathbf{G}_t$ in (23) can be sped up by replacing the pseudo-inverse with the transpose operator:

$$\Delta \mathbf{G}_t = (\mathbf{U}_t^{(1)})^\top \Delta \mathbf{X}_t^{(1)} \mathbf{Z}_t^\top. \quad (27)$$

We refer to this variant of ATD as ATD-O.

4.2.2 Adaptive CP Decomposition

It is well known that CP decomposition is viewed as a special case of Tucker decomposition when the core tensor is an identity tensor \mathcal{I} of size $r \times r \times \cdots \times r$, thanks to the following relation:

$$\mathcal{I} \prod_{n=1}^N \times_n \mathbf{U}^{(n)} = \sum_{i=1}^r \mathbf{u}_i^{(1)} \circ \mathbf{u}_i^{(2)} \circ \cdots \circ \mathbf{u}_i^{(N)}. \quad (28)$$

Therefore, we can derive a new adaptive CP decomposition from ATD. Particularly in Step 1 and Step 2 of ATD, we recast the design matrix \mathbf{H}_t in (13) into $\mathbf{H}_t = \odot_{n=1}^N \mathbf{U}_{t-1}^{(n)}$ and the coefficient matrix $\mathbf{W}_k^{(n)}$ in (16) into $\mathbf{W}_k^{(n)} = (\odot_{i=1, i \neq n}^N \mathbf{U}_{t-1}^{(i)}) \odot \mathbf{u}_k^\top$. Meanwhile, Step 3 of ATD is no longer required as we set the core tensor \mathbf{G}_t to \mathcal{I} . This modification of ATD is referred to as ACP which stands for adaptive CP decomposition.

4.3 Performance Analysis

4.3.1 Memory Storage and Computational Complexity

We assume that the fixed dimensions of the streaming tensor are equal to I and the desired Tucker rank is $\mathbf{r}_{\text{TD}} = [r, r, \dots, r]$. In terms of memory storage, ATD requires $\mathcal{O}(r^{N+1})$ and $\mathcal{O}(NIr)$ words of memory for saving the core tensor \mathcal{G} and N tensor factors $\{\mathbf{U}^{(n)}\}_{n=1}^N$, respectively. In addition, the cost to save N matrices $\mathbf{S}_t^{(n)}$ is $\mathcal{O}(Nr^2)$ words of memory. In total, ATD requires $\mathcal{O}(Nr(I+r) + r^{N+1})$ words of memory at each time t . Meanwhile, ACP costs a lower space complexity of $\mathcal{O}(Nr(I+r))$ as it does not need to save the core tensor at each time t .

In terms of computational complexity, the computation of ATD comes from three main estimations: (i) the weight vector \mathbf{u}_t , (ii) the tensor factors $\{\mathbf{U}^{(n)}\}_{n=1}^N$, and (iii) the core tensor \mathcal{G} . The two former estimations are similar to that of ACP, so they require a cost of $\mathcal{O}(|\Omega_t|r + (I^{N-1} + |\mathcal{S}_1|)r^2)$ flops in a parallel scheme where \mathcal{S}_1 denotes the size of the sampling set of (13). The latter estimation costs $\mathcal{O}(|\Omega_t|r + I^{N-1}r^{2(N+1)})$ flops for computing $\Delta\mathcal{X}$ and $\Delta\mathcal{G}$. If using the randomize technique in this stage, the complexity is reduced to $\mathcal{O}(|\Omega_t|r + |\mathcal{S}_2|r^{2(N+1)})$ flops where \mathcal{S}_2 is the set of samples selected from (25). Therefore, the overall computational complexity of ATD is $\mathcal{O}(|\Omega_t|r + (I^{N-1} + |\mathcal{S}_1|)r^2 + |\mathcal{S}_2|r^{2(N+1)})$ flops in parallel scheme. For ACP, the overall computational complexity is $\mathcal{O}(|\Omega_t|r + (NI^{N-1} + |\mathcal{S}_1|)r^2)$ flops and reduces to $\mathcal{O}(|\Omega_t|r + (I^{N-1} + |\mathcal{S}_1|)r^2)$ flops in a parallel scheme. Note that when a preconditioning step (e.g., SRHT) is needed to guarantee the incoherence of \mathbf{H}_{Ω_t} , ATD and ACP require an additional cost of $\mathcal{O}(|\Omega_t|r \log r^2)$ flops.

4.3.2 Convergence Guarantee

Our main theoretical result is stated in the following lemma.

Lemma 1. *Given assumptions (A1)-(A4), $\lambda = 1$, the true \mathcal{G} and \mathcal{U} are fixed, the solutions $\{\mathcal{G}_t, \mathcal{U}_t\}_{t=1}^\infty$ generated by ATD converge to a stationary point of f_t when $t \rightarrow +\infty$, i.e., $\nabla f_t(\mathcal{G}_t, \mathcal{U}_t) \rightarrow 0$ as $t \rightarrow +\infty$.*

Proof of Lemma 1 can be obtained by applying the same framework to derive the asymptotic convergence of adaptive algorithms for problems of online matrix and tensor factorization [8, 9, 16, 55, 56, 70]. In particular, the analysis consists of the following three main stages: (S1) the surrogate function $g_t(\mathcal{G}, \mathcal{U})$ is strongly bi-convex in the sense that \mathcal{G} and \mathcal{U} are seen as multivariate variables. Solutions $\{\mathcal{G}_t, \mathcal{U}_t\}_{t=1}^\infty$ generated by ATD are bounded and their variations between two successive time instances satisfy $\|\mathbf{U}_{t+1}^{(n)} - \mathbf{U}_t^{(n)}\|_F \rightarrow \mathcal{O}(1/t)$ a.s. (S2) The nonnegative sequence $\{g_t(\mathcal{G}_t, \mathcal{U}_t)\}_{t=1}^\infty$ is quasi-martingale and hence convergent almost surely. Furthermore, $g_t(\mathcal{G}_t, \mathcal{U}_t) - f_t(\mathcal{G}_t, \mathcal{U}_t) \rightarrow 0$ a.s. (S3) The empirical cost function $f_t(\mathcal{G}, \mathcal{U})$ is continuously differentiable

and Lipschitz. The sequence of solutions $\{\mathbf{g}_t, \mathbf{u}_t\}_{t=1}^{\infty}$ converges to a stationary point of $f_t(\mathbf{g}, \mathbf{u})$, i.e., when $t \rightarrow \infty$, the gradient $\nabla f_t(\mathbf{g}_t, \mathbf{u}_t) \rightarrow 0$ a.s. We refer the readers to our companion work on robust tensor tracking in [16] for details on this proof framework.

5 Experimental Procedures

5.1 Resource Availability

Lead Contact: Further information and requests for resources should be directed to and will be fulfilled by the lead contact, Dr. Nguyen Linh Trung (linhtrung@vnu.edu.vn).

Materials Availability: No new materials were generated in this study.

Data and Code Availability: All original code has been deposited at https://github.com/thanhtbt/tensor_tracking and archived in Figshare under the <https://doi.org/10.6084/m9.figshare.22276105>. Any additional information required to reanalyze the data reported in this paper is available from the lead contact upon request.

5.2 Evaluations

In this section, experiments are conducted to evaluate the performance of the two proposed algorithms (ACP and ATD) on both synthetic and real data. We also compare them with several state-of-the-art algorithms to provide practical evidences of their effectiveness and efficiency. All experiments are implemented on a windows computer with an Intel Core i5-8300H and 16GB of RAM.

5.2.1 Performance of ACP

We assess the performance of ACP w.r.t. the following aspects: (i) impact of algorithm parameters on its tracking ability; (ii) performance of ACP in non-stationary and time-varying environments; (iii) effectiveness and efficiency of ACP as compared with other adaptive CP algorithms.

Experiment Setup: According to the setup of OLSTEC [9], a time-varying model for streaming tensors is constructed as follows. At $t = 0$, the loading factor $\mathbf{U}_t^{(n)}$ is generated at random whose entries are i.i.d. drawn from the Gaussian distribution $\mathcal{N}(0, 1)$. At time $t > 0$, $\mathbf{U}_t^{(n)} \in \mathbb{R}^{I_n \times r}$ is varied under the model $\mathbf{U}_t^{(n)} = \mathbf{U}_{t-1}^{(n)} \mathbf{Q}_t$, where $\mathbf{Q}_t \in \mathbb{R}^{r \times r}$ is a rotation matrix to control the variation of $\mathbf{U}^{(n)}$ between t and $t - 1$, which is defined by

$$\mathbf{Q}_t = \left[\begin{array}{c|cc|c} \mathbf{I}_{p_t-1} & 0 & 0 & 0 \\ 0 & \cos(\alpha_t) & -\sin(\alpha_t) & 0 \\ 0 & \sin(\alpha_t) & \cos(\alpha_t) & 0 \\ \hline 0 & 0 & 0 & \mathbf{I}_{r-p_t-1} \end{array} \right], \quad (29)$$

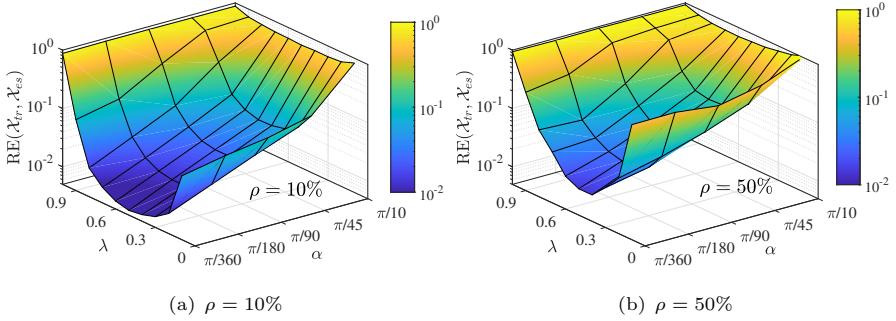


Fig. 2: Effect of the forgetting factor λ on the performance of ACP versus the rotation angle α .

where $p_t = \text{mod}(t + r - 2, r - 1) + 1$ and α_t is the rotation angle with $0 \leq \alpha_t \leq \pi/2$. Specifically, the higher value of α_t is, the faster the loading factor $\mathbf{U}^{(n)}$ changes. The t -th slice \mathcal{X}_t with missing entries is derived from the model:

$$\mathcal{X}_t = \mathcal{P}_t \circledast \left(\mathcal{I} \prod_{n=1}^N \times_n \mathbf{U}_t^{(n)} \times_{N+1} \mathbf{u}_t^\top + \sigma \mathcal{N}_t \right), \quad (30)$$

where \mathcal{P}_t is a binary mask tensor whose entries are generated randomly using the Bernoulli model with the probability ρ , i.e., ρ represents the missing density in the measurement; \mathcal{N}_t is a Gaussian noise tensor (with zero-mean, unit power entries) of the same size of \mathcal{X}_t and the factor σ is to control the noise level; and the weight vector \mathbf{u}_t is a Gaussian random vector living on \mathbb{R}^r space. In order to assess the adaptability of tensor algorithms to unforeseen events and unexpected corruption, we also create abrupt (significant) changes at some time instances during the tracking process by setting the rotation angle at this time instant to $\pi/2$.

To evaluate estimation accuracy, we measure the relative error (RE) metric defined by

$$\text{RE}(\mathbf{A}_{tr}, \mathbf{A}_{es}) = \frac{\|\mathbf{A}_{tr} - \mathbf{A}_{es}\|_F}{\|\mathbf{A}_{tr}\|_F}, \quad (31)$$

where \mathbf{A}_{tr} (resp. \mathbf{A}_{es}) refers to the ground truth (resp. estimation).⁸

Effect of Forgetting Factor λ : The choice of λ plays a central role in how effective and efficient ACP can be in nonstationary environments. In order to investigate the effect of the forgetting factor, we vary the value of λ from 0 to 1 and measure estimation accuracy of ACP in different tests with regard to

⁸Due to the permutation and scaling indeterminacy of the CP decomposition, we can find \mathbf{U}_{es} which is matched with \mathbf{U}_{tr} from \mathbf{U}_t , as follows: $\mathbf{U}_{es} = \mathbf{U}_t \mathbf{P}^\top \mathbf{D}^{-1}$, where the permutation matrix $\mathbf{P} \in \mathbb{R}^{r \times r}$ and the diagonal matrix $\mathbf{D} \in \mathbb{R}^{r \times r}$ are derived by minimizing the optimization $\arg\min_{\mathbf{D}, \mathbf{P}} \|\mathbf{U}_t - \mathbf{U}_{tr} \mathbf{D} \mathbf{P}\|_F^2$.

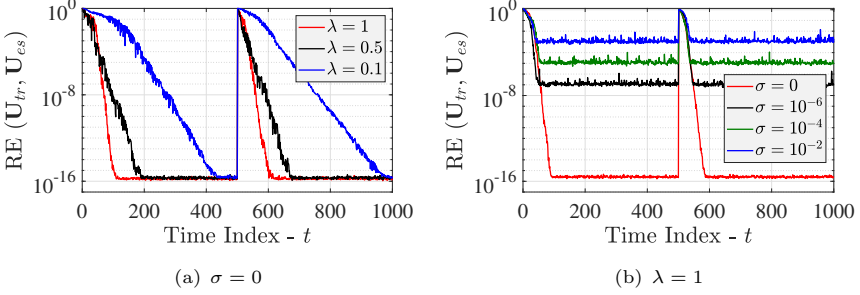


Fig. 3: Performance of ACP in stationary environments: $\mathcal{X}_t \in \mathbb{R}^{20 \times 20 \times 20 \times 1000}$, the true rank $r = 5$, an abrupt change at $t = 500$.

the rotational angle α . Fig. 2 illustrates the experimental results of applying ACP to a synthetic 4-order tensor whose size is $20 \times 20 \times 20 \times 500$ and its rank $r = 5$. The noise level σ is set at 10^{-3} , while the sketching parameter m is fixed at 10. It is clear that the optimal value of λ depends not only on the rotation angle α , but also on the missing density ρ . When λ increases from 0 to 1, the performance of ACP first increases and then drops. Particularly when λ is close to 1, most of the observations are taken into estimation of the underlying low-rank approximation of streaming tensors. However, the properties of old data may be very different from those of the latest observations, especially in fast time-varying environments. Therefore, the performance of ACP degrades significantly in such a case. When the value of λ is small, ACP discounts exponentially the influence of old observations, including the very recent ones. As a result, its convergence rate is slow in stationary or slowly time-varying environments. Accordingly, the forgetting factor λ should be neither too small nor too large. As can be seen in Fig. 2, the value of λ should be around 0.5 for reasonable performance. Thus, we fix $\lambda = 0.5$ in the next experiments. It is worth noting that in stationary environments, we can set the value of $\lambda = 1$ to achieve the best performance, please see Fig. 3 for an illustration.

Asymptotic Convergence Behavior: We next illustrate the convergence behavior of ACP in terms of the variation $\|\mathbf{U}_{t+1} - \mathbf{U}_t\|_F$ and the objective value $f_t(\mathbf{U}_t)$. We use the same 4-order tensor as above but with 1000 tensor slices. Two noise levels are considered (including $\sigma = 0$ and $\sigma = 10^{-3}$), while the missing density ρ is chosen among $\{10\%, 30\%, 50\%\}$. The experimental results are shown as in Fig. 4 and Fig. S1 (in the supplementary data). We can see that the convergence results agree with those stated in Lemma 1.

Noisy and Dynamic Environments: First, the robustness of ACP is investigated against the noise variance. We test ACP's tracking ability on the same static 4-order tensor as above with different values of the noise level σ . Fig. 5(a) shows that the value of σ does not affect the convergence rate of ACP, but only its estimation error. Specifically, when we increase the noise level σ ,

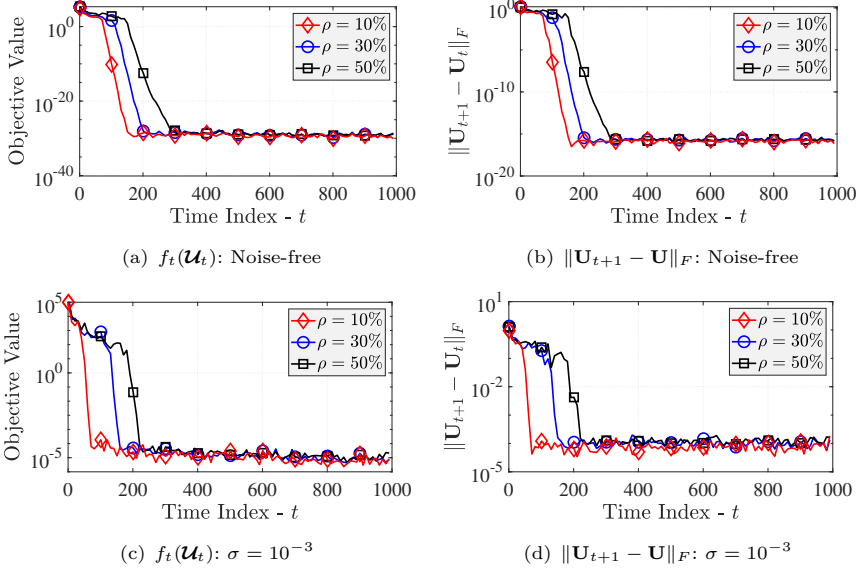


Fig. 4: Convergence behavior of ACP in terms of $f_t(\mathbf{U}_t)$ and $\|\mathbf{U}_{t+1} - \mathbf{U}_t\|_F$.

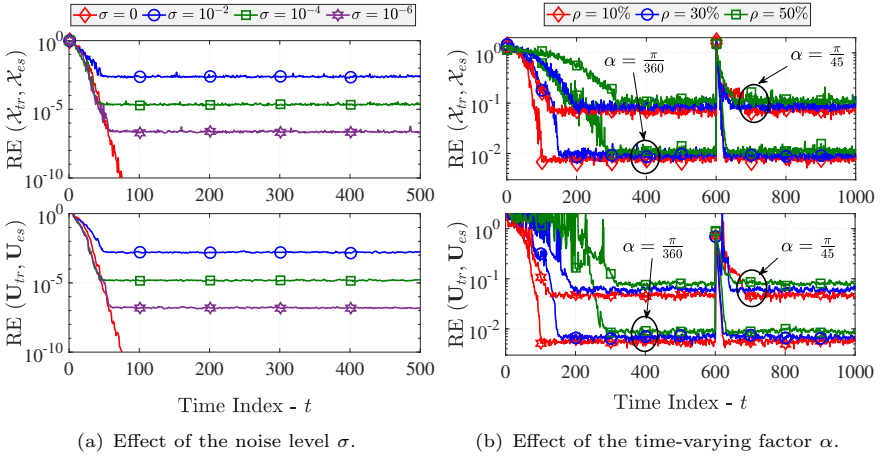


Fig. 5: Performance of ACP in noisy and time-varying environments.

the relative error (RE) between the ground truth and estimation gradually increases, but towards an error bound.

Next, we use the same tensor, but the number of slices is double in order to illustrate the robustness of ACP against time-varying environments. In particular, the proposed algorithm is evaluated in two scenarios, including a

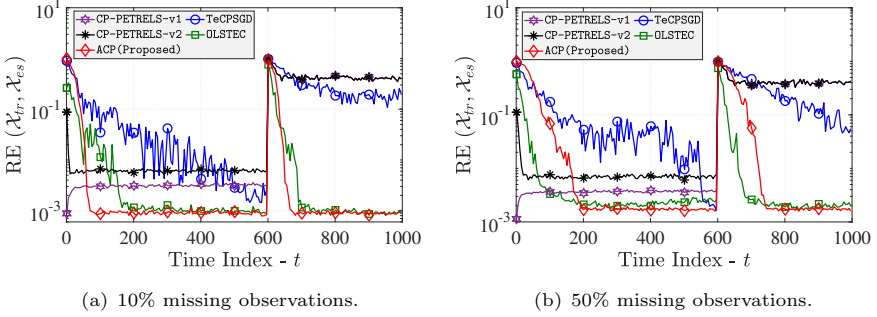


Fig. 6: Tracking performance of adaptive CP algorithms on synthetic 3rd-order tensors: Tensor: size $20 \times 20 \times 1000$, rank $r_{CP} = 5$, the noise level $\sigma_n = 10^{-3}$, and the rotation angle $\alpha = \pi/360$.

slow time-varying case (i.e., $\alpha = \pi/360$) and a fast time-varying case (i.e., $\alpha = \pi/45$). Also, at time $t = 600$, we make an abrupt change in these models. In addition, the missing density ρ is chosen among $\{10\%, 30\%, 50\%\}$. Experimental results indicate that ACP is capable of tracking t-LRA in dynamic environments, as shown in Fig. 5(b). In both scenarios, the relative error (RE) between the ground truth and estimation always converges towards a steady state error bound. The missing density ρ has only an influence on the convergence rate of ACP. Specifically, the lower the missing density ρ is, the faster ACP converges.

Evaluation of Effectiveness and Efficiency: To demonstrate the effectiveness and efficiency of our algorithm, we compare the performance of ACP in terms of estimation accuracy and running time with the state-of-the-art adaptive CP decompositions for incomplete tensors, including OLSTEC [9], CP-PETRELS [10], TeCPSGD [8]. For a fair comparison, the parameters of these algorithms were carefully fine-tuned to achieve good performance. Particularly, the forgetting factor λ is set at 0.7 and 0.98, respectively, for OLSTEC and CP-PETRELS. Moreover, OLSTEC and TeCPSGD are also dependent on a penalty parameter which is set at 10^{-3} in both cases. As CP-PETRELS is sensitive to initialization, we use a set of L training data samples and apply the batch CP method to obtain a good starting point for it. Here, we consider two versions of CP-PETRELS with $L = 100$ and $L = 20$, denoted by CP-PETRELS-v1 and CP-PETRELS-v2, respectively. We use random initialization for ACP, OLSTEC, and TeCPSGD.

Since these algorithms are capable of tracking 3-order tensors only, we use synthetic streaming tensors of size $N \times N \times 1000$ in this task. The noise level is fixed at $\sigma = 10^{-3}$ while we set the rotation angle at $\alpha = \pi/360$. The performance of these algorithms is evaluated on a small tensor $20 \times 20 \times 1000$ and a big tensor $200 \times 200 \times 1000$. Results are shown in Fig. 6 and Fig. 7. We can see that ACP and OLSTEC provide the best tracking performance in

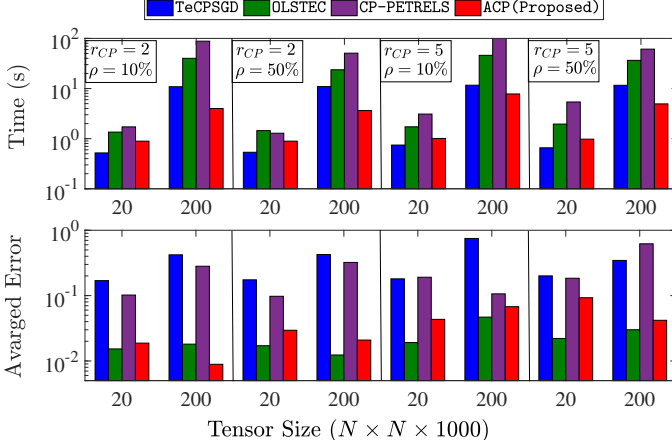


Fig. 7: Running time and averaged error of adaptive CP algorithms.

terms of estimation accuracy. When the number of missing data is small (e.g., $\rho = 10\%$), ACP converges faster than OLSTEC. TeCPSGD’s convergence rate is much slower than that of ACP and OLSTEC and its estimation accuracy is also worse than them. Thanks to the second-order estimation, CP-PETRELS has the fastest convergence rate but its relative errors are higher than those of the others. With respect to the running time, ACP is several times faster than OLSTEC, especially in big tensor tests. TeCPSGD is also a fast adaptive algorithm, thanks to SGD, while the running time of CP-PETRELS is high. For additional performance comparison results, we refer the readers to the supplemental information (i.e., Figs. S2–S4).

5.2.2 Performance of ATD

Experimental Setup: Follow the setup above, the incomplete slice \mathcal{X}_t at time t is generated randomly using the following model:

$$\mathcal{X}_t = \mathcal{P}_t \circledast \left(\mathcal{G}_t \prod_{n=1}^N \times_n \mathbf{U}_t^{(n)} \times \mathbf{u}_t^\top + \sigma \mathcal{N}_t \right), \quad (32)$$

where the loading factor $\mathbf{U}_t^{(n)}$ and the core tensor \mathcal{G}_t are updated by the following rules

$$\mathbf{U}_t^{(n)} = \mathbf{U}_{t-1}^{(n)} + \varepsilon \mathbf{N}_t^{(n)} \text{ and } \mathcal{G}_t = \mathcal{G}_{t-1} + \varepsilon \mathcal{V}_t, \quad (33)$$

where $\mathbf{U}_0^{(n)}, \mathbf{N}_t^{(n)} \in \mathbb{R}^{I_n \times r_n}$ and $\mathcal{V}_t \in \mathbb{R}^{r_1 \times r_2 \times \dots \times r_{(N+1)}}$ are the Gaussian noises whose entries are distributed i.i.d from $\mathcal{N}(0, 1)$ and the time-varying factor ε is to control their variation.

Besides the relative error (RE) metric, we also use the subspace estimation performance (SEP) [71] metric to evaluate the subspace estimation accuracy,

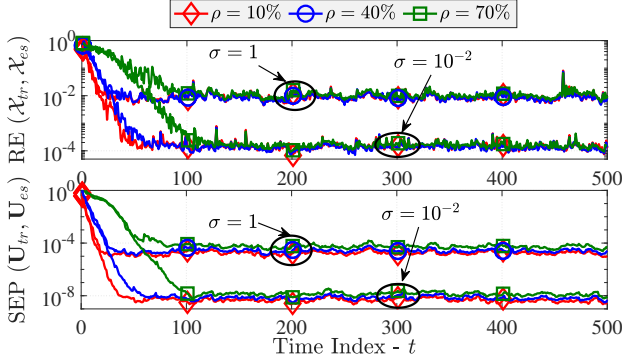


Fig. 8: Performance of ATD versus the missing density ρ and the noise level σ : On the 4-order tensor of size $20 \times 20 \times 20 \times 500$ and its Tucker rank $\mathbf{r}_{\text{TD}} = [3, 3, 3, 3]$.

which is defined by

$$\text{SEP}(\mathbf{U}_{tr}, \mathbf{U}_{es}) = \frac{\text{tr}(\mathbf{U}_{es}^\# (\mathbf{I} - \mathbf{U}_{tr} \mathbf{U}_{tr}^\#) \mathbf{U}_{es})}{\text{tr}(\mathbf{U}_{es}^\# (\mathbf{U}_{tr} \mathbf{U}_{tr}^\#) \mathbf{U}_{es})}, \quad (34)$$

where \mathbf{U}_{tr} (resp. \mathbf{U}_{es}) refers to the true loading factor (resp. estimated factor). The lower the value of SEP is, the more accurate the algorithm is.

Robustness of ATD: In order to demonstrate the robustness of ATD against data corruption, we change the number of missing entries in the measurement by varying the value of ρ and evaluate its performance on different noise levels. We also compare ATD with three well-known batch Tucker algorithms for tensor completion, including iHOOI [72], ALSaS [72], and WTucker [73]. These algorithms are iterative, so their procedure will be stopped when the accuracy tolerance tol ⁹ or the maximum iteration ITER_{\max} has been met. As a convergence guarantee, we fix the value of tol at 10^{-4} , while the value of ITER_{\max} is set at 500, 500, and 100, respectively, for iHOOI, ALSaS and WTucker. For ATD, the forgetting factor λ is fixed at 0.7 in the following experiments.

In this task, we use a static tensor of size $20 \times 20 \times 20 \times 500$ (i.e., the time-varying factor $\varepsilon = 0$) whose core is generated at random from a Gaussian distribution of zero-mean and unit variance. Under the Tucker model with $\mathbf{r}_{\text{TD}} = [3, 3, 3, 3]$, the performance of ATD on the tensor is illustrated in Fig. 8. Results show that ATD can successfully track the multilinear low-rank model in all test cases. Similar to ACP, the missing density ρ has influence only on the convergence rate of ATD, i.e., the higher the value of ρ is, the more slowly ATD converges. Performance comparison among Tucker algorithms is reported

⁹The tolerance tol is used to bound the fitting/estimation error of iterative algorithms which is defined as $\|\mathcal{P} \circledast (\mathcal{X} - \mathcal{G} \prod_{n=1}^N \times_n \mathbf{U}^{(n)})\|_F / \|\mathcal{P} \circledast \mathcal{X}\|_F$.

Table 2: Performance of Tucker algorithms on a static 4-order tensor of size $20 \times 20 \times 20 \times 500$ and the noise level $\sigma = 10^{-2}$.

Missing	$\rho = 50\%$			$\rho = 70\%$			Rank
Metric	RE(\mathcal{X})	SEP(\mathbf{U})	Time(s)	RE(\mathcal{X})	SEP(\mathbf{U})	Time(s)	
iHOOI	3.0e-4	4.2e-8	88.1	8.1e-4	4.7e-7	345.3	[3, 3, 3, 3]
ALSaS	3.1e-4	4.3e-8	109.9	7.8e-4	4.9e-7	539.5	
WTucker	2.1e-4	2.4e-8	209.1	3.5e-4	1.3e-7	597.4	
ATD	6.4e-5	7.6e-9	2.5	1.8e-4	1.4e-8	5.7	
iHOOI	9.1e-2	5.1e-4	192.9	3.5e-1	1.3e-2	571.5	[10, 10, 10, 10]
ALSaS	1.0e-4	2.8e-9	719.1	8.3e-4	3.4e-8	3754.6	
WTucker	3.7e-5	2.8e-10	241.2	5.0e-5	3.3e-10	631.7	
ATD	1.7e-5	6.8e-11	21.7	3.2e-5	2.5e-10	58.2	

in Tab. 2, and shown in Figs. 9 and S5 (in the supplemental information document). We can see that ATD achieves the similar estimation accuracy to iHOOI, ALSaS, and WTucker. As ATD starts from a random point, it needs a certain number of observations to converge during the early stage of tracking. By contrast, iHOOI, ALSaS, and WTucker process data samples in a group or batch. Therefore, their relative error often remains almost the same when the time-varying factor and the noise level are fixed over time. Experimental results also indicate that ATD is the fastest algorithm, much faster than the state-of-the-art Tucker algorithms, thanks to the adaptive scheme and the randomized technique. For example, when dealing with the case of 50% missing observations and $\mathbf{r}_{\text{TD}} = [3, 3, 3, 3]$, the running time of ATD is only 2.51 seconds, which is 35 times faster than the second-fastest algorithm, iHOOI.

Tracking Ability in Dynamic Environments: We continue to investigate the tracking ability of ATD in nonstationary and time-varying environments by changing the time-varying factor ε in the range $[10^{-4}, 10^{-1}]$. We use the same tensor dimensions as in the previous task. We also create a significant subspace change at time $t = 300$ to see how fast ATD can converge. Fig. 10 shows the convergence behavior of ATD versus the time-varying factor ε . We can see that the convergence rate of ATD is not affected by ε but only its error. The performance of the orthogonal ATD (ATD-O) is illustrated in Fig. S6 (in the supplemental information document). Both ATD and ATD-O converge towards the same steady-state error bound. However, the convergence rate of ATD-O is slightly better than that of ATD.

5.2.3 Real Data

To demonstrate the effectiveness of our algorithms on real datasets, we consider two related applications: video completion and multichannel EEG analysis.

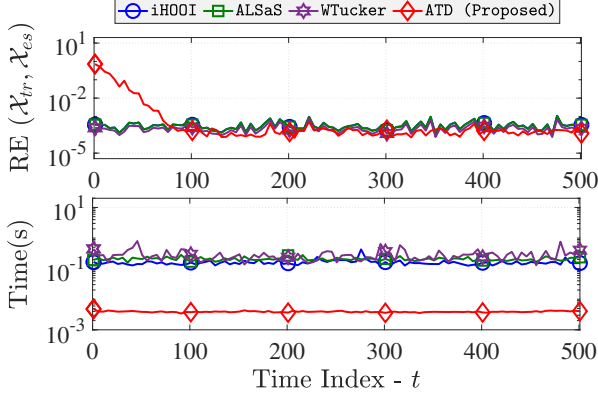


Fig. 9: Performance of Tucker algorithms in the case where 50% entries are observed and Tucker rank $\mathbf{r}_{\text{TD}} = [3, 3, 3, 3]$, and the noise level $\sigma = 10^{-2}$.

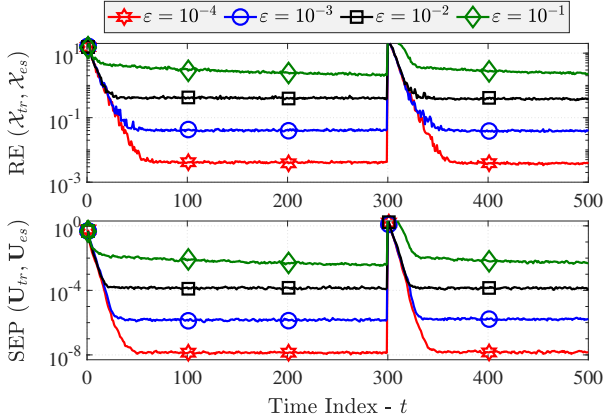


Fig. 10: Effect of the time-varying factor ε on the performance of ATD: Tucker rank $[3, 3, 3, 3]$, 90% entries are observed, the noise is $\sigma = 10^{-2}$ and an abrupt change at $t = 300$.

Video Completion: In this task, four real video surveillance sequences are used, including **Highway**, **Hall**, **Lobby**, and **Park** (video sequences: <http://jacarini.dinf.usherbrooke.ca/>). Specifically, **Highway** contains 1700 frames of size 320×240 pixels showing the two-lane traffic surveillance of vehicles on a highway. **Hall**, which has 3584 frames of size 174×144 pixels, shows an airport hall with many people coming in and out. **Lobby** consists of 1546 frames of size 128×160 pixels captured in an office lobby where the background was changed by switching lights on/off. **Park** includes 600 frames of size 288×352 pixels shot by a thermal camera. To create missing pixels in

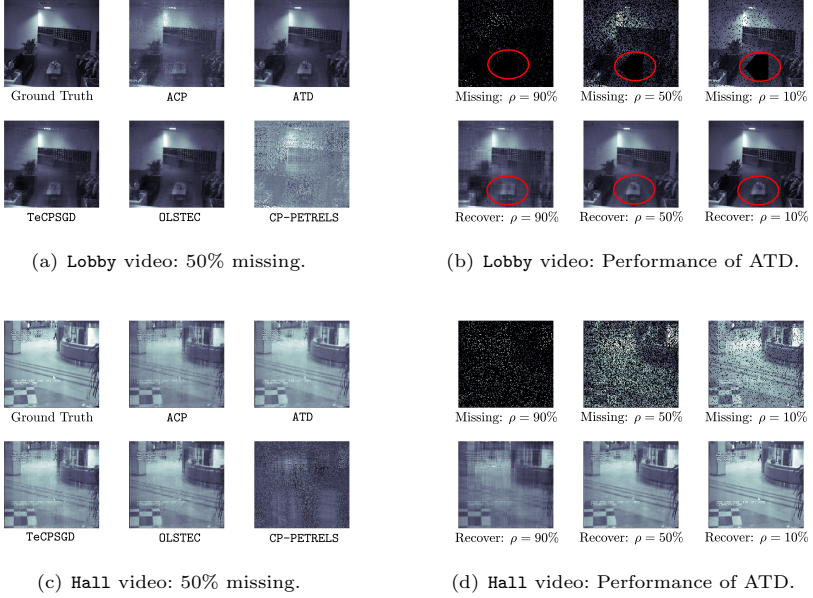


Fig. 11: Performance of adaptive tensor completion algorithms on video datasets.

video, we use a set of binary masks of the same size as the video frames and their entries are derived from a Bernoulli model with probability ρ (i.e., ρ indicates the missing density). Then, we apply the Hadamard product to multiply these masks with all video frames (each mask corresponds to each frame) in order to create zero entries, as in the previous tasks on synthetic data.

We compare the proposed algorithms with OLSTEC [9], TeCPSGD [8] and CP-PETRELS [10]. We set the value of λ at 0.7 and 0.999, respectively, for OLSTEC and CP-PETRELS. Besides, OLSTEC and TeCPSGD also depend on the regularization parameter μ whose value is fixed at 0.1. The performance of these algorithms is shown statistically in Tab. 3 and graphically in Fig. 11. We can see that ATD outperforms adaptive CP algorithms in almost all tests. ACP also provides reasonable estimation accuracy on these data as compared to others. CP-PETRELS seems to fail to track the video background and thus fails to recover missing data. With respect to the running time, experimental results indicate that ACP is the fastest adaptive CP decomposition. We refer readers to the supplemental information document for more experimental results (see Figs. S7-S10).

Multichannel EEG Analysis: We follow the experimental framework in [74, 75] to illustrate the use of ACP for analyzing multichannel EEG signals. In this task, we use a public EEG dataset collected on 14 subjects during the stimulation of hands (EEG data: <http://www.erpwavelab.org/index.htm>).

Table 3: Performance of adaptive tensor decompositions on video datasets.

Methods			TeCPSGD		OLSTEC		CP-PETRELS		ACP		ATD	
Dataset	Size	Missing	RE(\mathcal{X})	Time(s)	RE(\mathcal{X})	Time(s)	RE(\mathcal{X})	Time(s)	RE(\mathcal{X})	Time(s)	RE(\mathcal{X})	Time(s)
Highway	$320 \times 240 \times 1700$	10%	0.2057	36.582	0.1693	132.02	0.9250	451.41	0.2178	14.437	0.1484	36.587
		50%	0.2111	35.252	0.1709	95.188	0.9346	273.98	0.2251	13.295	0.1526	33.269
		90%	0.2256	27.103	0.1849	54.246	0.9224	107.79	0.2725	13.017	0.1964	26.996
Hall	$174 \times 144 \times 3584$	10%	0.1456	15.060	0.1247	83.789	0.9819	339.10	0.1457	11.852	0.1006	36.293
		50%	0.1450	14.916	0.1260	74.869	0.9269	188.15	0.1602	11.808	0.1045	31.576
		90%	0.1614	12.532	0.1497	47.235	0.9281	71.576	0.2341	11.897	0.1426	25.047
Lobby	$128 \times 160 \times 1546$	10%	0.1324	5.672	0.1213	29.490	0.9161	107.44	0.1258	4.613	0.0868	14.590
		50%	0.1452	4.920	0.1228	21.940	0.8596	61.051	0.1881	4.711	0.0884	10.630
		90%	0.1733	4.022	0.1530	14.701	0.9736	22.150	0.2602	3.811	0.1333	9.245
Park	$288 \times 352 \times 600$	10%	0.1057	10.303	0.0905	49.213	0.9945	186.28	0.1270	6.458	0.0686	16.157
		50%	0.1246	9.940	0.0916	33.660	0.9892	127.30	0.1441	5.825	0.0759	14.052
		90%	0.1369	8.497	0.1006	22.031	0.9627	50.435	0.2001	5.179	0.1122	10.966

The EEG signals are recorded using a system of 64 channels (electrodes) and we have 28 measurements per subject in total.

We construct a 3rd-order EEG tensor of measurement \times channel \times time-frequency by taking a continuous wavelet transform of each EEG channel, as in our past works [16, 76]. Note that, the resulting time-frequency representations are reshaped into vectors of length 4392 and are hence streamed. In a nutshell, we have the EEG tensor whose size is $28 \times 64 \times 4392$ and its rank is set at 3 as provided in [74, 75]. At each time, data of 20 (and 40) channels are assumed to be discarded randomly for our missing observation purpose.

We evaluate the performance of ACP by comparing it with the adaptive NL-PETRELS algorithm in [75] and the batch CP-WOPT algorithm in [74]. To have a good initialization for NL-PETRELS, the 1500 first slices are used to construct the training tensor. Also, the forgetting factor λ is set at 0.999. By contrast, ACP is not as sensitive to initialization conditions, so it is initialized at random. We consider results obtained by using the batch algorithm as our ground truth.

Under the CP model with $r_{CP} = 3$, taking the tensor decomposition of the EEG tensor results in three loading factors $\mathbf{A} \in \mathbb{R}^{28 \times 3}$, $\mathbf{B} \in \mathbb{R}^{64 \times 3}$ and $\mathbf{C} \in \mathbb{R}^{4392 \times 3}$ corresponding to, respectively, the measurement, channel and time-frequency modes. Fig. 12 and Fig. 13 illustrate the estimation of \mathbf{A} , \mathbf{B} and \mathbf{C} using CP-WOPT, NL-PETRELS and ACP. In particular, we use bar plots and 3D head plots to represent the column vectors of \mathbf{A} and \mathbf{B} , while the time-frequency representations corresponding to the columns of \mathbf{C} are plotted as matrices. We can see from Fig. 12 that both adaptive algorithms are capable of tracking three EEG loading factors. Our ACP provides a slightly better estimation as compared to that of CP-WOPT. However, in the presence of

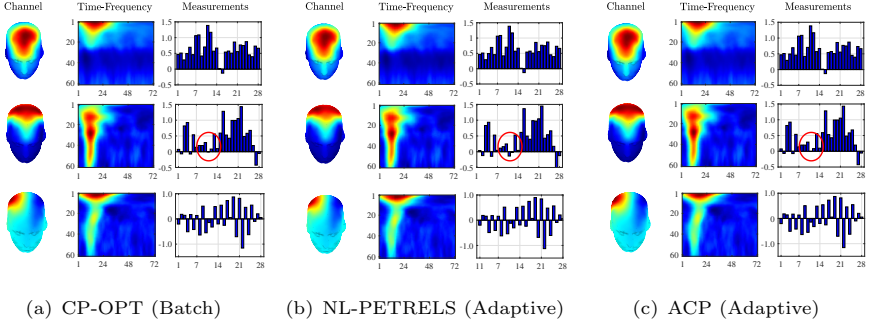


Fig. 12: Waveform-preserving character of ACP on the EEG tensor: 20 channels are missing.

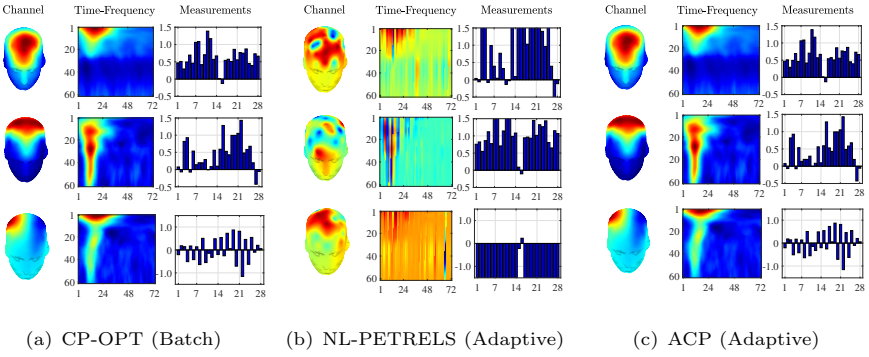


Fig. 13: Waveform-preserving character of ACP on the EEG tensor: 40 channels are missing.

Table 4: Multichannel EEG Analysis: Running times of tensor algorithms.

Missing Data	20 Channels	40 Channels	60 Channels
CP-WOPT	87.51(s)	90.15(s)	95.29(s)
NL-PETRELS	59.72(s)	39.83(s)	37.12(s)
ACP	1.84(s)	1.75(s)	1.42(s)

highly incomplete observations (e.g. 40 channels are missing), NL-PETRELS fails to estimate the EEG loading factors while our ACP algorithms still works well, as shown in Fig 13. The running time of the three algorithms is reported in Tab. 4. In particular, ACP is much faster than NL-PETRELS and CP-WOPT.

6 Conclusions

In this paper, we have proposed two new low-complexity algorithms, namely ACP and ATD, for adaptive decomposition of higher-order incomplete and streaming tensors. Developed based on CP decomposition, ACP estimates a multilinear LRA of streaming tensors from noisy and high-dimensional data with high accuracy, even when the decomposition model may change slowly with time. In parallel, developed based on Tucker decomposition, ATD is a fast randomized tracker, able to recover missing entries from highly incomplete observations. Leveraging the stochastic approximation and the uniform sampling technique, ATD has been shown to be one of the fastest Tucker algorithms, much faster than the batch algorithms while providing good estimation accuracy. Experimental results have indicated that ATD outperforms state-of-the-art adaptive CP algorithms as well as the proposed ACP one in the application of video completion. Thanks to the Tucker format, ATD is more stable than ACP, but with more parameters and is thus slower than ACP.

Acknowledgments. This research was conducted under the research project QG.22.62 “Multi-dimensional data analysis and application to Alzheimer’s disease diagnosis” of Vietnam National University, Hanoi.

Author Contributions. Conceptualization, L.T. Thanh; Methodology, L.T. Thanh, K. Abed-Meraim; Formal Analysis, L.T. Thanh; Validation, K. Abed-Meraim, N.L. Trung, A. Hafiane; Software, L.T. Thanh; Writing — Original Draft, L.T. Thanh; Writing — Review and Editing, L.T. Thanh, K. Abed-Meraim, N.L. Trung, A. Hafiane; Visualization, L.T. Thanh; Supervision, K. Abed-Meraim, N.L. Trung, A. Hafiane; Funding Acquisition, N.L. Trung.

Declaration of Interests. The authors declare no competing interests.

Appendix A Coherence of Kronecker and Khatri-Rao Structure Matrices

We first revisit the definition of the leverage scores and coherence of a matrix.

Definition 2. (Leverage Scores & Coherence [44, Definition 2.1]). *Given a matrix $\mathbf{A} = [\mathbf{a}_1^\top; \mathbf{a}_2^\top; \dots; \mathbf{a}_m^\top] \in \mathbb{R}^{m \times r}$ with $m > r$, its i -th leverage score is defined as*

$$\tau_i(\mathbf{A}) \triangleq \mathbf{a}_i^\top (\mathbf{A}^\top \mathbf{A})^\# \mathbf{a}_i = \|\mathbf{U}_A(i, :)\|_2^2, \quad i = 1, 2, \dots, m. \quad (\text{A1})$$

Here, $\mathbf{U}_A \in \mathbb{R}^{m \times r}$ is the left singular vector matrix of \mathbf{A} . The coherence of \mathbf{A} is the largest leverage score $\mu(\mathbf{A}) = \max_i \tau_i(\mathbf{A})$.

The leverage score $\tau_i(\mathbf{A})$ evaluates the contribution of \mathbf{a}_i in constituting \mathbf{A} 's row space. Accordingly, if the value of $\mu(\mathbf{A})$ is high, \mathbf{A} contains at least one “strong” row whose removal would have a pernicious effect on its row space. When the value of $\mu(\mathbf{A})$ is small (e.g., $\mu(\mathbf{A}) \approx r/m \ll 1$), no specific row is more important than others, i.e., information is approximately uniformized across all rows. In such a case, the matrix \mathbf{A} is called incoherent. The following proposition indicates that the Kronecker structure and Khatri-Rao structure may increase the incoherence from their factors.

Proposition 1. *Given $\mathbf{B} = \odot_{n=1}^N \mathbf{A}_n$, $\mathbf{C} = \otimes_{n=1}^N \mathbf{A}_n$ and let $\bar{\mu}_{\mathbf{A}} = \frac{1}{N} \sum_{n=1}^N \mu(\mathbf{A}_n)$. We always have*

$$\mu(\mathbf{B}) \stackrel{(i)}{\leq} \prod_{n=1}^N \mu(\mathbf{A}_n) \stackrel{(ii)}{=} \mu(\mathbf{C}) \stackrel{(iii)}{\leq} \bar{\mu}_{\mathbf{A}}^N < 1.$$

The first inequality (i) is indeed a corollary of Lemma 4 in [47] which shows that $\mu(\mathbf{A}_1 \odot \mathbf{A}_2) \leq \mu(\mathbf{A}_1)\mu(\mathbf{A}_2)$ for any \mathbf{A}_1 and \mathbf{A}_2 of suitable sizes. The equality (ii) follows Lemma 3 in [47] which indicates that $\mu(\mathbf{A}_1 \otimes \mathbf{A}_2) = \mu(\mathbf{A}_1)\mu(\mathbf{A}_2)$ for any matrices \mathbf{A}_1 and \mathbf{A}_2 . The second inequality (iii) is obtained by applying the AM-GM inequality to the set of N positive numbers $\{\mu(\mathbf{A}_n)\}_{n=1}^N$. Accordingly, when dealing with a large N and/or with some incoherent factors, $\mu(\mathbf{B}) \leq \bar{\mu}_{t-1}^N \ll \bar{\mu}_{t-1} < 1$, i.e., \mathbf{B} and \mathbf{C} have low coherence.

References

- [1] Chen, M., Mao, S., Liu, Y.: Big data: A survey. *Mobile Netw. Appl.* **19**(2), 171–209 (2014). <https://doi.org/10.1007/s11036-013-0489-0>
- [2] Sidiropoulos, N.D., De Lathauwer, L., Fu, X., Huang, K., Papalexakis, E.E., Faloutsos, C.: Tensor decomposition for signal processing and machine learning. *IEEE Trans. Signal Process.* **65**(13), 3551–3582 (2017). <https://doi.org/10.1109/TSP.2017.2690524>
- [3] Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM Rev.* **51**(3), 455–500 (2009). <https://doi.org/10.1137/07070111X>,
- [4] Harshman, R.A.: Foundations of the PARAFAC procedure: Models and conditions for an explanatory multimodal factor analysis. *UCLA Working Papers Phon.* **16**, 1–84 (1970). <https://www.psychology.uwo.ca/faculty/harshman/wpppfac0.pdf>
- [5] De Lathauwer, L., De Moor, B., Vandewalle, J.: A multilinear singular value decomposition. *SIAM J. Matrix Anal. Appl.* **21**(4), 1253–1278 (2000). <https://doi.org/10.1137/S0895479896305696>

- [6] Nion, D., Sidiropoulos, N.D.: Adaptive algorithms to track the PARAFAC decomposition of a third-order tensor. *IEEE Trans. Signal Process.* **57**(6), 2299–2310 (2009). <https://doi.org/10.1109/TSP.2009.2016885>
- [7] Thanh, L.T., Abed-Meraim, K., Trung, N.L., Hafiane, A.: A contemporary and comprehensive survey on streaming tensor decomposition. *IEEE Trans. Knowl. Data Eng.* (2022). early access, <https://doi.org/10.1109/TKDE.2022.3230874>
- [8] Mardani, M., Mateos, G., Giannakis, G.B.: Subspace learning and imputation for streaming matrices and tensors. *IEEE Trans. Signal Process.* **63**(10), 2663–2677 (2015). <https://doi.org/10.1109/TSP.2015.2417491>
- [9] Kasai, H.: Fast online low-rank tensor subspace tracking by CP decomposition using recursive least squares from incomplete observations. *Neurocomput.* **347**, 177–190 (2019). <https://doi.org/10.1016/j.neucom.2018.11.030>
- [10] Chinh, T.M., Nguyen, V.D., Trung, N.L., Abed-Meraim, K.: Adaptive PARAFAC decomposition for third-order tensor completion. In: *IEEE Int. Conf. Commun. Elect.*, pp. 297–301 (2016). <https://doi.org/10.1109/CCE.2016.7562652>
- [11] Recht, B., Fazel, M., Parrilo, P.A.: Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization. *SIAM Rev.* **52**(3), 471–501 (2010). <https://doi.org/10.1137/070697835>
- [12] Hu, Z., Nie, F., Wang, R., Li, X.: Low rank regularization: A review. *Neural Netw.* **136**, 218–232 (2021). <https://doi.org/10.1016/j.neunet.2020.09.021>
- [13] Ahn, D., Kim, S., Kang, U.: Accurate online tensor factorization for temporal tensor streams with missing values. In: *ACM Int. Conf. Inf. Knowl. Manag.*, pp. 2822–2826 (2021). <https://doi.org/10.1145/3459637.3482048>
- [14] Zhang, Z., Hawkins, C.: Variational Bayesian inference for robust streaming tensor factorization and completion. In: *IEEE Int. Conf. Data Min.*, pp. 1446–1451 (2018). <https://doi.org/10.1109/ICDM.2018.00200>
- [15] Dongjin, L., Kijung, S.: Robust factorization of real-world tensor streams with patterns, missing values, and outliers. In: *IEEE Int. Conf. Data Eng.*, pp. 840–851 (2021). <https://doi.org/10.1109/ICDE51399.2021.00078>
- [16] Thanh, L.T., Abed-Meraim, K., Trung, N.L., Hafiane, A.: Robust tensor tracking with missing data and outliers: Novel adaptive CP decomposition and convergence analysis. *IEEE Trans. Signal Process.* **70**, 4305–4320 (2022). <https://doi.org/10.1109/TSP.2022.3201640>

- [17] Zhou, S., Vinh, N.X., Bailey, J., Jia, Y., Davidson, I.: Accelerating online CP decompositions for higher order tensors. In: ACM Int. Conf. Knowl. Discover. Data Min., pp. 1375–1384 (2016). <https://doi.org/10.1145/2939672.2939763>
- [18] Smith, S., Huang, K., Sidiropoulos, N.D., Karypis, G.: Streaming tensor factorization for infinite data sources. In: SIAM Int. Conf. Data Min., pp. 81–89 (2018). <https://doi.org/10.1137/1.9781611975321.10>
- [19] Thanh, L.T., Abed-Meraim, K., Linh-Trung, N., Hafiane, A.: A Fast Randomized Adaptive CP Decomposition for Streaming Tensors. In: IEEE Int. Conf. Acoust. Speech Signal Process., pp. 2910–2914 (2021). <https://doi.org/10.1109/ICASSP39728.2021.9413554>
- [20] Zeng, C., Ng, M.K.: Incremental CP tensor decomposition by alternating minimization method. SIAM J. Matrix Anal. Appl. **42**(2), 832–858 (2021). <https://doi.org/10.1137/20M1319097>
- [21] Lyu, H., Strohmeier, C., Needell, D.: Online nonnegative CP-dictionary learning for Markovian data. J. Mach. Learn. Res. **23**(148), 1–50 (2022). <https://dl.acm.org/doi/pdf/10.5555/3586589.3586737>
- [22] Kasai, H., Mishra, B.: Low-rank tensor completion: a Riemannian manifold preconditioning approach. In: Int. Conf. Mach. Learn., pp. 1012–1021 (2016). <https://proceedings.mlr.press/v48/kasai16.html>
- [23] Fang, S., Kirby, R.M., Zhe, S.: Bayesian streaming sparse Tucker decomposition. In: Conf. Uncertain. Artif. Intell., pp. 558–567 (2021). <https://proceedings.mlr.press/v161/fang21b.html>
- [24] Zdunek, R., Fonal, K.: Incremental nonnegative Tucker decomposition with block-coordinate descent and recursive approaches. Symmetry **14**(1), 113 (2022). <https://doi.org/10.3390/sym14010113>
- [25] Jang, J.-G., Kang, U.: Static and streaming Tucker decomposition for dense tensors. ACM Trans. Knowl. Discov. Data **17**(5), 1–34 (2023). <https://doi.org/10.1145/3568682>
- [26] Sun, J., Tao, D., Papadimitriou, S., Yu, P.S., Faloutsos, C.: Incremental tensor analysis: Theory and applications. ACM Trans. Knowl. Discov. Data **2**(3), 1–37 (2008). <https://doi.org/10.1145/1409620.1409621>
- [27] Traore, A., Berar, M., Rakotomamonjy, A.: Online multimodal dictionary learning. Neurocomput. **368**, 163–179 (2019). <https://doi.org/10.1016/j.neucom.2019.08.053>
- [28] Li, P., Feng, J., Jin, X., Zhang, L., Xu, X., Yan, S.: Online robust low-rank

- tensor modeling for streaming data analysis. *IEEE Trans. Neural Netw. Learn. Syst.* **30**(4), 1061–1075 (2019). <https://doi.org/10.1109/TNNLS.2018.2860964>
- [29] Chachlakis, D.G., Dhanaraj, M., Prater-Bennette, A., Markopoulos, P.P.: Dynamic L1-norm Tucker tensor decomposition. *IEEE J. Sel. Topics Signal Process.* **15**(3), 587–602 (2021). <https://doi.org/10.1109/JSTSP.2021.3058846>
- [30] Thanh, L.T., Duy, T.T., Abed-Meraim, K., Linh Trung, N., Hafiane, A.: Robust online Tucker dictionary learning from multidimensional data streams. In: *IEEE Asia-Pacific Signal Inf. Process. Assoc. Annu. Conf.*, pp. 1815–1820 (2022). <https://doi.org/10.23919/APSIPAASC55919.2022.9980029>
- [31] Gilman, K., Tarzanagh, D.A., Balzano, L.: Grassmannian optimization for online tensor completion and tracking with the t-SVD. *IEEE Trans. Signal Process.* **70**, 2152–2167 (2022). <https://doi.org/10.1109/TSP.2022.3164837>
- [32] Martin, C.D., Shafer, R., LaRue, B.: An order- p tensor factorization with applications in imaging. *SIAM J. Sci. Comput.* **35**(1), 474–490 (2013). <https://doi.org/10.1137/110841229>
- [33] Zhang, Z., Aeron, S.: Exact tensor completion using t-SVD. *IEEE Trans. Signal Process.* **65**(6), 1511–1526 (2017). <https://doi.org/10.1109/TSP.2016.2639466>
- [34] Jiang, F., Liu, X.-Y., Lu, H., Shen, R.: Efficient multi-dimensional tensor sparse coding using t-linear combination. In: *AAAI Conf. Artif. Intell.*, pp. 3326–3333 (2018). <https://doi.org/10.1609/aaai.v32i1.11620>
- [35] De Lathauwer, L.: Decompositions of a higher-order tensor in block terms – Part II: Definitions and uniqueness. *SIAM J. Matrix Anal. Appl.* **30**(3), 1033–1066 (2008). <https://doi.org/10.1137/070690729>
- [36] Gujral, E., Papalexakis, E.E.: OnlineBTD: Streaming algorithms to track the block term decomposition of large tensors. In: *IEEE Int. Conf. Data Sci. Adv. Anal.*, pp. 168–177 (2020). <https://doi.org/10.1109/DSAA49011.2020.00029>
- [37] Rontogiannis, A.A., Kofidis, E., Giampouras, P.V.: Online rank-revealing block-term tensor decomposition. In: *Asilomar Conf. Signals Syst. Comput.*, pp. 1678–1682 (2021). <https://doi.org/10.1109/ICASSP39728.2021.9415104>
- [38] Thanh, L.T., Abed-Meraim, K., Linh-Trung, N., Boyer, R.: Adaptive

- Algorithms for Tracking Tensor-Train Decomposition of Streaming Tensors. In: Eur. Signal Process. Conf., pp. 995–999 (2020). <https://doi.org/10.23919/Eusipco47968.2020.9287780>
- [39] Thanh, L.T., Abed-Meraim, K., Linh Trung, N., Hafiane, A.: Robust tensor tracking with missing data under tensor-train format. In: Eur. Signal. Process. Conf., pp. 832–836 (2022). <https://doi.org/10.23919/EUSIPCO55093.2022.9909702>
- [40] Yu, J., Zou, T., Zhou, G.: Online subspace learning and imputation by tensor-ring decomposition. Neural Netw. **153**, 314–324 (2022). <https://doi.org/10.1016/j.neunet.2022.05.023>
- [41] Song, Q., Huang, X., Ge, H., Caverlee, J., Hu, X.: Multi-aspect streaming tensor completion. In: ACM Int. Conf. Knowl. Disc. Data Min., pp. 435–443 (2017). <https://doi.org/10.1145/3097983.3098007>
- [42] Najafi, M., He, L., Yu, P.S.: Outlier-robust multi-aspect streaming tensor completion and factorization. In: Int. Joint Conf. Artificial Intell., pp. 3187–3194 (2019). <https://doi.org/10.24963/ijcai.2019/442>
- [43] Nimishakavi, M., Mishra, B., Gupta, M., Talukdar, P.: Inductive framework for multi-aspect streaming tensor completion with side information. In: ACM Int. Conf. Inf. Knowl. Manag., pp. 307–316 (2018). <https://doi.org/10.1145/3269206.3271713>
- [44] Mahoney, M.W.: Randomized algorithms for matrices and data. Found. Trends Mach. Learn. **3**(2), 123–224 (2011). <http://dx.doi.org/10.1561/22000000035>
- [45] Wang, Y., Tung, H.-Y., Smola, A.J., Anandkumar, A.: Fast and guaranteed tensor decomposition via sketching. In: Adv. Neural Inf. Process. Syst., pp. 991–999 (2015). <https://dl.acm.org/doi/10.5555/2969239.2969350>
- [46] Song, Z., Woodruff, D., Zhang, H.: Sublinear time orthogonal tensor decomposition. In: Adv. Neural Inf. Process. Syst., pp. 793–801 (2016). <https://dl.acm.org/doi/10.5555/3157096.3157185>
- [47] Battaglino, C., Ballard, G., Kolda, T.G.: A practical randomized CP tensor decomposition. SIAM J. Matrix Anal. Appl. **39**(2), 876–901 (2018). <https://doi.org/10.1137/17M1112303>
- [48] Malik, O.A., Becker, S.: Low-rank Tucker decomposition of large tensors using Tensorsketch. In: Adv. Neural Inf. Process. Syst., pp. 10096–10106 (2018). <https://dl.acm.org/doi/10.5555/3327546.3327674>

- [49] Che, M., Wei, Y.: Randomized algorithms for the approximations of Tucker and the tensor train decompositions. *Adv. Comput. Math.* **45**(1), 395–428 (2019). <https://doi.org/10.1007/s10444-018-9622-8>
- [50] Che, M., Wei, Y., Yan, H.: Randomized algorithms for the low multilinear rank approximations of tensors. *J. Comput. Appl. Math.* **390**, 113380 (2021). <https://doi.org/10.1016/j.cam.2020.113380>
- [51] Minster, R., Saibaba, A.K., Kilmer, M.E.: Randomized algorithms for low-rank tensor decompositions in the Tucker format. *SIAM J. Math. Data Sci.* **2**(1), 189–215 (2020). <https://doi.org/10.1137/19M1261043>
- [52] Ahmadi-Asl, S., Abukhovich, S., Asante-Mensah, M.G., Cichocki, A., Phan, A.H., Tanaka, T., Oseledets, I.: Randomized algorithms for computation of Tucker decomposition and higher-order SVD (HOSVD). *IEEE Access* **9**, 28684–28706 (2021). <https://doi.org/10.1109/ACCESS.2021.3058103>
- [53] Cichocki, A., Lee, N., Oseledets, I.V., Phan, A.-H., Zhao, Q., Mandic, D.P.: Tensor networks for dimensionality reduction and large-scale optimization: Part I low-rank tensor decompositions. *Found. Trends Mach. Learn.* **9**(4-5), 249–429 (2016). <http://dx.doi.org/10.1561/22000000059>
- [54] De Silva, V., Lim, L.-H.: Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM J. Matrix Anal. Appl.* **30**(3), 1084–1127 (2008). <https://doi.org/10.1137/06066518X>
- [55] Mairal, J., Bach, F., Ponce, J., Sapiro, G.: Online learning for matrix factorization and sparse coding. *J. Mach. Learn. Res.* **11**(2), 19–60 (2010). <https://dl.acm.org/doi/10.5555/1756006.1756008>
- [56] Thanh, L.T., Nguyen, V.D., Trung, N.L., Abed-Meraim, K.: Robust subspace tracking with missing data and outliers: Novel algorithm with convergence guarantee. *IEEE Trans. Signal Process.* **69**, 2070–2085 (2021). <https://doi.org/10.1109/TSP.2021.3066795>
- [57] Thanh, L.T., Abed-Meraim, K., Hafiane, A., Trung, N.L.: Sparse subspace tracking in high dimensions. In: *IEEE Int. Conf. Acoust. Speech Signal Process.*, pp. 5892–5896 (2022). <https://doi.org/10.1109/ICASSP43922.2022.9746546>
- [58] Chatterjee, S.: A deterministic theory of low rank matrix completion. *IEEE Trans. Inf. Theory* **66**(12), 8046–8055 (2020). <https://doi.org/10.1109/TIT.2020.3019569>
- [59] Candes, E.J., Tao, T.: The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inf. Theory* **56**(5), 2053–2080 (2010).

<https://doi.org/10.1109/TIT.2010.2044061>

- [60] Mairal, J.: Incremental majorization-minimization optimization with application to large-scale machine learning. *SIAM J. Optim.* **25**(2), 829–855 (2015). <https://doi.org/10.1137/140957639>
- [61] Raskutti, G., Mahoney, M.W.: A statistical perspective on randomized sketching for ordinary least-squares. *J. Mach. Learn. Res.* **17**(1), 7508–7538 (2016). <https://dl.acm.org/doi/10.5555/2946645.3053495>
- [62] Farrar, D.E., Glauber, R.R.: Multicollinearity in regression analysis: the problem revisited. *Rev. Econ. Stat.* **49**(1), 92–107 (1967). <https://doi.org/10.2307/1937887>
- [63] Allen, M.P.: The problem of multicollinearity. *Understanding Regression Analysis*, 176–180 (1997). <https://doi.org/10.1007/978-0-585-25657-3.37>
- [64] Tropp, J.A.: Improved analysis of the subsampled randomized Hadamard transform. *Adv. Adapt. Data Anal.* **3**(01n02), 115–126 (2011). <https://doi.org/10.1142/S1793536911000787>
- [65] Balzano, L., Chi, Y., Lu, Y.M.: Streaming PCA and subspace tracking: The missing data case. *Proc. IEEE* **106**(8), 1293–1310 (2018). <https://doi.org/10.1109/JPROC.2018.2847041>
- [66] Chi, Y., Eldar, Y.C., Robert, C.: PETRELS: Parallel subspace estimation and tracking by recursive least squares from partial observations. *IEEE Trans. Signal Process.* **61**(23), 5947–5959 (2013). <https://doi.org/10.1109/TSP.2013.2282910>
- [67] Spall, J.C.: *Introduction to Stochastic Search and Optimization*, (2005). (John Wiley & Sons)
- [68] Langville, A.N., Stewart, W.J.: The Kronecker product and stochastic automata networks. *J. Comput. Appl. Math.* **167**(2), 429–447 (2004). <https://doi.org/10.1016/j.cam.2003.10.010>
- [69] Diao, H., Jayaram, R., Song, Z., Sun, W., Woodruff, D.: Optimal sketching for Kronecker product regression and low rank approximation. In: *Adv. Neural Inf. Process. Syst.*, pp. 4739–4750 (2019). <https://dl.acm.org/doi/10.5555/3454287.3454713>
- [70] Feng, J., Xu, H., Yan, S.: Online robust PCA via stochastic optimization. In: *Adv. Neural Inf. Process. Syst.*, pp. 404–412 (2013). <https://dl.acm.org/doi/abs/10.5555/2999611.2999657>
- [71] Trung, N.L., Nguyen, V.D., Thameri, M., Chinh, T.M., Abed-Meraim, K.:

- Low-complexity adaptive algorithms for robust subspace tracking. *IEEE J. Sel. Topics Signal Process.* **12**(6), 1197–1212 (2018). <https://doi.org/10.1109/JSTSP.2018.2876626>
- [72] Xu, Y.: Fast algorithms for higher-order singular value decomposition from incomplete data. *J. Comput. Math.* **35**(4), 395–420 (2017). <https://doi.org/10.4208/jcm.1608-m2016-0641>
- [73] Filipović, M., Jukić, A.: Tucker factorization with missing data with application to low- n -rank tensor completion. *Multidim. Syst. Signal Process.* **26**(3), 677–692 (2015). <https://doi.org/10.1007/s11045-013-0269-9>
- [74] Acar, E., Dunlavy, D.M., Kolda, T.G., Mørup, M.: Scalable tensor factorizations for incomplete data. *Chemometr. Intell. Lab.* **106**(1), 41–56 (2011). <https://doi.org/10.1016/j.chemolab.2010.08.004>
- [75] Trung, N.L., Chinh, T.M., Nguyen, V.D., Abed-Meraim, K.: A non-linear tensor tracking algorithm for analysis of incomplete multi-channel EEG data. In: *IEEE Int. Symp. Medical Inf. Commun. Tech.*, pp. 114–119 (2018). <https://doi.org/10.1109/ISMICT.2018.8573711>
- [76] Thanh, L.T., Dao, N.T.A., Dung, N.V., Trung, N.L., Abed-Meraim, K.: Multi-channel EEG epileptic spike detection by a new method of tensor decomposition. *J. Neural Eng.* **17**(1), 016023 (2020). <https://doi.org/10.1088/1741-2552/ab5247>