

1 **LOW-RANK TUCKER APPROXIMATION OF A TENSOR FROM**
2 **STREAMING DATA ***

3 YIMING SUN[†], YANG GUO[‡], CHARLENE LUO[§], JOEL TROPP[¶], AND
4 MADELEINE UDELL[†]

5 **Abstract.** This paper describes a new algorithm for computing a low-Tucker-rank approximation
6 of a tensor. The method applies a randomized linear map to the tensor to obtain a *sketch* that
7 captures the important directions within each mode, as well as the interactions among the modes.
8 The sketch can be extracted from streaming or distributed data or with a single pass over the tensor,
9 and it uses storage proportional to the degrees of freedom in the output Tucker approximation. The
10 algorithm does not require a second pass over the tensor, although it can exploit another view to
11 compute a superior approximation. The paper provides a rigorous theoretical guarantee on the
12 approximation error. Extensive numerical experiments show that that the algorithm produces useful
13 results that improve on the state of the art for streaming Tucker decomposition.

14 **Key words.** Tucker decomposition, tensor compression, dimension reduction, sketching method,
15 randomized algorithm, streaming algorithm

16 **AMS subject classifications.** 68Q25, 68R10, 68U05

17 **1. Introduction.** Large-scale datasets with natural tensor (multidimensional
18 array) structure arise in a wide variety of applications including computer vision [39],
19 neuroscience [9], scientific simulation [3], sensor networks [31], and data mining [21].
20 In many cases, these tensors are too large to manipulate, to transmit, or even to store
21 in a single machine. Luckily, tensors often exhibit a low-rank structure, and can be
22 approximated by a low-rank tensor factorization, such as CANDECOMP/PARAFAC
23 (CP), tensor train, or Tucker factorization [20]. These factorizations reduce the
24 storage costs by exposing the latent structure. Sufficiently low rank tensors can be
25 compressed by several orders of magnitude with negligible loss. However, computing
26 these factorizations can require substantial computational resources. One challenge is
27 that these large tensors may not fit in main memory on our computer.

28 In this paper, we develop a new algorithm to compute a low-rank Tucker approxi-
29 mation of a tensor from streaming data, using working storage proportional to the
30 degrees of freedom in the output Tucker approximation. The algorithm forms a linear
31 sketch of the tensor, and it operates on the sketch to compute a low-rank Tucker
32 approximation. The main computational work is all performed on a small tensor whose
33 size is proportional to the core tensor in the Tucker factorization. We derive detailed
34 probabilistic error bounds on the quality of the approximation in terms of the tail
35 energy of any matricization of the target tensor.

36 This algorithm is useful in at least three concrete problem settings:

- 37 1. **Streaming:** Data about the tensor is received sequentially. At each time, we
38 observe a low-dimensional slice, an individual entry, or an additive update to
39 the tensor (the “turnstile” model [28]). For example, each slice of the tensor
40 may represent one time step in a computer simulation or the measurements
41 from a sensor array at a particular time. In the streaming setting, the complete
42 tensor is not stored; indeed, it may be much larger than available computing

*Submitted to the editors DATE.

[†]Cornell University, Ithaca (ys784@cornell.edu, udell@cornell.edu).

[‡]University of Wisconsin-Madison, Madison, WI (yguo@cs.wisc.edu).

[§]Columbia University, New York, NY (cl3788@columbia.edu).

[¶]California Institute of Technology, Pasadena, CA (jtropp@cms.caltech.edu)

resources.

Our algorithm can approximate a tensor, presented as a data stream, by sketching the updates and storing the sketch. The linearity of the sketching operation guarantees that sketching commutes with slice, entrywise, or additive updates. Our method forms an approximation of the tensor only after all the data has been observed, rather than approximating the tensor-observed-so-far at any time. This protocol allows for offline data analysis, including many scientific applications. Conversely, this protocol is not suitable for real-time monitoring.

2. **Limited memory:** Data describing the tensor is stored on the hard disk of a computer with much smaller RAM. This setting reduces to the streaming setting by streaming the data from disk.
3. **Distributed:** Data describing the tensor may be stored on many different machines. Communicating data among these machines may be costly due to low network bandwidth or high latency. Our algorithm can approximate tensors stored in a distributed computing environment by sketching the data on each slave machine and transmitting the sketch to a master, which computes the sum of the sketches. Linearity of the sketch guarantees that the sum of the sketches is the sketch of the full tensor.

In the streaming setting, the tensor is not stored, so we require an algorithm that can compute an approximation from a single pass over the data. In contrast, multiple passes over the data are possible in the memory-limited or distributed settings.

This paper presents algorithms for all these settings, among other contributions:

- We present a new linear sketch for higher order tensors that we call the *Tucker sketch*. This sketch captures the principal subspace of the tensor along each mode (corresponding to factor matrices in a Tucker decomposition) and the action of the tensor that links these subspaces (corresponding to the core). The sketch is linear, so it naturally handles streaming or distributed data. The Tucker sketch can be constructed from any dimension reduction map, and it can be used directly to, e.g., cluster the fibers of the tensor along some mode. It also can be used to approximate the original tensor.
- We develop a practical algorithm to compute a low-rank Tucker approximation from the Tucker sketch. This algorithm requires a single pass over the data to form the sketch, and does not require further data access. A variant of this algorithm, using the truncated QR decomposition, yields a quasi-optimal method for tensor approximation that (in expectation) matches the guarantees for HOSVD or ST-HOSVD up to constants.
- We show how to efficiently compress the output of our low-rank Tucker approximation to any fixed rank, without further data access. This method exploits the spectral decay of the original tensor, and it often produces results that are superior to truncated QR. It can also be used to adaptively choose the final size of the Tucker decomposition sufficient to achieve a desired approximation quality.
- We propose a two-pass algorithm that uses additional data access to improve on the one-pass method. This two-pass algorithm was also proposed in the simultaneous work [27]. Both the one-pass and two-pass methods are appropriate for limited memory or distributed data settings.
- We develop provable probabilistic guarantees on the performance of both the one-pass and two-pass algorithms when the tensor sketch is composed of Gaussian dimension reduction maps.

- We exhibit several random maps that can be used to sketch the tensor. Compared to the Gaussian map, these alternatives are cheaper to store, easier to apply, and experimentally deliver similar performance as measured by the tensor approximation error. In particular, we demonstrate the benefits of a Khatri–Rao product of random matrices, which we call the Tensor Random Projection (TRP), which uses exceedingly low storage.
- We perform a comprehensive simulation study with synthetic data, and we consider applications to several real datasets. These results demonstrate the practical performance of our method. In comparison to the only existing one-pass Tucker approximation algorithm [26], our methods reduce the approximation error by more than an order of magnitude given the same storage budget.
- We have developed and released an open-source package in python, available at <https://github.com/udellgroup/tensorsketch>, that implements our algorithms.

2. Background and Related Work. We begin with a short review of tensor notation and some related work on low-rank matrix and tensor approximation.

2.1. Notation. Our paper follows the notation of [20]. We denote *scalar*, *vector*, *matrix*, and *tensor* variables respectively by lowercase letters (x), boldface lowercase letters (\mathbf{x}), boldface capital letters (\mathbf{X}), and boldface Euler script letters (\mathcal{X}). For two vectors \mathbf{x} and \mathbf{y} , we write $\mathbf{x} > \mathbf{y}$ if \mathbf{x} is greater than \mathbf{y} elementwise.

Define $[N] := \{1, \dots, N\}$. For a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$, we respectively denote its i th row, j th column, and (i, j) -th element by $\mathbf{X}(i, \cdot)$, $\mathbf{X}(\cdot, j)$, and $\mathbf{X}(i, j)$ for each $i \in [m]$, $j \in [n]$. We write $\mathbf{X}^\dagger \in \mathbb{R}^{n \times m}$ for the *Moore–Penrose pseudoinverse* of a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$. In particular, $\mathbf{X}^\dagger = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ if $m \geq n$ and \mathbf{X} has full column rank; $\mathbf{X}^\dagger = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1}$, if $m < n$ and \mathbf{X} has full row rank.

2.1.1. Kronecker and Khatri–Rao product. For two matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and $\mathbf{B} \in \mathbb{R}^{K \times L}$, we define the *Kronecker product* $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{IK \times JL}$ as

$$(2.1) \quad \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} \mathbf{A}(1, 1)\mathbf{B} & \cdots & \mathbf{A}(1, J)\mathbf{B} \\ \vdots & \ddots & \vdots \\ \mathbf{A}(I, 1)\mathbf{B} & \cdots & \mathbf{A}(I, J)\mathbf{B} \end{bmatrix}.$$

For $J = L$, we define the *Khatri–Rao product* as $\mathbf{A} \odot \mathbf{B}$, that is, the “matching column-wise” Kronecker product. The resulting matrix of size $(IK) \times J$ is defined as

$$\mathbf{A} \odot \mathbf{B} = [\mathbf{A}(\cdot, 1) \otimes \mathbf{B}(\cdot, 1) \cdots \mathbf{A}(\cdot, J) \otimes \mathbf{B}(\cdot, J)]$$

2.1.2. Tensor basics. For a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, its *mode* or *order* is the number N of dimensions. If $I = I_1 = \cdots = I_N$, we denote $\mathbb{R}^{I_1 \times \cdots \times I_N}$ as \mathbb{R}^{I^N} . The inner product of two tensors \mathcal{X}, \mathcal{Y} is defined as $\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{i_1=1}^{I_1} \cdots \sum_{i_N=1}^{I_N} \mathcal{X}_{i_1 \dots i_N} \mathcal{Y}_{i_1 \dots i_N}$. The *Frobenius norm* of \mathcal{X} is $\|\mathcal{X}\|_F = \sqrt{\langle \mathcal{X}, \mathcal{X} \rangle}$.

2.1.3. Tensor unfoldings. Let $\bar{I} = \prod_{j=1}^N I_j$ and $I_{(-n)} = \prod_{j \neq n} I_j$, and let $\text{vec}(\mathcal{X})$ denote the vectorization of \mathcal{X} . The *mode- n unfolding* of \mathcal{X} is the matrix $\mathbf{X}^{(n)} \in \mathbb{R}^{I_n \times I_{(-n)}}$. The inner product for tensors matches that of any mode- n unfolding:

$$(2.2) \quad \langle \mathcal{X}, \mathcal{Y} \rangle = \langle \mathbf{X}^{(n)}, \mathbf{Y}^{(n)} \rangle = \text{Tr}((\mathbf{X}^{(n)})^\top \mathbf{Y}^{(n)}).$$

2.1.4. Mode n -Rank of A Tensor. The *mode- n rank* is the rank of the mode- n unfolding. We say a tensor \mathcal{X} has (multilinear) rank $\mathbf{r}(\mathcal{X}) = (r_1, \dots, r_N)$ if its *mode- n rank* is r_n for each $n \in [N]$.

135 **2.1.5. Tensor contractions.** Write $\mathcal{G} = \mathcal{X} \times_n \mathbf{U}$ for the *mode- n (matrix) product*
 136 of \mathcal{X} with $\mathbf{U} \in \mathbb{R}^{J \times I_n}$. That is, $\mathcal{G} = \mathcal{X} \times_n \mathbf{U} \iff \mathbf{G}^{(n)} = \mathbf{U} \mathbf{X}^{(n)}$. The tensor \mathcal{G}
 137 has dimension $I_1 \times \cdots \times I_{n-1} \times J \times I_{n+1} \times \cdots \times I_N$. Mode products with respect to
 138 different modes commute: for $\mathbf{U} \in \mathbb{R}^{J_1 \times I_n}$, $\mathbf{V} \in \mathbb{R}^{J_2 \times I_m}$,

$$139 \quad \mathcal{X} \times_n \mathbf{U} \times_m \mathbf{V} = \mathcal{X} \times_m \mathbf{V} \times_n \mathbf{U} \quad \text{if } n \neq m.$$

140 Mode products obey the associative rule. This rule simplifies mode products with
 141 matrices along the same mode: for $\mathbf{A} \in \mathbb{R}^{J_1 \times I_n}$, $\mathbf{B} \in \mathbb{R}^{J_2 \times J_1}$,

$$142 \quad \mathcal{X} \times_n \mathbf{A} \times_n \mathbf{B} = \mathcal{X} \times_n (\mathbf{B}\mathbf{A}).$$

143 **2.1.6. Tail energy.** To state our results, we will need a tensor equivalent for the
 144 decay in the spectrum of a matrix. For each unfolding $\mathbf{X}^{(n)}$, define the ρ th tail energy

$$145 \quad (\tau_\rho^{(n)})^2 := \sum_{k > \rho}^{\min(I_n, I_{(-n)})} \sigma_k^2(\mathbf{X}^{(n)}),$$

146 where $\sigma_k(\mathbf{X}^{(n)})$ is the k th largest singular value of $\mathbf{X}^{(n)}$.

147 **2.2. Tucker Approximation.** Given a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$ and target rank
 148 $\mathbf{r} = (r_1, \dots, r_N)$, the goal of multilinear approximation is to approximate \mathcal{X} by a
 149 tensor of multilinear rank \mathbf{r} . Concretely, we search over Tucker decompositions
 150 of the approximating tensor with *core tensor* $\mathcal{G} \in \mathbb{R}^{r_1 \times \cdots \times r_N}$ and *factor matrices*
 151 $\mathbf{U}_n \in \mathbb{R}^{I_n \times r_n}$ for $n \in [N]$ with each \mathbf{U}_n satisfying $\mathbf{U}_n^\top \mathbf{U}_n = \mathbf{I}$. For brevity, we define
 152 $\llbracket \mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N \rrbracket = \mathcal{G} \times_1 \mathbf{U}_1 \times_2 \cdots \times_N \mathbf{U}_N$. Any best rank- \mathbf{r} Tucker approximation is
 153 of the form $\llbracket \mathcal{G}^*; \mathbf{U}_1^*, \dots, \mathbf{U}_N^* \rrbracket$, where \mathcal{G}^* , \mathbf{U}_n^* solve the *Tucker approximation* problem

$$154 \quad (2.3) \quad \begin{aligned} & \text{minimize} && \|\mathcal{X} - \mathcal{G} \times_1 \times \cdots \times_N \mathbf{U}_{n+1} \times_N \mathbf{U}_N\|_F^2 \\ & \text{subject to} && \mathbf{U}_n^\top \mathbf{U}_n = \mathbf{I}. \end{aligned}$$

155 The problem (2.3) is a challenging nonconvex optimization problem. Moreover, the
 156 solution is not unique [20]. We use the notation $\llbracket \mathcal{X} \rrbracket_{\mathbf{r}}$ to represent a best rank- \mathbf{r} Tucker
 157 approximation of the tensor \mathcal{X} , which in general we cannot compute.

158 **2.2.1. HOSVD.** The standard approach to computing a rank $\mathbf{r} = (r_1, \dots, r_N)$
 159 Tucker approximation for a tensor \mathcal{X} begins with the higher order singular value
 decomposition (HOSVD) [10, 37]. ([Algorithm 2.1](#)):

Algorithm 2.1 Higher order singular value decomposition (HOSVD) [10, 37]

Given: tensor \mathcal{X} , target rank $\mathbf{r} = (r_1, \dots, r_N)$

1. *Factors.* For $n \in [N]$, compute the top r_n left singular vectors \mathbf{U}_n of $\mathbf{X}^{(n)}$.
2. *Core.* Contract these with \mathcal{X} to form the core

$$\mathcal{G} = \mathcal{X} \times_1 \mathbf{U}_1^T \cdots \times_N \mathbf{U}_N^T.$$

Return: Tucker approximation $\mathcal{X}_{\text{HOSVD}} = \llbracket \mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N \rrbracket$

160

161 The HOSVD can be computed in two passes over the tensor [43, 7]. We describe
 162 this method briefly here, and in more detail in the next section. In the first pass,
 163 sketch each matricization $\mathbf{X}^{(n)}$, $n \in [N]$, and use randomized linear algebra (e.g., the
 164 randomized range finder of [16]) to (approximately) recover its range \mathbf{U}_n . To form the

165 core $\mathcal{X} \times_1 \mathbf{U}_1^T \cdots \times_N \mathbf{U}_N^T$ requires a second pass over \mathcal{X} , since the factor matrices \mathbf{U}_n
 166 depend on \mathcal{X} . The main algorithmic contribution of this paper is to develop a method
 167 to approximate both the factor matrices and the core in just one pass over \mathcal{X} .

168 It is possible to improve the accuracy of the resulting approximation. The higher
 169 order orthogonal iteration (HOOI) [11], for example, uses the HOSVD to initialize an
 170 alternating minimization method, and sequentially minimizes over each of the factor
 171 matrices and the core tensor. However, this method is rarely used in practice due to
 172 the memory and computation required.

173 **2.2.2. ST-HOSVD.** The sequentially truncated higher-order singular value de-
 174 composition (ST-HOSVD) modifies the HOSVD to reduce the computational burden
 175 [38]. This method compresses the target tensor after extracting each factor matrix.
 176 The resulting algorithm can be accelerated using randomized matrix approximations
 177 [27], but seems to require N passes over the tensor. Hence the method is difficult to
 178 implement when the data is too large to store locally.

Algorithm 2.2 Sequentially truncated HOSVD (ST-HOSVD) [38]

Given: tensor \mathcal{X} , target rank $\mathbf{r} = (r_1, \dots, r_N)$

1. $\mathcal{G} = \mathcal{X}$
2. For $n = 1$ to N
 - Compute a best rank r_n approximation of the mode n unfolding of \mathcal{G} :

$$\mathbf{U}_n, \Sigma_n, \mathbf{V}_n = \text{TruncatedSVD}(\mathcal{G}^{(n)}, \mathbf{r}_n).$$

- Form the updated tensor \mathcal{G} from its mode n unfolding $\mathbf{G}^{(n)} \leftarrow \Sigma_n \mathbf{V}_n^T$.

Return: Tucker approximation $\mathcal{X}_{\text{ST-HOSVD}} = \llbracket \mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N \rrbracket$

179 **2.2.3. Quasi-optimality.** A method for tensor approximation is called *quasi-*
 180 *optimal* if the error of the resulting approximation is comparable to the best possible:
 181 more precisely, we say an approximation method is quasi-optimal with factor d if
 182 for any \mathcal{X} and any multilinear rank \mathbf{r} , the rank- \mathbf{r} approximation $\hat{\mathcal{X}}$ produced by the
 183 method satisfies

$$\|\mathcal{X} - \hat{\mathcal{X}}\|_F \leq d \|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F.$$

184
 185 We call a randomized tensor approximation method quasi-optimal if this inequality
 186 holds in expectation. This definition shows the advantage of a quasi-optimal approxi-
 187 mation method: the method finds a good approximation of the tensor whenever a good
 188 rank- \mathbf{r} approximation exists. Moreover, it exactly recovers a rank- \mathbf{r} decomposition of
 189 a tensor that is exactly rank \mathbf{r} .

190 Both the HOSVD and the ST-HOSVD are quasi-optimal with factor \sqrt{N} [38,
 191 15, 13]. This paper demonstrates the first known quasi-optimal streaming Tucker
 192 approximations.

193 **2.3. Previous Work.** The only previous work on streaming Tucker approxima-
 194 tion is [26], which develops a streaming method called Tucker TensorSketch (T.-TS)
 195 [26, Algorithm 2]. T.-TS improves on the HOOI by sketching the data matrix in the
 196 least squares problems. However, the success of the approach depends on the quality
 197 of the initial core and factor matrices, and the alternating least squares algorithm
 198 takes several iterations to converge.

199 In contrast, our work is motivated by the HOSVD (not HOOI), and requires no
 200 initialization or iteration. We treat the tensor as a *multilinear* operator. The sketch

201 identifies a low-dimensional subspace *for each mode of the tensor* that captures the
 202 action of the operator along that mode. The reconstruction produces a low-Tucker-rank
 203 multilinear operator with the same action on this low-dimensional tensor product space.
 204 This linear algebraic view allows us to develop the first guarantees on approximation
 205 error for this class of problems¹. Moreover, we show numerically that our algorithm
 206 achieves a better approximation of the original tensor given the same storage budget.

207 More generally, there is a large literature on randomized algorithms for matrix
 208 factorizations and for solving optimization problems; for example, see the review
 209 articles [16, 41]. In particular, our method is strongly motivated by the recent papers
 210 [34, 35], which provide methods for one-pass matrix approximation. The novelty
 211 of this paper is in our design of a core sketch (and reconstruction) for the Tucker
 212 decomposition, together with provable performance guarantees. The proof requires
 213 a careful accounting of the errors resulting from the factor sketches and from the
 214 core sketch. The structure of the Tucker sketch guarantees that these errors are
 215 independent.

216 Many researchers have used randomized algorithms to compute tensor decomposi-
 217 tions. For example, [40, 6] apply sketching techniques to the CP decomposition, while
 218 [36] suggests sparsifying the tensor. Several papers aim to make Tucker decomposition
 219 efficient in the limited-memory or distributed settings [5, 43, 3, 19, 24, 7].

220 **3. Dimension Reduction Maps.** In this section, we first introduce some com-
 221 monly used randomized dimension reduction maps together with some mathematical
 222 background, and explain how to calculate and update sketches.

223 **3.1. Dimension Reduction Map.** Dimension reduction maps (DRMs) take a
 224 collection of high-dimensional objects to a lower-dimensional space while maintaining
 225 certain geometric properties [29]. For example, we may wish to preserve the pairwise
 226 distances between vectors, or to preserve the column space of matrices. We call the
 227 output of a DRM on an object x a *sketch* of x .

228 Common DRMs include matrices with i.i.d. Gaussian entries or i.i.d. Rademacher
 229 entries (uniform on $\{\pm 1\}$). The Scrambled Subsampled Randomized Fourier Transform
 230 (SSRFT) [42] and sparse random projections [1, 25] can achieve similar performance
 231 with fewer computational and storage requirements; see supplement for details.

232 Our theoretical bounds rely on properties of the Gaussian DRM. However, our
 233 numerical experiments indicate that many other DRMs yield qualitatively similar
 234 results; see supplement for details.

235 **3.2. Tensor Random Projection.** Here we present a strategy for reducing the
 236 storage of the random map that makes use of the tensor random projection (TRP), and
 237 extremely low storage structured dimension reduction map proposed in [32]. The *tensor*
 238 *random projection (TRP)* $\Omega : \prod_{n=1}^N I_n \rightarrow \mathbb{R}^k$ is defined as the iterated Khatri–Rao
 239 product of DRMs $\mathbf{A}_n \in \mathbb{R}^{I_n \times k}$, $n \in [N]$:

$$240 \quad (3.1) \quad \Omega = \mathbf{A}_1 \odot \cdots \odot \mathbf{A}_N.$$

241 Each $\mathbf{A}_n \in \mathbb{R}^{I_n \times k}$ can be a Gaussian map, a Rademacher matrix, an SSRFT, etc. The
 242 number of constituent maps N and their dimensions I_n for $n \in [N]$ are parameters of

¹ The guarantees in [26] hold only when a new sketch is applied for each subsequent least squares solve; the resulting algorithm cannot be used in a streaming setting. In contrast, the practical streaming method T-TS fixes the sketch for each mode, and so has no known guarantees. Interestingly, experiments in [26] show that the method achieves lower error using a fixed sketch (with no guarantees) than using fresh sketches at each iteration.

243 the TRP, and control the quality of the map; see [32] for details. The TRP map is
 244 a row-product random matrix, which behaves like a Gaussian map in many respects
 245 [30]. Our experimental results confirm this behavior.

246 For simplicity, suppose I_n is the same for each $n \in [N]$. Then the TRP can be
 247 formed (and stored) using only kNI random variables, while standard dimension
 248 reduction maps use randomness (and storage) that grows as I^N when applied to a
 249 generic (dense) tensor. Table 1 compares the computational and storage costs for
 different DRMs.

DRM	Storage	Computation
Gaussian	kI^N	kI^N
Sparse	μkI^N	μkI^N
SSRFT	I^N	$\log(k)I^N$
TRP	kNI	kI^N

Table 1: Performance of Different Dimension Reduction Maps: We compare the storage and the computational cost of applying a DRM mapping \mathbb{R}^{I^N} to \mathbb{R}^k to a dense tensor in \mathbb{R}^{I^N} . Here μ is the fraction of nonzero entries in the sparse DRMs. The TRP considered here is composed of Gaussian DRMs.

250

251 We do not need to explicitly form or store the TRP map Ω . Instead, we can
 252 store the constituent DRMs $\mathbf{A}_1, \dots, \mathbf{A}_N$ and compute the action of the map on the
 253 matricized tensor using the definition of the TRP. The additional computation required
 254 is minimal, and it empirically incurs almost no performance loss.

255 **4. Algorithms for Tucker approximation.** In this section, we present our
 256 proposed tensor sketch, our algorithms for one- and two-pass Tucker approximation,
 257 and we discuss the computational complexity and storage required for both sparse and
 258 dense input tensors. We present guarantees for these methods in section 5.

259 **4.1. Tensor compression via sketching.** Our Tucker sketch generalizes the
 260 matrix sketch of [35] to higher order tensors. To compute a Tucker sketch for tensor
 261 $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ with sketch size parameters \mathbf{k} and \mathbf{s} , draw independent, random DRMs

262 (4.1)
$$\Omega_1, \Omega_2, \dots, \Omega_N \quad \text{and} \quad \Phi_1, \Phi_2, \dots, \Phi_N,$$

263 with $\Omega_n \in \mathbb{R}^{I_{(-n)} \times k_n}$ and $\Phi_n \in \mathbb{R}^{I_n \times s_n}$ for $n \in [N]$. Use these DRMs to compute

264
$$\mathbf{V}_n = \mathbf{X}^{(n)} \Omega_n \quad \in \mathbb{R}^{I_n \times k_n}, \quad n \in [N],$$

265
$$\mathcal{H} = \mathcal{X} \times_1 \Phi_1^\top \cdots \times_N \Phi_N^\top \quad \in \mathbb{R}^{s_1 \times \dots \times s_N}.$$

267 The *factor sketch* \mathbf{V}_n captures the span of the mode- n fibers of \mathcal{X} for each $n \in [N]$,
 268 while the *core sketch* \mathcal{H} contains information about the interaction between different
 269 modes. See Algorithm 4.1 for pseudocode.

270 To produce a rank $\mathbf{r} = \{r_1, \dots, r_N\}$ Tucker approximation of \mathcal{X} , choose sketch size
 271 parameters $\mathbf{k} = (k_1, \dots, k_N) \geq \mathbf{r}$ and $\mathbf{s} = (s_1, \dots, s_N) \geq \mathbf{k}$. (Vector inequalities hold
 272 elementwise.) Our approximation guarantees depend closely on the parameters \mathbf{k} and
 273 \mathbf{s} . As a rule of thumb, we suggest selecting $\mathbf{s} = 2\mathbf{k} + 1$, as the theory requires $\mathbf{s} > 2\mathbf{k}$,
 274 and choosing \mathbf{k} as large as possible given storage limitations.

275 The sketches \mathbf{V}_n and \mathcal{H} are linear functions of the original tensor \mathcal{X} so can be
 276 computed in a single pass over \mathcal{X} . Linearity enables easy computation of the sketch

277 even in the streaming model or distributed model (see supplement for details). Storing
 278 the sketches requires memory $\sum_{n=1}^N I_n \cdot k_n + \Pi_{i=1}^N s_n$: much less than the full tensor.

Algorithm 4.1 Tucker Sketch

Given: RDRM (a function that generates a random DRM)

- 1: **function** TUCKERSKETCH($\mathcal{X}; \mathbf{k}, \mathbf{s}$)
 - 2: Form DRMs $\Omega_n = \text{RDRM}(I_{(-n)}, k_n)$ and $\Phi_n = \text{RDRM}(I_n, s_n)$, $n \in [N]$
 - 3: Compute factor sketches $\mathbf{V}_n = \mathbf{X}^{(n)} \Omega_n$, $n \in [N]$
 - 4: Compute core sketch $\mathcal{H} = \mathcal{X} \times_1 \Phi_1^\top \times \cdots \times_N \Phi_N^\top$
 - 5: **return** $(\mathcal{H}, \mathbf{V}_1, \dots, \mathbf{V}_N, \{\Phi_n, \Omega_n\}_{n \in [N]})$
 - 6: **end function**
-

279 *Remark 4.1.* The DRMs $\Omega_n \in \mathbb{R}^{I_{(-n)} \times k_n}$ are large—much larger than the size of
 280 the Tucker factorization we seek! Even using a low memory mapping such as the
 281 SSRFT and sparse random map, the storage required grows as $\mathcal{O}(I_{(-n)})$. However, we
 282 do not need to store these matrices. Instead, we can generate (and regenerate) them
 283 as needed using a (stored) random seed.²

284 *Remark 4.2.* Alternatively, the TRP (subsection 3.2) can be used to limit the
 285 storage required for Ω_n . The Khatri–Rao structure in the sketch need not match the
 286 structure in the matricized tensor. However, we can take advantage of the structure
 287 of our problem to reduce storage even further. We generate DRMs $\mathbf{A}_n \in \mathbb{R}^{I_n \times k}$ for
 288 $n \in [N]$ and define $\Omega_n = \mathbf{A}_1 \odot \cdots \odot \mathbf{A}_{n-1} \odot \mathbf{A}_{n+1} \odot \cdots \odot \mathbf{A}_N$ for each $n \in [N]$. Hence
 289 we need not store the maps Ω_n , but only the small matrices \mathbf{A}_n . The storage required
 290 is thereby reduced from $\mathcal{O}(N(\prod_{n=1}^N I_n)k)$ to $\mathcal{O}((\sum_{n=1}^N I_n)k)$, while the approximation
 291 error is essentially unchanged. We use this method in our experiments.

292 **4.2. Low-Rank Approximation.** Now we explain how to construct a Tucker
 293 decomposition of \mathcal{X} with target multilinear rank \mathbf{k} from the factor and core sketches.

294 We first present a simple two-pass algorithm, Algorithm 4.2, that uses only the
 295 factor sketches by projecting the unfolded matrix of original tensor \mathcal{X} to the column
 296 space of each factor sketch. (Notice that Algorithm 4.2 does not use the core sketch,
 297 so the core sketch parameter \mathbf{s} of the Tucker Sketch is set to 0.)

298 To project to the column space of each factor matrix, we calculate the QR
 299 decomposition of each factor sketch:

$$300 \quad (4.2) \quad \mathbf{V}_n = \mathbf{Q}_n \mathbf{R}_n \quad \text{for } n \in [N],$$

301 where $\mathbf{Q}_n \in \mathbb{R}^{I_n \times k_n}$ has orthonormal columns and $\mathbf{R}_n \in \mathbb{R}^{k_n \times k_n}$ is upper triangular.
 302 Consider the tensor approximation

$$303 \quad (4.3) \quad \hat{\mathcal{X}}_2 = \mathcal{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top.$$

304 This approximation admits the guarantees stated in Theorem 5.1. Using the commu-
 305 tativity of the mode product between different modes, we can rewrite $\hat{\mathcal{X}}$ as

$$306 \quad (4.4) \quad \hat{\mathcal{X}}_2 = \underbrace{[\mathcal{X} \times \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top]}_{\mathbf{w}_2} \times_1 \mathbf{Q}_1 \times_2 \cdots \times_N \mathbf{Q}_N = \llbracket \mathbf{w}_2; \mathbf{Q}_1, \dots, \mathbf{Q}_N \rrbracket,$$

² Our theory assumes the DRMs are random, whereas our experiments use pseudorandom numbers. In fact, for many pseudorandom number generators it is NP-hard to determine whether the output is random or pseudorandom [2]. In particular, we expect both to perform similarly for tensor approximation.

307 which gives an explicit Tucker approximation $\tilde{\mathcal{X}}$ of our original tensor. The core
 308 approximation $\mathcal{W}_2 \in \mathbb{R}^{k_1 \times \dots \times k_N}$ is much smaller than the original tensor \mathcal{X} . To
 309 compute this approximation, we need access to \mathcal{X} twice: once to compute $\mathbf{Q}_1, \dots, \mathbf{Q}_N$,
 310 and again to apply them to \mathcal{X} in order to form \mathcal{W}_2 .

Algorithm 4.2 Two-Pass Sketch and Low-Rank Recovery

Given: tensor \mathcal{X} , sketch parameters \mathbf{k}

1. *Sketch.* $(\mathcal{H}, \mathbf{V}_1, \dots, \mathbf{V}_N, \{\Phi_n, \Omega_n\}_{n \in [N]}) = \text{TUCKERSKETCH}(\mathcal{X}; \mathbf{k}, 0)$
2. *Recover factor matrices.* For $n \in [N]$, $(\mathbf{Q}_n, \sim) = \text{QR}(\mathbf{V}_n)$
3. *Recover core.* $\mathcal{W}_2 = \mathcal{X} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N$

Return: Tucker approximation $\hat{\mathcal{X}}_2 = \llbracket \mathcal{W}_2; \mathbf{Q}_1, \dots, \mathbf{Q}_N \rrbracket$ with rank $\leq \mathbf{k}$

311 This two-pass algorithm, [Algorithm 4.2](#), has a simple motivation. In the first step
 312 of the HOSVD, [Algorithm 2.1](#), we approximately compute the top r_n eigenvectors
 313 of each matricization $\mathbf{X}^{(n)}$ using the randomized SVD [16]. Indeed, the same idea
 314 was proposed in concurrent work [27], which extends the idea to the ST-HOSVD
 315 and provides an error analysis. The error analyses of the two papers essentially
 316 coincide for [Algorithm 4.2](#). One major difference is that the authors of [27] focus on
 317 the computational benefits gained using the randomized SVD, while here we focus
 318 primarily on the benefits due to reduced storage.

319 To find an algorithm for streaming data, when it is impossible to store the full
 320 tensor, we require a one-pass algorithm.

321 *One-Pass Approximation.* To develop an one-pass method, we must use the core
 322 sketch \mathcal{H} (the compression of \mathcal{X} using the random projections Φ_n) to approximate \mathcal{W}_2
 323 (the compression of \mathcal{X} using random projections \mathbf{Q}_n). To develop intuition, consider
 324 the following calculation: if the factor matrix approximations \mathbf{Q}_n capture the range of
 325 \mathcal{X} well, then projection onto their ranges in each mode approximately preserves the
 326 action of \mathcal{X} :

$$327 \quad \mathcal{X} \approx \mathcal{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top.$$

328 Recall that for tensor \mathcal{A} and matrices \mathbf{B} and \mathbf{C} with compatible sizes, $\mathcal{A} \times_n (\mathbf{BC}) =$
 329 $(\mathcal{A} \times_n \mathbf{C}) \times_n \mathbf{B}$. Use this rule to recognize the two-pass core approximation \mathcal{W}_2 :

$$330 \quad \mathcal{X} \approx (\mathcal{X} \times_1 \mathbf{Q}_1^\top \times \cdots \times_N \mathbf{Q}_N^\top) \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N = \mathcal{W}_2 \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N$$

331 Now contract both sides of this approximate equality with the DRMs Φ_n to identify
 332 the core sketch \mathcal{H} :

$$333 \quad \mathcal{H} := \mathcal{X} \times_1 \Phi_1^\top \cdots \times_N \Phi_N^\top \approx \mathcal{W}_2 \times_1 \Phi_1^\top \mathbf{Q}_1 \times \cdots \times_N \Phi_N^\top \mathbf{Q}_N.$$

334 We have chosen $s > \mathbf{k}$ so each $\Phi_n^\top \mathbf{Q}_n$ has a left inverse with high probability. Hence,
 335 we can solve the approximate equality for \mathcal{W}_2 :

$$336 \quad \mathcal{W}_2 \approx \mathcal{H} \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger =: \mathcal{W}_1.$$

337 The right-hand side of the approximation defines the one-pass core approximation \mathcal{W}_1 .
 338 [Lemma A.2](#) controls the error in this approximation.

339 [Algorithm 4.3](#) summarizes the resulting one-pass algorithm. One (streaming) pass
 340 over the tensor can be used to sketch the tensor; to recover the tensor, we only access
 341 the sketches. [Theorem 5.3](#) (below) bounds the overall quality of the approximation.

Algorithm 4.3 One-Pass Sketch and Low-Rank Recovery**Given:** tensor \mathcal{X} , sketch parameters \mathbf{k} and $\mathbf{s} > \mathbf{k}$

1. *Sketch.* $(\mathcal{H}, \mathbf{V}_1, \dots, \mathbf{V}_N, \{\Phi_n, \Omega_n\}_{n \in [N]}) = \text{TUCKERSKETCH}(\mathcal{X}; \mathbf{k}, \mathbf{s})$
2. *Recover factor matrices.* For $n \in [N]$, $(\mathbf{Q}_n, \sim) = \text{QR}(\mathbf{V}_n)$
3. *Recover core.* $\mathcal{W}_1 = \mathcal{H} \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times \dots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger$

Return: Tucker approximation $\hat{\mathcal{X}}_1 = \llbracket \mathcal{W}_1; \mathbf{Q}_1, \dots, \mathbf{Q}_N \rrbracket$ with rank $\leq \mathbf{k}$

342 The computational complexity and storage required by [Algorithm 4.3](#) is presented
 343 in [Table 2](#). These requirements compare favorably to the only previous method for
 344 streaming Tucker approximation [\[26\]](#); see the supplement for details.

Stage	Computation	Storage
Sketching	$\mathcal{O}(((1 - (s/I)^N)/(1 - (s/I)) + Nk)I^N)$	$kNI + s^N$
Recovery	$\mathcal{O}((k^2 s^N (1 - (k/s)^N))/(1 - k/s) + k^2 NI)$	
Total	$\mathcal{O}((s(1 - (s/I)^N))/(1 - s/I) + Nk)I^N)$	

Table 2: Computational Complexity of one-pass approximation ([Algorithm 4.3](#)) on tensor $\mathcal{X} \in \mathbb{R}^{I \times \dots \times I}$ with parameters (k, s) , using a TRP composed of Gaussian DRMs inside the Tucker sketch. Most of the time is spent sketching the tensor \mathcal{X} .

345 **4.3. Fixed-Rank Approximation.** The low-rank approximation methods [Al-](#)
 346 [gorithms 4.2](#) and [4.3](#) of the previous section produce approximations with rank no
 347 more than \mathbf{k} . It is often valuable to truncate this approximation to a user-specified
 348 target rank $\mathbf{r} \leq \mathbf{k}$ [\[35, Figure 4\]](#). Increasing \mathbf{k} relative to \mathbf{r} can improve the quality
 349 of the final approximation by using more intermediate storage, without changing the
 350 storage requirements of the final approximation to \mathcal{X} . In this section, we introduce a
 351 few methods to compute fixed-rank approximations with rank no more than \mathbf{r} by way
 352 of a sketch with parameter $\mathbf{k} \geq \mathbf{r}$.

353 **4.3.1. Truncated QR.** One simple fix to [Algorithms 4.2](#) and [4.3](#) results in a
 354 final approximation with rank \mathbf{r} rather than \mathbf{k} : simply replace the QR decomposition
 355 with a truncated QR decomposition [\[14\]](#). Indeed, we will show that this simple change
 356 results in a one-pass algorithm that achieves quasi-optimality with factor $2\sqrt{N}$, nearly
 357 matching the guarantee for the HOSVD and ST-HOSVD. This approach is best for
 358 tensors with many modes that are (almost) exactly rank \mathbf{r} . However, for tensors with
 359 few modes and slower spectral decay, a more sophisticated fixed-rank approximation
 360 method outperforms this naïve approach.

361 **4.3.2. Optimal Fixed-Rank Approximation..** For tensors with few modes,
 362 we recommend computing a rank- \mathbf{r} approximation to \mathcal{X} by forming an initial ap-
 363 proximation with rank $\mathbf{k} \geq \mathbf{r}$ using a randomized method such as [Algorithm 4.3](#) or
 364 [Algorithm 4.2](#) and then truncating it to rank \mathbf{r} using a deterministic method such
 365 as ST-HOSVD. In previous work, we have found that rank truncation is essential
 366 to ensure that the final approximation is fully reliable: for matrices, we find the
 367 top singular values and vectors of the approximation are accurate when $\mathbf{k} \gtrsim 4\mathbf{r}$ [\[35\]](#).
 368 Rank truncation can also be used to choose the final size of the Tucker decomposition
 369 adaptively to achieve a desired approximation quality.

370 For moderate \mathbf{k} , it is computationally easy to truncate an initial rank- \mathbf{k} approxi-
 371 mation to rank \mathbf{r} , thanks to the following lemma.

372 LEMMA 4.1 (Core truncation). *Let $\mathcal{W} \in \mathbb{R}^{k_1 \times \dots \times k_N}$ be a tensor with $\mathbf{k} \geq \mathbf{r}$, and*
 373 *let $\mathbf{Q}_n \in \mathbb{R}^{I_n \times k_n}$ be orthogonal matrices for each $n \in [N]$. Then*

374
$$\llbracket \mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N \rrbracket_{\mathbf{r}} = \llbracket \mathcal{W} \rrbracket_{\mathbf{r}} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N.$$

375 Lemma 4.1 shows that we can compute the optimal rank- \mathbf{r} approximation of the (large)
 376 tensor $\hat{\mathcal{X}} = \mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N$ by calculating the optimal rank- \mathbf{r} approximation of
 377 the (small) core \mathcal{W} . Interestingly, the same result holds if we replace the best rank- \mathbf{r}
 378 Tucker approximation $\llbracket \cdot \rrbracket$ by the HOOI [8, Lemma A.1].

379 *Proof of Lemma 4.1.* The target tensor to be approximated, $\mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N$,
 380 lies in the subspace spanned by the \mathbf{Q}_n , $\{\mathcal{X} : \mathcal{X}^{(n)} \in \text{range}(\mathbf{Q}_n)\}$. By the Pythagorean
 381 theorem, any optimal Tucker decomposition also lies in this subspace.

382 Suppose $\llbracket \mathcal{W}' ; \mathbf{V}_1, \dots, \mathbf{V}_N \rrbracket$ is an optimal Tucker decomposition. Since its n th
 383 unfolding is in $\text{range}(\mathbf{Q}_n)$, each \mathbf{V}_n can be written as $\mathbf{Q}_n \mathbf{U}_n$ for some orthogonal
 384 $\mathbf{U}_n \in \mathbb{R}^{k_n \times r_n}$. Then, using the orthogonal invariance of the Frobenius norm,

385
$$\begin{aligned} & \|\mathcal{W} \times_1 \mathbf{Q}_1 \times \cdots \times_N \mathbf{Q}_N - \mathcal{W}' \times_1 \mathbf{Q}_1 \mathbf{U}_1 \times \cdots \times_N \mathbf{Q}_N \mathbf{U}_N\|_F \\ 386 &= \|\mathcal{W} - \mathcal{W}' \times_1 \mathbf{U}_1 \times \cdots \times_N \mathbf{U}_N\|_F \geq \|\mathcal{W} - \llbracket \mathcal{W} \rrbracket_{\mathbf{r}}\|_F \\ 387 &= \|\mathcal{W} \times_1 \mathbf{Q}_1 \times \cdots \times_N \mathbf{Q}_N - \llbracket \mathcal{W} \rrbracket_{\mathbf{r}} \times_1 \mathbf{Q}_1 \times \cdots \times_N \mathbf{Q}_N\|_F. \quad \square \end{aligned}$$

389 Motivated by this lemma, to produce a fixed rank- \mathbf{r} approximation of \mathcal{X} , we
 390 compress the core tensor approximation from Algorithm 4.2 or Algorithm 4.3 to rank
 391 \mathbf{r} . This compression is cheap because the core approximation $\mathcal{W} \in \mathbb{R}^{k_1 \times \dots \times k_N}$ is small.

392 We present this method (using ST-HOSVD as the compression algorithm) as
 393 Algorithm 4.4. One convenient aspect of this scheme is that the rank- \mathbf{k} approximation
 394 can be stored in memory, allowing users to experiment with different desired final
 395 ranks \mathbf{r} and with different algorithms \mathcal{A} to compress the core to rank \mathbf{r} . Reasonable
 396 choices to compress the core include the HOSVD, the ST-HOSVD, or TTHRESH
 397 [4]. It is possible to use these strategies to adaptively compute a core approximation
 398 that achieves some target approximation error. For example, for the HOSVD, one
 399 can successively truncate the core of the HOSVD (using the ordering property [10,
 400 Theorem 2]); for the ST-HOSVD, one can use the error tolerance strategy of [38];
 401 and the iterative strategy of TTHRESH naturally terminates upon reaching a target
 402 error approximation. Both HOSVD and ST-HOSVD are quasi-optimal [38], while
 403 ST-HOSVD requires less storage and computation.

Algorithm 4.4 Fixed-rank approximation

Given: Tucker approximation $\llbracket \mathcal{W}; \mathbf{Q}_1, \dots, \mathbf{Q}_N \rrbracket$ of tensor \mathcal{X} , rank target \mathbf{r} , algorithm $\mathcal{A}(\mathcal{W}, \mathbf{r})$ to compute rank \mathbf{r} approximation to \mathcal{W} (e.g., ST-HOSVD).

1. *Approximate core with fixed rank.* $\mathcal{G}, \mathbf{U}_1, \dots, \mathbf{U}_N = \mathcal{A}(\mathcal{W}, \mathbf{r})$
2. *Compute factor matrices.* For $n \in [N]$, $\mathbf{P}_n = \mathbf{Q}_n \mathbf{U}_n$

Return: Tucker approximation $\hat{\mathcal{X}}_{\mathbf{r}} = \llbracket \mathcal{G}; \mathbf{P}_1, \dots, \mathbf{P}_N \rrbracket$ with rank $\leq \mathbf{r}$

404 **5. Guarantees.** In this section, we present probabilistic guarantees on the pre-
 405 ceding algorithms. We show that the approximation error for the one-pass algorithm
 406 is the sum of the error from the two-pass algorithm and the error resulting from the
 407 core approximation. We present most of the proofs in this section, and defer some
 408 more technical parts to the appendix.

409 **5.1. Low-rank approximation.** [Theorem 5.1](#) guarantees the performance of
410 the two-pass method ([Algorithm 4.2](#)).

411 **THEOREM 5.1** (Two-pass low-rank approximation). *Sketch the tensor \mathbf{X} using*
412 *a Tucker sketch with parameters \mathbf{k} using DRMs with i.i.d. standard normal entries.*
413 *Then the approximation $\hat{\mathbf{X}}_2$ computed by the two-pass method ([Algorithm 4.2](#)) satisfies*

$$414 \quad \mathbb{E}\|\mathbf{X} - \hat{\mathbf{X}}_2\|_F^2 \leq \min_{1 \leq \rho \leq k-1} \sum_{n=1}^N \left(1 + \frac{\rho_n}{k_n - \rho_n - 1}\right) (\tau_{\rho_n}^{(n)})^2.$$

415 This theorem shows that the proposed randomized tensor approximation works best for
416 tensors whose unfoldings exhibit spectral decay. As a simple consequence, we see that
417 the two-pass method with $\mathbf{k} > \mathbf{r} + 1$ perfectly recovers a tensor with exact (multilinear)
418 rank \mathbf{r} , since in that case $\tau_{r_n}^{(n)} = 0$ for each $n \in [N]$. However, the theorem states a
419 stronger bound: the method exploits decay in the spectrum, wherever (in the first k_n
420 singular values of each mode n unfolding) it occurs.

421 Another useful consequence shows that the rank- \mathbf{k} approximation computed with
422 this two-pass method competes with the best rank- \mathbf{r} approximation.

423 **COROLLARY 5.2.** *Suppose $\mathbf{k} \geq 2\mathbf{r} + 1$. Then the approximation $\hat{\mathbf{X}}_2$ computed by*
424 *the two-pass method ([Algorithm 4.2](#)) satisfies*

$$425 \quad \mathbb{E}\|\mathbf{X} - \hat{\mathbf{X}}_2\|_F^2 \leq 2 \sum_{n=1}^N (\tau_{r_n}^{(n)})^2 \leq 2N\|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_{\mathbf{r}}\|_F^2.$$

426 *Proof of [Theorem 5.1](#).* Suppose $\hat{\mathbf{X}}_2$ is the low-rank approximation from [Algo-](#)
427 [rithm 4.2](#). Use the definition of the mode- n product and the commutativity of the
428 mode product between different modes to see

$$429 \quad \begin{aligned} \hat{\mathbf{X}}_2 &= [\mathbf{X} \times_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top] \times_1 \mathbf{Q}_1 \times_1 \cdots \times_N \mathbf{Q}_N \\ &= \mathbf{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top. \end{aligned}$$

430 Here we see that $\hat{\mathbf{X}}_2$ is a *multi-linear orthogonal projection* of \mathbf{X} onto the subspace
431 spanned by the \mathbf{Q}_n , $\{\mathbf{X} : \mathbf{X}^{(n)} \in \text{range}(\mathbf{Q}_n)\} = \{\llbracket \mathbf{W}; \mathbf{Q}_1, \dots, \mathbf{Q}_N \rrbracket : \mathbf{W} \in \mathbb{R}^{k_1 \times \cdots \times k_N}\}$.
432 (See [\[12\]](#) for more background on multi-linear orthogonal projection.) As in [\[38\]](#),
433 [Theorem 5.1](#)], we sequentially apply the Pythagorean theorem to each mode to show

$$434 \quad (5.1) \quad \|\hat{\mathbf{X}}_2 - \mathbf{X}\|_F^2 \leq \sum_{n=1}^N \left\| (\mathbf{I} - \mathbf{Q}_n \mathbf{Q}_n^\top) \mathbf{X}^{(n)} \right\|_F^2.$$

435 We then take the expectation over \mathbf{Q}_n for each term in the sum and use [Corollary B.1](#)
436 to show this expectation is bounded by the corresponding tail energy,

$$437 \quad \mathbb{E} \left\| (\mathbf{I} - \mathbf{Q}_n \mathbf{Q}_n^\top) \mathbf{X}^{(n)} \right\|_F^2 \leq \min_{1 \leq \rho_n \leq k_n - 1} \left(1 + \frac{\rho_n}{k_n - \rho_n - 1}\right) (\tau_{\rho_n}^{(n)})^2. \quad \square$$

438 [Theorem 5.3](#) guarantees the performance of the one pass method [Algorithm 4.3](#).

439 **THEOREM 5.3** (One-pass low-rank approximation). *Sketch the tensor \mathbf{X} using*
440 *a Tucker sketch with parameters \mathbf{k} and $\mathbf{s} \geq 2\mathbf{k}$ using DRMs with i.i.d. standard*
441 *normal entries. Then the approximation $\hat{\mathbf{X}}_1$ computed with the one-pass method*

442 (Algorithm 4.3) satisfies the bound

$$443 \quad \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_1\|_F^2 \leq (1 + \Delta) \min_{1 \leq \rho \leq k-1} \sum_{n=1}^N \left(1 + \frac{\rho_n}{k_n - \rho_n - 1}\right) (\tau_{\rho_n}^{(n)})^2,$$

444 where $\Delta := \max_{n=1}^N k_n / (s_n - k_n - 1)$.

445 The resulting rank- \mathbf{k} approximation, computed in a single pass, is nearly optimal.

446 COROLLARY 5.4. Suppose $\mathbf{k} \geq 2\mathbf{r} + 1$ and $\mathbf{s} \geq 2\mathbf{k}$. Then the approximation $\hat{\mathcal{X}}_2$
447 computed by the one-pass method (Algorithm 4.3) satisfies

$$448 \quad \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_1\|_F^2 \leq 4 \sum_{n=1}^N (\tau_{r_n}^{(n)})^2 \leq 4N \|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_r\|_F^2.$$

449 *Proof of Theorem 5.3.* We decompose the approximation error into the error due
450 to the factor matrix approximations and the error due to the core approximation.
451 Recall that $\hat{\mathcal{X}}_1$ is the one-pass approximation from Algorithm 4.3, and

$$452 \quad \hat{\mathcal{X}}_2 = \mathcal{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top,$$

453 is the two-pass approximation from Algorithm 4.2. The the one-pass and two-pass
454 approximations differ only in the core approximation:

$$455 \quad (5.2) \quad \hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2 = (\mathcal{W} - \mathcal{X} \times_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top) \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N.$$

456 Thus $\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2$ is in the subspace spanned by the \mathbf{Q}_n , $\{\mathcal{X} : \mathcal{X}^{(n)} \in \text{range}(\mathbf{Q}_n)\}$, while
457 $\hat{\mathcal{X}}_2 - \mathcal{X}$ is orthogonal to that subspace. Therefore,

$$458 \quad \langle \hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2, \hat{\mathcal{X}}_2 - \mathcal{X} \rangle = 0.$$

459 Now, use the Pythagorean theorem to bound the error of the one-pass approximation:

$$460 \quad (5.3) \quad \|\hat{\mathcal{X}}_1 - \mathcal{X}\|_F^2 = \|\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2\|_F^2 + \|\hat{\mathcal{X}}_2 - \mathcal{X}\|_F^2.$$

461 Consider the first term. Using (5.2), we see that

$$462 \quad \begin{aligned} \|\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2\|_F^2 &= \|(\mathcal{W}_1 - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top) \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N\|_F^2 \\ 463 \quad &= \|(\mathcal{W}_1 - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top)\|_F^2, \end{aligned}$$

465 where we use the orthogonal invariance of the Frobenius norm for the second equality.

466 Next, we use Lemma A.2 to bound the expected error from the core approximation as

$$467 \quad \mathbb{E}\|\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2\|_F^2 \leq \Delta \|\mathcal{X} - \hat{\mathcal{X}}_2\|_F^2.$$

468 Taking the expectation of (5.3) and using this bound on the core error, we find

$$469 \quad \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_1\|_F^2 \leq (1 + \Delta) \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_2\|_F^2$$

470 Finally, we use the two-pass approximation error bound Theorem 5.1:

$$471 \quad \mathbb{E}\|\hat{\mathcal{X}}_2 - \mathcal{X}\|_F^2 \leq \min_{1 \leq \rho_n < k_n - 1} \left[\sum_{n=1}^N \left(1 + \frac{\rho_n}{k_n - \rho_n - 1}\right) (\tau_{\rho_n}^{(n)})^2 \right]. \quad \square$$

472 We see that the additional error due to sketching the core is a multiplicative factor
 473 Δ more than the error due to sketching the factor matrices. This factor Δ decreases
 474 as the size of the core sketch \mathbf{s} increases.

475 **Theorem 5.3** also offers guidance on how to select the sketch size parameters \mathbf{s} and
 476 \mathbf{k} . In particular, suppose that the mode- n unfolding has a good rank r_n approximation
 477 for each mode n . Then the choices $k_n = 2r_n + 1$ and $s_n = 2k_n + 1$ ensure that

$$478 \quad \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_1\|_F^2 \leq 4 \sum_{n=1}^N (\tau_{r_n}^{(n)})^2.$$

479 More generally, as k_n/r_n and s_n/k_n increase, the leading constant in the approximation
 480 error tends to one.

481 **5.2. Fixed-rank approximation.** We now present bounds on the error of the
 482 fixed rank- \mathbf{r} approximations produced either using the truncated QR method in
 483 [Algorithms 4.2](#) and [4.3](#) or by using the fixed-rank approximation on the output of
 484 [Algorithms 4.2](#) and [4.3](#). The former method produces algorithms that are quasi-optimal
 485 with factor $\sqrt{2N}$ (two pass) or $2\sqrt{N}$ (one pass), matching the rate of the HOSVD
 486 and the ST-HOSVD. The latter method produces algorithms that are quasi-optimal
 487 with factor that grows linearly in the number of modes N , but which adapts better to
 488 spectral decay. For tensors with few modes and high dimension, such as those that
 489 appear in our numerical experiments, the latter methods substantially outperform the
 490 former.

491 The resulting bounds show that the best rank- \mathbf{r} approximation of the output
 492 from the one- or two-pass algorithms is comparable in quality to a true best rank- \mathbf{r}
 493 approximation of the input tensor. An important insight is that the sketch size
 494 parameters \mathbf{s} and \mathbf{k} that guarantee a good low-rank approximation also guarantee a
 495 good fixed-rank approximation: the error due to sketching depends only on the sketch
 496 size parameters \mathbf{k} and \mathbf{s} , and not on the target rank \mathbf{r} .

497 **5.2.1. Truncated QR.** We can modify the argument in the proofs of [Theorems](#)
 498 [5.1](#) and [5.3](#) to provide an error bound for the rank- \mathbf{r} approximation obtained by
 499 truncating the QR decomposition in [Algorithms 4.2](#) and [4.3](#) to rank \mathbf{r} as in [section 4.3.1](#).
 500 This bound will allow us to show quasi-optimality of the resulting algorithm.

501 **THEOREM 5.5** (Fixed-rank approximation via truncated QR). *Sketch the ten-*
 502 *tor \mathcal{X} using a Tucker sketch with parameters \mathbf{k} using DRMs with i.i.d. standard*
 503 *normal entries. The rank- \mathbf{r} approximation $\hat{\mathcal{X}}_2$ computed with the two-pass method*
 504 *([Algorithm 4.2](#)), using a rank- \mathbf{r} truncated QR in step 2 of the algorithm, satisfies*

$$505 \quad \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_2\|_F^2 \leq \sum_{n=1}^N \left(1 + \frac{r_n}{k_n - r_n - 1}\right) (\tau_{r_n}^{(n)})^2.$$

506 *Similarly, the rank- \mathbf{r} approximation $\hat{\mathcal{X}}_1$ computed with one-pass method ([Algorithm 4.3](#)),*
 507 *using a rank- \mathbf{r} truncated QR in step 2 of the algorithm, satisfies*

$$508 \quad \mathbb{E}\|\mathcal{X} - \hat{\mathcal{X}}_1\|_F^2 \leq (1 + \Delta) \sum_{n=1}^N \left(1 + \frac{r_n}{k_n - r_n - 1}\right) (\tau_{r_n}^{(n)})^2,$$

509 *where $\Delta := \max_{n=1}^N r_n / (s_n - r_n - 1)$.*

510 *Proof.* For the two-pass error, in the proof of [Theorem 5.1](#) use the tail bound
 511 from [Lemma B.2](#) to bound the error when $\mathbf{Q}_n \in \mathbb{R}^{I_n \times r_n}$ is chosen by the truncated
 512 QR algorithm [14]. For the one-pass error, use [Lemma A.2](#) to show the error can be
 513 no more than a factor $(1 + \Delta)$ times the error bound for the two-pass approximation
 514 with truncated QR, where $\Delta := \max_{n=1}^N r_n / (s_n - r_n - 1)$.

515 **COROLLARY 5.6** (Quasi-optimality with truncated QR). *For a given target rank*
 516 *\mathbf{r} , choose the sketch size parameters $\mathbf{k} = 2\mathbf{r} + 1$ and $\mathbf{s} = 2\mathbf{r} + 1$. Replace step*
 517 *2 of [Algorithms 4.2](#) and [4.3](#) by a rank- \mathbf{r} truncated QR. The resulting algorithms*
 518 *produce quasi-optimal rank- \mathbf{r} approximations. Specifically, the two-pass approximation*
 519 *([Algorithm 4.2](#)) with truncated QR is quasi-optimal with factor $\sqrt{2N}$ and the one-pass*
 520 *approximation ([Algorithm 4.3](#)) with truncated QR is quasi-optimal with factor $2\sqrt{N}$.*

521 In simultaneous work, [27] also proves the two-pass approximation with truncated
 522 QR is quasi-optimal.

523 **5.2.2. Optimal fixed-rank approximation.** We now bound the error induced
 524 by applying the fixed-rank approximation method [Algorithm 4.4](#) to a given (random)
 525 low-rank approximation.

526 Recall that $\llbracket \mathbf{X} \rrbracket_{\mathbf{r}}$ returns a best rank- \mathbf{r} approximation to \mathbf{X} .

527 **LEMMA 5.7.** *For any tensor \mathbf{X} and random approximation $\hat{\mathbf{X}}$ of the same size,*

$$528 \quad \mathbb{E} \|\mathbf{X} - \llbracket \hat{\mathbf{X}} \rrbracket_{\mathbf{r}}\|_F \leq \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_{\mathbf{r}}\|_F + 2\sqrt{\mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}\|_F^2}.$$

529 *Proof of [Lemma 5.7](#).* Our argument follows the proof of [33, Proposition 6.1]:

$$\begin{aligned} 530 \quad \|\mathbf{X} - \llbracket \hat{\mathbf{X}} \rrbracket_{\mathbf{r}}\|_F &\leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\hat{\mathbf{X}} - \llbracket \hat{\mathbf{X}} \rrbracket_{\mathbf{r}}\|_F \\ 531 &\leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\hat{\mathbf{X}} - \llbracket \mathbf{X} \rrbracket_{\mathbf{r}}\|_F \\ 532 &\leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\hat{\mathbf{X}} - \mathbf{X} + \mathbf{X} - \llbracket \mathbf{X} \rrbracket_{\mathbf{r}}\|_F \\ 533 &\leq 2\|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_{\mathbf{r}}\|_F. \end{aligned}$$

535 The first and the third lines use the triangle inequality, and the second line follows
 536 from the definition of the best rank- r approximation. Take the expectation and use
 537 Lyapunov's inequality to finish the proof. \square

538 **COROLLARY 5.8** (Quasi-optimality with truncated core). *Suppose $\mathbf{k} \geq 2\mathbf{r} + 1$ and*
 539 *$\mathbf{s} \geq 2\mathbf{k}$ and the core approximation \mathcal{A} in [Algorithm 4.4](#) computes an optimal rank- \mathbf{r}*
 540 *approximation to its input \mathbf{W} . That is, $\mathcal{A}(\mathbf{W}, \mathbf{r}) = \llbracket \mathbf{W} \rrbracket_{\mathbf{r}}$. Then the rank- \mathbf{r} approxi-*
 541 *mation algorithm produced by composing the two-pass approximation ([Algorithm 4.2](#)),*
 542 *resp. the one-pass approximation ([Algorithm 4.3](#)), with [Algorithm 4.4](#) is quasi-optimal*
 543 *with factor $\sqrt{2N}$, resp. $2\sqrt{N}$ for one pass.*

544 *Proof.* Use [Lemma 5.7](#) together with [Corollary 5.2](#) or [Corollary 5.4](#). \square

545 Of course, in general it is not possible to compute an optimal rank- \mathbf{r} approximation
 546 for the core. We can still bound the error of the resulting approximation if the
 547 approximation algorithm \mathcal{A} is quasi-optimal using the following lemma.

548 **LEMMA 5.9.** *Suppose the tensor \mathbf{X} and random approximation $\hat{\mathbf{X}}$ satisfy*

$$549 \quad \mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq C(N) \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_{\mathbf{r}}\|_F.$$

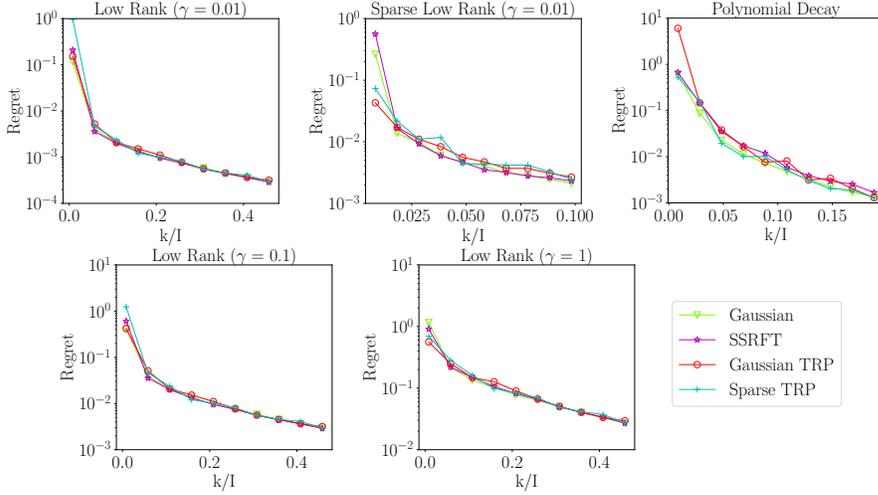


Fig. 1: *Different DRMs perform similarly.* We approximate 3D synthetic tensors (see subsection 6.3) with $I = 600$, using our one-pass algorithm with $r = 5$ and varying k ($s = 2k + 1$), using different DRMs in the Tucker sketch.

550 Further suppose algorithm $\mathcal{A}(\mathcal{W}, \mathbf{r})$ computes a quasi-optimal rank \mathbf{r} approximation to
 551 \mathcal{W} with factor $C'(N)$. Then

$$552 \quad (5.4) \quad \mathbb{E}\|\mathcal{X} - \mathcal{A}(\hat{\mathcal{X}}, \mathbf{r})\|_F \leq (C(N)C'(N) + C(N) + C'(N))\|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F.$$

553 *Proof.* We calculate that

$$\begin{aligned}
 554 \quad \mathbb{E}\|\mathcal{X} - \mathcal{A}(\hat{\mathcal{X}}, \mathbf{r})\|_F &\leq \mathbb{E}\left[\|\mathcal{X} - \hat{\mathcal{X}}\|_F + \|\hat{\mathcal{X}} - \mathcal{A}(\hat{\mathcal{X}}, \mathbf{r})\|_F\right] \\
 555 &\leq C(N)\|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F + C'(N)\mathbb{E}\|\hat{\mathcal{X}} - \llbracket \hat{\mathcal{X}} \rrbracket_{\mathbf{r}}\|_F \\
 556 &\leq C(N)\|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F + C'(N)\mathbb{E}\|\hat{\mathcal{X}} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F \\
 557 &\leq C(N)\|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F + C'(N)\mathbb{E}\left(\|\mathcal{X} - \hat{\mathcal{X}}\|_F + \|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F\right) \\
 558 &\leq (C'(N) + C(N) + C'(N)C(N))\|\mathcal{X} - \llbracket \mathcal{X} \rrbracket_{\mathbf{r}}\|_F. \quad \square
 \end{aligned}$$

560 **COROLLARY 5.10.** Suppose $\mathbf{k} \geq 2\mathbf{r} + 1$ and $\mathbf{s} \geq 2\mathbf{k}$ and the core approximation \mathcal{A}
 561 in Algorithm 4.4 is quasi-optimal with factor \sqrt{N} (such as the ST-HOSVD). Then the
 562 rank- \mathbf{r} approximation algorithm produced by composing the two-pass approximation
 563 (Algorithm 4.2), resp. the one-pass approximation (Algorithm 4.3) with Algorithm 4.4
 564 is quasi-optimal with factor $(1 + \sqrt{2})\sqrt{N} + \sqrt{2}N$, resp. $2N + 3\sqrt{N}$ for one pass.

565 *Proof.* Use Lemma 5.7 together with Corollary 5.2 or Corollary 5.4. \square

566 **6. Numerical Experiments.** In this section, we study the performance of our
 567 streaming Tucker approximation methods. We compare the performance using various
 568 different DRMs, including the TRP. We also compare our method with the algorithm
 569 proposed by [26] to show that, for the same storage budget, our method produces
 570 better approximations. Our two-pass algorithm outperforms the one-pass version, as
 571 expected. (Contrast this to [26], where the multi-pass method performs less well than
 572 the one-pass version.)

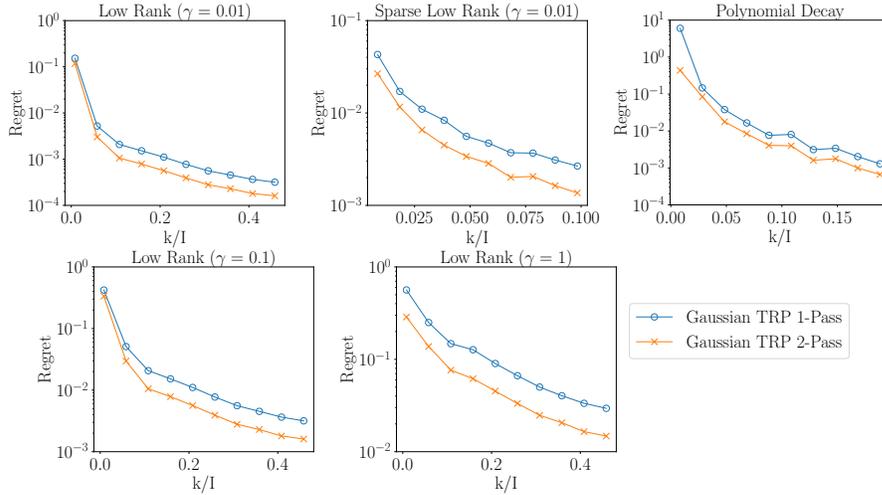


Fig. 2: *Two-pass improves on one-pass.* We approximate 3D synthetic tensors (see subsection 6.3) with $I = 600$, using our one-pass and two-pass algorithms with $r = 5$ and varying k ($s = 2k + 1$), using the Gaussian TRP in the Tucker sketch.

573 **6.1. Error metrics.** We measure the quality of an approximation $\hat{\mathcal{X}}$ to the
 574 original tensor \mathcal{X} using two different metrics. One is the relative error:

575 relative error:
$$\|\mathcal{X} - \hat{\mathcal{X}}\|_F / \|\mathcal{X}\|_F.$$

576 However, many tensors are not close to low rank, in which case every low-rank
 577 approximation will incur high relative error.

578 Our methods cannot solve the problem that many tensors are not low rank;
 579 instead, our goal is just to propose a faster, cheaper, more memory-efficient method
 580 to compute a low-rank approximation to the tensor that is *almost* as good as one
 581 computed using a more expensive method like the HOOI, HOSVD, or ST-HOSVD.
 582 To facilitate comparisons among approximation algorithms, we define another metric
 583 that we call *regret*. We found that the HOOI performs marginally better than the
 584 ST-HOSVD approximation on the examples featured in this section. To simplify plots
 585 and interpretations, we treat the HOOI as the gold standard, and we define the *regret*
 586 of an approximation relative to the HOOI as

587
$$\left(\|\mathcal{X} - \hat{\mathcal{X}}\|_F - \|\mathcal{X} - \mathcal{X}_{\text{HOOI}}\|_F \right) / \|\mathcal{X}\|_F.$$

588 The regret measures the increase in error incurred by using the approximation $\hat{\mathcal{X}}$
 589 rather than $\mathcal{X}_{\text{HOOI}}$. The regret of HOOI is 0. An approximation with a regret of .01
 590 is only 1% worse than HOOI, relative to the norm of the target tensor \mathcal{X} .

591 **6.2. Computational platform.** We ran all experiments on a server with 128
 592 Intel[®] Xeon[®] E7-4850 v4 2.10GHz CPU cores and 1056GB memory. All experiments
 593 are implemented in python. We use the default implementations available in the
 594 python package *tensorly* [22] for tensor algorithms such as the HOOI and ST-HOSVD.
 595 Code for the one and two pass approximation algorithms is available on Github
 596 at <https://github.com/udellgroup/tensorsketch>, as is the code that generates the
 597 experiments in this paper.

598 **6.3. Synthetic experiments.** All synthetic experiments use an input tensor
 599 with equal side lengths I . We consider three different data generation schemes:

- 600 • *Low-rank noise.* Generate a core tensor $\mathbf{C} \in \mathbb{R}^{r^N}$ with entries drawn i.i.d. from
 601 the uniform distribution $U(0, 1)$. Generate N random matrices $\mathbf{B}_1, \dots, \mathbf{B}_N \in$
 602 $\mathbb{R}^{r \times I}$ with i.i.d. $\mathcal{N}(0, 1)$ entries, and let $\mathbf{A}_1, \dots, \mathbf{A}_N \in \mathbb{R}^{r \times I}$ be orthonormal
 603 bases for their respective column spaces. Define $\mathbf{X}^\natural = \mathbf{C} \times_1 \mathbf{A}_1 \cdots \times_N \mathbf{A}_N$
 604 and the noise parameter $\gamma > 0$. Generate an input tensor as $\mathbf{X} = \mathbf{X}^\natural +$
 605 $(\gamma \|\mathbf{X}^\natural\|_F / I^{N/2}) \epsilon$ where the noise ϵ has i.i.d. $\mathcal{N}(0, 1)$ entries.
 606 • *Sparse low-rank noise.* We construct the input tensor \mathbf{X} as above (low-rank
 607 noise), but with sparse factor matrices \mathbf{A}_n : If δ_n is the sparsity (proportion
 608 of nonzero elements) of \mathbf{A}_n , then the sparsity of the true signal \mathbf{X}^\natural scales as
 609 $r^N \prod_{n=1}^N \delta_n$. We use $\delta_n = 0.2$ unless otherwise specified.
 610 • *Polynomial decay.* We construct the input tensor \mathbf{X} as

611
$$\mathbf{X} = \text{superdiag}(1, \dots, 1, 2^{-t}, 3^{-t}, \dots, (I - r)^{-t}).$$

612 The first r entries are 1. Recall **superdiag** converts a vector to N dimensional
 613 superdiagonal tensor. Our experiments use $t = 1$.

614 Our goal in including the polynomial and sparse setups is to demonstrate that the
 615 method performs robustly and reliably even when the distribution of the data is far
 616 from ideal for the method. In the polynomial decay setup, the original tensor is not
 617 particularly low rank, so even a rather expensive and accurate method (the HOOI)
 618 cannot achieve low error; yet **Figures 1 to 3** tell us that the penalty from using our
 619 cheaper methods is essentially the same regardless of the data distribution.

620 **6.3.1. Different dimension reduction maps perform similarly.** We first
 621 investigate the performance of our one-pass fixed-rank algorithm as the sketch size
 622 (hence, the required storage) varies, for several types of dimension reductions maps.
 623 We generate synthetic data as described above with $\mathbf{r} = (5, 5, 5)$, $I = 600$. **Figure 1**
 624 shows the error of the rank- r approximation as a function of the compression factor
 625 k/I . (Results for other input tensors appear in the supplement.) In general, the
 626 performance for different maps are similar, although our theory only guarantees results
 627 for the Gaussian map. We see that for all input tensors, the performance of our
 628 one-pass algorithm converges to that of HOOI as k increases.

629 **6.3.2. A second pass reduces error.** The second experiment compares our
 630 two-pass and one-pass algorithm. The design is similar to the first experiment. **Figure 2**
 631 shows that the two-pass algorithm typically outperforms the one-pass algorithm,
 632 especially in the high-noise, sparse, or rank-decay case. Both converge at the same
 633 asymptotic rate. (Results for other input tensors are available in the supplement.)

634 **6.3.3. Improvement on state-of-the-art.** The third experiment compares the
 635 performance of our two-pass and one-pass algorithms and Tucker TensorSketch (T-
 636 TS), as described in [26], the only extant one-pass algorithm. For a fair comparison,
 637 we allocate the same storage budget to each algorithm and compare the relative error
 638 of the resulting fixed-rank approximations. We approximate synthetic 3D tensors
 639 with equal side lengths $I_1 = I_2 = I_3 = I = 300$ and of equal multilinear rank
 640 $\mathbf{r} = (r, r, r)$ with $r = 10$. We use the suggested parameter settings for each algorithm:
 641 $\mathbf{k} = 2\mathbf{r} + 1$ and $\mathbf{s} = 2\mathbf{k} + 1$ for our methods; $K = 10$ for T.-TS. Our one-pass algorithm
 642 (with the Gaussian TRP) uses $((4k + 3)^N + (2r + 1)IN)$ storage, whereas T.-TS uses
 643 $(Kr^{2N} + Kr^{2N-2})$ storage (see supplement).

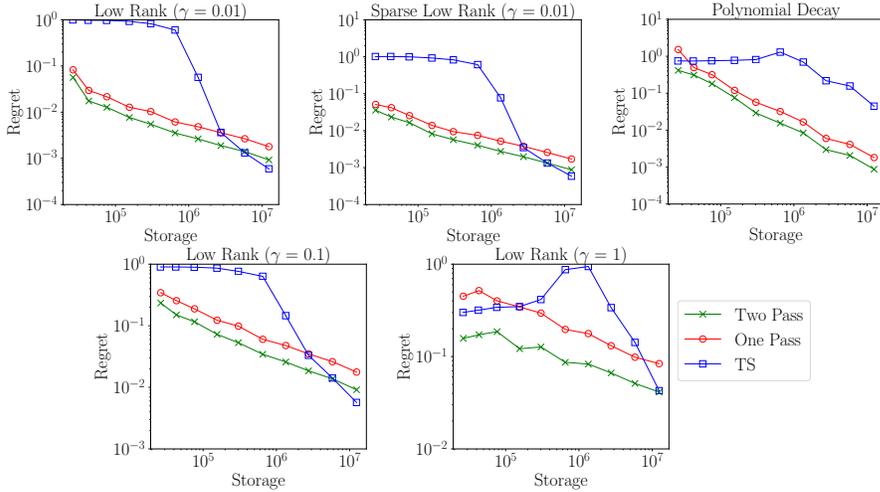


Fig. 3: *Approximations improve with more memory: synthetic data.* We approximate 3D synthetic tensors (see subsection 6.3) with $I = 300$, using T.-TS and our one-pass and two-pass algorithms with the Gaussian TRP to produce approximations with equal ranks $r = 10$. Notice every marker on the plot corresponds to a $2700\times$ compression!

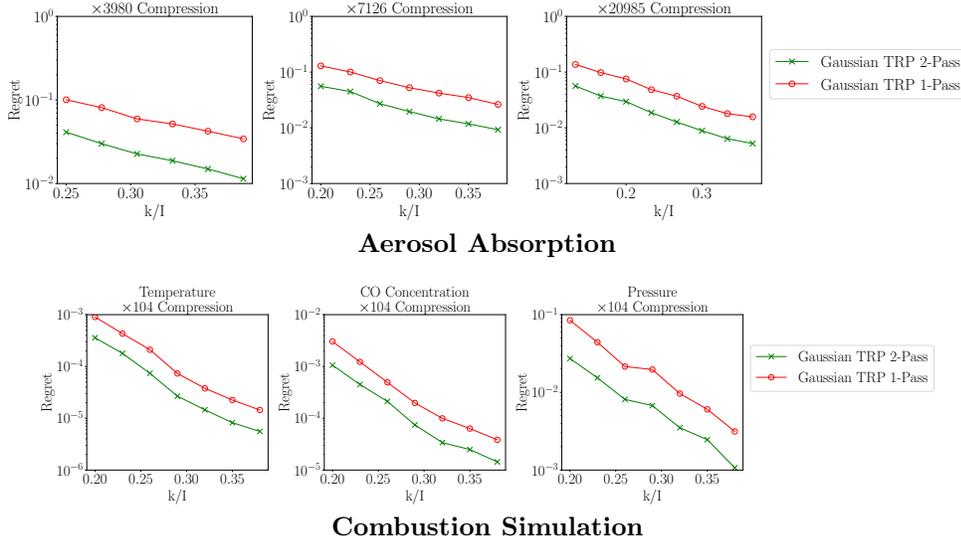
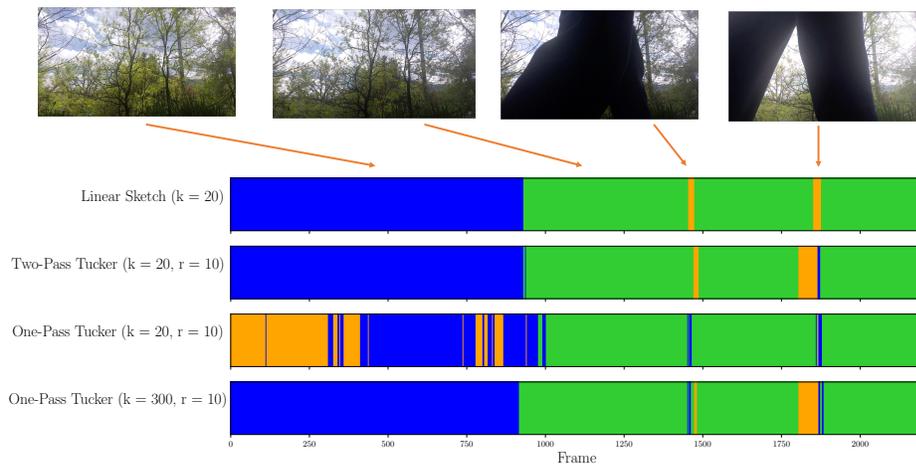


Fig. 4: *Approximations improve with more memory: real data.* We approximate aerosol absorption and combustion data using our one-pass and two-pass algorithms with the Gaussian TRP. We compare three target ranks ($r/I = 0.125, 0.1, 0.067$) for the former, and use the same target rank ($r/I = 0.1$) for each measured quantity in the combustion dataset. Notice $r/I = 0.1$ gives a hundred-fold compression. For reference, on the aerosol data, the HOOI gives an approximation with relative errors .23, .26, and .33 for each of the three ranks, respectively; on the combustion data, the relative error of HOOI is .0063, .032, and .28 for temperature, CO, and pressure, respectively.



Video Scene Classification

Fig. 5: *Video Scene Classification* ($2200 \times 1080 \times 1980$): We classify frames from the video data from [26] (collected as a third order tensor with size $2200 \times 1080 \times 1980$) using K -means with $K=3$ on vectors computed using four different methods. $s = 2k + 1$ throughout. 1) The linear sketch along the time dimension (Row 1). 2-3) the Tucker factor along the time dimension, computed via our two-pass (Row 2) and one-pass (Row 3) algorithms. 4) The Tucker factor along the time dimension, computed via our one-pass (Row 4) algorithm

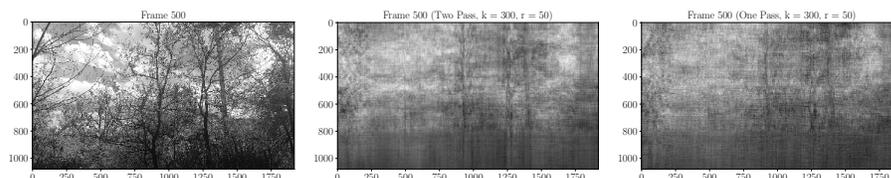


Fig. 6: *Visualizing Video Recovery*: Original frame (left); approximation by two-pass sketch (middle); approximation by one-pass sketch (right).

644 **Figure 3** shows that our algorithms generally perform as well as T.-TS, and
 645 dramatically outperforms for small storage budgets. One nice property of our method
 646 is that the regret consistently decreases with increasing storage. In contrast, the tensor
 647 sketch method behaves unpredictably as storage increases: there are wide plateaus
 648 where increasing storage hardly helps at all, and occasionally, increasing storage hurts
 649 performance. The performance of T.-TS is comparable with that of the algorithms
 650 presented in this paper only when the storage budget is large.

651 *Remark 6.1.* The paper [26] proposes a multi-pass method, Tucker Tensor-Times-
 652 Matrix-TensorSketch (TTMTS) that is dominated by the one-pass method Tucker
 653 TensorSketch (TS) in all numerical experiments; hence we compare only with T.-TS.

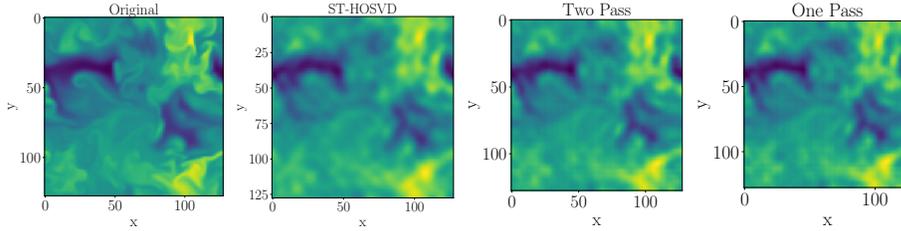


Fig. 7: *Visualizing Combustion Simulation:* All four figures show a slice of the temperature data along the first dimension. The approximation uses $\mathbf{r} = (281, 25, 25)$, $\mathbf{k} = (562, 50, 50)$, $\mathbf{s} = (1125, 101, 101)$, with the Gaussian TRP in the Tucker sketch.

654 **6.4. Applications.** We also apply our method to datasets drawn from three
 655 application domains: climate, combustion, and video.

- 656 • *Climate data.* We consider global climate simulation datasets from the Com-
 657 munity Earth System Model (CESM) Community Atmosphere Model (CAM)
 658 5.0 [17, 18]. The dataset on aerosol absorption has four dimensions: times,
 659 altitudes, longitudes, and latitudes ($240 \times 30 \times 192 \times 288$). The data on net
 660 radiative flux at surface and dust aerosol burden have three dimensions: times,
 661 longitudes, and latitudes ($1200 \times 192 \times 288$). Each of these quantities has a
 662 strong impact on the absorption of solar radiation and on cloud formation.
- 663 • *Combustion data.* We consider combustion simulation data from [23]. The
 664 data consists of three measured quantities (pressure, CO concentration, and
 665 temperature) each observed on a $1408 \times 128 \times 128$ spatial grid.
- 666 • *Video data.* Consider the 3D tensor from [26]: each slice of the tensor is a
 667 video frame. A low frame rate camera is mounted in a fixed position as people
 668 walk by to form the video, which consists of 2493 frames, each of size 1080 by
 669 1980. Stored as a `numpy.array`, the video data is 41.4 GB in total.

670 **6.4.1. Data compression.** We show that our proposed algorithms are able to
 671 successfully compress climate and combustion data even when the full data does not fit
 672 in memory. Since the multilinear rank of the original tensor is unknown, we perform
 673 experiments for three different target ranks. In this experiment, we hope to understand
 674 the effect of different choices of storage budget k to achieve the same compression
 675 ratio. We define the compression ratio as the ratio in size between the original input
 676 tensor and the output Tucker factors, i.e. $\frac{\prod_{i=1}^N J_i}{\sum_{i=1}^N r_i I_i + \prod_{i=1}^N r_i}$. As in our experiments
 677 on simulated data, Figure 4 shows that the two-pass algorithm outperforms the
 678 one-pass algorithm as expected. However, as the storage budget k increases, both
 679 methods converge to the performance of HOOI. The rate of convergence is faster for
 680 smaller target ranks. Performance of our algorithms on the combustion simulation
 681 is qualitatively similar, but converges faster to the performance of HOOI. Figure 7
 682 visualizes the recovery of the temperature data in combustion simulation for a slice
 683 along the first dimension. We could observe that the recovery for both two-pass and
 684 one-pass algorithm approximate the recovery from HOOI. Similar results on other
 685 datasets appear in the supplement.)

686 **6.4.2. Video scene classification.** We show how to use our single-pass method
 687 to classify scenes in the video data described above. The goal is to identify frames
 688 in which people appear. We remove the first 100 frames and last 193 frames where

689 the camera setup happened, as in [26]. We stream over the tensor and sketch it using
 690 parameters $k = 300, s = 601$. Finally, we compute a fixed-rank approximation with
 691 $\mathbf{r} = (10, 10, 10)$ and $(20, 20, 20)$. We apply K-means clustering to the resulting 10 or
 692 20 dimensional vectors corresponding to each of the remaining 2200 frames.

693 We experimented with clustering vectors found in three ways: from the unfolding
 694 along time direction after two-pass or one-pass Tucker approximation, or directly from
 695 the factor sketch along the time which we call it linear sketch. When matching the
 696 video frames with the classification result, we can see that the background light is
 697 relatively dark at the beginning, thus classified into `Class 0`. After a change in the
 698 background light, most other frames of the video are classified into `Class 1`. When
 699 a person passes by the camera, the frames are classified into `Class 2`. Right after
 700 the person passed by, the frames are classified into `Class 0`, the brighter background
 701 scene, due to the light adjustment.

702 Our classification results (using the linear sketch or approximation) are similar
 703 to those in [26] while using only 1/500 as much storage; the one-pass approximation
 704 requires more storage (but still less than [26]) to achieve similar performance. In
 705 particular, using the sketch itself, rather than the Tucker approximation, to summarize
 706 the data enables very efficient video scene classification. Interestingly, classification
 707 works well even though the video is not very low rank along the spatial dimensions.
 708 Figure 6 shows that the scene is poorly approximated even with $\mathbf{s} = 601, 601, 601$,
 709 $\mathbf{k} = (300, 300, 300)$, and $\mathbf{r} = (50, 50, 50)$.

710 **Acknowledgments.** MU, YS, and YG were supported in part by DARPA Award
 711 FA8750-17-2-0101 and NSF CCF-1740822. JAT gratefully acknowledges support from
 712 ONR Awards N00014-11-10025, N00014-17-12146, and N00014-18-12363. The authors
 713 wish to thank Osman Asif Malik and Stephen Becker for their help in understanding
 714 and implementing Tucker TensorSketch, and Tamara Kolda for insightful comments
 715 on an early draft.

716 References.

- 717 [1] D. ACHLIOPTAS, *Database-friendly random projections: Johnson-Lindenstrauss*
 718 *with binary coins*, Journal of computer and System Sciences, 66 (2003), pp. 671–
 719 687.
- 720 [2] S. ARORA AND B. BARAK, *Computational complexity: a modern approach*,
 721 Cambridge University Press, 2009.
- 722 [3] W. AUSTIN, G. BALLARD, AND T. G. KOLDA, *Parallel tensor compression for*
 723 *large-scale scientific data*, in Parallel and Distributed Processing Symposium, 2016
 724 IEEE International, IEEE, 2016, pp. 912–922.
- 725 [4] R. BALLESTER-RIPOLL, P. LINDSTROM, AND R. PAJAROLA, *Tthresh: Tensor*
 726 *compression for multidimensional visual data*, IEEE transactions on visualization
 727 and computer graphics, (2019).
- 728 [5] M. BASKARAN, B. MEISTER, N. VASILACHE, AND R. LETHIN, *Efficient and scal-*
 729 *able computations with sparse tensors*, in High Performance Extreme Computing
 730 (HPEC), 2012 IEEE Conference on, IEEE, 2012, pp. 1–6.
- 731 [6] C. BATTAGLINO, G. BALLARD, AND T. G. KOLDA, *A practical randomized cp*
 732 *tensor decomposition*, SIAM Journal on Matrix Analysis and Applications, 39
 733 (2018), pp. 876–901.
- 734 [7] C. BATTAGLINO, G. BALLARD, AND T. G. KOLDA, *Faster parallel tucker tensor*
 735 *decomposition using randomization*, (2019).
- 736 [8] P. BREIDING AND N. VANNIEUWENHOVEN, *A riemannian trust region method for*
 737 *the canonical tensor rank approximation problem*, SIAM Journal on Optimization,

- 738 28 (2018), pp. 2435–2465.
- 739 [9] A. CICHOCKI, *Tensor decompositions: a new concept in brain data analysis?*,
740 arXiv preprint arXiv:1305.0395, (2013).
- 741 [10] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *A multilinear singular*
742 *value decomposition*, SIAM journal on Matrix Analysis and Applications, 21
743 (2000), pp. 1253–1278.
- 744 [11] L. DE LATHAUWER, B. DE MOOR, AND J. VANDEWALLE, *On the best rank-1*
745 *and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors*, SIAM Journal
746 on Matrix Analysis and Applications, 21 (2000), pp. 1324–1342, <https://doi.org/10.1137/S0895479898346995>,
747 <https://doi.org/10.1137/S0895479898346995>,
748 <https://arxiv.org/abs/https://doi.org/10.1137/S0895479898346995>.
- 749 [12] V. DE SILVA AND L.-H. LIM, *Tensor rank and the ill-posedness of the best low-*
750 *rank approximation problem*, SIAM Journal on Matrix Analysis and Applications,
751 30 (2008), pp. 1084–1127.
- 752 [13] L. GRASEDYCK, *Hierarchical singular value decomposition of tensors*, SIAM
753 Journal on Matrix Analysis and Applications, 31 (2010), pp. 2029–2054.
- 754 [14] M. GU AND S. C. EISENSTAT, *Efficient algorithms for computing a strong rank-*
755 *revealing qr factorization*, SIAM Journal on Scientific Computing, 17 (1996),
756 pp. 848–869.
- 757 [15] W. HACKBUSCH, *Tensor spaces and numerical tensor calculus*, vol. 42, Springer
758 Science & Business Media, 2012.
- 759 [16] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with*
760 *randomness: Probabilistic algorithms for constructing approximate matrix decom-*
761 *positions*, SIAM review, 53 (2011), pp. 217–288.
- 762 [17] J. W. HURRELL, M. M. HOLLAND, P. R. GENT, S. GHAN, J. E. KAY,
763 P. J. KUSHNER, J.-F. LAMARQUE, W. G. LARGE, D. LAWRENCE, K. LINDSAY,
764 ET AL., *The community earth system model: a framework for collaborative research*,
765 Bulletin of the American Meteorological Society, 94 (2013), pp. 1339–1360.
- 766 [18] J. KAY, C. DESER, A. PHILLIPS, A. MAI, C. HANNAY, G. STRAND, J. AR-
767 BLASTER, S. BATES, G. DANABASOGLU, J. EDWARDS, ET AL., *The community*
768 *earth system model (cesm) large ensemble project: A community resource for*
769 *studying climate change in the presence of internal climate variability*, Bulletin of
770 the American Meteorological Society, 96 (2015), pp. 1333–1349.
- 771 [19] O. KAYA AND B. UÇAR, *High performance parallel algorithms for the tucker*
772 *decomposition of sparse tensors*, in Parallel Processing (ICPP), 2016 45th Interna-
773 tional Conference on, IEEE, 2016, pp. 103–112.
- 774 [20] T. G. KOLDA AND B. W. BADER, *Tensor decompositions and applications*, SIAM
775 review, 51 (2009), pp. 455–500.
- 776 [21] T. G. KOLDA AND J. SUN, *Scalable tensor decompositions for multi-aspect data*
777 *mining*, in 2008 Eighth IEEE International Conference on Data Mining, IEEE,
778 2008, pp. 363–372.
- 779 [22] J. KOSSAIFI, Y. PANAGAKIS, A. ANANDKUMAR, AND M. PANTIC, *Tensorly:*
780 *Tensor learning in python*, The Journal of Machine Learning Research, 20 (2019),
781 pp. 925–930.
- 782 [23] S. LAPOINTE, B. SAVARD, AND G. BLANQUART, *Differential diffusion effects,*
783 *distributed burning, and local extinctions in high karlovitz premixed flames*, Com-
784 bustion and flame, 162 (2015), pp. 3341–3355.
- 785 [24] J. LI, C. BATTAGLINO, I. PERROS, J. SUN, AND R. VUDUC, *An input-adaptive*
786 *and in-place approach to dense tensor-times-matrix multiply*, in High Performance
787 Computing, Networking, Storage and Analysis, 2015 SC-International Conference

- 788 for, IEEE, 2015, pp. 1–12.
- 789 [25] P. LI, T. J. HASTIE, AND K. W. CHURCH, *Very sparse random projections*, in
790 Proceedings of the 12th ACM SIGKDD international conference on Knowledge
791 discovery and data mining, ACM, 2006, pp. 287–296.
- 792 [26] O. A. MALIK AND S. BECKER, *Low-rank tucker decomposition of large tensors*
793 *using tensorsketch*, in Advances in Neural Information Processing Systems, 2018,
794 pp. 10116–10126.
- 795 [27] R. MINSTER, A. K. SAIBABA, AND M. E. KILMER, *Randomized algo-*
796 *rithms for low-rank tensor decompositions in the tucker format*, arXiv preprint
797 arXiv:1905.07311, (2019).
- 798 [28] S. MUTHUKRISHNAN ET AL., *Data streams: Algorithms and applications*, Foun-
799 dations and Trends® in Theoretical Computer Science, 1 (2005), pp. 117–236.
- 800 [29] S. OYMAK AND J. A. TROPP, *Universality laws for randomized dimension*
801 *reduction, with applications*, Information and Inference: A Journal of the IMA,
802 (2015).
- 803 [30] M. RUDELSON, *Row products of random matrices*, Advances in Mathematics, 231
804 (2012), pp. 3199–3231.
- 805 [31] J. SUN, D. TAO, S. PAPANIMITRIOU, P. S. YU, AND C. FALOUTSOS, *Incremental*
806 *tensor analysis: Theory and applications*, ACM Transactions on Knowledge
807 Discovery from Data (TKDD), 2 (2008), p. 11.
- 808 [32] Y. SUN, Y. GUO, J. A. TROPP, AND M. UDELL, *Tensor random projec-*
809 *tion for low memory dimension reduction*, in NeurIPS Workshop on Rela-
810 tional Representation Learning, 2018, [https://r2learning.github.io/assets/papers/](https://r2learning.github.io/assets/papers/CameraReadySubmission%2041.pdf)
811 [CameraReadySubmission%2041.pdf](https://r2learning.github.io/assets/papers/CameraReadySubmission%2041.pdf).
- 812 [33] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching*
813 *algorithms for low-rank matrix approximation*, SIAM Journal on Matrix Analysis
814 and Applications, 38 (2017), pp. 1454–1485.
- 815 [34] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *More practical*
816 *sketching algorithms for low-rank matrix approximation*, Tech. Report 2018-01,
817 California Institute of Technology, Pasadena, California, 2018.
- 818 [35] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Streaming low-rank*
819 *matrix approximation with an application to scientific simulation*, SIAM Journal
820 on Scientific Computing (SISC), (2019), <https://arxiv.org/abs/1902.08651>.
- 821 [36] C. E. TSOURAKAKIS, *Mach: Fast randomized tensor decompositions*, in Proceed-
822 ings of the 2010 SIAM International Conference on Data Mining, SIAM, 2010,
823 pp. 689–700.
- 824 [37] L. R. TUCKER, *Some mathematical notes on three-mode factor analysis*, Psy-
825 chometrika, 31 (1966), pp. 279–311.
- 826 [38] N. VANNIEUWENHOVEN, R. VANDEBRIL, AND K. MEERBERGEN, *A new trunca-*
827 *tion strategy for the higher-order singular value decomposition*, SIAM Journal on
828 Scientific Computing, 34 (2012), pp. A1027–A1052.
- 829 [39] M. A. O. VASILESCU AND D. TERZOPOULOS, *Multilinear analysis of image*
830 *ensembles: Tensorfaces*, in European Conference on Computer Vision, Springer,
831 2002, pp. 447–460.
- 832 [40] Y. WANG, H.-Y. TUNG, A. J. SMOLA, AND A. ANANDKUMAR, *Fast and*
833 *guaranteed tensor decomposition via sketching*, in Advances in Neural Information
834 Processing Systems, 2015, pp. 991–999.
- 835 [41] D. P. WOODRUFF ET AL., *Sketching as a tool for numerical linear algebra*,
836 Foundations and Trends® in Theoretical Computer Science, 10 (2014), pp. 1–157.
- 837 [42] F. WOOLFE, E. LIBERTY, V. ROKHLIN, AND M. TYGERT, *A fast randomized*

- 838 *algorithm for the approximation of matrices*, Applied and Computational Harmonic
839 Analysis, 25 (2008), pp. 335–366.
- 840 [43] G. ZHOU, A. CICHOCKI, AND S. XIE, *Decomposition of big tensors with low*
841 *multilinear rank*, arXiv preprint arXiv:1412.1885, (2014).

842 **Appendix A. Probabilistic Analysis of Core Sketch Error.** This section
 843 contains the most technical part of our proof. We provide a probabilistic error bound
 844 for the difference between the two-pass core approximation \mathcal{W}_2 from [Algorithm 4.2](#)
 845 and the one-pass core approximation \mathcal{W}_1 from [Algorithm 4.3](#).

846 Introduce for each $n \in [N]$ the orthonormal matrix \mathbf{Q}_n^\perp that forms a basis for the
 847 subspace orthogonal to \mathbf{Q}_n , so that $\mathbf{Q}_n^\perp (\mathbf{Q}_n^\perp)^\top = \mathbf{I} - \mathbf{Q}_n \mathbf{Q}_n^\top$. Next, define

$$848 \quad (\text{A.1}) \quad \mathbf{\Phi}_n^Q = \mathbf{\Phi}_n^\top \mathbf{Q}_n, \quad \mathbf{\Phi}_n^{Q^\perp} = \mathbf{\Phi}_n^\top \mathbf{Q}_n^\perp.$$

849 Recall that the DRMs $\mathbf{\Phi}_n$ are i.i.d. Gaussian. Thus, conditional on \mathbf{Q}_n , the random
 850 matrices $\mathbf{\Phi}_n^Q$ and $\mathbf{\Phi}_n^{Q^\perp}$ are statistically independent.

851 **A.1. Decomposition of Core Approximation Error.** In this section, we
 852 characterize the difference between the one- and two-pass core approximations $\mathcal{W}_1 -$
 853 $\mathcal{W}_2 = \mathcal{W}_1 - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top$.

854 **LEMMA A.1.** *Suppose that $\mathbf{\Phi}_n$ has full column rank for each $n \in [N]$. We define*
 855 $\mathbb{1}_{a=b} = 1$ if $a = b$ and 0 otherwise. Then

$$856 \quad \mathcal{W}_1 - \mathcal{W}_2 = \mathcal{W}_1 - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top = \sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum_{j=1}^N i_j \geq 1} \mathbf{y}_{i_1 \dots i_N},$$

857 where

$$858 \quad (\text{A.2}) \quad \mathbf{y}_{i_1 \dots i_N} = \mathcal{X} \times_1 \left(\mathbb{1}_{i_1=0} \mathbf{Q}_1^\top + \mathbb{1}_{i_1=1} (\mathbf{\Phi}_1^{Q_1})^\dagger \mathbf{\Phi}_1^{Q_1^\perp} (\mathbf{Q}_1^\perp)^\top \right) \\ \times_2 \cdots \times_N \left(\mathbb{1}_{i_N=0} \mathbf{Q}_N^\top + \mathbb{1}_{i_N=1} (\mathbf{\Phi}_N^{Q_N})^\dagger \mathbf{\Phi}_N^{Q_N^\perp} (\mathbf{Q}_N^\perp)^\top \right).$$

859 *Proof.* Let \mathcal{H} be the core sketch from [Algorithm 4.1](#). Write \mathcal{W}_1 as

$$\begin{aligned} 860 \quad \mathcal{W}_1 &= \mathcal{H} \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \\ &= (\mathcal{X} - \hat{\mathcal{X}}_2) \times_1 \mathbf{\Phi}_1^\top \times_2 \cdots \times_N \mathbf{\Phi}_N^\top \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \\ &\quad + \hat{\mathcal{X}}_2 \times_1 \mathbf{\Phi}_1^\top \times_2 \cdots \times_N \mathbf{\Phi}_N^\top \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger. \end{aligned}$$

861 Using the fact that $(\mathbf{\Phi}_n^\top \mathbf{Q}_n)^\dagger (\mathbf{\Phi}_n^\top \mathbf{Q}_n) = \mathbf{I}$, we can simplify the second term as

$$\begin{aligned} 862 \quad &\hat{\mathcal{X}}_2 \times_1 \mathbf{\Phi}_1^\top \times_2 \cdots \times_N \mathbf{\Phi}_N^\top \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \\ &= \mathcal{X} \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \mathbf{\Phi}_1^\top \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \mathbf{\Phi}_N^\top \mathbf{Q}_N \mathbf{Q}_N^\top \\ &= \mathcal{X} \times_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top, \end{aligned}$$

863 which is exactly the two-pass core approximation \mathcal{W}_2 . Therefore

$$864 \quad \mathcal{W}_1 - \mathcal{W}_2 = (\mathcal{X} - \hat{\mathcal{X}}_2) \times_1 \mathbf{\Phi}_1^\top \times_2 \cdots \times_N \mathbf{\Phi}_N^\top \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger.$$

865 We continue to simplify this difference:

$$\begin{aligned} 866 \quad (\text{A.3}) \quad &(\mathcal{X} - \hat{\mathcal{X}}) \times_1 \mathbf{\Phi}_1^\top \times_2 \cdots \times_N \mathbf{\Phi}_N^\top \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \\ &= (\mathcal{X} - \tilde{\mathcal{X}}) \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \mathbf{\Phi}_1^\top \times_2 \cdots \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \mathbf{\Phi}_N^\top \\ &= (\mathcal{X} - \tilde{\mathcal{X}}) \times_1 (\mathbf{\Phi}_1^\top \mathbf{Q}_1)^\dagger \mathbf{\Phi}_1^\top (\mathbf{Q}_1 \mathbf{Q}_1^\top + \mathbf{Q}_1^\perp (\mathbf{Q}_1^\perp)^\top) \cdots \\ &\quad \times_N (\mathbf{\Phi}_N^\top \mathbf{Q}_N)^\dagger \mathbf{\Phi}_N^\top (\mathbf{Q}_N \mathbf{Q}_N^\top + \mathbf{Q}_N^\perp (\mathbf{Q}_N^\perp)^\top) \\ &= (\mathcal{X} - \tilde{\mathcal{X}}) \times_1 (\mathbf{Q}_1^\top + (\mathbf{\Phi}_1^Q)^\dagger \mathbf{\Phi}_1^{Q^\perp} (\mathbf{Q}_1^\perp)^\top) \times_2 \cdots \\ &\quad \times_N (\mathbf{Q}_N^\top + (\mathbf{\Phi}_N^{Q_N})^\dagger \mathbf{\Phi}_N^{Q_N^\perp} (\mathbf{Q}_N^\perp)^\top). \end{aligned}$$

867 Many terms in this sum are zero. We use the following two facts:

- 868 1. $(\mathbf{X} - \hat{\mathbf{X}}) \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top = 0$.
 869 2. For each $n \in [N]$, $\tilde{\mathbf{X}} \times_n (\Phi_n^{Q_n})^\dagger \Phi_n^{Q_n^\perp} (\mathbf{Q}_n^\perp)^\top = 0$.

870 Here, 0 denotes a tensor with all zero elements. These facts can be obtained from the
 871 exchange rule of the mode product and the orthogonality between \mathbf{Q}_n^\perp and \mathbf{Q}_n . Using
 872 these two facts, we find that only the terms $\mathbf{y}_{i_1 \dots i_N}$ (defined in (A.2)) remain in the
 873 expression. Therefore, to complete the proof, we write (A.3) as

$$874 \quad \sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum_{n=1}^N i_n \neq 0} \mathbf{y}_{i_1 \dots i_N}.$$

875 **A.2. Probabilistic Core Error Bound.** In this section, we derive a probabilis-
 876 tic error bound based on the core error decomposition from Lemma A.1.

877 LEMMA A.2. *Sketch the tensor \mathbf{X} using a Tucker sketch with parameters \mathbf{k} and*
 878 *$\mathbf{s} > 2\mathbf{k}$ with i.i.d. Gaussian $\mathcal{N}(0, 1)$ DRMs. Define $\Delta = \max_{n=1}^N \frac{k_n}{s_n - k_n - 1}$. Let $\hat{\mathbf{X}}_2$ be*
 879 *the output from the two-pass low-rank approximation method (Algorithm 4.2). Then*

$$880 \quad (\text{A.4}) \quad \mathbb{E} \|\mathbf{W}_1 - \mathbf{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top\|_F^2 \leq \Delta \|\mathbf{X} - \hat{\mathbf{X}}_2\|$$

881 *Proof.* We use the fact that the core DRMs $\{\Omega_n\}_{n \in [N]}$ are independent of the
 882 factor matrix DRMs $\{\Phi_n\}_{n \in [N]}$, and that the randomness in each factor matrix
 883 approximation \mathbf{Q}_n comes solely from Ω_n .

884 For $i \in \{0, 1\}^N$, define $\mathcal{B}_{i_1 \dots i_N} =$

$$885 \quad \mathbf{X} \times_1 (\mathbb{1}_{i_1=0} \mathbf{Q}_1 \mathbf{Q}_1^\top + \mathbb{1}_{i_1=1} \mathbf{Q}_1^\perp (\mathbf{Q}_1^\perp)^\top) \cdots \times_N (\mathbb{1}_{i_N=0} \mathbf{Q}_N \mathbf{Q}_N^\top + \mathbb{1}_{i_N=1} \mathbf{Q}_N^\perp (\mathbf{Q}_N^\perp)^\top).$$

886 Lemma A.1 decomposes the core error as the sum of $\mathbf{y}_{i_1 \dots i_N}$ where $\sum_{n=1}^N i_n \geq 1$.
 887 Applying Lemma B.1 and using the orthogonal invariance of the Frobenius norm, we
 888 observe

$$889 \quad \mathbb{E} [\|\mathbf{y}_{i_1 \dots i_N}\|_F^2 \mid \Omega_1 \cdots \Omega_N] = \left(\prod_{n=1}^N \Delta_n^{i_n} \right) \|\mathcal{B}_{i_1 \dots i_N}\|_F^2 \leq \Delta \|\mathcal{B}_{i_1 \dots i_N}\|_F^2$$

890 when $\sum_{n=1}^N i_n \geq 1$, where $\Delta_n = \frac{k_n}{s_n - k_n - 1} < 1$ and $\Delta = \max_{n=1}^N \Delta_n$.

891 Suppose $\mathbf{q}_1, \mathbf{q}_2 \in \{0, 1\}^N$ are index (binary) vectors of length N . For different
 892 indices \mathbf{q}_1 and \mathbf{q}_2 , there exists some $1 \leq r \leq N$ such that their r th element is different.
 893 Without loss of generality, assume $\mathbf{q}_1(r) = 0$ and $\mathbf{q}_2(r) = 1$ to see

$$894 \quad (\text{A.5}) \quad \langle \mathcal{B}_{\mathbf{q}_1}, \mathcal{B}_{\mathbf{q}_2} \rangle = \langle \dots \mathbf{Q}_r^\top \mathbf{Q}_r^\perp \dots \rangle = 0.$$

895 Similarly we can show that the inner product between $\mathcal{B}_{\mathbf{q}_1}$ and $\mathcal{B}_{\mathbf{q}_2}$ is zero with different
 896 $\mathbf{q}_1, \mathbf{q}_2$. Noticing that $\mathcal{B}_{0, \dots, 0} = \hat{\mathbf{X}}_2$, we have

$$897 \quad \|\mathbf{X} - \hat{\mathbf{X}}_2\|_F^2 = \left\| \sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum_{n=1}^N i_n \geq 1} \mathcal{B}_{i_1 \dots i_N} \right\|_F^2 = \sum_{\substack{(i_1, \dots, i_N) \in \{0,1\}^N, \\ \sum_{n=1}^N i_n \geq 1}} \|\mathcal{B}_{i_1 \dots i_N}\|_F^2.$$

898

899 Put these together and use the Pythagorean theorem to finish the proof:

$$\begin{aligned}
900 & \mathbb{E} [\|\mathbf{W} - \mathbf{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top\|_F^2 \mid \boldsymbol{\Omega}_1, \dots, \boldsymbol{\Omega}_N] \\
901 & = \sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum_{n=1}^N i_n \geq 1} \mathbb{E} [\|\mathbf{y}_{i_1 \dots i_N}\|_F^2 \mid \boldsymbol{\Omega}_1, \dots, \boldsymbol{\Omega}_N] \\
902 & \leq \Delta \left(\sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum_{n=1}^N i_n \geq 1} \|\mathbf{B}_{i_1 \dots i_N}\|_F^2 \right) = \Delta \|\mathbf{X} - \hat{\mathbf{X}}_2\|_F^2. \quad \square \\
903 &
\end{aligned}$$

904 **Appendix B. Random matrix projections.**

905 Proofs for the lemmas in this section can be found in [16, sections 9 and 10].

906 LEMMA B.1. Assume that $t > q$. Suppose $\mathbf{G}_1 \in \mathbb{R}^{t \times q}$ and $\mathbf{G}_2 \in \mathbb{R}^{t \times p}$ have i.i.d.
907 standard normal entries. For any matrix \mathbf{B} with conforming dimensions,

$$908 \quad \mathbb{E} \|\mathbf{G}_1^\dagger \mathbf{G}_2 \mathbf{B}\|_F^2 = \frac{q}{t - q - 1} \|\mathbf{B}\|_F^2.$$

909 LEMMA B.2. Given a fixed $\mathbf{A} \in \mathbb{R}^{m \times n}$ and random $\boldsymbol{\Omega} \in \mathbb{R}^{n \times k}$ with i.i.d. standard
910 normal entries, let $\mathbf{Q}_\rho = \text{TruncatedQR}(\mathbf{A}\boldsymbol{\Omega}, \rho) \in \mathbb{R}^{n \times \rho}$ for $\rho < k - 1$ [14]. Then

$$911 \quad (\text{B.1}) \quad \mathbb{E} \|(\mathbf{I} - \mathbf{Q}_\rho \mathbf{Q}_\rho^\top) \mathbf{A}\|_F^2 \leq \frac{\rho}{k - \rho - 1} \tau_\rho.$$

912 COROLLARY B.1. Under the same conditions as in Lemma B.2, suppose $\mathbf{Q} =$
913 $\text{QR}(\mathbf{A}\boldsymbol{\Omega})$ is an orthogonal matrix spanning the column space of $\mathbf{A}\boldsymbol{\Omega}$. Then

$$914 \quad (\text{B.2}) \quad \mathbb{E} \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^\top) \mathbf{A}\|_F^2 \leq \min_{1 \leq \rho < k-1} \frac{\rho}{k - \rho - 1} \tau_\rho.$$

915 *Proof.* For each $\rho < k - 1$,

$$916 \quad \|(\mathbf{I} - \mathbf{Q}_\rho \mathbf{Q}_\rho^\top) \mathbf{A}\|_F^2 \leq \|(\mathbf{I} - \mathbf{Q}_\rho \mathbf{Q}_\rho^\top) \mathbf{A}\|_F^2 \leq \frac{\rho}{k - \rho - 1} \tau_\rho,$$

917 using (B.1) for the second inequality. Minimize over $\rho < k - 1$ to reach the result. \square