

ACCELERATING ALTERNATING LEAST SQUARES FOR TENSOR DECOMPOSITION BY PAIRWISE PERTURBATION

LINJIAN MA* AND EDGAR SOLOMONIK*

Abstract. The alternating least squares algorithm for CP and Tucker decomposition is dominated in cost by the tensor contractions necessary to set up the quadratic optimization subproblems. We introduce a novel family of algorithms that uses perturbative corrections to the subproblems rather than recomputing the tensor contractions. This approximation is accurate when the factor matrices are changing little across iterations, which occurs when alternating least squares approaches convergence. We provide a theoretical analysis to bound the approximation error. Our numerical experiments demonstrate that the proposed pairwise perturbation algorithms are easy to control and converge to minima that are as good as alternating least squares. The experimental results show improvements of up to 3.1X with respect to state-of-the-art alternating least squares approaches for various model tensor problems and real datasets.

Key words. tensor, CP decomposition, Tucker decomposition, alternating least squares

AMS subject classifications. 15A69, 15A72, 65F35, 65K10, 65Y20, 65Y04, 65Y05, 68W25

1. Introduction. Tensor decompositions provide general techniques for approximation and modeling of high dimensional data [12, 15, 19, 21, 37, 50]. They are fundamental in methods for computational chemistry [10, 27, 30], physics [47], and quantum information [29, 47]. Tensor decompositions are performed on tensors arising both in the context of numerical-PDEs (e.g. as part of preconditioners [49]) as well as in data-driven statistical modeling [4, 37, 41, 44]. The alternating least squares (ALS) method, which is most commonly used to compute many of these tensor decompositions, has become a target for parallelization [24, 32], performance optimization [13, 54], and acceleration by randomization [9]. We propose a new algorithm, *pairwise perturbation*, that asymptotically accelerates ALS iteration complexity for CP and Tucker decomposition by leveraging an approximation that is provably accurate for well-conditioned problems and is effective when the algorithm approaches the optimization local minima.

Each iteration of ALS is a sweep over quadratic optimization subproblems for each individual factor matrix composing the decomposition. For both CP and Tucker decomposition, computational cost of each sweep is dominated by the tensor contractions needed to setup the quadratic optimization subproblem for every factor matrix. These contractions are redone at every ALS sweep since they involve the factor matrices, all of which change after each sweep. We propose to circumvent these contractions in the scenario when the factor matrices are changing only slightly at each sweep, which is expected when ALS approaches a local minima. Our method approximates the setup of each quadratic optimization subproblem by computing perturbative corrections to the right-hand side due to the change in each factor matrix since a previous ALS sweep. To do so, pairwise perturbative operators are computed that propagate the change to each factor matrix to the subproblem needed to update each other factor matrix. Computing these operators costs slightly more than a typical ALS sweep. These operators are then reused to *approximately* perform more ALS sweeps until the changes to the factor matrices are deemed large, at which point, regular ALS sweeps are performed. Once the updates performed in these regular sweeps are again small, the pairwise operators are recomputed. Each sweep computed approximately in this way costs asymptotically less than a regular ALS sweep.

For CP decomposition, CP-ALS [12, 23] is widely used as it is robust and makes a relatively large amount of progress for the amount of computation required [37] (although alternatives based on gradient and subgradient descent are also competitive [1]). Within CP-ALS, the computational bottleneck of each sweep involves an operation called *the matricized tensor-times Khatri-Rao product* (MTTKRP). Similarly, the costliest operation in the ALS-based Tucker decomposition (Tucker-ALS) method is called *the tensor times matrix-chain* (TTMc) product. For an order N tensor with modes of dimension s , approximated computation of ALS sweeps via pairwise perturbation reduces the cost of that sweep from $O(s^N R)$ to $O(s^2 R + s R^2)$ for a rank- R CP decomposition and from $O(s^N R)$ to $O(s^2 R^{N-1})$ for a rank- R Tucker decomposition.

To quantify the accuracy of the pairwise perturbation algorithm, in Section 4, we provide an error analysis for both MTTKRP and TTMc operations. For both operations, we first view the ALS procedure in terms of pairwise updates, pushing updates to least-squares problems of all factor matrices as soon as any one of them is updated. This reformulation is algebraically equivalent to the original ALS procedure. If the relative

*Department of Computer Science, University of Illinois at Urbana-Champaign, Urbana, IL, 61801 (lma16@illinois.edu, solomon2@illinois.edu).

change to each factor matrix since pairwise perturbation operators were constructed is bounded by $O(\epsilon)$, we can bound the absolute error of the way pairwise perturbation propagates updates in MTTKRP/TTMc calculations due to changes in any one of the other factor matrices. For order three tensors, this absolute error bound yields a relative error bound that depends on a matrix condition number. For the TTMc operation in Tucker decomposition, we derive a 2-norm relative error bound for the overall TTMc calculations (as opposed to updates thereof) of $O(\epsilon^2)$ that holds when the residual of the Tucker decomposition is somewhat less than the norm of the original tensor. We also derive a Frobenius norm error bound of $O(\epsilon^2(s/R)^{N/2})$ for TTMc, which only assumes that *higher-order singular value decomposition* (HOSVD) [16, 58] is performed to initialize Tucker-ALS (which is typical). In addition, in the appendix, we show that for the CP decomposition, if the factor matrices have changed by $O(\epsilon)$ in norm, the relative error in pairwise perturbation for the overall MTTKRP calculation is bounded by a term that scales with $O(\epsilon^2)$ and a tensor condition number. However, we demonstrate that in the worst case scenario, for decomposition of any large tensor, this tensor condition number can be infinite.

In order to evaluate the performance benefit of pairwise perturbation, in Section 5, we compare per ALS sweep and full decomposition performance using a NumPy-based [22] sequential implementation. Our microbenchmark results compare the performance of one CP-ALS sweep with different input tensor sizes. We consider the initialization sweep, in which the pairwise perturbation operators are calculated, as well as the approximated sweep, in which the operators are not recalculated, of the pairwise perturbation algorithm. These results show that the approximated pairwise perturbation sweeps are up to 6.3X faster than one ALS sweep with the dimension tree algorithm [7, 13, 33, 34, 38, 51, 59] for an order three tensor with dimension size 960, and up to 33.0X faster than one ALS sweep for an order six tensor. We then study the performance and numerical behavior of pairwise perturbation for the decomposition of synthetic tensors and application datasets. Our experimental results show that pairwise perturbation achieves fitness as high as standard ALS, and achieves speed-ups of up to 3.1X for CP decomposition and up to 1.13X for Tucker decomposition with respect to state of the art ALS algorithms.

We also evaluate the performance of pairwise perturbation based on a distributed-memory parallel implementation on many nodes of an Intel KNL system (Stampede2) using Cyclops Tensor Framework [55] and ScaLAPACK [11] libraries. Our experimental results show that pairwise perturbation achieves fitness as high as standard ALS, and achieves speed-ups of up to 1.75X with respect to a standard ALS implementation on top of the Cyclops library on Stampede2.

2. Background. This section first outlines the notation used throughout this paper, then outlines the basic alternating least square algorithms for both CP and Tucker decomposition.

2.1. Notation and Definitions. Our analysis makes use of tensor algebra in both element-wise equations and specialized notation for tensor operations [37]. For vectors, bold lowercase Roman letters are used, e.g., \mathbf{x} . For matrices, bold uppercase Roman letters are used, e.g., \mathbf{X} . For tensors, bold calligraphic fonts are used, e.g., \mathcal{X} . An order N tensor corresponds to an N -dimensional array with dimensions $s_1 \times \cdots \times s_N$. Elements of vectors, matrices, and tensors are denoted in parentheses, e.g., $\mathbf{x}(i)$ for a vector \mathbf{x} , $\mathbf{X}(i, j)$ for a matrix \mathbf{X} , and $\mathcal{X}(i, j, k, l)$ for an order 4 tensor \mathcal{X} . Columns of a matrix \mathbf{X} are denoted by $\mathbf{x}_i = \mathbf{X}(:, i)$. The mode- n matrix product of an order N tensor $\mathcal{X} \in \mathbb{R}^{s_1 \times \cdots \times s_N}$ with a matrix $\mathbf{A} \in \mathbb{R}^{J \times s_n}$ is denoted by $\mathcal{X} \times_n \mathbf{A}$, with the result having dimensions $s_1 \times \cdots \times s_{n-1} \times J \times s_{n+1} \times \cdots \times s_N$. The mode- n vector product of \mathcal{X} with a vector $\mathbf{v} \in \mathbb{R}^{s_n}$ is denoted by $\mathcal{X} \times_n \mathbf{v}^T$, with the result having dimensions $s_1 \times \cdots \times s_{n-1} \times s_{n+1} \times \cdots \times s_N$. Matricization is the process of unfolding a tensor into a matrix. Given a tensor \mathcal{X} the mode- n matricized version is denoted by $\mathbf{X}_{(n)} \in \mathbb{R}^{s_n \times K}$ where $K = \prod_{m=1, m \neq n}^N s_m$. We generalize this notation to define the unfoldings of a tensor \mathcal{X} with dimensions $s_1 \times \cdots \times s_N$ into an order $M+1$ tensor, $\mathcal{X}_{(i_1, \dots, i_M)} \in \mathbb{R}^{s_{i_1} \times \cdots \times s_{i_M} \times K}$, where $K = \prod_{i \in \{1, \dots, N\} \setminus \{i_1, \dots, i_M\}} s_i$, e.g., $\mathcal{X}(j, k, l, m) = \mathcal{X}_{(1,3)}(j, l, k + (m-1)s_2)$. We use parenthesized superscripts as labels for different tensors, e.g., $\mathcal{X}^{(1)}$ and $\mathcal{X}^{(2)}$ are generally unrelated tensors.

The Hadamard product of two matrices $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{I \times J}$ resulting in matrix $\mathbf{W} \in \mathbb{R}^{I \times J}$ is denoted by $\mathbf{W} = \mathbf{U} * \mathbf{V}$, where $\mathbf{W}(i, j) = \mathbf{U}(i, j)\mathbf{V}(i, j)$. The outer product of K vectors $\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(K)}$ of corresponding sizes s_1, \dots, s_K is denoted by $\mathcal{X} = \mathbf{u}^{(1)} \circ \cdots \circ \mathbf{u}^{(K)}$ where $\mathcal{X} \in \mathbb{R}^{s_1 \times \cdots \times s_K}$ is an order K tensor. The Kronecker product of vectors $\mathbf{u} \in \mathbb{R}^I$ and $\mathbf{v} \in \mathbb{R}^J$ is denoted by $\mathbf{w} = \mathbf{u} \otimes \mathbf{v}$ where $\mathbf{w} \in \mathbb{R}^{IJ}$. For matrices $\mathbf{A} \in \mathbb{R}^{I \times K}$ and $\mathbf{B} \in \mathbb{R}^{J \times K}$, their Khatri-Rao product results in a matrix of size $(IJ) \times K$ defined by $\mathbf{A} \odot \mathbf{B} = [\mathbf{a}_1 \otimes \mathbf{b}_1, \dots, \mathbf{a}_K \otimes \mathbf{b}_K]$.

2.2. Tensor Norm. The spectral norm of any tensor $\mathcal{T} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ is

$$\|\mathcal{T}\|_2 = \max_{\substack{\forall i \in \{2, \dots, N\}, \mathbf{x}^{(i)} \in \mathbb{R}^{s_i} \\ \|\mathbf{x}^{(2)}\|_2 = \dots = \|\mathbf{x}^{(N)}\|_2 = 1}} \left\| \mathcal{T} \times_{i \in \{2, \dots, N\}} \mathbf{x}^{(i)T} \right\|_2,$$

where \mathcal{T} is contracted with $\mathbf{x}^{(i)}$ along its i th mode. The spectral tensor norm corresponds to the magnitude of the largest tensor singular value [40]. Computing the spectral norm is NP-hard [25], but can usually be done in practice by specialized variants of ALS [18]. The spectral norm is invariant under reordering of modes of \mathcal{T} . Lemma 2.1 shows submultiplicativity of this norm for the multilinear multiplication.

LEMMA 2.1. *Given any tensor $\mathcal{T} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ and matrix $\mathbf{M} \in \mathbb{R}^{s_N \times R}$, if $\mathcal{V} = \mathcal{T} \times_N \mathbf{M}^T$ then $\|\mathcal{V}\|_2 \leq \|\mathcal{T}\|_2 \|\mathbf{M}\|_2$.*

Proof. There exist unit vectors $\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$ such that

$$\|\mathcal{V}\|_2 = \left\| \mathcal{V} \times_{i \in \{2, \dots, N\}} \mathbf{x}^{(i)T} \right\|_2 = \left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{x}^{(i)T} \times_N (\mathbf{M} \mathbf{x}^{(N)})^T \right\|_2.$$

Let $\mathbf{z} = \mathbf{M} \mathbf{x}^{(N)}$, so $\|\mathbf{z}\|_2 \leq \|\mathbf{M}\|_2$. If $\|\mathbf{z}\|_2 = 0$, then $\|\mathcal{V}\|_2 = 0$, the inequality holds. Otherwise, since

$$\left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{x}^{(i)T} \times_N \mathbf{z}^T \right\|_2 \leq \left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{x}^{(i)T} \times_N \mathbf{z}^T \right\|_2 \frac{\|\mathbf{M}\|_2}{\|\mathbf{z}\|_2} \leq \|\mathcal{T}\|_2 \|\mathbf{M}\|_2,$$

the inequality still holds. \square

2.3. CP Decomposition with ALS. The CP tensor decomposition [23, 26] is a higher-order generalization of the matrix singular value decomposition (SVD). The CP decomposition is denoted by

$$\mathcal{X} \approx \left[\left[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \right] \right], \quad \text{where } \mathbf{A}^{(i)} = \left[\mathbf{a}_1^{(i)}, \dots, \mathbf{a}_R^{(i)} \right],$$

and serves to approximate a tensor by a sum of R tensor products of vectors,

$$\mathcal{X} \approx \sum_{r=1}^R \mathbf{a}_r^{(1)} \circ \dots \circ \mathbf{a}_r^{(N)}.$$

The CP-ALS method alternates among quadratic optimization problems for each of the factor matrices $\mathbf{A}^{(n)}$, resulting in linear least squares problems for each row,

$$\mathbf{A}_{\text{new}}^{(n)} \mathbf{P}^{(n)T} \cong \mathbf{X}_{(n)},$$

where the matrix $\mathbf{P}^{(n)} \in \mathbb{R}^{I_n \times R}$, where $I_n = s_1 \times \dots \times s_{n-1} \times s_{n+1} \times \dots \times s_N$, is formed by Khatri-Rao products of the other factor matrices,

$$\mathbf{P}^{(n)} = \mathbf{A}^{(1)} \odot \dots \odot \mathbf{A}^{(n-1)} \odot \mathbf{A}^{(n+1)} \odot \dots \odot \mathbf{A}^{(N)}.$$

These linear least squares problems are often solved via the normal equations [37]. We also adopt this strategy here to devise the pairwise perturbation method. The normal equations for the n th factor matrix are

$$\mathbf{A}_{\text{new}}^{(n)} \mathbf{\Gamma}^{(n)} = \mathbf{X}_{(n)} \mathbf{P}^{(n)},$$

where $\mathbf{\Gamma} \in \mathbb{R}^{R \times R}$ can be computed via

$$\mathbf{\Gamma}^{(n)} = \mathbf{S}^{(1)} * \dots * \mathbf{S}^{(n-1)} * \mathbf{S}^{(n+1)} * \dots * \mathbf{S}^{(N)}, \quad \text{with each } \mathbf{S}^{(i)} = \mathbf{A}^{(i)T} \mathbf{A}^{(i)}.$$

These equations also give the n th component of the optimality conditions for the unconstrained minimization of the nonlinear objective function,

$$f(\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}) = \frac{1}{2} \left\| \mathcal{X} - \left[\left[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \right] \right] \right\|_F^2,$$

Algorithm 2.1 CP-ALS: ALS procedure for CP decomposition

```

1: Input: Tensor  $\mathcal{X} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ , stopping criteria  $\Delta$ 
2: Initialize  $[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$  as uniformly distributed random matrices within  $[0, 1]$ , initialize  $\mathbf{G}^{(n)} \leftarrow \mathbf{A}^{(n)}$ ,  $\mathbf{S}^{(n)} \leftarrow \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$  for  $n \in \{1, \dots, N\}$ 
3: while  $\sum_{i=1}^N \|\mathbf{G}^{(i)}\|_F > \Delta \|\mathcal{X}\|_F$  do
4:   for  $n \in \{1, \dots, N\}$  do
5:      $\mathbf{\Gamma}^{(n)} \leftarrow \mathbf{S}^{(1)} * \dots * \mathbf{S}^{(n-1)} * \mathbf{S}^{(n+1)} * \dots * \mathbf{S}^{(N)}$ 
6:     Update  $\mathbf{M}^{(n)}$  based on the dimension tree algorithm shown in Figure 1
7:      $\mathbf{A}_{\text{new}}^{(n)} \leftarrow \mathbf{M}^{(n)} \mathbf{\Gamma}^{(n)\dagger}$ 
8:      $\mathbf{G}^{(n)} \leftarrow (\mathbf{A}^{(n)} - \mathbf{A}_{\text{new}}^{(n)}) \mathbf{\Gamma}^{(n)}$ 
9:      $\mathbf{A}^{(n)} \leftarrow \mathbf{A}_{\text{new}}^{(n)}$ 
10:     $\mathbf{S}^{(n)} \leftarrow \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$ 
11:   end for
12: end while
13: return  $[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$ 

```

for which the n th component of the gradient is

$$\frac{\partial f}{\partial \mathbf{A}^{(n)}} = \mathbf{G}^{(n)} = \mathbf{A}^{(n)} \mathbf{\Gamma}^{(n)} - \mathbf{X}_{(n)} \mathbf{P}^{(n)} = \left(\mathbf{A}^{(n)} - \mathbf{A}_{\text{new}}^{(n)} \right) \mathbf{\Gamma}^{(n)}.$$

Algorithm 2.1 presents the basic ALS method described above, keeping track of the Frobenius norm of the N components of the overall gradient to ascertain convergence.

The *Matricized Tensor Times Khatri-Rao Product* or MTTKRP computation, $\mathbf{M}^{(n)} = \mathbf{X}_{(n)} \mathbf{P}^{(n)}$, is the main computational bottleneck of CP-ALS [8]. The computational cost of MTTKRP is $\Theta(s^N R)$ if $s_n = s$ for all $n \in \{1, \dots, N\}$. With the dimension tree algorithm, which will be detailed in Section 2.5, the computational complexity for all the MTTKRP calculations in one ALS sweep is $4s^N R$ to leading order in s . The normal equations worsen the conditioning, but are advantageous for CP-ALS, since $\mathbf{\Gamma}^{(n)}$ can be computed and inverted in just $O(s^2 R + R^3)$ cost and the MTTKRP can be amortized by dimension trees. If QR is used instead of the normal equations, the product of \mathbf{Q} with the right-hand sides would have the cost $2s^N R$ and would need to be done for each linear least squares problem, increasing the overall leading order cost by a factor of $N/2$.

Algorithm 2.2 Tucker-ALS: ALS procedure for Tucker decomposition

```

1: Input: Tensor  $\mathcal{X} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ , decomposition ranks  $\{R_1, \dots, R_N\}$ , stopping criteria  $\Delta$ 
2: Initialize  $[\mathcal{G}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$  using HOSVD, initialize  $\mathcal{F} \leftarrow \mathcal{G}$ 
3: while  $\|\mathcal{F}\|_F > \Delta \|\mathcal{X}\|_F$  do
4:   for  $n \in \{1, \dots, N\}$  do
5:     Update  $\mathcal{Y}^{(n)}$  based on the dimension tree algorithm
6:      $\mathbf{A}^{(n)} \leftarrow R_n$  leading left singular vectors of  $\mathcal{Y}_{(n)}^{(n)}$ 
7:   end for
8:    $\mathcal{G}_{\text{new}} \leftarrow \mathcal{Y}^{(N)} \times_N \mathbf{A}^{(N)T}$ 
9:    $\mathcal{F} \leftarrow \mathcal{G}_{\text{new}} - \mathcal{G}$ 
10:   $\mathcal{G} \leftarrow \mathcal{G}_{\text{new}}$ 
11: end while
12: return  $[\mathcal{G}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$ 

```

2.4. Tucker Decomposition with ALS. In this section we review the ALS method for computing a low-rank Tucker decomposition of a tensor [58]. Tucker decomposition approximates a tensor by a core tensor contracted by matrices with orthonormal columns along each mode. The Tucker decomposition is given by

$$\mathcal{X} \approx \left[[\mathcal{G}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}] \right] = \mathcal{G} \times_1 \mathbf{A}^{(1)} \times_2 \mathbf{A}^{(2)} \dots \times_N \mathbf{A}^{(N)}.$$

The corresponding element-wise expression is

$$\mathcal{X}(x_1, \dots, x_N) \approx \sum_{\{z_1, \dots, z_N\}} \mathcal{G}(z_1, \dots, z_N) \prod_{r \in \{1, \dots, N\}} \mathbf{A}^{(r)}(x_r, z_r).$$

The core tensor \mathcal{G} is of order N with dimensions (Tucker ranks) $R_1 \times \cdots \times R_N$ (throughout error and cost analysis we assume each $R_n = R$ for $n \in \{1, \dots, N\}$). The matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{s_n \times R_n}$ have orthonormal columns.

The *higher-order singular value decomposition* (HOSVD) [16, 58] computes the leading left singular vectors of each one-mode unfolding of \mathcal{X} , providing a good starting point for the Tucker-ALS algorithm. The classical HOSVD computes the truncated SVD of $\mathbf{X}_{(n)} \approx \mathbf{U}^{(n)} \mathbf{\Sigma}^{(n)} \mathbf{V}^{(n)T}$ and sets $\mathbf{A}^{(n)} = \mathbf{U}^{(n)}$ for $n \in \{1, \dots, N\}$. The interlaced HOSVD [20, 60] instead computes the truncated SVD of

$$\mathbf{Z}_{(n)}^{(n)} = \mathbf{U}^{(n)} \mathbf{\Sigma}^{(n)} \mathbf{V}^{(n)T} \quad \text{where} \quad \mathbf{Z}^{(1)} = \mathcal{X} \quad \text{and} \quad \mathbf{Z}_{(n)}^{(n+1)} = \mathbf{\Sigma}^{(n)} \mathbf{V}^{(n)T}.$$

The interlaced HOSVD is cheaper, since the size of each $\mathbf{Z}^{(n)}$ is $s^{N-n+1} R^{n-1}$.

The ALS method for Tucker decomposition [5, 17, 37], which is also called the *higher-order orthogonal iteration* (HOOI), then proceeds by fixing all except one factor matrix, and computing a low-rank matrix factorization to update that factor matrix and the core tensor. To update the n th factor matrix, Tucker-ALS factorizes

$$\mathbf{Y}^{(n)} = \mathcal{X} \times_1 \mathbf{A}^{(1)T} \cdots \times_{n-1} \mathbf{A}^{(n-1)T} \times_{n+1} \mathbf{A}^{(n+1)T} \cdots \times_N \mathbf{A}^{(N)T},$$

which is called the *Tensor Times Matrix-chain* or TTMc, into a product of an matrix with orthonormal columns $\mathbf{A}^{(n)}$ and the core tensor \mathcal{G} , so that $\mathbf{Y}_{(n)}^{(n)} \approx \mathbf{A}^{(n)} \mathbf{G}_{(n)}$. This factorization can be done by taking $\mathbf{A}^{(n)}$ to be the R_n leading left singular vectors of $\mathbf{Y}_{(n)}^{(n)}$. This Tucker-ALS procedure is given in Algorithm 2.2.

As in previous work [14, 48], our implementation computes these singular vectors by finding the left eigenvectors of the Gram matrix $\mathbf{W} = \mathbf{Y}_{(n)}^{(n)} \mathbf{Y}_{(n)}^{(n)T}$. Computing the Gram matrix sacrifices some numerical stability, but avoids a large SVD and provides consistency of the signs of the singular vectors across ALS sweeps.

2.5. The Dimension Tree Algorithm. For CP-ALS, the tensor contractions for MTTKRP can be amortized across the linear least squares problems necessary for a given ALS sweep (for loop iteration in Algorithm 2.1). Such amortization techniques are referred to as dimension tree algorithms and a variety of dimension trees have been studied to minimize costs [7, 13, 33, 34, 38, 51, 59]. As our analysis focuses on leading order cost in s , simple binary dimension trees are an optimal choice. These dimension trees for $N = 3, 4$ are illustrated in Figure 1a,1b. We define the partially contracted MTTKRP intermediates $\mathcal{M}^{(i_1, i_2, \dots, i_m)}$ therein as follows,

$$(2.1) \quad \mathcal{M}^{(i_1, i_2, \dots, i_m)} = \mathcal{X}_{(i_1, i_2, \dots, i_m)} \bigcirc_{j \in \{1, \dots, N\} \setminus \{i_1, i_2, \dots, i_m\}} \mathbf{A}^{(j)}.$$

Elementwise,

$$\mathcal{M}^{(i_1, i_2, \dots, i_m)}(x_{i_1}, x_{i_2}, \dots, x_{i_m}, k) = \sum_{\{x_1, \dots, x_N\} \setminus \{x_{i_1}, x_{i_2}, \dots, x_{i_m}\}} \mathcal{X}(x_1, \dots, x_N) \prod_{r \in \{1, \dots, N\} \setminus \{i_1, i_2, \dots, i_m\}} \mathbf{A}^{(r)}(x_r, k),$$

where $\mathcal{M}^{(1, \dots, N)}$ is the input tensor \mathcal{X} . The first level contractions (contractions between the input tensor and one factor matrix) can be done via matrix multiplications between the reshaped input tensor and the factor matrix. These contractions have a cost of $O(s^N R)$ and are generally the most time-consuming part of ALS. Other contractions (transforming one intermediate into another intermediate) can be done via batched matrix-vector products, and the complexity of an i th level contraction is $O(s^{N+1-i} R)$. Because two first level contractions are necessary for the construction of tree dimension tree, as is illustrated in Figure 1a,1b, to calculate all the $\mathcal{M}^{(n)}$ in one ALS sweep, to leading order in s , the computational complexity is $4s^N R$.

For Tucker-ALS, The *Tensor Times Matrix-chain* or TTMc that computes each $\mathcal{Y}^{(n)}$ is the main computational bottleneck of Tucker-ALS [35] and can also be amortized by the dimension tree. The intermediates for Tucker dimension tree are the partially contracted TTMc, $\mathcal{Y}^{(i_1, i_2, \dots, i_m)}$, defined as follows,

$$\mathcal{Y}^{(i_1, i_2, \dots, i_m)} = \mathcal{X} \times_{j \in \{1, \dots, N\} \setminus \{i_1, i_2, \dots, i_m\}} \mathbf{A}^{(j)T},$$

where \mathcal{X} is contracted with all the matrices $\mathbf{A}^{(j)}$ except $\mathbf{A}^{(i_1)}, \dots, \mathbf{A}^{(i_m)}$. Each contraction can be done via matrix multiplications, and the complexity of an i th level contraction is $O(s^{N+1-i} R^i)$. Similar to CP-ALS, to calculate all the $\mathcal{Y}^{(n)}$ in one ALS sweep, to leading order in s , the computational complexity is $4s^N R$.

3. Pairwise Perturbation Algorithms. We now introduce a pairwise perturbation (PP) algorithm to accelerate the ALS procedure when the iterative optimization steps are approaching a local minimum. We first derive the approximation for order three tensors, then generalize the algorithm to order N tensors. The key idea of the pairwise perturbation method is to compute *pairwise perturbation operators*, which correlate a pair of factor matrices. These tensors are then used to repeatedly update the quadratic subproblems for each tensor. As we will show, these updates are provably accurate if the factor matrices do not change significantly since their state at the time of formation of the pairwise perturbation operators.

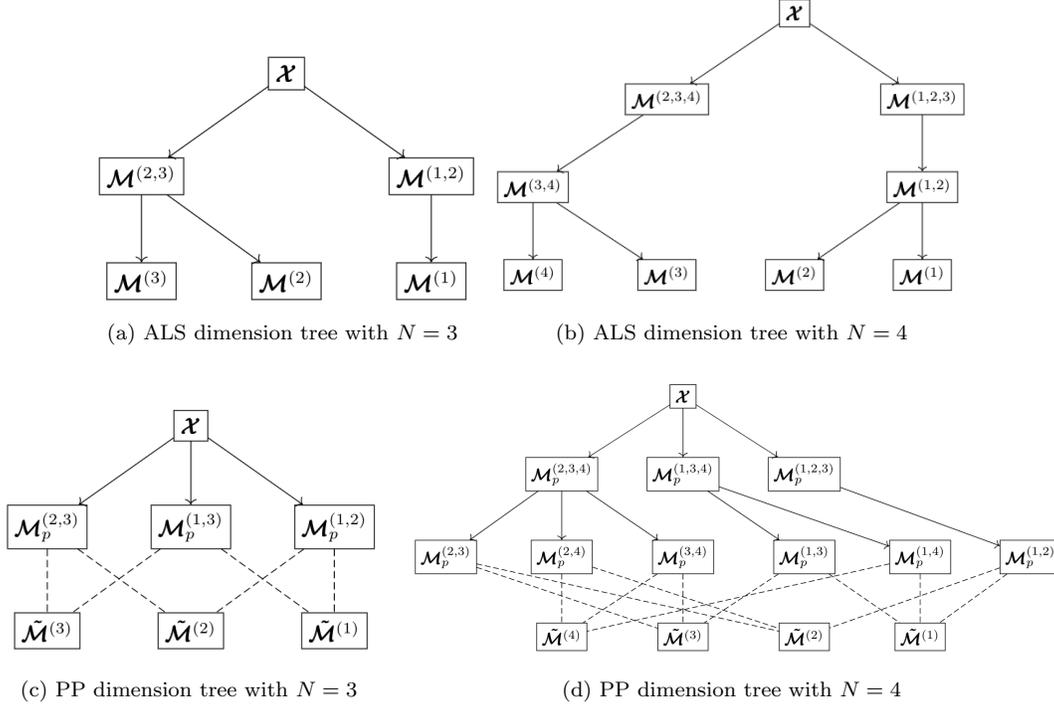


FIGURE 1. Dimension trees for ALS and pairwise perturbation. In (c)(d), the solid arrows denote the data dependencies in building pairwise perturbation operators, and is calculated in the PP initialization step. The dashed lines denote the data dependencies in the PP approximated step calculations.

3.1. Pairwise Perturbation for Order Three Tensors.

3.1.1. CP-ALS. The pairwise perturbation procedure for CP-ALS approximates the MTTKRP outputs. Consider an order three equi-dimensional tensor with size in each mode s and CP rank R , the first mode MTTKRP can be expressed as $\mathbf{M}^{(1)} = \mathbf{X}_{(1)} (\mathbf{A}^{(2)} \odot \mathbf{A}^{(3)})$. Let $\mathbf{A}_p^{(n)}$ denote the $\mathbf{A}^{(n)}$ calculated with regular ALS at some number of sweeps prior to the current one. Then $\mathbf{A}^{(n)}$ at the current sweep can be expressed as

$$\mathbf{A}^{(n)} = \mathbf{A}_p^{(n)} + d\mathbf{A}^{(n)},$$

and $\mathbf{M}^{(1)}$ can be expressed as

$$(3.1) \quad \mathbf{M}^{(1)} = \underbrace{\mathbf{X}_{(1)} (\mathbf{A}_p^{(2)} \odot \mathbf{A}_p^{(3)}) + \mathbf{X}_{(1)} (\mathbf{A}_p^{(2)} \odot d\mathbf{A}^{(3)}) + \mathbf{X}_{(1)} (d\mathbf{A}^{(2)} \odot \mathbf{A}_p^{(3)})}_{\mathbf{U}^{(1)}} + \mathbf{X}_{(1)} (d\mathbf{A}^{(2)} \odot d\mathbf{A}^{(3)}).$$

The pairwise perturbation procedure for CP-ALS approximates $\mathbf{M}^{(1)}$ with $\tilde{\mathbf{M}}^{(1)} = \mathbf{U}^{(1)} + \mathbf{V}^{(1)}$, where $\mathbf{U}^{(1)}$ is the first three terms in Equation 3.1 and $\mathbf{V}^{(1)}$ approximates the final term through approximating the input tensor \mathcal{X} by its approximate CP decomposition,

$$\mathbf{X}_{(1)} (d\mathbf{A}^{(2)} \odot d\mathbf{A}^{(3)}) \approx \mathbf{V}^{(1)} = \left(\left[\left[\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)} \right] \right]_{(1)} \right) (d\mathbf{A}^{(2)} \odot d\mathbf{A}^{(3)}) = \mathbf{A}^{(1)} \left((\mathbf{A}^{(2)T} d\mathbf{A}^{(2)}) * (\mathbf{A}^{(3)T} d\mathbf{A}^{(3)}) \right),$$

which can be calculated with the cost of $O(sR^2)$. The remaining error term is

$$\left(\boldsymbol{\mathcal{X}} - \left[\left[\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}\right]\right]_{(1)}\right) \left(d\mathbf{A}^{(2)} \odot d\mathbf{A}^{(3)}\right).$$

Therefore, the norm of the error scales as $O(C\epsilon^2)$ if each $\|d\mathbf{A}^{(i)}\|_2 \leq \epsilon$ and the decomposition residual norm is bounded by C .

The approximated MTTKRP, $\tilde{\mathbf{M}}^{(1)}$, can be rewritten as a function of $\mathcal{M}_p^{(i_1, i_2, \dots, i_m)}$, which is defined in the same way as $\mathcal{M}^{(i_1, i_2, \dots, i_m)}$ in Equation 2.1 except that $\boldsymbol{\mathcal{X}}$ is contracted with $\mathbf{A}_p^{(j)}$ for $j \in \{1, \dots, N\} \setminus \{i_1, i_2, \dots, i_m\}$,

$$\tilde{\mathbf{M}}^{(1)}(x, k) = \mathbf{M}_p^{(1)}(x, k) + \sum_{y=1}^s \mathcal{M}_p^{(1,2)}(x, y, k) d\mathbf{A}^{(2)}(y, k) + \sum_{y=1}^s \mathcal{M}_p^{(1,3)}(x, y, k) d\mathbf{A}^{(3)}(y, k) + \mathbf{V}^{(1)}(x, k).$$

PP has two steps: the initialization step, where the terms $\mathbf{M}_p^{(1)}$ and pairwise perturbation operators $\mathcal{M}_p^{(1,2)}$, $\mathcal{M}_p^{(1,3)}$ are calculated, and the approximated step, where these terms are used in the equation above to calculate $\tilde{\mathbf{M}}^{(1)}$. Using the dimension tree structure shown in Figure 1c, the initialization step for all the three modes can be done with the leading order cost of $6s^3R$, 1.5X the cost of the ALS dimension tree. Each approximated step for all the modes can be done with the leading order cost of $3(4s^2R + 6sR^2)$ overall.

3.1.2. Tucker-ALS. We derive a similar pairwise perturbation algorithm for order three Tucker-ALS. The first mode of TTMc can be expressed as $\boldsymbol{\mathcal{Y}}^{(1)} = \boldsymbol{\mathcal{Y}} \times_2 \mathbf{A}^{(2)T} \times_3 \mathbf{A}^{(3)T}$. PP approximates $\boldsymbol{\mathcal{Y}}^{(1)}$ with

$$\tilde{\boldsymbol{\mathcal{Y}}}^{(1)} = \boldsymbol{\mathcal{X}} \times_2 \mathbf{A}_p^{(2)T} \times_3 \mathbf{A}_p^{(3)T} + \boldsymbol{\mathcal{X}} \times_2 \mathbf{A}_p^{(2)T} \times_3 d\mathbf{A}^{(3)T} + \boldsymbol{\mathcal{X}} \times_2 d\mathbf{A}^{(2)T} \times_3 \mathbf{A}_p^{(3)T},$$

and the error term is $\boldsymbol{\mathcal{X}} \times_2 d\mathbf{A}^{(2)T} \times_3 d\mathbf{A}^{(3)T}$. The expression above can be rewritten as a function of $\boldsymbol{\mathcal{Y}}_p^{(i_1, i_2, \dots, i_m)}$, which is defined in the same way as $\boldsymbol{\mathcal{Y}}^{(i_1, i_2, \dots, i_m)}$ except that $\boldsymbol{\mathcal{X}}$ is contracted with $\mathbf{A}_p^{(j)}$ for $\boldsymbol{\mathcal{Y}}_p^{(i_1, i_2, \dots, i_m)}$,

$$\tilde{\boldsymbol{\mathcal{Y}}}^{(1)} = \boldsymbol{\mathcal{Y}}_p^{(1)} + \boldsymbol{\mathcal{Y}}_p^{(1,2)} \times_2 d\mathbf{A}^{(2)T} + \boldsymbol{\mathcal{Y}}_p^{(1,3)} \times_3 d\mathbf{A}^{(3)T}.$$

Using the dimension tree structure, the initialization step for all the three modes can be done with the leading order cost of $6s^3R$, 1.5X the cost of the ALS dimension tree. Each approximated step for all the modes can be done with the leading order cost of $12s^2R^2$ overall.

3.2. General Pairwise Perturbation Algorithm. We now generalize PP to order N tensors.

3.2.1. CP-ALS. The MTTKRP in n th mode, $\mathbf{M}^{(n)}$, can be expressed as

$$\mathbf{M}^{(n)} = \mathbf{X}_{(n)} \bigodot_{i=1, i \neq n}^N \left(\mathbf{A}_p^{(i)} + d\mathbf{A}^{(i)}\right).$$

$\mathbf{M}^{(n)}$ can be expressed as a function of $\mathcal{M}_p^{(i_1, i_2, \dots, i_m)}$ as follows,

$$\begin{aligned} \mathbf{M}^{(n)}(y, k) = & \mathbf{M}_p^{(n)}(y, k) + \sum_{i=1, i \neq n}^N \sum_{x=1}^{s_i} \mathcal{M}_p^{(i,n)}(x, y, k) d\mathbf{A}^{(i)}(x, k) + \\ & \sum_{i=1, i \neq n}^N \sum_{j=i+1, j \neq n}^N \sum_{x=1}^{s_i} \sum_{z=1}^{s_j} \mathcal{M}_p^{(i,j,n)}(x, z, y, k) d\mathbf{A}^{(i)}(x, k) d\mathbf{A}^{(j)}(z, k) + \dots \end{aligned}$$

From the above expression we observe that, except the first two terms, all terms include the contraction between tensor $\mathcal{M}_p^{(i_1, i_2, \dots, i_m)}$ and at least two matrices $d\mathbf{A}^{(i)}$, so that their norm scales quadratically with the norm of the perturbative updates $d\mathbf{A}^{(i)}$. Therefore, their norm scales as $O(\epsilon^2)$ if $\|d\mathbf{A}^{(i)}\|_2 \leq \epsilon$. The pairwise perturbation algorithm obtains an effective approximation by keeping the first two terms (these terms are illustrated in Figure 1d for an order four tensor), and approximating the input tensor using its approximate CP decomposition in the third term to lower the error to a greater extent,

$$(3.2) \quad \tilde{\mathbf{M}}^{(n)}(y, k) = \mathbf{M}_p^{(n)}(y, k) + \sum_{i=1, i \neq n}^N \sum_{x=1}^{s_i} \mathcal{M}_p^{(i,n)}(x, y, k) d\mathbf{A}^{(i)}(x, k) + \sum_{i,j=1, i, j \neq n, i \neq j}^N \mathbf{V}^{(n,i,j)}(y, k),$$

$$\text{where } \mathbf{M}_p^{(n)} = \mathbf{X}_{(n)} \bigcirc_{i=1, i \neq n}^N \mathbf{A}_p^{(i)}, \quad \mathcal{M}_p^{(i,n)} = \mathcal{X}_{(i,n)} \bigcirc_{j \in \{1, \dots, N\} \setminus \{i, n\}}^N \mathbf{A}_p^{(j)},$$

$$\text{and } \mathbf{V}^{(n,i,j)} = \mathbf{A}^{(n)} \left(\left(\mathbf{A}^{(i)T} d\mathbf{A}^{(i)} \right) * \left(\mathbf{A}^{(j)T} d\mathbf{A}^{(j)} \right) * \bigstar_{k=1, k \neq i, j, n}^N \left(\mathbf{A}^{(k)T} \mathbf{A}^{(k)} \right) \right).$$

We evaluate the benefit of including the $\mathbf{V}^{(n,i,j)}$ correction in Section 5.1. Given $\mathcal{M}_p^{(i,n)}$ and $\mathbf{M}_p^{(n)}$, calculation of $\tilde{\mathbf{M}}^{(n)}$ for $n \in \{1, \dots, N\}$ requires $2N^2 (s^2 R + sR^2)$ operations overall. Further, we show in Section 4.1 that the column-wise relative approximation error of $\tilde{\mathbf{M}}^{(n)}$ with respect to $\mathbf{M}^{(n)}$ is small if each $\left\| d\mathbf{a}_k^{(n)} \right\|_2 / \left\| \mathbf{a}_k^{(n)} \right\|_2$ for $n \in \{1, \dots, N\}, k \in \{1, \dots, R\}$ is sufficiently small. Algorithm 3.1 presents the PP-CP-ALS method described above.

Algorithm 3.1 PP-CP-ALS: Pairwise perturbation procedure for CP-ALS

```

1: Input: tensor  $\mathcal{X} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ , stopping criteria  $\Delta$ , PP tolerance  $\epsilon < 1$ 
2: Initialize  $[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$  as uniformly distributed random matrices within  $[0, 1]$ , initialize  $\mathbf{G}^{(n)}, d\mathbf{A}^{(n)} \leftarrow \mathbf{A}^{(n)}$ ,  $\mathbf{S}^{(n)} \leftarrow \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$  for  $i \in \{1, \dots, N\}$ 
3: while  $\sum_{i=1}^N \left\| \mathbf{G}^{(i)} \right\|_F > \Delta \|\mathcal{X}\|_F$  do
4:   if  $\forall i \in \{1, \dots, N\}, \left\| d\mathbf{A}^{(i)} \right\|_F < \epsilon \left\| \mathbf{A}^{(i)} \right\|_F$  then
5:     Compute  $\mathcal{M}_p^{(i,n)}, \mathbf{M}_p^{(n)}$  for  $i, n \in \{1, \dots, N\}$  via dimension tree in Section 3.2.3
6:     for  $n \in \{1, \dots, N\}$  do
7:        $\mathbf{A}_p^{(n)} \leftarrow \mathbf{A}^{(n)}, d\mathbf{A}^{(n)} \leftarrow \mathbf{O}$ 
8:     end for
9:     while  $\sum_{i=1}^N \left\| \mathbf{G}^{(i)} \right\|_F > \Delta \|\mathcal{X}\|_F$  and  $\forall i \in \{1, \dots, N\}, \left\| d\mathbf{A}^{(i)} \right\|_F < \epsilon \left\| \mathbf{A}^{(i)} \right\|_F$  do
10:      for  $n \in \{1, \dots, N\}$  do
11:         $\Gamma^{(n)} \leftarrow \mathbf{S}^{(1)} * \dots * \mathbf{S}^{(n-1)} * \mathbf{S}^{(n+1)} * \dots * \mathbf{S}^{(N)}$ 
12:        Update  $\tilde{\mathbf{M}}^{(n)}$  based on Equation 3.2
13:         $\mathbf{A}_{\text{new}}^{(n)} \leftarrow \tilde{\mathbf{M}}^{(n)} \Gamma^{(n)\dagger}$ 
14:         $\mathbf{G}^{(n)} \leftarrow (\mathbf{A}^{(n)} - \mathbf{A}_{\text{new}}^{(n)}) \Gamma^{(n)}$ 
15:         $\mathbf{A}^{(n)} \leftarrow \mathbf{A}_{\text{new}}^{(n)}$ 
16:         $\mathbf{S}^{(n)} \leftarrow \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$ 
17:         $d\mathbf{A}^{(n)} = \mathbf{A}_{\text{new}}^{(n)} - \mathbf{A}_p^{(n)}$ 
18:      end for
19:    end while
20:   end if
21:   Perform regular ALS sweep as in Algorithm 2.1, taking  $d\mathbf{A}^{(n)} \leftarrow \mathbf{A}_{\text{new}}^{(n)} - \mathbf{A}^{(n)}$  for each  $n \in \{1, \dots, N\}$ 
22: end while
23: return  $[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$ 

```

3.2.2. Tucker-ALS. We derive a similar pairwise perturbation algorithm for Tucker-ALS. Similar to the expression for $\mathbf{M}^{(n)}$ in CP-ALS, $\mathcal{Y}^{(n)}$ can be expressed as

$$\mathcal{Y}^{(n)} = \mathcal{X} \bigtimes_{i=1, i \neq n}^N \left(\mathbf{A}_p^{(i)T} + d\mathbf{A}^{(i)T} \right).$$

The expression above can be rewritten as a function of $\mathcal{Y}_p^{(i_1, i_2, \dots, i_m)}$,

$$\mathcal{Y}^{(n)} = \mathcal{Y}_p^{(n)} + \sum_{i=1, i \neq n}^N \mathcal{Y}_p^{(i,n)} \times_i d\mathbf{A}^{(i)T} + \sum_{i=1, i \neq n}^N \sum_{j=i+1, j \neq n}^N \mathcal{Y}_p^{(i,j,n)} \times_i d\mathbf{A}^{(i)T} \times_j d\mathbf{A}^{(j)T} + \dots$$

The pairwise perturbation algorithm again takes only the first order terms in $d\mathbf{A}^{(i)}$, computing

$$\tilde{\mathcal{Y}}^{(n)} = \mathcal{Y}_p^{(n)} + \sum_{i=1, i \neq n}^N \mathcal{Y}_p^{(i,n)} \times_i d\mathbf{A}^{(i)T}, \quad \text{where } \mathcal{Y}_p^{(n)} = \mathcal{X} \bigtimes_{l=1, l \neq n}^N \mathbf{A}_p^{(l)T} \quad \text{and} \quad \mathcal{Y}_p^{(i,n)} = \mathcal{X} \bigtimes_{j \in \{1, \dots, N\} \setminus \{i, n\}} \mathbf{A}_p^{(j)T}.$$

Given $\mathcal{Y}_p^{(i,n)}$ and $\mathcal{Y}_p^{(n)}, \tilde{\mathcal{Y}}^{(n)}$ for $n \in \{1, \dots, N\}$ can be calculated with $2N^2 s^2 R^{N-1}$ cost overall. In Section 4.2, we show that the relative Frobenius norm approximation error of $\tilde{\mathcal{Y}}^{(n)}$ with respect to $\mathcal{Y}^{(n)}$ is small, so long as each $\|d\mathbf{A}^{(n)}\|_F / \|\mathbf{A}^{(n)}\|_F$ is sufficiently small. Algorithm 3.2 presents the PP-Tucker-ALS method described above.

Algorithm 3.2 PP-Tucker-ALS: Pairwise perturbation procedure for Tucker-ALS

```

1: Input: tensor  $\mathcal{X} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ , decomposition ranks  $\{R_1, \dots, R_N\}$ , stopping criteria  $\Delta$ , PP tolerance  $\epsilon$ 
2: Initialize  $[\mathcal{G}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$  using HOSVD, initialize  $d\mathbf{A}^{(n)} \leftarrow \mathbf{A}^{(n)}$  for  $i \in \{1, \dots, N\}$ , initialize  $\mathcal{F} \leftarrow \mathcal{G}$ 
3: while  $\|\mathcal{F}\|_F > \Delta \|\mathcal{X}\|_F$  do
4:   if  $\forall i \in \{1, \dots, N\}, \|d\mathbf{A}^{(i)}\|_F < \epsilon \|\mathbf{A}^{(i)}\|_F$  then
5:     Compute  $\mathcal{Y}_p^{(i,n)}, \mathcal{Y}_p^{(n)}$  for  $i, n \in \{1, \dots, N\}$  via dimension tree in Section 3.2.3
6:     for  $n \in \{1, \dots, N\}$  do
7:        $\mathbf{A}_p^{(n)} \leftarrow \mathbf{A}^{(n)}, d\mathbf{A}^{(n)} \leftarrow \mathbf{O}$ 
8:     end for
9:     while  $\|\mathcal{F}\|_F > \Delta \|\mathcal{X}\|_F$  and  $\|\mathcal{F}\|_F < \epsilon \|\mathcal{X}\|_F$  do
10:      for  $n \in \{1, \dots, N\}$  do
11:         $\mathcal{Y}^{(n)} \leftarrow \mathcal{Y}_p^{(n)} + \sum_{i=1, i \neq n}^N \mathcal{Y}_p^{(i,n)} \times_i d\mathbf{A}^{(i)}$ 
12:         $\mathbf{A}^{(n)} \leftarrow R_n$  leading left singular vectors of  $\mathbf{Y}_{(n)}^{(n)}$ 
13:         $d\mathbf{A}^{(n)} \leftarrow \mathbf{A}^{(n)} - \mathbf{A}_p^{(n)}$ 
14:      end for
15:       $\mathcal{G}_{\text{new}} \leftarrow \mathcal{Y}^{(N)} \times_N \mathbf{A}^{(N)T}$ 
16:       $\mathcal{F} \leftarrow \mathcal{G}_{\text{new}} - \mathcal{G}$ 
17:       $\mathcal{G} \leftarrow \mathcal{G}_{\text{new}}$ 
18:    end while
19:   end if
20:   Perform regular ALS sweep as in Algorithm 2.2, taking  $d\mathbf{A}^{(n)} \leftarrow \mathbf{A}_{\text{new}}^{(n)} - \mathbf{A}^{(n)}$  for each  $n \in \{1, \dots, N\}$ 
21:    $\mathcal{G}_{\text{new}} \leftarrow \mathcal{Y}^{(N)} \times_N \mathbf{A}^{(N)T}$ 
22:    $\mathcal{F} \leftarrow \mathcal{G}_{\text{new}} - \mathcal{G}$ 
23:    $\mathcal{G} \leftarrow \mathcal{G}_{\text{new}}$ 
24: end while
25: return  $[\mathcal{G}; \mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$ 

```

3.2.3. Dimension Trees for Pairwise Perturbation Operators. Computation of the pairwise perturbation operators $\mathcal{M}_p^{(i,n)}$ and of $\mathcal{M}_p^{(n)}$ can benefit from amortization of common tensor contraction (Khatri-Rao product or multilinear multiplication) subexpressions. In the context of ALS, this technique is known as dimension trees and has been successfully employed to accelerate TTMC and MTTKRP. The same trees can be used for both CP and Tucker, although the tensor intermediates and contraction operations are different (Khatri-Rao products for CP and multilinear multiplication for Tucker). We describe the trees for CP decomposition, computing each $\mathcal{M}_p^{(i,n)}$ and $\mathcal{M}_p^{(n)}$. Figure 1c,1d describes the dimension tree for $N = 3, 4$. Our tree constructions assume that the tensors are equidimensional, if this is not the case, the largest dimensions should be contracted first.

The main goal of the dimension tree is to perform a minimal number of contractions to obtain each $\mathcal{M}_p^{(i,n)}$. Each matrix $\mathcal{M}_p^{(n)}$ can be simply obtained by a contraction with $\mathcal{M}_p^{(i,n)}$ for any $i \neq n$. Each level of the tree for $l = 1, \dots, N-1$ should contain intermediate tensors containing $N-l+1$ uncontracted modes belonging to the original tensor (the root is the original tensor $\mathcal{X} = \mathcal{M}^{(1, \dots, N)}$). For any pair of the original tensor modes, each level should contain an intermediate for which these modes are uncontracted. Since the leaves at level $l = N-1$ have two uncontracted modes, they will include each $\mathcal{M}_p^{(i,n)}$ for $i < n$ and have $\binom{N}{2}$ tensors overall. At level l it then suffices to compute $\binom{l+1}{2}$ tensors $\mathcal{M}^{(i,j,l+2,l+3,\dots,N)}, \forall i, j \in \{1, \dots, l+1\}, i < j$. Each $\mathcal{M}^{(i,j,l+2,l+3,\dots,N)}$ can be computed by contraction of $\mathcal{M}^{(s,t,v,l+2,l+3,\dots,N)}$ and $\mathbf{A}^{(w)}$ where $\{s, t, v\} = \{i, j, w\}$ with $w = \max_{w \in \{l-1, l, l+1\} \setminus \{i, j\}}(w)$ and $s < t < v$.

TABLE 1

Cost comparison between pairwise perturbation algorithm and ALS dimension tree algorithm for CP and Tucker decompositions.

	DT ALS	PP initialization step	PP approximated step
CP	$4s^N R$	$6s^N R$	$2N^2(s^2 R + sR^2)$
Tucker	$4s^N R$	$6s^N R$	$2N^2 s^2 R^{N-1}$

The construction of pairwise perturbation operators for CP decomposition costs

$$2R \sum_{l=2}^{N-1} \binom{l+1}{2} s^{N-l+2} = 6s^N R + 12s^{N-1} R + O(s^{N-2} R^2).$$

The cost to form pairwise perturbation operators for Tucker decomposition is

$$2 \sum_{l=2}^{N-1} \binom{l+1}{2} s^{N-l+2} R^{l-1} = 6s^N R + 12s^{N-1} R^2 + O(s^{N-2} R^3).$$

We summarize the leading order computational costs for both algorithms in Table 1. The PP initialization step, which involves the PP operator construction and does one more first level contraction, is computationally 1.5X more expensive than the ALS algorithm.

As for the memory footprint, ALS with the best choice of dimension tree requires intermediate tensors of size $O(s^{\lceil N/2 \rceil} R)$. As an example, for the order four case shown in Figure 1b, the first and second level contractions are combined to save memory, so that $\mathbf{M}^{(3,4)}$ and $\mathbf{M}^{(1,2)}$ are stored, both of size $O(s^2 R)$. The PP dimension tree described above and in Figure 1d needs at least $O(s^{N-1} R)$ auxiliary memory to store the first level contraction results. The memory needed for PP can be reduced similar to ALS. For example, when calculating the PP operator $\mathbf{M}_p^{(1,3)}$ for an order four tensor, we can bypass the first level contraction and save its memory via directly performing a contraction between the input tensor and the Khatri-Rao product output $\mathbf{A}^{(1)} \odot \mathbf{A}^{(3)}$. Combining the first $l \leq N - 2$ levels of contractions requires $O(s^{N-l} R + N^2 s^2 R)$ auxiliary memory, but incurs a cost of $O(l^2 s^{N-1} R)$.

4. Error Analysis. In this section, we formally bound the approximation error of the pairwise perturbation algorithm relative to ALS. We show that quadratic optimization problems computed by pairwise perturbation differ only slightly from ALS so long as the factor matrices have not changed significantly since the construction of the pairwise perturbation operators.

4.1. CP-ALS. To bound the error of pairwise perturbation, we view the ALS procedure for CP decomposition in terms of pairwise updates (Algorithm 4.1), pushing updates to least-squares problems of all tensors as soon as any one of them is updated. This reformulation is algebraically equivalent to Algorithm 2.1, but makes oracle-like use of $\mathbf{M}^{(m,n)}$ (Equation 2.1), recomputing which would increase the computational cost. We can bound the error of the way pairwise perturbation propagates updates to any right-hand side $\mathbf{M}^{(m)}$ due to changes in any one of the other factor matrices $\delta \mathbf{A}^{(n)}$. We define the update $\mathbf{H}^{(m,n)}$ in terms of its columns,

$$\mathbf{h}_k^{(m,n)}(x) = \sum_{y=1}^{s_n} \mathbf{M}^{(m,n)}(x, y, k) \delta \mathbf{A}^{(n)}(y, k), \quad \text{where } \delta \mathbf{A}^{(n)} = \mathbf{A}_{\text{new}}^{(n)} - \mathbf{A}^{(n)}.$$

Note that $\delta \mathbf{A}^{(n)}$ denotes the update of n th factor between two neighboring sweeps, which should be distinguished from $d\mathbf{A}^{(n)}$, denoting the perturbation of n th factor in PP. Based on the definition, the update of each $\mathbf{M}^{(m)}$ after an ALS sweep is the summation of $\mathbf{H}^{(m,n)}$ expressed as $\delta \mathbf{M}^{(m)} = \sum_{n=1, n \neq m}^N \mathbf{H}^{(m,n)}$.

For simplicity, we first perform an error analysis for the case where the second order correction terms $\mathbf{V}^{(n,i,j)}$ are not included in PP. In Theorem 4.1, we prove that when the column-wise norm of $d\mathbf{A}^{(n)} = \mathbf{A}^{(n)} - \mathbf{A}_p^{(n)}$ relative to the norm of $\mathbf{A}^{(n)}$ for $n \in \{1, \dots, N\}$ is small, the absolute error of column-wise results for $\mathbf{H}^{(m,n)}$ calculated from pairwise perturbation with respect to that calculated from exact ALS is also small. Corollary 4.2 provides a simple relative error bound for third-order tensors. Overall, these bounds demonstrate that pairwise perturbation should generally compute updates with small relative error with

Algorithm 4.1 CP-ALS: Reinterpreted ALS procedure for CP decomposition

```

1: Input: Tensor  $\mathcal{X} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ , stopping criteria  $\Delta$ 
2: Initialize  $[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$  as uniformly distributed random matrices within  $[0, 1]$ , initialize  $\mathbf{G}^{(n)}, \delta \mathbf{A}^{(n)} \leftarrow \mathbf{A}^{(n)}$ ,
    $\mathbf{S}^{(n)} \leftarrow \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$  for  $i \in \{1, \dots, N\}$ 
3: for  $n \in \{1, \dots, N\}$  do
4:   Update  $\mathbf{M}^{(n)}$  based on the dimension tree algorithm shown in Figure 1
5: end for
6: while  $\sum_{i=1}^N \|\mathbf{G}^{(i)}\|_F > \Delta \|\mathcal{X}\|_F$  do
7:   for  $n \in \{1, \dots, N\}$  do
8:      $\mathbf{\Gamma}^{(n)} \leftarrow \mathbf{S}^{(1)} * \dots * \mathbf{S}^{(n-1)} * \mathbf{S}^{(n+1)} * \dots * \mathbf{S}^{(N)}$ 
9:      $\mathbf{A}_{\text{new}}^{(n)} \leftarrow \mathbf{M}^{(n)} \mathbf{\Gamma}^{(n)\dagger}$ 
10:     $\delta \mathbf{A}^{(n)} = \mathbf{A}_{\text{new}}^{(n)} - \mathbf{A}^{(n)}$ 
11:     $\mathbf{G}^{(n)} \leftarrow -\delta \mathbf{A}^{(n)} \mathbf{\Gamma}^{(n)}$ 
12:     $\mathbf{A}^{(n)} \leftarrow \mathbf{A}_{\text{new}}^{(n)}$ 
13:     $\mathbf{S}^{(n)} \leftarrow \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$ 
14:    for  $m \in \{1, \dots, N\}, m \neq n$  do
15:      Update  $\mathbf{M}^{(m)}$  as  $\mathbf{M}^{(m)}(x, k) = \mathbf{M}^{(m)}(x, k) + \sum_{y=1}^{s_n} \mathcal{M}^{(m,n)}(x, y, k) \delta \mathbf{A}^{(n)}(y, k)$ 
16:    end for
17:   end for
18: end while
19: return  $[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]$ 

```

respect to the magnitude of the perturbation of the factor matrices since the setup of the pairwise operators. However, this relative error can be amplified during other steps of ALS, which are ill-conditioned, i.e., can suffer from catastrophic cancellation (the same would hold for round-off error).

We then perform an error analysis for the case where the second order correction terms $\mathbf{V}^{(n,i,j)}$ are included in PP in Theorem 4.3. We show that the second order corrections can tighten the leading order error by a factor related to the CP decomposition accuracy.

THEOREM 4.1. *For $k \in \{1, \dots, R\}$, if $\|\mathbf{d}\mathbf{a}_k^{(l)}\|_2 / \|\mathbf{a}_k^{(l)}\|_2 \leq \epsilon < 1$ for all $l \in \{1, \dots, N\}$, the pairwise perturbation algorithm without second order corrections computes the update $\tilde{\mathbf{H}}^{(1,N)}$ with columnwise error,*

$$\|\tilde{\mathbf{h}}_k^{(1,N)} - \mathbf{h}_k^{(1,N)}\|_2 = O(N\epsilon) \|\hat{\mathbf{T}}\|_2 \prod_{j=2}^{N-1} \|\mathbf{a}_k^{(j)}\|_2,$$

where $\mathbf{H}^{(1,N)}$ is the update to the matrix $\mathbf{M}^{(1)}$ due to the change $\delta \mathbf{A}^{(N)}$ performed by a regular ALS sweep, and $\hat{\mathbf{T}} = \mathcal{X} \times_N \delta \mathbf{a}_k^{(N)T}$. Analogous bounds hold for $\mathbf{H}^{(m,n)}$ for any $m, n \in \{1, \dots, N\}, m \neq n$.

Proof. The ALS update and approximated update are

$$(4.1) \quad \mathbf{h}_k^{(1,N)} = \hat{\mathbf{T}} \times_{i \in \{2, \dots, N-1\}} \mathbf{a}_k^{(i)T} \quad \text{and} \quad \tilde{\mathbf{h}}_k^{(1,N)} = \hat{\mathbf{T}} \times_{i \in \{2, \dots, N-1\}} \left(\mathbf{a}_k^{(i)T} - \mathbf{d}\mathbf{a}_k^{(i)T} \right).$$

We can expand the error as

$$(4.2) \quad \tilde{\mathbf{h}}_k^{(1,N)} - \mathbf{h}_k^{(1,N)} = \sum_{S \subset \{2, \dots, N-1\}, S \neq \emptyset} \hat{\mathbf{T}} \times_{i \in \{2, \dots, N-1\}} \mathbf{v}_k^{(i)T}, \quad \text{where } \mathbf{v}_k^{(i)} = \begin{cases} -\mathbf{d}\mathbf{a}_k^{(i)} & : i \in S \\ \mathbf{a}_k^{(i)} & : i \notin S \end{cases}.$$

Consequently, we can upper-bound the error due to terms with $|S| = d$ by

$$\binom{N-2}{d} \epsilon^d \|\hat{\mathbf{T}}\|_2 \prod_{j=2}^{N-1} \|\mathbf{a}_k^{(j)}\|_2 = O(N\epsilon)^d \|\hat{\mathbf{T}}\|_2 \prod_{j=2}^{N-1} \|\mathbf{a}_k^{(j)}\|_2.$$

Therefore, the error bound when $|S| = d$ scales as $O(N\epsilon)^d$, and the leading order error is $O(N\epsilon)$. \square

Note that this error bound involves $\hat{\mathbf{T}}$, which is small in norm due to being constructed from contraction with $\delta \mathbf{a}_k^{(N)}$. Thus, the error norm generally scales as $O(\epsilon^2)$ relative to the norm of the original tensor \mathcal{X} , since $O(N\epsilon) \|\hat{\mathbf{T}}\|_2 \prod_{j=2}^{N-1} \|\mathbf{a}_k^{(j)}\|_2 = O(N\epsilon^2) \|\mathcal{X}\|_2 \prod_{j=2}^{N-1} \|\mathbf{a}_k^{(j)}\|_2$.

COROLLARY 4.2. For $N = 3$, using the bounds from the proof of Theorem 4.1, under the same assumptions, we obtain the absolute error bound,

$$\left\| \tilde{\mathbf{h}}_k^{(1,3)} - \mathbf{h}_k^{(1,3)} \right\|_2 \leq \left\| \hat{\mathbf{T}} \right\|_2 \left\| \mathbf{a}_k^{(2)} \right\|_2 \epsilon,$$

where $\hat{\mathbf{T}} = \mathcal{X} \times_3 \delta \mathbf{a}_k^{(3)T}$. Further, since $\mathbf{h}_k^{(1,3)} = \hat{\mathbf{T}} \mathbf{a}_k^{(2)}$, the relative error is bounded by

$$\frac{\left\| \tilde{\mathbf{h}}_k^{(1,3)} - \mathbf{h}_k^{(1,3)} \right\|_2}{\left\| \mathbf{h}_k^{(1,3)} \right\|_2} \leq \kappa(\hat{\mathbf{T}}) \epsilon.$$

From Theorem 4.1, we can conclude that the relative error in computing any column update $\mathbf{h}_k^{(i,j)}$ is $O(\epsilon)$ when $\epsilon \ll 1$ and the correct update is sufficiently large, e.g., for $i = 1$ and $j = N$, $\left\| \mathbf{h}_k^{(1,N)} \right\|_2 = \Omega\left(\left\| \hat{\mathbf{T}} \right\|_2 \prod_{i=2}^{N-1} \left\| \mathbf{a}_k^{(i)} \right\|_2\right)$. When this is the case, we can also bound the error of the update to the columns of the right-hand sides $\delta \mathbf{M}^{(n)}$ formed in ALS, so long as the sum of the updates $\mathbf{H}^{(n,m)}$ for $m \neq n$ is not too small in norm relative to each update matrix.

We now perform analysis for the case where the second order corrections $\mathbf{V}^{(n,i,j)}$ are included in PP.

THEOREM 4.3. For $k \in \{1, \dots, R\}$, if $\left\| d\mathbf{a}_k^{(l)} \right\|_2 / \left\| \mathbf{a}_k^{(l)} \right\|_2 \leq \epsilon < 1$ for all $l \in \{1, \dots, N\}$, the pairwise perturbation algorithm with second order correction terms computes the update term $\tilde{\mathbf{H}}^{(1,N)}$ with columnwise error,

$$\left\| \tilde{\mathbf{h}}_k^{(1,N)} - \mathbf{h}_k^{(1,N)} \right\|_2 = O(N\epsilon) \left\| \hat{\mathbf{P}} - \hat{\mathbf{T}} \right\|_2 \prod_{j=2}^{N-1} \left\| \mathbf{a}_k^{(j)} \right\|_2 + O((N\epsilon)^2) \left\| \hat{\mathbf{T}} \right\|_2 \prod_{j=2}^{N-1} \left\| \mathbf{a}_k^{(j)} \right\|_2,$$

where $\hat{\mathbf{P}} = \mathcal{Z} \times_N \delta \mathbf{a}_k^{(N)T}$, and \mathcal{Z} denotes the approximate CP decomposition of \mathcal{X} , $[[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}]]$. $\mathbf{H}^{(1,N)}$ is the update to the matrix $\mathbf{M}^{(1)}$ due to the change $\delta \mathbf{A}^{(N)}$ performed by a regular ALS sweep, and $\hat{\mathbf{T}} = \mathcal{X} \times_N \delta \mathbf{a}_k^{(N)T}$. Analogous bounds hold for $\mathbf{H}^{(m,n)}$ for any $m, n \in \{1, \dots, N\}$, $m \neq n$.

Proof. The ALS approximated update is

$$(4.3) \quad \tilde{\mathbf{h}}_k^{(1,N)} = \hat{\mathbf{T}} \times_{i \in \{2, \dots, N-1\}} \left(\mathbf{a}_k^{(i)T} - d\mathbf{a}_k^{(i)T} \right) + \sum_{i \in \{2, \dots, N-1\}} \hat{\mathbf{P}} \times_i d\mathbf{a}_k^{(i)T} \times_{j \in \{2, \dots, N-1\}, j \neq i} \mathbf{a}_k^{(j)T}.$$

We can expand the error as

$$\tilde{\mathbf{h}}_k^{(1,N)} - \mathbf{h}_k^{(1,N)} = \sum_{S \subset \{2, \dots, N-1\}, |S| \geq 2} \hat{\mathbf{T}} \times_{i \in \{2, \dots, N-1\}} \mathbf{v}_k^{(i)T} + \sum_{i \in \{2, \dots, N-1\}} \left(\hat{\mathbf{P}} - \hat{\mathbf{T}} \right) \times_i d\mathbf{a}_k^{(i)T} \times_{j \in \{2, \dots, N-1\}, j \neq i} \mathbf{a}_k^{(j)T},$$

where $\mathbf{v}_k^{(i)} = -d\mathbf{a}_k^{(i)}$ if $i \in S$ and $\mathbf{v}_k^{(i)} = \mathbf{a}_k^{(i)}$ otherwise. By the same analysis as in Theorem 4.1, the error due to each term with $|S| = d$, $d \geq 2$ can be bounded as $O(N\epsilon)^d \left\| \hat{\mathbf{T}} \right\|_2 \prod_{j=2}^{N-1} \left\| \mathbf{a}_k^{(j)} \right\|_2$. We can then upper-bound the error due to the second term by

$$(4.4) \quad O(N\epsilon) \left\| \hat{\mathbf{P}} - \hat{\mathbf{T}} \right\|_2 \prod_{j=2}^{N-1} \left\| \mathbf{a}_k^{(j)} \right\|_2,$$

thus completing the proof. \square

From Theorem 4.3, we can conclude that when the approximate CP decomposition is close to \mathcal{X} , the term expressed in Equation 4.4 will have small magnitude, making the absolute error second order accurate in terms of ϵ .

In Appendix 8.4, we also obtain relative error bounds on MTTKRPs (the right-hand sides in the linear least squares subproblems). However, this error bound is relative to the condition number of \mathcal{X} (defined in Appendix 8.1), which is infinite for sufficiently large tensors.

4.2. Tucker-ALS. For Tucker decomposition, the pairwise perturbation approximation satisfies better bounds than for CP decomposition, due to the orthogonality of the factor matrices. We can not only obtain the similar bound as Theorem 4.1, but also obtain stronger results assuming that either the residual of the Tucker decomposition is bounded (it suffices that the decomposition achieves one digit of accuracy in residual) or that the ratio of rank to dimension is not too large. We demonstrate that

- similar to Algorithm 4.1 and Theorem 4.1, when we view the ALS procedure for Tucker decomposition of equidimensional tensors in terms of pairwise updates, we can bound the error of updates to any right-hand side $\mathbf{Y}^{(m)}$ due to changes in any one of the other factor matrices $\delta\mathbf{A}^{(n)}$. We define the update $\mathcal{J}^{(m,n)}$ as

$$\mathcal{J}^{(m,n)} = \mathbf{Y}^{(m,n)} \times_n \delta\mathbf{A}^{(n)T}, \quad \text{where} \quad \delta\mathbf{A}^{(n)} = \mathbf{A}_{\text{new}}^{(n)} - \mathbf{A}^{(n)}.$$

The columnwise absolute error bound for MTTKRP holds for $\mathcal{J}^{(m,n)}$ when the column-wise 2-norm relative perturbations of the input matrices are bounded by $O(\epsilon)$ (Theorem 4.4),

- the relative error of $\mathbf{Y}^{(m)}$ for $m \in \{1, \dots, N\}$ satisfies the bound of $O(\epsilon^2)$, so long as the residual of Tucker decomposition is small (Theorem 4.6),
- the relative error of $\mathbf{Y}^{(m)}$ for $m \in \{1, \dots, N\}$ is bounded in Frobenius norm by $O(\epsilon^2)$ for a fixed problem size assuming that HOSVD is performed to initialize Tucker-ALS (Theorem 4.9).

THEOREM 4.4. *For an order N tensor \mathcal{X} with dimension sizes s , if $\|\mathbf{d}\mathbf{a}_k^{(n)}\|_2 / \|\mathbf{a}_k^{(n)}\|_2 \leq \epsilon < 1$ for all $n \in \{1, \dots, N\}, k \in \{1, \dots, R\}$, the pairwise perturbation algorithm computes update $\mathcal{J}^{(1,N)}$ with error,*

$$\left\| \tilde{\mathbf{j}}_{i_2, \dots, i_N}^{(1,N)} - \mathbf{j}_{i_2, \dots, i_N}^{(1,N)} \right\|_2 = O(N\epsilon) \left\| \hat{\mathbf{T}} \right\|_2 \prod_{j=2}^{N-1} \left\| \mathbf{a}_k^{(j)} \right\|_2,$$

where $\hat{\mathbf{T}} = \mathcal{X} \times_N \delta\mathbf{a}_{i_N}^{(N)T}$ and $\mathbf{j}_{i_2, \dots, i_N}^{(1,N)}(x) = \mathcal{J}^{(1,N)}(x, i_2, \dots, i_N)$.

Proof. The proof is similar to that of Theorem 4.1. The ALS update and approximated update after a change $\delta\mathbf{A}^{(N)}$ are

$$\mathbf{j}_{i_2, \dots, i_N}^{(1,N)} = \hat{\mathbf{T}} \times_{j=2}^{N-1} \mathbf{a}_{i_j}^{(j)T} \quad \text{and} \quad \tilde{\mathbf{j}}_{i_2, \dots, i_N}^{(1,N)} = \hat{\mathbf{T}} \times_{j=2}^{N-1} \left(\mathbf{a}_{i_j}^{(j)T} - \mathbf{d}\mathbf{a}_{i_j}^{(j)T} \right).$$

The error bound proceeds by analogy to the proof of Theorem 4.1. \square

Using Lemma 2.1, we prove in Lemma 4.5 that after contracting a tensor with a matrix with orthonormal columns, whose row length is higher or equal to the column length, the contracted tensor norm is the same as the original tensor norm.

LEMMA 4.5. *Given tensor $\mathcal{G} \in \mathbb{R}^{r_1 \times \dots \times r_N}$, the mode- n product for any $n \in \{1, \dots, N\}$, with a matrix with orthonormal columns $\mathbf{M} \in \mathbb{R}^{s \times r_n}$, $r_n \leq s$, satisfies $\|\mathcal{G}\|_2 = \|\mathcal{G} \times_n \mathbf{M}\|_2$.*

Proof. Based on the submultiplicative property of the tensor norm (Lemma 2.1),

$$\|\mathcal{G}\|_2 = \|\mathcal{G} \times_n (\mathbf{M}^T \mathbf{M})\|_2 = \|\mathcal{G} \times_n \mathbf{M} \times_n \mathbf{M}^T\|_2 \leq \|\mathcal{G} \times_n \mathbf{M}\|_2 \|\mathbf{M}^T\|_2 = \|\mathcal{G} \times_n \mathbf{M}\|_2,$$

and simultaneously, $\|\mathcal{G} \times_n \mathbf{M}\|_2 \leq \|\mathcal{G}\|_2 \|\mathbf{M}\|_2 = \|\mathcal{G}\|_2$. \square

Using Lemma 4.5, we prove in Theorem 4.6 that when the relative error of the matrices $\mathbf{A}^{(n)}$ for $n \in \{1, \dots, N\}$ is small and the residual of the Tucker decomposition is loosely bounded, the relative error bound for the $\mathbf{Y}^{(n)}$ is independent of the tensor condition number defined in Section 8.

THEOREM 4.6. *Given tensor $\mathcal{X} \in \mathbb{R}^{s_1 \times \dots \times s_N}$, if $\|\mathbf{d}\mathbf{A}^{(n)}\|_2 \leq \epsilon \ll 1$ for $n \in \{1, \dots, N\}$ and $\|\mathcal{X} - [\mathcal{G}; \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)}]\|_2 \leq \frac{1}{3} \|\mathcal{X}\|_2$, $\tilde{\mathbf{Y}}^{(n)}$ is constructed with error,*

$$\frac{\|\tilde{\mathbf{Y}}^{(n)} - \mathbf{Y}^{(n)}\|_2}{\|\mathbf{Y}^{(n)}\|_2} = O(\epsilon^2).$$

Proof.

$$\frac{\|\tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)}\|_2}{\|\mathbf{y}^{(n)}\|_2} \leq \binom{N}{2} \max_{i,j} \frac{\|\mathbf{y}_p^{(i,j,n)} \times_i d\mathbf{A}^{(i)T} \times_j d\mathbf{A}^{(j)T}\|_2}{\|\mathbf{y}^{(n)}\|_2} \leq \binom{N}{2} \max_{i,j} \frac{\|\mathbf{y}_p^{(i,j,n)}\|_2 \|d\mathbf{A}^{(i)}\|_2 \|d\mathbf{A}^{(j)}\|_2}{\|\mathbf{y}^{(n)}\|_2}.$$

Let $\tilde{\mathbf{x}} = \llbracket \mathbf{g}; \mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(N)} \rrbracket$, $\mathcal{R} = \mathbf{x} - \tilde{\mathbf{x}}$. Define the tensors $\mathcal{Z}^{(i,j,n)}$ by contraction of \mathcal{R} with all except three factor matrices,

$$\mathcal{Z}^{(i,j,n)} = \mathcal{R} \times_{r \in \{1, \dots, N\} \setminus \{i, j, n\}} \mathbf{A}^{(r)T}.$$

For $\|\mathbf{x} - \tilde{\mathbf{x}}\|_2 = \|\mathcal{R}\|_2 \leq \frac{1}{3}\|\mathbf{x}\|_2$, we have $\frac{2}{3}\|\mathbf{x}\|_2 \leq \|\tilde{\mathbf{x}}\|_2 \leq \frac{4}{3}\|\mathbf{x}\|_2$. Based on Lemma 4.5,

$$\begin{aligned} \|\mathbf{y}^{(n)}\|_2 &= \|\mathbf{g} \times_n \mathbf{A}^{(n)} + \mathcal{Z}^{(i,j,n)} \times_i \mathbf{A}^{(i)T} \times_j \mathbf{A}^{(j)T}\|_2 \geq \|\mathbf{g}\|_2 - \|\mathcal{Z}^{(i,j,n)}\|_2 \|\mathbf{A}^{(i)T}\|_2 \|\mathbf{A}^{(j)T}\|_2 \\ &\geq \|\mathbf{g}\|_2 - \|\mathcal{R}\|_2 \geq \frac{1}{3}\|\mathbf{x}\|_2. \end{aligned}$$

Additionally,

$$\|\mathbf{y}^{(i,j,n)}\|_2 = \|\mathbf{g} \times_i \mathbf{A}^{(i)} \times_j \mathbf{A}^{(j)} \times_n \mathbf{A}^{(n)} + \mathcal{Z}^{(i,j,n)}\|_2 \leq \|\mathbf{g}\|_2 + \|\mathcal{R}\|_2 \leq \frac{5}{3}\|\mathbf{x}\|_2.$$

Therefore,

$$\frac{\|\tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)}\|_2}{\|\mathbf{y}^{(n)}\|_2} \leq \binom{N}{2} \max_{i,j} \frac{\|\mathbf{y}_p^{(i,j,n)}\|_2 \|d\mathbf{A}^{(i)}\|_2 \|d\mathbf{A}^{(j)}\|_2}{\|\mathbf{y}^{(n)}\|_2} \leq \binom{N}{2} \frac{5}{3} \frac{\|\mathbf{x}\|_2 \epsilon^2}{\frac{1}{3}\|\mathbf{x}\|_2} = O(\epsilon^2). \quad \square$$

We now derive a Frobenius norm error bound that is independent of residual norm and tensor condition number, and is based the ratio of the tensor dimensions and the Tucker rank. We arrive at this result (Theorem 4.9) by obtaining a lower bound on the residual achieved by the HOSVD (Lemmas 4.7 and 4.8).

LEMMA 4.7. *Given tensor $\mathbf{x} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ and matrix $\mathbf{A} \in \mathbb{R}^{R \times s_n}$, where $R < \max\{s_n, \prod_{i=1, i \neq n}^N s_i\}$ and \mathbf{A} consists of R leading left singular vectors of $\mathbf{X}_{(n)}$. Let $\mathcal{Z} = \mathbf{x} \times_n \mathbf{A}$, $\|\mathbf{x}\|_F \geq \|\mathcal{Z}\|_F \geq \sqrt{\frac{R}{s_n}} \|\mathbf{x}\|_F$.*

Proof. The singular values of $\mathbf{A}\mathbf{X}_{(n)}$ are the first R singular values of $\mathbf{X}_{(n)}$. Since the square of the Frobenius norm of a matrix is the sum of the squares of the singular values, $\|\mathcal{Z}\|_F^2 = \|\mathbf{A}\mathbf{X}_{(n)}\|_F^2 \geq (R/s_n) \|\mathbf{X}_{(n)}\|_F^2 = (R/s_n) \|\mathbf{x}\|_F^2$ and $\|\mathcal{Z}\|_F \leq \|\mathbf{x}\|_F$. \square

LEMMA 4.8. *For any equidimensional order N tensor \mathbf{x} with size s , $\|\mathbf{y}^{(n)}\|_F \geq (\frac{R}{s})^{N/2} \|\mathbf{x}\|_F$ if Tucker-ALS starts from an interlaced HOSVD.*

Proof. In Tucker-ALS, $\|\mathbf{g}\|_F$ is strictly increasing after each Tucker iteration, where \mathbf{g} is \mathbf{x} 's HOSVD core tensor. Since the interlaced SVD computes each $\mathbf{A}^{(n)}$ from the truncated SVD of the product of \mathbf{x} and the first $n-1$ factor matrices, we can apply Lemma 4.7 N times,

$$\begin{aligned} \|\mathbf{x} \times_1 \mathbf{A}^{(1)T} \dots \times_{N-1} \mathbf{A}^{(N-1)T}\|_F &\geq \|\mathbf{g}\|_F \geq \sqrt{\frac{R}{s}} \|\mathbf{x} \times_1 \mathbf{A}^{(1)T} \dots \times_{N-1} \mathbf{A}^{(N-1)T}\|_F, \\ &\vdots \\ \|\mathbf{x}\|_F &\geq \|\mathbf{g}\|_F \geq (R/s)^{N/2} \|\mathbf{x}\|_F. \end{aligned} \quad \square$$

THEOREM 4.9. *Given any equidimensional order N tensor \mathbf{x} with size s , if $\|d\mathbf{A}^{(n)}\|_F \leq \epsilon$ for $n \in [1, N]$, $\tilde{\mathbf{y}}^{(n)}$ is constructed with error,*

$$\frac{\|\tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)}\|_F}{\|\mathbf{y}^{(n)}\|_F} = O\left(\epsilon^2 \left(\frac{s}{R}\right)^{N/2}\right),$$

assuming that HOSVD is used to initialize Tucker-ALS and the residual associated with factor matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(n)}$ is no higher than that attained by HOSVD.

Proof.

$$\frac{\|\tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)}\|_F}{\|\mathbf{y}^{(n)}\|_F} \leq \binom{N}{2} \max_{i,j} \frac{\|\mathbf{y}_p^{(i,j,n)} \times_i d\mathbf{A}^{(i)T} \times_j d\mathbf{A}^{(j)T}\|_F}{\|\mathbf{y}^{(n)}\|_F}.$$

From Lemma 4.8, we have

$$\frac{\|\mathbf{y}_p^{(i,j,n)} \times_i d\mathbf{A}^{(i)T} \times_j d\mathbf{A}^{(j)T}\|_F}{\|\mathbf{y}^{(n)}\|_F} \leq \frac{\|\mathbf{X}\|_F \|d\mathbf{A}^{(i)}\|_F \|d\mathbf{A}^{(j)}\|_F}{\left(\frac{R}{s}\right)^{N/2} \|\mathbf{X}\|_F}.$$

Consequently, we can bound the relative error by

$$\frac{\|\tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)}\|_F}{\|\mathbf{y}^{(n)}\|_F} \leq \binom{N}{2} (s/R)^{N/2} \max_{i,j} \|d\mathbf{A}^{(i)}\|_F \|d\mathbf{A}^{(j)}\|_F = O\left(\epsilon^2 \left(\frac{s}{R}\right)^{N/2}\right). \quad \square$$

5. Experiments. We evaluate the performance of the pairwise perturbation algorithms on both synthetic tensors and application datasets. The synthetic experiments enable us to test tensors with known factors and to measure how effectively the algorithm works across many problem instances. We also consider publicly available tensor datasets as well as tensors of interest for quantum chemistry calculations and demonstrate the effectiveness of our algorithms on practical problems. We focus on the experiments on CP decomposition, because for many cases in Tucker decomposition, performing HOSVD and running CP decomposition on a much smaller core tensor is sufficient for getting accurate results.

We use the metrics *relative residual* and *fitness* to evaluate the convergence of the decomposition. Let $\tilde{\mathbf{X}}$ denote the tensor reconstructed by the factor matrices and the core tensor, the relative residual and fitness are defined as follows,

$$r = \frac{\|\mathbf{X} - \tilde{\mathbf{X}}\|_F}{\|\mathbf{X}\|_F}, \quad f = 1 - r.$$

We compare the performance of our own implementations of regular ALS with dimension trees to the pairwise perturbation algorithm. Both algorithms are implemented in Python with NumPy for sequential calculation and with Cyclops Tensor Framework (v1.5.5) [55], which is a distributed-memory library for matrix/tensor contractions that uses MPI for interprocessor communication and OpenMP for threading. We also make use of a wrapper Cyclops provides for ScaLAPACK [11] routines to solve symmetric positive definite linear systems of equations and compute the SVD¹.

The experimental results are collected on the Stampede2 supercomputer located at the University of Texas at Austin. We leverage the Knight's Landing (KNL) nodes exclusively, each of which consists of 68 cores, 96 GB of DDR RAM, and 16 GB of MCDRAM. These nodes are connected via a 100 Gb/sec fat-tree Omni-Path interconnect. For both NumPy and Cyclops implementations, we use Intel compilers and the MKL library for threaded BLAS routines, including batched BLAS routines, which are efficient for Khatri-Rao products arising in MTTKRP in CP decomposition, and employ the HPTT library [56] for high-performance tensor transposition. All storage and computation assumes the tensors are dense.

5.1. Sequential Experimental Results. We collect the sequential results on one KNL node on Stampede2, leveraging 64 threads for MKL and HPTT routines.

We compare the per-sweep time of the ALS dimension tree to the pairwise perturbation initialization and approximated sweep in Figure 2. Each initialization sweep constructs the PP operators and updates all the factor matrices, while an approximated sweep computes approximate updates to all the factor matrices using the PP operators constructed in the last initialization sweep. We also provide the reference per-sweep time of the ALS implementation from MATLAB Tensor Toolbox [36]. As can be seen, both ALS sweep times on top of NumPy and Cyclops are comparable to the Tensor Toolbox. For both decompositions and all the configurations, the time of an PP initialization sweep is 1.5-2.0X the time of a dimension tree based ALS sweep, while the approximated steps can have up to 6.3X speed-up for an order three tensor and 33.0X speed-up for an order six tensor for CP, and up to 10.6X speed-up for an order 6 tensor for Tucker. In

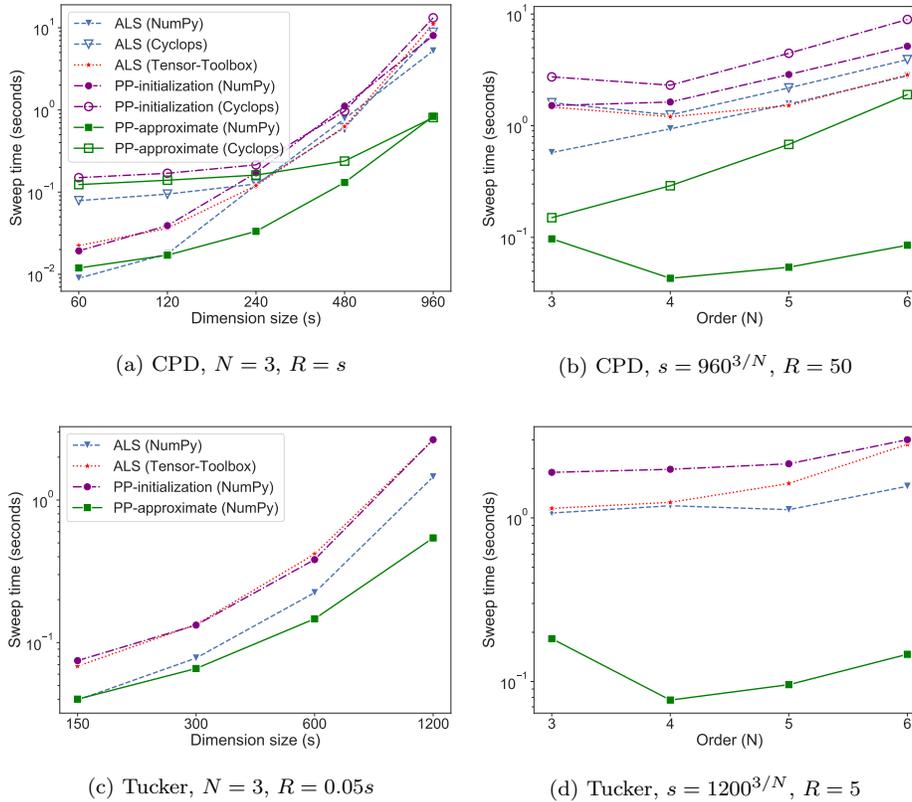


FIGURE 2. Sequential ALS sweep time comparison for both CP and Tucker decompositions. Results are taken as the mean time across 5 sweeps. The line label of (b) is the same as (a), of (d) is the same as (c). In (a)(c), we vary the dimension size and the decomposition rank, and fix the input tensor order. In (b)(d), we vary the input tensor order, and fix the input tensor size and the decomposition rank.

addition, larger speed-up can be achieved with the increase of dimension size s and the tensor order N , which is consistent with Table 1.

We use five different tensors to test the sequential performance of pairwise perturbation. Sequential performance results are collected using NumPy, as NumPy has better sequential performance than Cyclops, as shown in Figures 2a and 2b. For all the experiments, the pairwise perturbation tolerance is set as 0.1 for CP decomposition, and set as 0.3 for Tucker decomposition.

1. **Tensors with random collinearity** [9]. We create tensors based on known randomly-generated factor matrices $\mathbf{A}^{(n)}$. The factor matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{s \times R}$ are randomly generated so that the columns have collinearity defined based on a scalar C (selected randomly for the tensor from a given interval $[a, b)$), so that

$$\frac{\langle \mathbf{a}_i^{(n)}, \mathbf{a}_j^{(n)} \rangle}{\|\mathbf{a}_i^{(n)}\|_2 \|\mathbf{a}_j^{(n)}\|_2} = C, \quad \forall i, j \in \{1, \dots, R\}, i \neq j.$$

Higher collinearity corresponds to greater overlap between columns within each factor matrix, which makes the convergence of CP-ALS slower [53].

2. **Tensors made by random matrices**. We create tensors based on known uniformly distributed randomly generated factor matrices $\mathbf{A}^{(n)} \in [0, 1]^{s \times R}$,

$$\mathcal{X} = \left[\left[\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)} \right] \right].$$

In the experiments, we set R to be the same as the decomposition rank.

¹All of our code is available at https://github.com/LinjianMa/tensor_decompositions.

3. **Quantum chemistry tensor.** We also test on the density fitting intermediate tensor arising in quantum chemistry, which is the Cholesky factor of the two-electron integral tensor [27, 30]. For an order 4 two-electron integral tensor \mathcal{T} , its Cholesky factor is an order 3 tensor \mathcal{D} , with their relations shown as follows:

$$\mathcal{T}(a, b, c, d) = \sum_{s=1}^P \mathcal{D}(a, b, s) \mathcal{D}(c, d, s),$$

where P is the third mode dimension size of \mathcal{D} . CP decomposition can be performed on \mathcal{D} to provide the compressed form of the density fitting intermediate and can be used to speed up post Hartree-Fork calculations [28]. We generate the density fitting tensor via the PySCF library [57], which represents the compressed restricted Hartree-Fock wave function of an 8 water molecule chain system with a STO3G basis set. The generated tensor has size $904 \times 56 \times 56$. We set the CP rank to be 400.

4. **COIL dataset.** COIL-100 is an image-recognition data set that contains images of objects in different poses [46] and has been used previously as a tensor decomposition benchmark [9, 61]. There are 100 different object classes, each of which is imaged from 72 different angles. Each image has 128×128 pixels in three color channels. Transferring the data into tensor format, we have a $128 \times 128 \times 3 \times 7200$ tensor. We fix the CP decomposition rank to be 15 and the Tucker decomposition rank to be $10 \times 10 \times 3 \times 50$.
5. **Time-Lapse hyperspectral radiance images.** We consider the 3D hyperspectral imaging dataset called ‘‘Souto wood pile’’ [45]. The dataset is usually used on the benchmark of nonnegative tensor decomposition [7, 39]. The hyperspectral data consists of a tensor with dimensions $1024 \times 1344 \times 33 \times 9$. We fix the CP decomposition rank to be 50 and the Tucker decomposition rank to be $100 \times 100 \times 3 \times 3$.

The order three tensors are tested to justify the relative error bound shown in Section 4.1. The performance of PP on higher order CP decompositions is also considered. We focus on the high rank CP decomposition, because for the cases with rank $R < s$, Tucker decomposition or HOSVD can be used to effectively compress the input tensor from dimensions of size s to R , and then CP decomposition can be performed.

We test the synthetic tensors for CP decomposition. These tensors are all generated based on known factor matrices whose column sizes are equal to the decomposition rank, so these tensors have exact decompositions. For Tensor 1, we test on both order three tensors with both dimension sizes s and decomposition rank R equal to 400 and order four tensors with $s = R = 120$, and test the performance of pairwise perturbation on tensors with different collinearity for the exact input factor matrices. For Tensor 2, we test on order three tensors with $s = R$, and test the performance of pairwise perturbation with different dimension size and corresponding rank.

We display the speed-ups of pairwise perturbation compared to the dimension tree algorithm for synthetic tensors in Figure 3. Figures 3a and 3b show the speed-up distribution with different exact factor matrices collinearity. We stop the algorithm when the stopping tolerance (defined as the fitness difference between two neighboring sweeps) is reached. It can be seen that for both order three and order four tensors, PP achieves up to 2.0X speed-up, and high speed-up is achieved with tighter stopping tolerance. We find that the stricter stopping tolerance of 10^{-5} is valuable, as generally it permits about one more digit of accuracy to be achieved in fitness compared to a tolerance of 10^{-4} . In addition, experiments with a 10^{-4} stopping tolerance sometimes stop at transient swamps [43] with high decomposition residual, where ALS makes small progress for a period but the residual norm decreases more rapidly afterwards. In addition, PP tends to have higher speed-ups with relatively high collinearity. This is because tensors with high collinearity will converge in more sweeps, and more PP approximated sweeps are activated as can be seen in Table 2. PP starts working early for almost all the experiments, as can be observed in Figure 3c, where PP starts to have speed-up when the fitness is around 0.975 and the experiment time is less than 20 seconds, and in Table 2, where almost all the PP initialization steps start within 20 sweeps. In addition, the fitness increases monotonically in Figure 3c, indicating that PP controls the approximation error well.

Figure 3c also illustrates the importance of the second-order correction term, $\mathbf{V}^{(n,i,j)}$, in Equation 3.2. We set the PP tolerance to be 0.02 for the PP experiment without corrections, which results in more conservative use of PP approximate steps than with the 0.1 tolerance we use for PP with the second-order correction. As can be seen, without the correction, PP suffers from more instability and no speed-up is achieved for this experiment. Therefore, for all other experiments, the correction terms are included as part of PP.

Figure 3d shows the speed-up distribution with different dimension size for order three tensors made by random factor matrices. It can be seen from the figure that PP achieves up to 3.0X speed-up, and PP has larger speed-ups on larger tensors, consistent with the cost analysis.

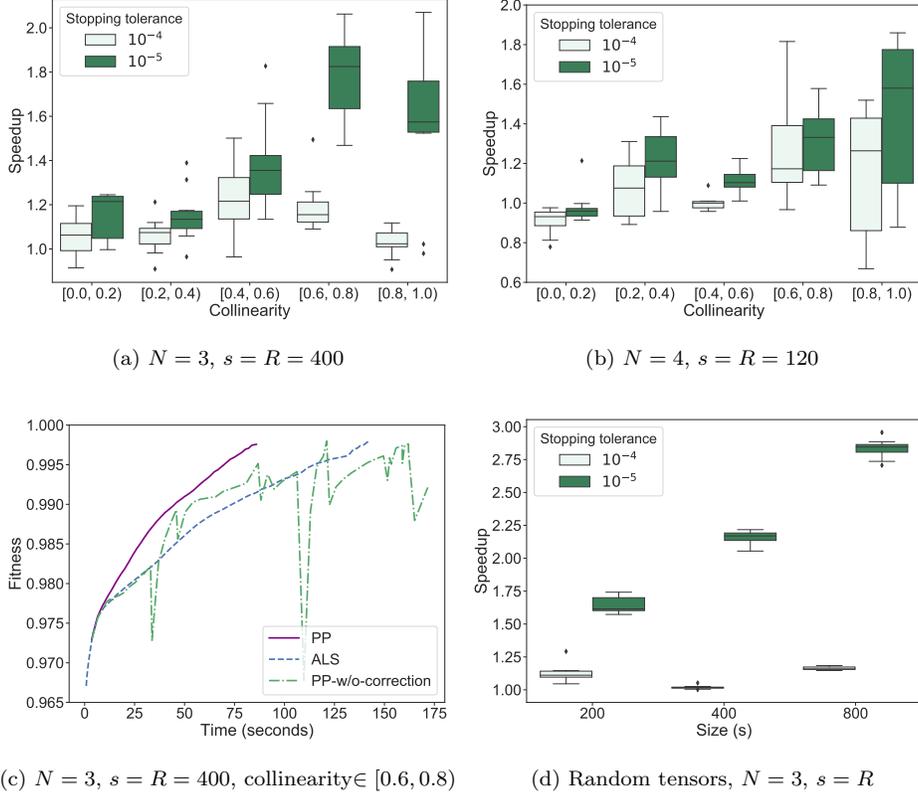


FIGURE 3. (a)(b) Box plot of the relation between PP speed-up and input collinearity ranges for tensors with specific collinearity. (c) Fitness-time relation for the decomposition of one tensor with specific collinearity. (d) Box plot of the relation between PP speed-up and size in each mode for order 3 random tensors. For all the box plots, each box is based on 10 experiments with different random seeds. Each box shows the 25th-75th quartiles, the median is indicated by a horizontal line inside the box, and outliers are displayed as dots.

TABLE 2

Detailed statistics of the results shown in Figure 3. From left to right: the tensor configuration (col stands for collinearity), number of exact ALS sweeps within the PP algorithm, number of PP initialization sweeps, number of PP approximated sweeps, index of sweep when PP is first initialized (approximation begins), the fitness when PP is first initialized, and the final fitness of the experiment. All the data are the average statistics from ten experiments.

Configuration	Num-ALS	Num-PP-init	Num-PP-approx	PP-init-sweep	PP-init-fit	Final-fit
$N=3, \text{col} \in [0.0, 0.2]$	19.9	2.5	11.4	12.7	0.8203	0.9330
$N=3, \text{col} \in [0.2, 0.4]$	49.1	18.4	35.3	7.7	0.7937	0.9991
$N=3, \text{col} \in [0.4, 0.6]$	60.8	52.9	149.1	8.8	0.9345	0.9999
$N=3, \text{col} \in [0.6, 0.8]$	54.8	50.1	252.1	5.7	0.9751	0.9962
$N=3, \text{col} \in [0.8, 1.0]$	12.8	9.4	51.1	4.3	0.9940	0.9966
$N=4, \text{col} \in [0.0, 0.2]$	20.1	3.3	2.4	13.7	0.6802	0.8235
$N=4, \text{col} \in [0.2, 0.4]$	15.4	1.9	5.6	14.0	0.9525	0.9945
$N=4, \text{col} \in [0.4, 0.6]$	34.0	7.5	13.5	22.6	0.9477	0.9935
$N=4, \text{col} \in [0.6, 0.8]$	46.1	29.3	73.3	9.1	0.9365	0.9990
$N=4, \text{col} \in [0.8, 1.0]$	47.5	26.4	62.4	6.2	0.9831	0.9963

We also test the performance of pairwise perturbation on CP decomposition of the quantum chemistry tensor, as is shown in Figure 4, with detailed statistics shown in Table 3. In addition to the original ALS algorithm, we consider two other ALS variants for this problem: the ALS algorithm with different update step size, and the ALS algorithm with a symmetry constraint [30]. The algorithm with different update step size updates the factor matrices $\mathbf{A}^{(n)}$ based on

$$\mathbf{A}_{new}^{(n)} = (1 - \lambda)\mathbf{A}^{(n)} + \lambda\mathbf{M}^{(n)}\mathbf{\Gamma}^{(n)\dagger},$$

where λ is the update step size. A good choice of λ can help achieving better convergence. The symmetry

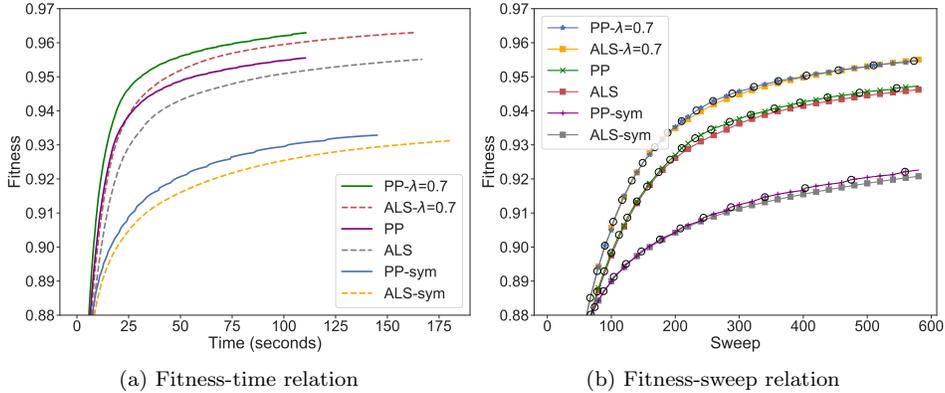


FIGURE 4. Comparison of PP and the dimension tree algorithm for CP decomposition on the quantum chemistry tensor with different variants. PP-sym/ALS-sym denotes the decomposition with symmetry constraint. PP- $\lambda=0.7$ /ALS- $\lambda=0.7$ denotes the decomposition with step size chosen to be 0.7. (b) shows detailed fitness-sweep relation for part of the sweeps. In (b), squares on the dimension tree lines represent the results per 20 sweeps (including all PP initialization, PP approximated and ALS sweeps), and the black circles on pairwise perturbation lines represent the time when pairwise perturbation re-initializes.

TABLE 3

Detailed statistics of different experiments. From left to right: the tensor type, number of ALS sweeps until PP experiments are finished, number of PP initialization sweeps, number of PP approximated sweeps, the average time of each ALS sweep, the average time of each PP initialization sweep, and average time of each PP approximated sweep.

Tensor	Num-ALS	Num-PP-init	Num-PP-approx	Time-ALS	Time-PP-init	Time-PP-approx
Chemistry (Figure 4)	44	40	1416	0.1116	0.1655	0.0703
Coil (Figure 5a)	31	22	147	2.357	3.660	0.0648
TimeLapse (Figure 5b)	23	16	161	0.4087	0.9236	0.0562
Chemistry (Figure 7)	88	54	1358	5.338	9.608	2.254

constrained algorithm considers the input tensor is symmetric in the two equidimensional modes and restricts the two factor matrices for these two modes to be the same: $\mathcal{X} = \llbracket \mathbf{A}, \mathbf{B}, \mathbf{B} \rrbracket$. We update \mathbf{A} the same as the original ALS step, and update \mathbf{B} with the update step size $\lambda = 0.8$ to avoid divergence.

As is shown in Figure 4a, for all the variants of ALS algorithms, PP performs better than the dimension tree algorithm, achieving 1.25-1.52X speed-up. All the experiments are stopped after 1500 sweeps. It can also be observed in Figure 4b that PP usually restarts once approximately every 40 sweeps, and for each sweep, the fitness of both ALS and PP are almost the same, indicating that PP controls the approximation error well.

We test the performance of pairwise perturbation on real image datasets with NumPy in Figure 5, with detailed statistics shown in Table 3. We display the fitness and execution time for CP decomposition of the two image datasets in Figure 5a, 5b. We observe that pairwise perturbation achieves a lower execution time for them. The speed-up for the Coil Dataset is 2.72X and for the Time-Lapse Dataset is 3.1X.

Pairwise perturbation is also used to speedup HOOI procedure in Tucker decomposition. However, as noted in other work [6], we observed that ALS sweeps do not significantly lower the residual beyond what is achieved by the first sweep (HOSVD). We display the fitness and the execution time for Tucker decomposition of the two real datasets in Figure 5c, 5d. The speed-up for the Coil Dataset is 1.05X and for the Time-Lapse Dataset is 1.13X. The reason for no obvious speed-up for the Coil Dataset is that the tensor is not equidimensional (one dimension is 7200, while others are all smaller or equal to 128). Therefore, when updating the factor matrix with a dimension of 7200, the number of operations necessary to construct the SVD input for PP are similar to that for the dimension tree Tucker algorithm. For the Time-Lapse Dataset, the tensor dimensions are more evenly distributed (two dimensions are greater than 1000), and we observe a greater speed-up. We conclude that the proposed Tucker PP algorithm performs better when used on the tensors whose dimensions are approximately equal.

5.2. Parallel Performance. We perform a parallel scaling analysis to compare the simulation time for one ALS sweep of the dimension tree algorithm to the initialization and the approximated step of the pairwise perturbation algorithm with Cyclops in Figure 6. Parallelism is used to accelerate the tensor contractions via calling Cyclops kernels as well as the linear system solve via calling ScaLAPACK kernels. The Cyclops library

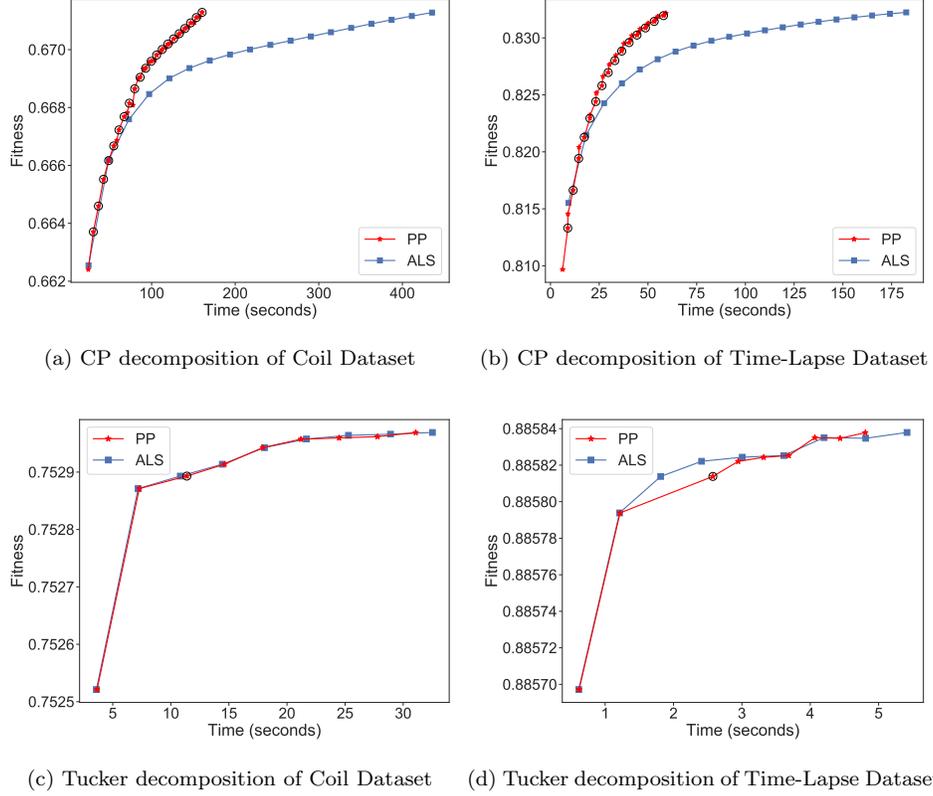


FIGURE 5. Experimental results on image datasets between pairwise perturbation and ALS for CP and Tucker decompositions. Each dot on the ALS/PP lines represents the results per 10 sweeps for CP and per sweep for Tucker decomposition (including all PP initialization, PP approximated and ALS sweeps), and the black circles on pairwise perturbation lines represent the time when pairwise perturbation restarts.

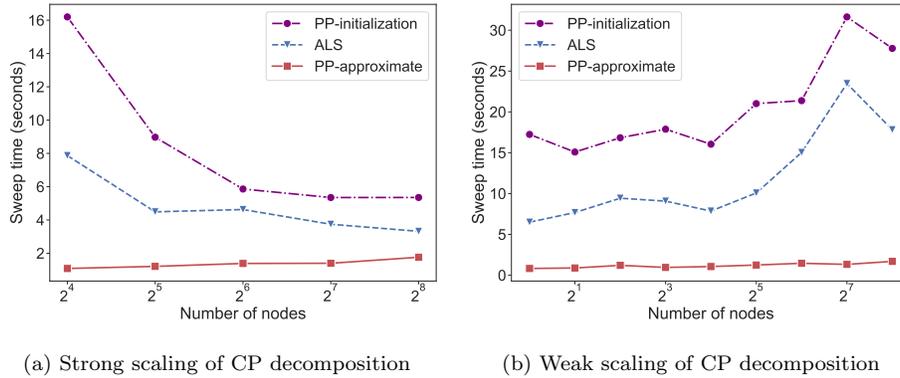


FIGURE 6. Benchmark results for ALS sweeps with Cyclops, taken as the mean time across 5 sweeps.

reduces each tensor contraction to a matrix multiplication. For the PP initialization step, this approach either keeps the input tensor in place, performs local multiplications, and afterwards performs a reduction on the output tensor when the rank R is small, or performs a general 3D parallel matrix multiplication when R is high. For the PP approximated step, this approach parallelizes small-sized batched matrix-vector products and result in over-parallelization. We direct readers to the reference [42] for a detailed communication cost analysis and a more communication efficient algorithm for parallel pairwise perturbation.

We use 8 processes per node and 8 threads per process for the benchmark experiments. The pairwise perturbation initialization step includes the construction of the pairwise perturbation operators, and is there-

fore much slower than the approximated steps. For strong scaling, we consider order $N = 6$ tensors with dimension $s = 50$ and rank $R = 6$ CP and Tucker decompositions. For weak scaling, on p processors, we consider order $N = 6$ tensors with dimension $s = \lfloor 32p^{1/6} \rfloor$ and rank $R = \lfloor 4p^{1/6} \rfloor$.

For weak scaling, Figure 6 shows that with the increase of number of nodes, the step time for all three steps increases. The approximated step time of pairwise perturbation is always much faster (7.8 and 10.5 times faster on 1 node and 256 nodes, respectively, compared to the dimension tree based ALS step time) than the two other steps, showing the good scalability of pairwise perturbation. For strong scaling, the figure shows that the approximated step time of pairwise perturbation increases with the number of nodes, while the two other step times decrease. The PP approximated step is much cheaper computationally and becomes dominated by communication with increasing node counts, thereby slowing down in step time. For the two other steps, the matrix calculation time will be decreased a lot with the increase of node number, thereby the step time is decreased. Overall, we observe that the potential performance benefit of pairwise perturbation is greater for weak scaling.

5.3. Parallel Experimental Results. We test the parallel performance of pairwise perturbation with Cyclops on a quantum chemistry tensor. Similar to Section 5.1, we generate the order three density fitting tensor representing the compressed restricted Hartree-Fock wave function of an 40 water molecule chain system with a STO3G basis set. The generated tensor has size $4520 \times 280 \times 280$. We set the CP rank to be 1800. We show the parallel performance with Cyclops for the quantum chemistry tensor in Figure 7, with

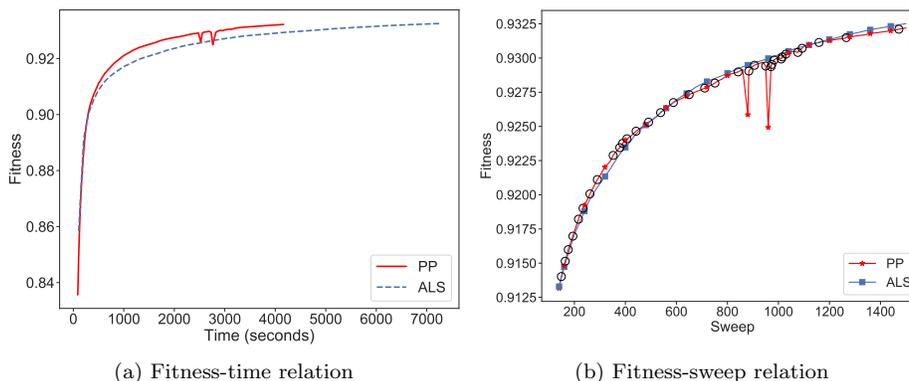


FIGURE 7. Comparison of PP and the dimension tree algorithm for CP decomposition on the quantum chemistry tensor with Cyclops. (b) shows detailed fitness-sweep relation for part of the sweeps. In (b), squares on the dimension tree lines represent the results per 20 sweeps (including all PP initialization, PP approximated and ALS sweeps), and the black circles on pairwise perturbation lines represent the time when pairwise perturbation re-initializes.

detailed statistics shown in Table 3. We perform experiments on 4 KNL nodes, leveraging 64 processors on each node. For the PP experiment, after first level contractions of the PP dimension tree, we redistribute the resulting tensor across all the processes so that it is partitioned in the rank mode, which makes the PP approximated steps faster. It can be seen that PP performs better than the dimension tree algorithm, achieving 1.75X speed-up to reach a fitness of 0.933. It can also be observed in Figure 7b that for most of the sweeps, the fitness of both the dimension tree algorithm and PP are almost the same, indicating that PP controls the approximation error well.

6. Conclusion. We have provided the pairwise perturbation algorithm for both CP and Tucker decompositions for dense tensors. The advantage of this algorithm is that it uses perturbative corrections rather than recomputing the tensor contractions to set up the quadratic optimization subproblems. Our error analysis demonstrates that it is accurate when the factor matrices change little. Specifically, our implementation of pairwise perturbation shows speed-ups for CP-ALS of up to 3.1X across all synthetic and application data with respect to the best known method for exact CP-ALS with the NumPy-based sequential implementation.

We leave analysis and benchmarking of pairwise perturbation with sparse tensors for future work. Since contraction between the input tensor and the first factor matrix requires fewer operations, there is less likely to be a benefit in using pairwise perturbation. Additionally, it is likely of interest to investigate more efficient adaptations of pairwise perturbation for non-equidimensional tensors and to experiment with alternative schemes for switching between regular ALS and pairwise perturbation.

7. Acknowledgments. The authors are grateful to Daniel Kressner for pointing out the connection to the Hurwitz problem, to Fan Huang for finding the $8 \times 8 \times 8$ perfectly conditioned tensor, and to Nick Vannieuwenhoven for helpful comments. The authors were supported by the US NSF OAC SSI program, award No. 1931258. This work used the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562. We used XSEDE to employ Stampede2 at the Texas Advanced Computing Center (TACC) through allocation TG-CCR180006.

REFERENCES

- [1] E. Acar, D. M. Dunlavy, and T. G. Kolda. A scalable optimization approach for fitting canonical tensor decompositions. *Journal of Chemometrics*, 25(2):67–86, 2011.
- [2] J. F. Adams. Vector fields on spheres. *Annals of Mathematics*, pages 603–632, 1962.
- [3] J. F. Adams, P. D. Lax, and R. S. Phillips. On matrices whose real linear combinations are nonsingular. *Proceedings of the American Mathematical Society*, 16(2):318–322, 1965.
- [4] A. Anandkumar, R. Ge, D. J. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15(1):2773–2832, 2014.
- [5] C. A. Andersson and R. Bro. Improving the speed of multi-way algorithms: Part I. Tucker3. *Chemometrics and intelligent laboratory systems*, 42(1-2):93–103, 1998.
- [6] W. Austin, G. Ballard, and T. G. Kolda. Parallel tensor compression for large-scale scientific data. In *Parallel and Distributed Processing Symposium, 2016 IEEE International*, pages 912–922. IEEE, 2016.
- [7] G. Ballard, K. Hayashi, and R. Kannan. Parallel nonnegative CP decomposition of dense tensors. *arXiv preprint arXiv:1806.07985*, 2018.
- [8] G. Ballard, N. Knight, and K. Rouse. Communication lower bounds for matricized tensor times Khatri-Rao product. In *2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 557–567. IEEE, 2018.
- [9] C. Battaglino, G. Ballard, and T. G. Kolda. A practical randomized CP tensor decomposition. *SIAM Journal on Matrix Analysis and Applications*, 39(2):876–901, 2018.
- [10] U. Benedikt, H. Auer, M. Espig, W. Hackbusch, and A. A. Auer. Tensor representation techniques in post-Hartree-Fock methods: matrix product state tensor format. *Molecular Physics*, 111(16-17):2398–2413, 2013.
- [11] L. S. Blackford, J. Choi, A. Cleary, E. D’Azevedo, J. Demmel, I. Dhillon, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK User’s Guide*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1997.
- [12] J. D. Carroll and J.-J. Chang. Analysis of individual differences in multidimensional scaling via an N-way generalization of Eckart-Young decomposition. *Psychometrika*, 35(3):283–319, 1970.
- [13] V. T. Chakaravarthy, J. W. Choi, D. J. Joseph, X. Liu, P. Murali, Y. Sabharwal, and D. Sreedhar. On optimizing distributed Tucker decomposition for dense tensors. In *Parallel and Distributed Processing Symposium (IPDPS), 2017 IEEE International*, pages 1038–1047. IEEE, 2017.
- [14] J. Choi, X. Liu, and V. Chakaravarthy. High-performance dense Tucker decomposition on GPU clusters. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis*, page 42. IEEE Press, 2018.
- [15] A. Cichocki, N. Lee, I. Oseledets, A.-H. Phan, Q. Zhao, and D. P. Mandic. Tensor networks for dimensionality reduction and large-scale optimization: Part 1 low-rank tensor decompositions. *Foundations and Trends in Machine Learning*, 9(4-5):249–429, 2016.
- [16] L. De Lathauwer, B. De Moor, and J. Vandewalle. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000.
- [17] L. De Lathauwer, B. De Moor, and J. Vandewalle. On the best rank-1 and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors. *SIAM journal on Matrix Analysis and Applications*, 21(4):1324–1342, 2000.
- [18] S. Friedland, V. Mehrmann, R. Pajarola, and S. K. Suter. On best rank one approximation of tensors. *Numerical Linear Algebra with Applications*, 20(6):942–955, 2013.
- [19] L. Grasedyck, D. Kressner, and C. Tobler. A literature survey of low-rank tensor approximation techniques. *GAMM-Mitteilungen*, 36(1):53–78, 2013.
- [20] W. Hackbusch. *Tensor spaces and numerical tensor calculus*, volume 42. Springer Science & Business Media, 2012.
- [21] N. Hao, L. Horesh, and M. Kilmer. Nonnegative tensor decomposition. In *Compressed Sensing & Sparse Filtering*, pages 123–148. Springer, 2014.
- [22] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del R’io, M. Wiebe, P. Peterson, P. G’erard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant. Array programming with NumPy. *Nature*, 585(7825):357–362, Sept. 2020.
- [23] R. A. Harshman. Foundations of the PARAFAC procedure: models and conditions for an explanatory multimodal factor analysis. 1970.
- [24] K. Hayashi, G. Ballard, J. Jiang, and M. Tobia. Shared memory parallelization of MTTKRP for dense tensors. *arXiv preprint arXiv:1708.08976*, 2017.
- [25] C. J. Hillar and L.-H. Lim. Most tensor problems are NP-hard. *J. ACM*, 60(6):45:1–45:39, Nov. 2013.
- [26] F. L. Hitchcock. The expression of a tensor or a polyadic as a sum of products. *Studies in Applied Mathematics*, 6(1-4):164–189, 1927.
- [27] E. G. Hohenstein, R. M. Parrish, and T. J. Martínez. Tensor hypercontraction density fitting. I. Quartic scaling second- and third-order Møller-Plesset perturbation theory. *The Journal of Chemical Physics*, 137(4):044103, 2012.
- [28] E. G. Hohenstein, R. M. Parrish, C. D. Sherrill, and T. J. Martínez. Communication: Tensor hypercontraction. III.

- Least-squares tensor hypercontraction for the determination of correlated wavefunctions, 2012.
- [29] T. Huckle, K. Waldherr, and T. Schulte-Herbrüggen. Computations in quantum tensor networks. *Linear Algebra and its Applications*, 438(2):750 – 781, 2013. Tensors and Multilinear Algebra.
- [30] F. Hummel, T. Tsatsoulis, and A. Grüneis. Low rank factorization of the Coulomb integrals for periodic coupled cluster theory. *The Journal of chemical physics*, 146(12):124105, 2017.
- [31] A. Hurwitz. Über die Composition der quadratischen Formen von beliebig vielen Variablen. *Nachrichten von der Gesellschaft der Wissenschaften zu Göttingen, Mathematisch-Physikalische Klasse*, 1898:309–316, 1898.
- [32] L. Karlsson, D. Kressner, and A. Uchmajew. Parallel algorithms for tensor completion in the CP format. *Parallel Computing*, 57:222–234, 2016.
- [33] O. Kaya. *High performance parallel algorithms for tensor decompositions*. PhD thesis, 2017.
- [34] O. Kaya and Y. Robert. Computing dense tensor decompositions with optimal dimension trees. *Algorithmica*, 81(5):2092–2121, 2019.
- [35] O. Kaya and B. Uçar. High performance parallel algorithms for the Tucker decomposition of sparse tensors. In *Parallel Processing (ICPP)*, 2016 45th International Conference on, pages 103–112. IEEE, 2016.
- [36] T. G. Kolda and B. W. Bader. Matlab tensor toolbox. Technical report, Sandia National Laboratories, 2006.
- [37] T. G. Kolda and B. W. Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.
- [38] J. Li, J. Choi, I. Perros, J. Sun, and R. Vuduc. Model-driven sparse CP decomposition for higher-order tensors. In *2017 IEEE international parallel and distributed processing symposium (IPDPS)*, pages 1048–1057. IEEE, 2017.
- [39] A. P. Liavas, G. Kostoulas, G. Lourakis, K. Huang, and N. D. Sidiropoulos. Nesterov-based alternating optimization for nonnegative tensor factorization: algorithm and parallel implementation. *IEEE Trans. Signal Process*, 66:944–953, 2017.
- [40] L.-H. Lim. Singular values and eigenvalues of tensors: a variational approach. In *Computational Advances in Multi-Sensor Adaptive Processing*, 2005 1st IEEE International Workshop on, pages 129–132. IEEE, 2005.
- [41] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(1):208–220, 2013.
- [42] L. Ma and E. Solomonik. Efficient parallel CP decomposition with pairwise perturbation and multi-sweep dimension tree. *arXiv preprint arXiv:2010.12056*, 2020.
- [43] M. J. Mohlenkamp. The dynamics of swamps in the canonical tensor approximation problem. *SIAM Journal on Applied Dynamical Systems*, 18(3):1293–1333, 2019.
- [44] J. G. Nagy and M. E. Kilmer. Kronecker product approximation for preconditioning in three-dimensional imaging applications. *IEEE Transactions on Image Processing*, 15(3):604–613, March 2006.
- [45] S. M. Nascimento, K. Amano, and D. H. Foster. Spatial distributions of local illumination color in natural scenes. *Vision Research*, 120:39–44, 2016.
- [46] S. A. Nene, S. K. Nayar, and H. Murase. Columbia object image library (coil-100).
- [47] R. Orús. A practical introduction to tensor networks: Matrix product states and projected entangled pair states. *Annals of Physics*, 349:117 – 158, 2014.
- [48] I. V. Oseledets and E. E. Tyrtyshnikov. Breaking the curse of dimensionality, or how to use SVD in many dimensions. *SIAM Journal on Scientific Computing*, 31(5):3744–3759, 2009.
- [49] W. Pazner and P.-O. Persson. Approximate tensor-product preconditioners for very high order discontinuous Galerkin methods. *Journal of Computational Physics*, 354:344–369, 2018.
- [50] I. Perros, R. Chen, R. Vuduc, and J. Sun. Sparse hierarchical Tucker factorization and its application to healthcare. In *Data Mining (ICDM)*, 2015 IEEE International Conference on, pages 943–948. IEEE, 2015.
- [51] A.-H. Phan, P. Tichavský, and A. Cichocki. Fast alternating LS algorithms for high order CANDECOMP/PARAFAC tensor factorizations. *IEEE Transactions on Signal Processing*, 61(19):4834–4846, 2013.
- [52] J. Radon. Lineare Scharen orthogonaler Matrizen. In *Abhandlungen aus dem Mathematischen Seminar der Universität Hamburg*, volume 1, pages 1–14. Springer, 1922.
- [53] M. Rajih, P. Comon, and R. A. Harshman. Enhanced line search: a novel method to accelerate PARAFAC. *SIAM journal on matrix analysis and applications*, 30(3):1128–1147, 2008.
- [54] M. D. Schatz, T. M. Low, R. A. van de Geijn, and T. G. Kolda. Exploiting symmetry in tensors for high performance: multiplication with symmetric tensors. *SIAM Journal on Scientific Computing*, 36(5):C453–C479, 2014.
- [55] E. Solomonik, D. Matthews, J. R. Hammond, J. F. Stanton, and J. Demmel. A massively parallel tensor contraction framework for coupled-cluster computations. *Journal of Parallel and Distributed Computing*, 74(12):3176–3190, 2014.
- [56] P. Springer, T. Su, and P. Bientinesi. HPTT: a high-performance tensor transposition C++ library. In *Proceedings of the 4th ACM SIGPLAN International Workshop on Libraries, Languages, and Compilers for Array Programming*, pages 56–62. ACM, 2017.
- [57] Q. Sun, T. C. Berkelbach, N. S. Blunt, G. H. Booth, S. Guo, Z. Li, J. Liu, J. D. McClain, E. R. Sayfutyarova, S. Sharma, et al. PySCF: the Python-based simulations of chemistry framework. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 8(1):e1340, 2018.
- [58] L. R. Tucker. Some mathematical notes on three-mode factor analysis. *Psychometrika*, 31(3):279–311, 1966.
- [59] N. Vannieuwenhoven, K. Meerbergen, and R. Vandebril. Computing the gradient in optimization algorithms for the CP decomposition in constant memory through tensor blocking. *SIAM Journal on Scientific Computing*, 37(3):C415–C438, 2015.
- [60] N. Vannieuwenhoven, R. Vandebril, and K. Meerbergen. A new truncation strategy for the higher-order singular value decomposition. *SIAM Journal on Scientific Computing*, 34(2):A1027–A1052, 2012.
- [61] G. Zhou, A. Cichocki, and S. Xie. Decomposition of big tensors with low multilinear rank. *arXiv preprint arXiv:1412.1885*, 2014.

8. Appendix: Error Bounds based on a Tensor Condition Number. We provide relative error bounds for the pairwise perturbation algorithm for both CP-ALS and Tucker-ALS for tensors that are ‘well-

conditioned', in a sense that is defined in this appendix. However, results related to the Hurwitz problem regarding multiplicative relations of quadratic forms [31], imply that equidimensional order three tensors can have a bounded condition number only if their dimension is $s \in \{1, 2, 4, 8\}$. We provide families of tensors with unit condition number with such dimensions. The results shed further light on the stability of MTTKRP as well as ALS, and yield nontrivial bounds for small tensors. For factorization of large tensors, the bounds proven in this section are not meaningful, since their condition number is necessarily infinite for at least one ordering of modes.

8.1. Tensor Condition Number. We consider a notion of tensor condition number that corresponds to a global bound on the conditioning of the multilinear vector-valued function, $g_{\mathcal{T}} : \otimes_{i=2}^N \mathbb{R}^{s_i} \rightarrow \mathbb{R}^{s_1}$ associated with the product of the tensor with vectors along all except the first mode,

$$g_{\mathcal{T}}(\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}) = \mathcal{T} \times_{i \in \{2, \dots, N\}} \mathbf{x}^{(i)T},$$

where \mathcal{T} is contracted with $\mathbf{x}^{(i)}$ along its i th mode. The norm and condition number are given by extrema of the norm amplification of $g_{\mathcal{T}}$, which are described by the amplification function $f_{\mathcal{T}} : \otimes_{i=2}^N \mathbb{R}^{s_i} \rightarrow \mathbb{R}$,

$$f_{\mathcal{T}}(\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}) = \frac{\|g_{\mathcal{T}}(\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})\|_2}{\|\mathbf{x}^{(2)}\|_2 \cdots \|\mathbf{x}^{(N)}\|_2}.$$

The spectral norm of the tensor corresponds to its supremum,

$$\|\mathcal{T}\|_2 = \sup\{f_{\mathcal{T}}\}.$$

The tensor condition number can be defined as

$$\kappa(\mathcal{T}) = \sup\{f_{\mathcal{T}}\} / \inf\{f_{\mathcal{T}}\},$$

which enables quantification of the worst-case relative amplification of error with respect to input for the product of a tensor with vectors along all except the first mode. In particular, $\kappa(\mathcal{T})$ provides an upper bound on the relative norm of the perturbation of $g_{\mathcal{T}}$ with respect to the relative norm of any perturbation to any input vector.

For a matrix $\mathbf{M} \in \mathbb{R}^{s_1 \times s_2}$, if $s_1 > s_2$ the above notion of condition number gives $\kappa(\mathbf{M}) = \frac{\sigma_{\max}(\mathbf{M})}{\sigma_{\min}(\mathbf{M})}$ where $\sigma_{\min}(\mathbf{M})$ is the smallest singular value of \mathbf{M} in the reduced SVD, while if $s_1 < s_2$, then $\kappa(\mathbf{M}) = \infty$. When tensor dimensions are unequal, the condition number is infinite if the first dimension is not the largest, so for some i , $s_i > s_1$. Aside from this condition, the ordering of modes of \mathcal{T} does not affect the condition number, since for any $m > 1$, the supremum/infimum of $f_{\mathcal{T}}$ over the domain of unit vectors are for some choice of $\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(m-1)}, \mathbf{x}^{(m+1)}, \dots, \mathbf{x}^{(N)}$ the maximum/minimum singular values of

$$\mathbf{K} = \mathcal{T} \times_{i \in \{2, \dots, m-1, m+1, \dots, N\}} \mathbf{x}^{(i)T}.$$

8.2. Well-Conditioned Tensors. We provide two examples of order three tensors that have unit condition number. Other perfectly conditioned tensors can be obtained by multiplying the above tensors by an orthogonal matrix along any mode (we prove below that such transformations preserve condition number). The first example has $s_i = 2$, and yields a Givens rotation when contracted with a vector along the last mode. It is composed of two slices:

$$\begin{bmatrix} 1 & \\ & 1 \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} & 1 \\ -1 & \end{bmatrix}.$$

The second example has $s_i = 4$ and is composed of four slices:

$$\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & -1 \end{bmatrix}, \begin{bmatrix} -1 & 1 & & \\ & & 1 & \\ & & & 1 \\ & & & & 1 \end{bmatrix}, \begin{bmatrix} & & 1 & \\ & -1 & & \\ 1 & & & \end{bmatrix}, \begin{bmatrix} & & & -1 \\ & & & & 1 \\ & 1 & & \\ & & & & 1 \end{bmatrix}.$$

Finally, for $s_i = 8$, we provide an example by giving matrices \mathbf{M} and \mathbf{N} , so that the tensor has nonzeros $\mathcal{T}(i, j, \mathbf{M}(i, j)) = \mathbf{N}(i, j)$ for each entry in \mathbf{M} ,

$$\mathbf{M} = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 \\ 2 & 1 & 4 & 3 & 6 & 5 & 8 & 7 \\ 3 & 4 & 1 & 2 & 7 & 8 & 5 & 6 \\ 4 & 3 & 2 & 1 & 8 & 7 & 6 & 5 \\ 5 & 6 & 7 & 8 & 1 & 2 & 3 & 4 \\ 6 & 5 & 8 & 7 & 2 & 1 & 4 & 3 \\ 7 & 8 & 5 & 6 & 3 & 4 & 1 & 2 \\ 8 & 7 & 6 & 5 & 4 & 3 & 2 & 1 \end{bmatrix}, \mathbf{N} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ -1 & 1 & 1 & -1 & 1 & -1 & -1 & 1 \\ -1 & -1 & 1 & 1 & 1 & 1 & -1 & -1 \\ -1 & 1 & -1 & 1 & 1 & -1 & 1 & -1 \\ -1 & -1 & -1 & -1 & 1 & 1 & 1 & 1 \\ -1 & 1 & -1 & 1 & -1 & 1 & -1 & 1 \\ -1 & 1 & 1 & -1 & -1 & 1 & 1 & -1 \\ -1 & -1 & 1 & 1 & -1 & -1 & 1 & 1 \end{bmatrix}.$$

The fact that the latter two tensors have unit condition number can be verified by symbolic algebraic manipulation or numerical tests.

These tensors provide solutions to special cases of the Hurwitz problem [31], which seeks bilinear forms z_1, \dots, z_n in variables x_1, \dots, x_l and y_1, \dots, y_m such that

$$(x_1^2 + \dots + x_l^2) (y_1^2 + \dots + y_m^2) = z_1^2 + \dots + z_n^2.$$

Consequently, if for \mathcal{T} and any vectors \mathbf{x}, \mathbf{y} ,

$$\frac{\|\mathcal{T} \times_2 \mathbf{x}^T \times_3 \mathbf{y}^T\|_2}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2} = 1 \quad \Rightarrow \quad \|\mathcal{T} \times_2 \mathbf{x}^T \times_3 \mathbf{y}^T\|_2^2 = \|\mathbf{x}\|_2^2 \|\mathbf{y}\|_2^2,$$

so we can define bilinear forms,

$$z_i = \sum_j \sum_k \mathcal{T}(i, j, k) x_j y_k,$$

that provide a solution to the Hurwitz problem. Such equidimensional tensors with unit condition number exist for dimension $s \in \{1, 2, 4, 8\}$ [52], corresponding to the Hurwitz problem with $l = m = n = s$. However, solutions to the Hurwitz problem with $l = m = n$ cannot exist for any other dimension. Furthermore, tight bounds exist on the dimension s_3 for a tensor of dimensions $s \times s \times s_3$ to have bounded condition number ($\inf\{f_{\mathcal{T}}\} > 0$). This problem is equivalent to finding s_3 matrices of dimension $s \times s$, such that any nonzero linear combination thereof is invertible. Factorizing $s = 2^{4a+b}c$, where $b \in \{0, 1, 2, 3\}$ and c is odd, $s_3 \leq 8a + 2^b$ [2, 3].

8.3. Properties of the Tensor Condition Number. In our analysis, we make use of the following submultiplicativity property of the tensor condition number with respect to multilinear multiplication (the property also generalizes to pairs of arbitrary order tensors contracted over one mode).

LEMMA 8.1. *For any $\mathcal{T} \in \mathbb{R}^{s_1 \times \dots \times s_N}$ and matrix \mathbf{M} , if $\mathcal{V} = \mathcal{T} \times_N \mathbf{M}^T$ then $\kappa(\mathcal{V}) \leq \kappa(\mathcal{T})\kappa(\mathbf{M})$.*

Proof. Assume $\kappa(\mathcal{V}) > \kappa(\mathcal{T})\kappa(\mathbf{M})$, then there exist unit vectors $\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$ and $\mathbf{y}^{(2)}, \dots, \mathbf{y}^{(N)}$ such that

$$\kappa(\mathcal{T})\kappa(\mathbf{M}) < \kappa(\mathcal{V}) = \frac{\left\| \mathcal{V} \times_{i \in \{2, \dots, N\}} \mathbf{x}^{(i)T} \right\|_2}{\left\| \mathcal{V} \times_{i \in \{2, \dots, N\}} \mathbf{y}^{(i)T} \right\|_2} = \frac{\left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{x}^{(i)T} \times_N (\mathbf{M} \mathbf{x}^{(N)})^T \right\|_2}{\left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{y}^{(i)T} \times_N (\mathbf{M} \mathbf{y}^{(N)})^T \right\|_2}.$$

Let $\mathbf{u} = \mathbf{M} \mathbf{x}^{(N)}$ and $\mathbf{v} = \mathbf{M} \mathbf{y}^{(N)}$, so $\|\mathbf{u}\|_2 / \|\mathbf{v}\|_2 \leq \kappa(\mathbf{M})$, yielding a contradiction,

$$\kappa(\mathcal{V}) \leq \frac{\left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{x}^{(i)T} \times_N (\mathbf{u} / \|\mathbf{u}\|_2)^T \right\|_2}{\left\| \mathcal{T} \times_{i \in \{2, \dots, N-1\}} \mathbf{y}^{(i)T} \times_N (\mathbf{v} / \|\mathbf{v}\|_2)^T \right\|_2} \kappa(\mathbf{M}) \leq \kappa(\mathcal{T})\kappa(\mathbf{M}). \quad \square$$

Applying Lemma 8.1 with a vector, i.e. when $\mathbf{M} \in \mathbb{R}^{s_N \times 1}$ and so has condition number $\kappa(\mathbf{M}) = 1$, implies $\kappa(\mathcal{T} \times_N \mathbf{M}^T) \leq \kappa(\mathcal{T})$. By an analogous argument to the proof of Lemma 8.1, we can also conclude that the norm and infimum of such a product of \mathcal{T} with unit vectors are bounded by those of \mathcal{T} , giving the following corollary.

COROLLARY 8.2. For any $\mathcal{T} \in \mathbb{R}^{s_1 \times \dots \times s_N}$, vector $\mathbf{u} \in \mathbb{R}^{s_n}$, and any $n \in \{1, \dots, N\}$ such that $\exists m \in \{1, \dots, N\}$ with $s_m \geq s_n$ and $m \neq n$, if $\mathcal{V} = \mathcal{T} \times_n \mathbf{u}^T$, then $\|\mathcal{V}\|_2 \leq \|\mathbf{u}\|_2 \|\mathcal{T}\|_2$, $\inf\{f_{\mathcal{V}}\} \geq \|\mathbf{u}\|_2 \inf\{f_{\mathcal{T}}\}$, and $\kappa(\mathcal{V}) \leq \kappa(\mathcal{T})$.

For an orthogonal matrix \mathbf{M} , Lemma 8.1 can be applied in both directions, namely for $\mathcal{V} = \mathcal{T} \times_N \mathbf{M}^T$ and $\mathcal{T} = \mathcal{V} \times_N \mathbf{M}$, so we observe that $\kappa(\mathcal{V}) = \kappa(\mathcal{T})$. Using this fact, we demonstrate in the following theorem that any tensor \mathcal{T} can be transformed by orthogonal matrices along each mode, so that one of its fibers has norm $\|\mathcal{T}\|_2/\kappa(\mathcal{T})$.

THEOREM 8.1. For any $\mathcal{T} \in \mathbb{R}^{s_1 \times \dots \times s_N}$, there exist orthogonal matrices $\mathbf{Q}^{(2)} \dots \mathbf{Q}^{(N)}$, with $\mathbf{Q}^{(i)} \in \mathbb{R}^{s_i \times s_i}$, such that $\mathcal{V} = \mathcal{T} \times_2 \mathbf{Q}^{(2)} \dots \times_N \mathbf{Q}^{(N)}$ satisfies $\kappa(\mathcal{V}) = \kappa(\mathcal{T})$, $\|\mathcal{V}\|_2 = \|\mathcal{T}\|_2$, and the first fiber of \mathcal{V} , i.e. the vector \mathbf{v} with $\mathbf{v}^{(i)} = \mathcal{V}(i, 0, \dots, 0)$, satisfies $\|\mathbf{v}\|_2 = \|\mathcal{T}\|_2/\kappa(\mathcal{T})$.

Proof. Given a tensor \mathcal{T} with infinite condition number, there must exist $N-1$ unit vectors $\mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}$, such that $\left\| \mathcal{T} \times_{i \in \{2, \dots, N\}} \mathbf{x}^{(i)T} \right\|_2 = \|\mathcal{T}\|_2/\kappa(\mathcal{T})$. We define $N-1$ orthogonal matrices $\mathbf{Q}^{(2)}, \dots, \mathbf{Q}^{(N)}$ such that $\mathbf{Q}^{(i)T} \mathbf{x}^{(i)} = \mathbf{e}_i$. We can then contract \mathcal{T} with these matrices along the last $N-1$ modes, resulting in \mathcal{V} , with the same condition number as \mathcal{T} (by Lemma 8.1) and the same norm (by a similar argument). Then, we have that the first fiber of \mathcal{V} is

$$\mathbf{v} = \mathcal{V} \times_{i \in \{2, \dots, N\}} \mathbf{e}_i^T = \mathcal{T} \times_{i \in \{2, \dots, N\}} \mathbf{x}^{(i)T},$$

and consequently $\|\mathbf{v}\|_2 = \|\mathcal{T}\|_2/\kappa(\mathcal{T})$. \square

By Theorem 8.1, the condition number of a tensor is infinity if and only if it can be transformed by products with orthogonal matrices along the last $N-1$ modes into a tensor with a zero fiber. Further, any tensor \mathcal{T} may be perturbed to have infinite condition number by adding to it some $\delta\mathcal{T}$ with relative norm $\|\delta\mathcal{T}\|_2/\|\mathcal{T}\|_2 = 1/\kappa(\mathcal{T})$.

8.4. PP-CP-ALS Error Bound using Tensor Condition Number. For CP decomposition, we obtain condition-number-dependent column-wise error bounds on $\mathbf{M}^{(n)}$ (the right-hand sides in the linear least squares subproblems), based on the magnitude of the relative perturbation to $\mathbf{A}^{(n)}$ since the formation of the pairwise perturbation operators.

THEOREM 8.3. If $\frac{\|\mathbf{d}\mathbf{a}_k^{(n)}\|_2}{\|\mathbf{a}_k^{(n)}\|_2} \leq \epsilon \ll 1$ for $n \in \{1, \dots, N\}$, $k \in \{1, \dots, R\}$ and $s_m \leq s_n$ for any $m \in \{1, \dots, N\}$, the pairwise perturbation algorithm without second order corrections computes $\tilde{\mathbf{M}}^{(n)}$ with column-wise error,

$$\frac{\|\tilde{\mathbf{m}}_k^{(n)} - \mathbf{m}_k^{(n)}\|_2}{\|\mathbf{m}_k^{(n)}\|_2} = O(\epsilon^2 \kappa(\mathcal{X})),$$

where $\mathbf{M}^{(n)}$ is the matrix given by a regular ALS sweep.

Proof. We bound the error due to second order perturbations in $d\mathbf{A}^{(1)}, \dots, d\mathbf{A}^{(n)}$, by similar analysis, higher-order perturbations would lead to errors smaller by a factor of $O(\text{poly}(N)\epsilon)$ and are consequently negligible if $\epsilon \ll 1$. Consider the order four tensors $\mathcal{M}^{(i,j,n)}$ (Equation 2.1) based on the current factor matrices $\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(N)}$ and the pairwise perturbation operators $\mathcal{M}_p^{(i,j,n)}$ based on past factor matrices $\mathbf{A}_p^{(1)}, \dots, \mathbf{A}_p^{(N)}$. The contribution of second order terms to the error is

$$\tilde{\mathbf{m}}_k^{(n)}(x) - \mathbf{m}_k^{(n)}(x) \approx \sum_{\substack{i,j \in \{1, \dots, n-1, n+1, \dots, N\} \\ i \neq j}} \sum_{y=1}^s \sum_{z=1}^s \mathcal{M}_p^{(i,j,n)}(x, y, z, k) d\mathbf{a}_k^{(i)}(y) d\mathbf{a}_k^{(j)}(z).$$

This absolute error has magnitude,

$$\left\| \tilde{\mathbf{m}}_k^{(n)} - \mathbf{m}_k^{(n)} \right\|_2 \leq \binom{N}{2} \max_{i,j} \left\| \mathcal{M}_p^{(i,j,n)}(\cdot, \cdot, \cdot, k) \right\|_2 \left\| d\mathbf{a}_k^{(i)} \right\|_2 \left\| d\mathbf{a}_k^{(j)} \right\|_2.$$

Using the fact that for any i, j we can express $\mathbf{m}_k^{(n)}$ as

$$\mathbf{m}_k^{(n)}(x) = \sum_{y=1}^s \sum_{z=1}^s \mathcal{M}^{(i,j,n)}(x, y, z, k) \mathbf{a}_k^{(i)}(y) \mathbf{a}_k^{(j)}(z),$$

we can lower bound the magnitude of the answer with respect to any $\mathcal{M}^{(i,j,n)}$,

$$\left\| \mathbf{m}_k^{(n)} \right\|_2 \geq \inf \left\{ \left\| f_{\mathcal{M}^{(i,j,n)}(\cdot, \cdot, \cdot, k)} \right\|_2 \right\} \left\| \mathbf{a}_k^{(i)} \right\|_2 \left\| \mathbf{a}_k^{(j)} \right\|_2.$$

Combining the upper bound on the absolute error with the lower bound on norm,

$$\frac{\left\| \tilde{\mathbf{m}}_k^{(n)} - \mathbf{m}_k^{(n)} \right\|_2}{\left\| \mathbf{m}_k^{(n)} \right\|_2} \leq \binom{N}{2} \max_{i,j} \frac{\left\| \mathcal{M}_p^{(i,j,n)}(\cdot, \cdot, \cdot, k) \right\|_2 \left\| d\mathbf{a}_k^{(i)} \right\|_2 \left\| d\mathbf{a}_k^{(j)} \right\|_2}{\inf \left\{ \left\| f_{\mathcal{M}^{(i,j,n)}(\cdot, \cdot, \cdot, k)} \right\|_2 \right\} \left\| \mathbf{a}_k^{(i)} \right\|_2 \left\| \mathbf{a}_k^{(j)} \right\|_2}.$$

Lemma 8.1 implies that for any i, j, k ,

$$\left\| \mathcal{M}_p^{(i,j,n)}(\cdot, \cdot, \cdot, k) \right\|_2 \leq \|\mathcal{X}\|_2 \prod_{l \in \{1, \dots, N\} \setminus \{i, j, n\}} \left\| \mathbf{A}_p^{(l)}(\cdot, k) \right\|_2$$

and that

$$\inf \left\{ \left\| f_{\mathcal{M}^{(i,j,n)}(\cdot, \cdot, \cdot, k)} \right\|_2 \right\} \geq \inf \{ \|\mathcal{X}\|_2 \} \prod_{l \in \{1, \dots, N\} \setminus \{i, j, n\}} \left\| \mathbf{A}^{(l)}(\cdot, k) \right\|_2.$$

Since, $\left\| \mathbf{A}_p^{(l)}(\cdot, k) \right\|_2 \leq (1 + \epsilon) \left\| \mathbf{A}^{(l)}(\cdot, k) \right\|_2$, we obtain the bound,

$$\frac{\left\| \tilde{\mathbf{m}}_k^{(n)} - \mathbf{m}_k^{(n)} \right\|_2}{\left\| \mathbf{m}_k^{(n)} \right\|_2} \leq \binom{N}{2} \kappa(\mathcal{X}) (1 + \epsilon)^{N-3} \epsilon^2 \approx \binom{N}{2} \kappa(\mathcal{X}) \epsilon^2. \quad \square$$

This error bound is relative to the condition number of \mathcal{X} , which means the bound is sensitive to the input tensor and that the error may be unbounded if \mathcal{X} has an exact CP decomposition of rank at most $\min_i s_i$.

8.5. PP-Tucker-ALS Error Bound using Tensor Condition Number. For Tucker decomposition, we again obtain bounds based on the perturbation to $\mathbf{A}^{(n)}$, this time for $\mathbf{y}^{(n)}$ (the tensors on whose matricizations a truncated SVD is performed). Using Lemma 8.1, we prove in Theorem 8.4 that when the tensor has same length in each mode and the relative error of the matrices $\mathbf{A}^{(n)}$ for $n \in \{1, \dots, N\}$ is small, the relative error for the $\tilde{\mathbf{y}}^{(n)}$ is also small.

THEOREM 8.4. *Given an order N equidimensional tensor \mathcal{X} with size s , if $\left\| d\mathbf{A}^{(n)} \right\|_2 \leq \epsilon \ll 1$ for $n \in \{1, \dots, N\}$, the pairwise perturbation algorithm computes $\tilde{\mathbf{y}}^{(n)}$ with error,*

$$\frac{\left\| \tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)} \right\|_2}{\left\| \mathbf{y}^{(n)} \right\|_2} = O(\epsilon^2 \kappa(\mathcal{X})).$$

Proof. As in Theorem 8.3, we bound the error due to second-order terms,

$$\frac{\left\| \tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)} \right\|_2}{\left\| \mathbf{y}^{(n)} \right\|_2} = \binom{N}{2} \max_{i,j} \frac{\left\| \mathbf{y}_p^{(i,j,n)} \times_i d\mathbf{A}^{(i)T} \times_j d\mathbf{A}^{(j)T} \right\|_2}{\left\| \mathbf{y}^{(i,j,n)} \times_i \mathbf{A}^{(i)T} \times_j \mathbf{A}^{(j)T} \right\|_2}.$$

From Lemma 8.1, we have

$$\frac{\left\| \mathbf{y}_p^{(i,j,n)} \times_i d\mathbf{A}^{(i)T} \times_j d\mathbf{A}^{(j)T} \right\|_2}{\left\| \mathbf{y}^{(i,j,n)} \times_i \mathbf{A}^{(i)T} \times_j \mathbf{A}^{(j)T} \right\|_2} \leq \frac{\left\| \mathbf{y}_p^{(i,j,n)} \right\|_2 \left\| d\mathbf{A}^{(i)} \right\|_2 \left\| d\mathbf{A}^{(j)} \right\|_2}{\inf \left\{ \left\| f_{\mathbf{y}^{(i,j,n)}} \right\|_2 \right\} \left\| \mathbf{A}^{(i)} \right\|_2 \left\| \mathbf{A}^{(j)} \right\|_2}.$$

Since $\mathbf{A}^{(i)}$ and $\mathbf{A}^{(j)}$ are both matrices with orthonormal columns,

$$\frac{\left\| \tilde{\mathbf{y}}^{(n)} - \mathbf{y}^{(n)} \right\|_2}{\left\| \mathbf{y}^{(n)} \right\|_2} \leq \binom{N}{2} \max_{i,j} \frac{\left\| \mathbf{y}_p^{(i,j,n)} \right\|_2 \left\| d\mathbf{A}^{(i)} \right\|_2 \left\| d\mathbf{A}^{(j)} \right\|_2}{\inf \left\{ \left\| f_{\mathbf{y}^{(i,j,n)}} \right\|_2 \right\}} = O(\epsilon^2 \kappa(\mathcal{X})). \quad \square$$