

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 [38],
19 neuroscience [10], scientific simulation [4], sensor networks [29], 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. Indeed, one
27 particular challenge is that these large tensors may not fit in main memory on our
28 computer.

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

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

- 38 **1. Streaming:** Data from the tensor is generated sequentially. At each time
39 stamp, we may observe a low dimensional slice, an individual entry, or an
40 additive update to the tensor (the so-called “turnstile” model [26]). For
41 example, each slice of the tensor may represent a subsequent time step in a
42 simulation, or sensor measurements at a particular time. In the streaming

*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)

43 setting, the complete tensor is not stored; indeed, it may be much larger than
 44 available computing resources.

45 Our algorithm can approximate tensors revealed via streaming updates by
 46 sketching the updates and storing the sketch. Linearity of the sketch guarantees
 47 that sketching commutes with slice, entrywise, or additive updates. Our
 48 method forms an approximation of the tensor only after all the data has
 49 been observed, rather than approximating the tensor-observed-so-far at any
 50 time. This protocol allows for offline data analysis, including many scientific
 51 applications. Conversely, this protocol is not suitable for real-time monitoring.

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

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

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

- 66 • We present a new method to form a linear sketch of an unknown tensor. This
 67 sketch captures both the principal subspaces of the tensor along each mode,
 68 and the action of the tensor that links these subspaces. This tensor sketch
 69 can be formed from any dimension reduction map. Those sketches themselves
 70 are useful in many application even without forming Tucker decomposition
 71 like in video clustering.
- 72 • We develop a practical one-pass algorithm to compute a low rank Tucker
 73 approximation from streaming data. The algorithm sketches the tensor and
 74 then recovers a low rank Tucker approximation from this sketch.
- 75 • We propose a two-pass algorithm that improves on the one-pass method. Both
 76 the one-pass and two-pass methods are appropriate in a limited memory or
 77 distributed data setting.
- 78 • We develop provable probabilistic guarantees on the performance of both
 79 the one-pass and two-pass algorithms when the tensor sketch is composed of
 80 Gaussian dimension reduction maps.
- 81 • We exhibit several random maps that can be used to sketch the tensor.
 82 Compared to the Gaussian map, others are cheaper to store, easier to apply,
 83 and deliver similar performance experimentally in tensor approximation error.
 84 In particular, we demonstrate the effective performance of a row-product
 85 random matrix, which we call the Tensor Random Projection (TRP), which
 86 uses exceedingly low storage.
- 87 • We perform a comprehensive simulation study with synthetic data, and
 88 consider applications to several real datasets, to demonstrate the practical
 89 performance of our method. Our methods reduce approximation error com-
 90 pared to the only existing one-pass Tucker approximation algorithm [25] by
 91 more than an order of magnitude given the same storage budget.
- 92 • We have developed and released an open-source package in python that

93 implements our algorithms.

94 2. Background and Related Work.

95 **2.1. Notation.** Our paper follows the notation of [20]. We denote *scalar*, *vector*,
 96 *matrix*, and *tensor* variables respectively by lowercase letters (x), boldface lowercase
 97 letters (\mathbf{x}), boldface capital letters (\mathbf{X}), and boldface Euler script letters (\mathcal{X}). For two
 98 vectors \mathbf{x} and \mathbf{y} , we write $\mathbf{x} > \mathbf{y}$ if \mathbf{x} is greater than \mathbf{y} elementwise and $\mathbf{x} < \mathbf{y}$ as the
 99 opposite.

100 Define $[N] := \{1, \dots, N\}$. For a matrix $\mathbf{X} \in \mathbb{R}^{m \times n}$, we denote its i^{th} row, j^{th}
 101 column, and $(i, j)^{\text{th}}$ element as $\mathbf{X}(i, \cdot)$, $\mathbf{X}(:, j)$, and $\mathbf{X}(i, j)$, respectively, for $i \in [m]$,
 102 $j \in [n]$. We use $\mathbf{X}^\dagger \in \mathbb{R}^{n \times m}$ to denote the *Moore–Penrose pseudoinverse* of the matrix
 103 $\mathbf{X} \in \mathbb{R}^{m \times n}$. 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;
 104 $\mathbf{X}^\dagger = \mathbf{X}^\top (\mathbf{X} \mathbf{X}^\top)^{-1}$, if $m < n$ and \mathbf{X} has full row rank.

105 **2.1.1. Tail energy.** To state our results, we will need a tensor equivalent for the
 106 decay in the spectrum of a matrix. For each unfolding $\mathbf{X}^{(n)}$, define the ρ th *tail energy*

$$107 \quad (\tau_\rho^{(n)})^2 := \sum_{k>\rho}^{\min(I_n, I_{(-n)})} \sigma_k^2(\mathbf{X}^{(n)}),$$

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

109 **2.1.2. Kronecker and Khatri–Rao product.** For two matrices $\mathbf{A} \in \mathbb{R}^{I \times J}$ and
 110 $\mathbf{B} \in \mathbb{R}^{K \times L}$, we define the *Kronecker product* $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{IK \times JL}$ as

$$111 \quad (2.1) \quad \mathbf{A} \otimes \mathbf{B} = \begin{bmatrix} \mathbf{A}(1, 1)\mathbf{B} & \cdots & \mathbf{A}(1, J)\mathbf{B} \\ \vdots & \ddots & \vdots \\ \mathbf{A}(I, 1)\mathbf{B} & \cdots & \mathbf{A}(I, J)\mathbf{B} \end{bmatrix}.$$

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

$$114 \quad \mathbf{A} \odot \mathbf{B} = [\mathbf{A}(\cdot, 1) \otimes \mathbf{B}(\cdot, 1) \cdots \mathbf{A}(\cdot, K) \otimes \mathbf{B}(\cdot, K)]$$

115 **2.1.3. Tensor basics.** For a tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, its *mode* or *order* is the
 116 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
 117 product of two tensors \mathcal{X}, \mathcal{Y} is defined as $\langle \mathcal{X}, \mathcal{Y} \rangle = \sum_{i_1=1}^{I_1} \cdots \sum_{i_N=1}^{I_N} \mathcal{X}_{i_1 \dots i_N} \mathcal{Y}_{i_1 \dots i_N}$.
 118 The *Frobenius norm* of \mathcal{X} is $\|\mathcal{X}\|_F = \sqrt{\langle \mathcal{X}, \mathcal{X} \rangle}$.

119 **2.1.4. Tensor unfoldings.** Let $\bar{I} = \prod_{j=1}^N I_j$ and $I_{(-n)} = \prod_{j \neq n} I_j$, and let $\text{vec}(\mathcal{X})$
 120 denote the vectorization of \mathcal{X} . The *mode- n unfolding* of \mathcal{X} is the matrix $\mathbf{X}^{(n)} \in$
 121 $\mathbb{R}^{I_n \times I_{(-n)}}$. The inner product for tensors matches that of any mode- n unfolding:

$$122 \quad (2.2) \quad \langle \mathcal{X}, \mathcal{Y} \rangle = \langle \mathbf{X}^{(n)}, \mathbf{Y}^{(n)} \rangle = \text{Tr}((\mathbf{X}^{(n)})^\top \mathbf{Y}^{(n)}).$$

123 **2.1.5. Tensor rank.** The *mode- n rank* is the rank of the mode- n unfolding. We
 124 say the *rank* of \mathcal{X} is $\mathbf{r}(\mathcal{X}) = (r_1, \dots, r_N)$ if its *mode- n rank* is r_n for each $n \in [N]$.
 125 This notion of rank corresponds to the size of the core tensor in a Tucker factorization
 126 of \mathcal{X} . A *superdiagonal* tensor generalizes a diagonal matrix: all entries are zero except
 127 for the entries whose indices in each dimension are equal.

128 **2.1.6. Tensor contractions.** Write $\mathbf{G} = \mathbf{X} \times_n \mathbf{U}$ for the *mode-n (matrix) product*
 129 of \mathbf{X} with $\mathbf{U} \in \mathbb{R}^{J \times I_n}$. That is, $\mathbf{G} = \mathbf{X} \times_n \mathbf{U} \iff \mathbf{G}^{(n)} = \mathbf{U}\mathbf{X}^{(n)}$. The tensor \mathbf{G}
 130 has dimension $I_1 \times \dots \times I_{n-1} \times J \times I_{n+1} \times \dots \times I_N$. Mode products with respect to
 131 different modes commute: for $\mathbf{U} \in \mathbb{R}^{J_1 \times I_n}$, $\mathbf{V} \in \mathbb{R}^{J_2 \times I_m}$

$$132 \quad \mathbf{X} \times_n \mathbf{U} \times_m \mathbf{V} = \mathbf{X} \times_m \mathbf{V} \times_n \mathbf{U} \quad \text{if } n \neq m.$$

133 Mode products along the same mode simplify: for $\mathbf{A} \in \mathbb{R}^{J_1 \times I_n}$, $\mathbf{B} \in \mathbb{R}^{J_2 \times J_1}$,

$$134 \quad \mathbf{X} \times_n \mathbf{A} \times_n \mathbf{B} = \mathbf{X} \times_n (\mathbf{B}\mathbf{A}).$$

135 **2.2. Tucker Approximation.** Given a tensor $\mathbf{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ and target rank
 136 $\mathbf{r} = (r_1, \dots, r_N)$, the idea of Tucker approximation is finding a *core tensor* $\mathbf{G} \in$
 137 $\mathbb{R}^{r_1 \times \dots \times r_N}$ and matrices with orthogonal columns $\mathbf{U}_n \in \mathbb{R}^{I_n \times r_n}$ for $n \in [N]$, called
 138 *factor matrices*, so that

$$139 \quad \mathbf{X} \approx \mathbf{G} \times_1 \mathbf{U}_1 \times \dots \times_N \mathbf{U}_N.$$

140 For brevity, we define $\llbracket \mathbf{G}; \mathbf{U}_1, \dots, \mathbf{U}_N \rrbracket = \mathbf{G} \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N$. Any best rank- \mathbf{r}
 141 Tucker approximation is of the form $\llbracket \mathbf{G}^*; \mathbf{U}_1^*, \dots, \mathbf{U}_N^* \rrbracket$, where \mathbf{G}^* , \mathbf{U}_n^* solve the problem

$$142 \quad (2.3) \quad \begin{aligned} &\text{minimize} && \|\mathbf{X} - \mathbf{G} \times_1 \times \dots \times_{n+1} \times_N \mathbf{U}_N\|_F^2 \\ &\text{subject to} && \mathbf{U}_n^T \mathbf{U}_n = \mathbf{I}. \end{aligned}$$

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

146 **2.2.1. HOSVD.** The standard approach to computing a rank $\mathbf{r} = (r_1, \dots, r_N)$
 147 Tucker approximation for a tensor \mathbf{X} begins with the higher order singular value
 148 decomposition (HOSVD) [13, 36], (2.1):

Algorithm 2.1 Higher order singular value decomposition (HOSVD) [13, 36]

Given: tensor \mathbf{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 \mathbf{X} to form the core

$$\mathbf{G} = \mathbf{X} \times_1 \mathbf{U}_1 \times \dots \times_N \mathbf{U}_N.$$

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

148 The HOSVD can be computed in two passes over the tensor [42, 8]. We describe
 149 this method briefly here, and in more detail in the next section. In the first pass,
 150 sketch each matricization $\mathbf{X}^{(n)}$, $n \in [N]$, and use randomized linear algebra (e.g., the
 151 randomized range finder of [16]) to (approximately) recover its range \mathbf{U}_n . To form
 152 the core $\mathbf{G} = \mathbf{X} \times_1 \mathbf{U}_1 \times \dots \times_N \mathbf{U}_N$ requires a second pass over \mathbf{X} , since the factor matrices
 153 \mathbf{U}_n depend on \mathbf{X} . The main algorithmic contribution of this paper is to develop a
 154 method to approximate both the factor matrices and the core in just one pass over \mathbf{X} .

155 **2.2.2. ST-HOSVD.** [37] proposes a sequentially truncated higher-order singular
 156 value decomposition which enjoys the same approximation compared to HOSVD, but
 157 requires less operators. The idea, is to update the low rank approximation of the
 158 target tensor \mathbf{X} when we get factor matrix in Algorithm 2.1 for the following factor
 159 extraction. Readers can go to [37] for more details.

161 **2.2.3. HOOI.** The higher order orthogonal iteration (HOOI) [13] (2.2) improves
 162 on the resulting Tucker factorization by repeatedly minimizing the objective of 2.3
 over the the core and the factor matrices. Notice the core update (2.5) admits the

Algorithm 2.2 Higher order orthogonal iteration (HOOI) [13]

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

Initialize: compute $\mathcal{X} \approx [\mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N]$ using HOSVD

Repeat:

1. *Factors.* For each $n \in [N]$,

$$(2.4) \quad \mathbf{U}_n \leftarrow \arg \min_{\mathbf{U}_n} \|[\mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N] - \mathcal{X}\|_F^2,$$

2. *Core.*

$$(2.5) \quad \mathcal{G} \leftarrow \arg \min_{\mathcal{G}} \|[\mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N] - \mathcal{X}\|_F^2.$$

Return: Tucker approximation $\mathcal{X}_{\text{HOOI}} = [\mathcal{G}; \mathbf{U}_1, \dots, \mathbf{U}_N]$

163 closed form solution $\mathcal{G} \leftarrow \mathcal{X} \times_1 \mathbf{U}_1^\top \cdots \times_N \mathbf{U}_N^\top$, which motivates the second step of
 164 HOSVD.
 165

166 **2.3. Previous Work.** The only previous work on streaming Tucker approxima-
 167 tion is [25], which develops a streaming method called Tucker TensorSketch (T.-TS)
 168 [25, Algorithm 2]. T.-TS improves on HOOI by sketching the data matrix in the least
 169 squares problems. However, the success of the approach depends on the quality of
 170 the initial core and factor matrices, and the alternating least squares algorithm takes
 171 several iterations to converge.

172 In contrast, our work is motivated by HOSVD (not HOOI), and requires no
 173 initialization or iteration. We treat the tensor as a *multilinear* operator. The sketch
 174 identifies a low-dimensional subspace *for each input argument* that captures the action
 175 of the operator. The reconstruction produces a low-Tucker-rank multilinear operator
 176 with the same action on this low-dimensional tensor product space. This linear
 177 algebraic view allows us to develop the first guarantees on approximation error for this
 178 class of problems¹. Moreover, we show in numerical experiments that our algorithm
 179 achieves a better approximation of the original tensor given the same memory resources.

180 More generally, there is a large literature on randomized algorithms for matrix
 181 factorizations and for solving optimization problems; see e.g. the review articles [16, 40].
 182 In particular, our method is strongly motivated by the recent papers [33, 34], which
 183 provide methods for one-pass matrix approximation. The novelty of this paper is in
 184 our design of a core sketch (and reconstruction) for the Tucker decomposition, together
 185 with provable performance guarantees. The proof requires a careful accounting of the
 186 errors resulting from the factor sketches and from the core sketch. The structure of
 187 the Tucker sketch guarantees that these errors are independent.

188 Many researchers have used randomized algorithms to compute tensor decomposi-
 189 tions. For example, [39, 7] apply sketching techniques to the CP decomposition, while

¹ The guarantees in [25] 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 [25] show that the method achieves lower error using a fixed sketch (with no guarantees) than using fresh sketches at each iteration.

[35] suggests sparsifying the tensor. Several papers aim to make Tucker decomposition efficient in the limited-memory or distributed settings [6, 42, 4, 19, 23, 8].

3. Dimension Reduction Maps. In this section, we first introduce some commonly used randomized dimension reduction maps together with some mathematical background, and explain how to calculate and update sketches.

3.1. Dimension Reduction Map. Dimension reduction maps (DRMs) take a collection of high dimensional objects to a lower dimensional space while preserving certain geometric properties [27]. For example, we may wish to preserve the pairwise distances between vectors, or to preserve the column space of matrices. We call the output of a DRM on an object x a *sketch* of x .

Some common DRMs include matrices with i.i.d. Gaussian entries or i.i.d. ± 1 entries. The Scrambled Subsampled Randomized Fourier Transform (SSRFT) [41] and sparse random projections [1, 24] can achieve similar performance with fewer computational and storage requirements; see F for details.

Our theoretical bounds rely on properties of the Gaussian DRM. However, our numerical experiments indicate that many other DRMs yield qualitatively similar results; see, e.g., Figure 1, Figure 10 and Figure 9) in Figure I.

3.2. Tensor Random Projection. Here we present a strategy for reducing the storage of the random map that makes use of the tensor random projection (TRP), and extremely low storage structured dimension reduction map proposed in [30]. The *tensor random projection (TRP)* $\Omega : \prod_{n=1}^N I_n \rightarrow \mathbb{R}^k$ is defined as the iterated Khatri-Rao product of DRMs $\mathbf{A}_n \in \mathbb{R}^{I_n \times k}$, $n \in [N]$:

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

Each $\mathbf{A}_n \in \mathbb{R}^{I_n \times k}$ can be a Gaussian map, sign random projection, SSRFT, etc. The number of constituent maps N and their dimensions I_n for $n \in [N]$ are parameters of the TRP, and control the quality of the map; see [30] for details. The TRP map is a row-product random matrix, which behaves like a Gaussian map in many respects [28]. Our experimental results confirm this behavior.

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

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

Table 1: Performance of Different Dimension Reduction Maps: We compare the storage cost 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 sparse factor for sparse random projection. The TRP considered here is composed of Gaussian DRMs.

We do not need to explicitly form or store the TRP map Ω . Instead, we can store its constituent DRMs $\mathbf{A}_1, \dots, \mathbf{A}_N$ and compute the action of the map on the

224 matricized tensor using the definition of the TRP. The additional computation required
 225 is minimal and empirically incurs almost no performance loss.

226 **4. Algorithms for Tucker approximation.** In this section, we present our
 227 proposed tensor sketch and algorithms for one- and two-pass Tucker approximation,
 228 and discuss the computational complexity and storage cost of these methods for both
 229 sparse and dense input tensors. We present guarantees for these methods in 5.

230 **4.1. Tensor compression via sketching.**

231 *The Tucker sketch.* Our Tucker sketch generalizes the matrix sketch of [33] to
 232 higher order tensors. To compute a Tucker sketch for tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \dots \times I_N}$ with
 233 sketch size parameters \mathbf{k} and \mathbf{s} , draw independent, random DRMs

234 (4.1) $\Omega_1, \Omega_2, \dots, \Omega_N$ and $\Phi_1, \Phi_2, \dots, \Phi_N$,

235 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

236 $\mathbf{V}_n = \mathbf{X}^{(n)} \Omega_n \in \mathbb{R}^{I_n \times k_n}, \quad n \in [N],$

237 $\mathcal{H} = \mathcal{X} \times_1 \Phi_1^\top \cdots \times_N \Phi_N^\top \in \mathbb{R}^{s_1 \times \dots \times s_N}.$

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

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

246 The sketches \mathbf{V}_n and \mathcal{H} are linear functions of the original tensor, so we can
 247 compute the sketches in a single pass over the tensor \mathcal{X} . Linearity enables easy
 248 computation of the sketch even in the streaming model (E.1) or distributed model
 249 (E.2). Storing the sketches requires memory $\sum_{n=1}^N I_n \cdot k_n + \prod_{i=1}^N s_n$: much less than
 250 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  $\mathbf{V}_n \leftarrow \mathbf{X}^{(n)} \Omega_n$ ,  $n \in [N]$ 
4:   Compute core sketch  $\mathcal{H} \leftarrow \mathcal{X} \times_1 \Phi_1^\top \times \dots \times_N \Phi_N^\top$ 
5:   return ( $\mathcal{H}, \mathbf{V}_1, \dots, \mathbf{V}_N, \{\Phi_n, \Omega_n\}_{n \in [N]}$ )
6: end function

```

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

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 [3]. In particular, we expect both to perform similarly for tensor approximation.

257 *Remark 4.2.* Alternatively, the TRP (3.2) can be used to limit the storage of Ω_n
 258 required. The Khatri-Rao structure in the sketch need not match the structure in the
 259 matricized tensor. However, we can take advantage of the structure of our problem to
 260 reduce storage even further. We generate DRMs $\mathbf{A}_n \in \mathbb{R}^{I_n \times k}$ for $n \in [N]$ and define
 261 $\Omega_n = \mathbf{A}_1 \odot \cdots \odot \mathbf{A}_{n-1} \odot \mathbf{A}_{n+1} \odot \cdots \odot \mathbf{A}_N$ for each $n \in [N]$. Hence we need not store the
 262 maps Ω_n , but only the small matrices \mathbf{A}_n . The storage required is thereby reduced
 263 from $\mathcal{O}(N(\prod_{n=1}^N I_n)k)$ to $\mathcal{O}((\sum_{n=1}^N I_n)k)$, while the approximation error is essentially
 264 unchanged. We use this method in our experiments.

265 **4.2. Low-Rank Approximation.** Now we explain how to construct a Tucker
 266 decomposition of \mathfrak{X} with target Tucker rank \mathbf{k} from the factor and core sketches.

267 We first present a simple two-pass algorithm, 4.2, that uses only the factor sketches
 268 by projecting the unfolded matrix of original tensor \mathfrak{X} to the column space of each
 269 factor sketch. To project to the column space of each factor matrix, we calculate the
 270 QR decomposition of each factor sketch:

$$271 \quad (4.2) \quad \mathbf{V}_n = \mathbf{Q}_n \mathbf{R}_n \quad \text{for } n \in [N],$$

272 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.

273 Consider the tensor approximation

$$274 \quad (4.3) \quad \tilde{\mathfrak{X}} = \mathfrak{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top.$$

275 This approximation admits the guarantees stated in 5.1. Using the commutativity of
 276 the mode product between different modes, we can rewrite $\tilde{\mathfrak{X}}$ as

$$277 \quad (4.4) \quad \tilde{\mathfrak{X}} = \underbrace{[\mathfrak{X} \times \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top]}_{\mathcal{W}_2} \times_1 \mathbf{Q}_1 \times_2 \cdots \times_N \mathbf{Q}_N = [\mathcal{W}_2; \mathbf{Q}_1, \dots, \mathbf{Q}_N],$$

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

Algorithm 4.2 Two Pass Sketch and Low Rank Recovery

Given: tensor \mathfrak{X} , sketch parameters \mathbf{k} and $\mathbf{s} \geq \mathbf{k}$

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

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

282 *One-Pass Approximation.* To develop a one-pass method, we must use the core
 283 sketch \mathcal{H} — the compression of \mathfrak{X} using the random projections Φ_n — to approximate
 284 \mathcal{W}_2 — the compression of \mathfrak{X} using random projections \mathbf{Q}_n . To develop intuition,
 285 consider the following calculation: if the factor matrix approximations \mathbf{Q}_n capture
 286 the range of \mathfrak{X} well, then projection onto their ranges in each mode approximately
 287 preserves the action of \mathfrak{X} :

$$288 \quad \mathfrak{X} \approx \mathfrak{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top$$

289 Recall that for tensor \mathcal{A} , and matrix \mathbf{B} and \mathbf{C} with compatible sizes, $\mathcal{A} \times_n (\mathbf{BC}) =$
 290 $(\mathcal{A} \times_n \mathbf{C}) \times_n \mathbf{B}$. Use this rule to collect terms to recognize the two pass core approxi-
 291 mation \mathcal{W}_2 :

$$292 \quad \mathfrak{X} \approx (\mathfrak{X} \times_1 \mathbf{Q}_1^\top \times \cdots \times_N \mathbf{Q}_N^\top) \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N = \mathcal{W}_2 \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N$$

293 Now contract both sides of this approximate equality with the DRMs Φ_n and recognize
 294 the core sketch \mathcal{H} :

$$295 \quad \mathcal{H} := \mathfrak{X} \times_1 \Phi_1^\top \cdots \times_N \Phi_N^\top \approx \mathcal{W}_2 \times_1 \Phi_1^\top \mathbf{Q}_1 \times \cdots \times_N \Phi_N^\top \mathbf{Q}_N.$$

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

$$298 \quad \mathcal{W}_2 \approx \mathcal{H} \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger =: \mathcal{W}_1.$$

299 The right hand side of the approximation defines the one pass core approximation \mathcal{W}_1 .
 300 [B.2](#) controls the error in this approximation.

301 [4.3](#) summarizes the resulting one-pass algorithm. One (streaming) pass over the
 302 tensor can be used to sketch the tensor; to recover the tensor, we only access the
 303 sketches. [5.2](#) (below) bounds the overall quality of the approximation.

Algorithm 4.3 One Pass Sketch and Low Rank Recovery

Given: tensor \mathfrak{X} , sketch parameters k and $s \geq k$

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

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

304 The time and storage cost of [4.3](#) is given by [2](#). The time and storage complexity
 305 of these methods compare favorably to the only previous method for streaming Tucker
 306 approximation [25]; see [H](#) for details.

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

Table 2: Computational Complexity of [4.3](#) on tensor $\mathfrak{X} \in \mathbb{R}^{I \times \cdots \times I}$ with parameters (k, s) , using a TRP composed of Gaussian DRMs inside the Tucker sketch. By far the majority of the time is spent sketching the tensor \mathfrak{X} .

307 **4.3. Fixed-Rank Approximation.** Algorithm [4.2](#) and algorithm [4.3](#) produce
 308 a two-pass and one-pass rank- k tensor approximation respectively. It is often valuable
 309 to truncate this approximation to a user-specified target rank $r \leq k$ [34, Figure 4].

310 Our fixed rank approximation method is motivated by the following lemma:

311 **LEMMA 4.1.** *Let $\mathcal{W} \in \mathbb{R}^{k_1 \times \cdots \times k_N}$ be a tensor, and let $\mathbf{Q}_n \in \mathbb{R}^{I_n \times k_n}$ be orthogonal
 312 matrices with $k_n \geq r_n$ for $n \in [N]$. Then*

$$313 \quad [\mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N]_r = [\mathcal{W}]_r \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N.$$

314 (This lemma does not necessarily hold if the best rank- r Tucker approximation $[\cdot]$ is
 315 replaced by the output of any concrete algorithm such as HOSVD or HOOI.) The
 316 proof of [4.1](#) appears in [C](#).

Motivated by this lemma, to produce a fixed rank \mathbf{r} approximation of \mathfrak{X} , we compress the core tensor approximation from 4.2 or 4.3 to rank \mathbf{r} . This compression is cheap because the core approximation $\mathcal{W} \in \mathbb{R}^{k_1 \times \dots \times k_N}$ is small. We present this method (using HOOI as the the compression algorithm) as 4.4. Other compression algorithms can be used to trade off the quality of approximation with the difficulty of running the algorithm. Reasonable choices include the sequentially-truncated HOSVD (ST-HOSVD) [37] or TTHRESH [5]. Both HOSVD and ST-HOSVD are psuedual

Algorithm 4.4 Fixed rank approximation

Given: Tucker approximation $[\mathcal{W}; \mathbf{Q}_1, \dots, \mathbf{Q}_N]$ of tensor \mathfrak{X} , rank target \mathbf{r}

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

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

5. Guarantees. In this section, we present probabilistic guarantees on the preceding algorithms. We show that approximation error for the one-pass algorithm is the sum of the error from the two-pass algorithm and the error resulting from the core approximation. Proofs for the three theorems in this section can be found in the corresponding subsections of A.

5.1. Low rank approximation. 5.1 guarantees the performance of the two pass method 4.2.

THEOREM 5.1. *Sketch the tensor \mathfrak{X} using a Tucker sketch with parameters \mathbf{k} using DRMs with i.i.d. Gaussian $\mathcal{N}(0, 1)$ entries. Then the approximation $\hat{\mathfrak{X}}_2$ computed with the two pass method 4.2 satisfies*

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

The two pass method does not use the core sketch, so this result does not depend on \mathbf{s} . Theorem 5.2 guarantees the performance of one pass method 4.3.

THEOREM 5.2. *Sketch the tensor \mathfrak{X} using a Tucker sketch with parameters \mathbf{k} and $\mathbf{s} > 2\mathbf{k}$ using DRMs with i.i.d. Gaussian $\mathcal{N}(0, 1)$ entries. Then the approximation $\hat{\mathfrak{X}}_1$ computed with the one pass method 4.3 satisfies the bound*

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

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

The theorem shows that the method works best for tensors whose unfoldings exhibit spectral decay. As a simple consequence of this result, we see that the two pass method with $\mathbf{k} > \mathbf{r} + 1$ perfectly recovers a tensor with exact Tucker rank \mathbf{r} , since in that case $\tau_{\mathbf{r}_n}^{(n)} = 0$ for each $n \in [N]$. However, this theorem states a stronger bound: the method exploits decay in the spectrum, wherever (in the first k_n singular values of each mode n unfolding) it occurs.

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

351 *Remark 5.1.* Both HOSVD, ST-HOSVD achieves so called pseudo optimal with
 352 parameter \sqrt{N} . For example, for ST-HOSVD,

$$353 \quad (5.1) \quad \|[\mathcal{X} - [\mathcal{X}]_{\text{ST-}\mathbf{k}}]\|_F \leq \sqrt{N} \|[\mathcal{X} - [\mathcal{X}]_{\mathbf{k}}]\|_F \leq \sqrt{\sum_{n=1}^N (\tau_{k_n}^{(n)})^2},$$

354 where $[\mathcal{X}]_{\mathbf{k}}$ is the optimal Tucker rank \mathbf{k} approximation. Thus, the low rank approxi-
 355 mation for two pass algorithm is pseduo optimal with factor $\sqrt{2N}$ while for one pass
 356 algorithm, if we choose $s > 2\mathbf{k} + 1 (\Delta \leq 2)$, it is also pseduo optimal with factor $2\sqrt{N}$.

357 [5.2](#) also offers guidance on how to select the sketch size parameters s and \mathbf{k} . In
 358 particular, suppose that the mode- n unfolding has a good rank r_n approximation for
 359 each mode n . Then the choices $k_n = 2r_n + 1$ and $s_n = 2k_n + 1$ ensure that

$$360 \quad \mathbb{E} \|\mathcal{X} - \hat{\mathcal{X}}\|_F^2 \leq 4 \sum_{n=1}^N (\tau_{r_n}^{(n)})^2.$$

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

363 **5.2. Fixed rank approximation.** We now present a conditional analysis of the
 364 fixed rank approximation method given a low rank approximation. Recall that $[\cdot]_{\mathbf{r}}$
 365 returns a best rank- \mathbf{r} Tucker approximation.

366 **THEOREM 5.3.** *Suppose $\hat{\mathcal{X}} = [\mathcal{W}; \mathbf{Q}_1, \dots, \mathbf{Q}_N]$ approximates the target tensor
 367 \mathcal{X} , and let $\hat{\mathcal{X}}_{\mathbf{r}}$ denote some rank \mathbf{r} approximation to $\hat{\mathcal{X}}$ from some procedure like
 368 ST-HOSVD, Then for output any fix rank \mathbf{r}*

$$369 \quad \mathbb{E} \|\mathcal{X} - \hat{\mathcal{X}}_{\mathbf{r}}\|_F \leq \|\mathcal{X} - [\mathcal{X}]_{\mathbf{r}}\|_F + 2\sqrt{\mathbb{E} \|\mathcal{X} - \hat{\mathcal{X}}\|_F^2}.$$

370 The second term on the right-hand side of [5.3](#) is controlled by [5.1](#) and [5.2](#). Hence
 371 we can combine these results to provide guarantees for fixed rank approximation with
 372 either the two pass or one pass algorithms.

373 The resulting bound shows that the best rank- \mathbf{r} approximation of the output
 374 from the one or two pass algorithms is comparable in quality to a true best rank- \mathbf{r}
 375 approximation of the input tensor. An important insight is that the sketch size
 376 parameters s and \mathbf{k} that guarantee a good low rank approximation also guarantee a
 377 good fixed rank approximation: the error due to sketching depends only on the sketch
 378 size parameters \mathbf{k} and s , and not on the target rank \mathbf{r} .

379 In practice, one would truncate the rank of the approximation using HOOI [\(4.4\)](#),
 380 rather than the best rank \mathbf{r} approximation [\(??\)](#). Guarantees for resulting algorithm are
 381 beyond the scope of this paper, since there are no strong guarantees on the performance
 382 of HOOI; however, it is widely believed to produce an approximation that is usually
 383 quite close to the best rank \mathbf{r} approximation.

384 **5.3. Proof sketch.** To bound the approximation error of the algorithms pre-
 385 sented in the main body of this paper, we first develop several structural results showing
 386 an additive decomposition of the error. First, the total error is the sum of the error
 387 due to sketching and the error due to fixed rank approximation. Second, the sketching
 388 error is the sum of the error due to the factor matrix approximations and to the core
 389 approximation. Third, the error due to the factor matrix approximations is the sum of
 390 the error in each the modes, as the errors due to each mode are mutually orthogonal.

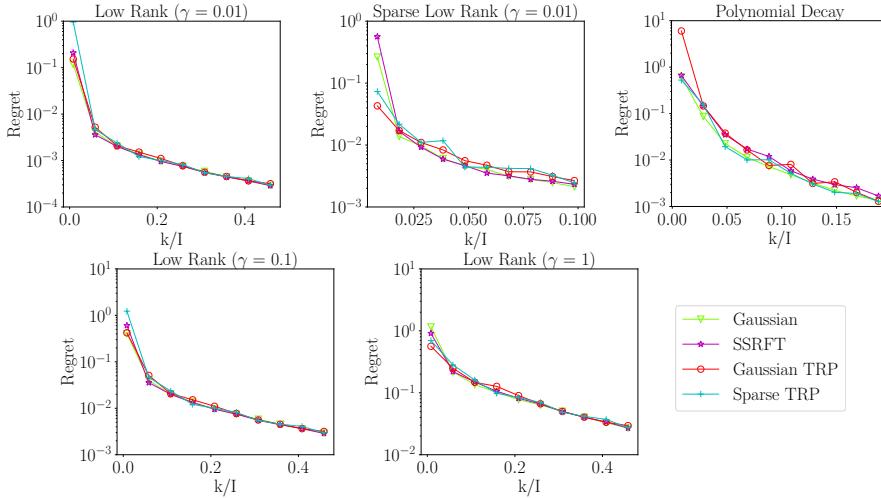


Fig. 1: *Different DRMs perform similarly.* We approximate 3D synthetic tensors (see 6.1) with $I = 600$, using our one-pass algorithm with $r = 5$ and varying k ($s = 2k + 1$), using a variety of DRMs in the Tucker sketch: Gaussian, SSRFT, Gaussian TRP, or Sparse TRP.

391 This finishes the approximation error bound for the two pass algorithm, 5.1. As for the
 392 error due to the core approximation, we rewrite the approximation error in the core
 393 tensor as a sum over each mode of errors that are mutually orthogonal. Indeed, these
 394 errors have the same form as the errors due to the factor matrix approximations, scaled
 395 down by a factor $\Delta(k, s)$ that depends on the sketch sizes \mathbf{k} and \mathbf{s} . This argument
 396 shows the error due to the core approximation is at most a factor $\Delta(k, s)$ times the
 397 error due to the factor matrix approximation.

398 **6. Numerical Experiments.** In this section, we study the performance of our
 399 method. We compare the performance of the method using various different DRMs,
 400 including TRP. We also compare our method with the algorithm proposed by [25] to
 401 show that for the same storage budget, our method produces better approximations.
 402 Our two-pass algorithm outperforms the one-pass version, as expected. (Contrast this
 403 to [25], where the multi-pass method performs less well than the one-pass version.)
 404 We evaluate the experimental results using two metrics:

$$\begin{aligned} \text{normalized error: } & \|\mathfrak{X} - \hat{\mathfrak{X}}\|_F / \|\mathfrak{X}\|_F \\ \text{regret: } & \left(\|\mathfrak{X} - \hat{\mathfrak{X}}\|_F - \|\mathfrak{X} - \mathfrak{X}_{\text{HOOI}}\|_F \right) / \|\mathfrak{X}\|_F. \end{aligned}$$

405 The normalized error measures the fraction of the energy in \mathfrak{X} captured by the
 406 approximation. The regret measures the increase in normalized error due to using the
 407 approximation $\hat{\mathfrak{X}}$ rather than using $\mathfrak{X}_{\text{HOOI}}$. The relative error measures the decrease in
 408 performance relative to HOOI. The normalized error of a rank \mathbf{r} Tucker approximation
 409 \hat{X} is always positive when \mathfrak{X} has a larger rank. In general, we find our proposed
 410 methods approaches the performance of HOOI for large enough storage budgets.

411 We ran all experiments on a server with 128 Intel® Xeon® E7-4850 v4 2.10GHz
 412 CPU cores and 1056GB memory. The code for our method is available at an anonymous
 413 Github repository <https://github.com/tensorsketch/tensorsketch>.

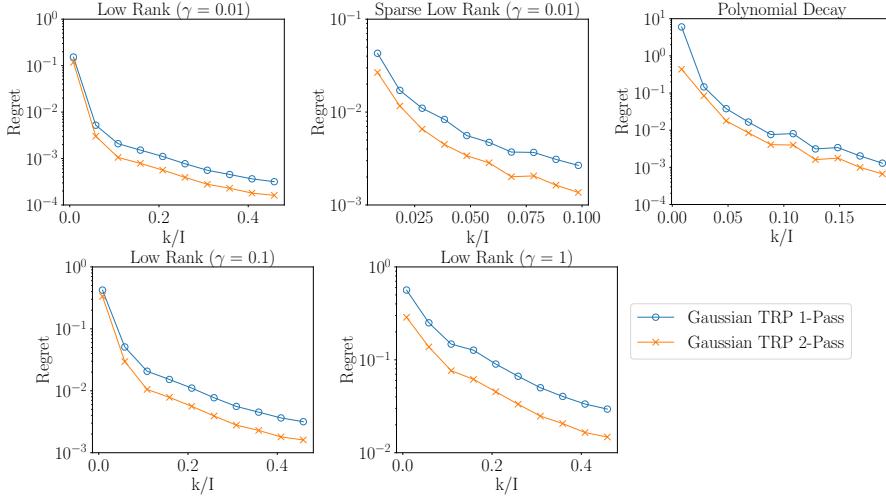


Fig. 2: *Two-pass improves on one-pass.* We approximate 3D synthetic tensors (see 6.1) with $I = 600$, using our one-pass and two-pass algorithms with $r = 5$ and varying k ($s = 2k + 1$), using the Gaussian TRP in the Tucker sketch.

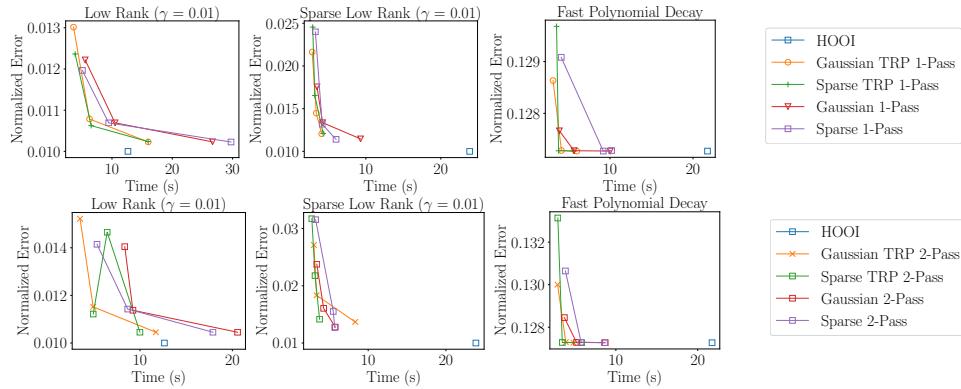


Fig. 3: *Faster approximations.* We approximate 3D synthetic tensors with $I = 600$ generated as described in 6.1, using HOOI and our one-pass and two-pass algorithms with $r = 5$ for a few different k ($s = 2k + 1$).

415 **6.1. Synthetic experiments.** All synthetic experiments use an input tensor
416 with equal side lengths I . We consider three different data generation schemes:

- 417 • *Low rank + noise.* Generate a core tensor $\mathcal{C} \in \mathbb{R}^{r^N}$ with entries drawn
418 from $\text{Unif}([0, 1])$. Independently generate N orthogonal factor matrices
419 $\mathbf{A}_1, \dots, \mathbf{A}_N \in \mathbb{R}^{r \times I}$. Define $\mathcal{X}^\natural = \mathcal{C} \times_1 \mathbf{A}_1 \cdots \times_N \mathbf{A}_N$ and the noise
420 parameter $\gamma > 0$. Generate an input tensor as $\mathcal{X} = \mathcal{X}^\natural + (\gamma \|\mathcal{X}^\natural\|_F / I^{N/2}) \boldsymbol{\epsilon}$ where
421 the noise $\boldsymbol{\epsilon}$ has i.i.d. $\mathcal{N}(0, 1)$ entries.
422 • *Sparse low rank + noise.* We construct the input tensor \mathcal{X} as above (Low
423 Rank + Noise), but with sparse factor matrices \mathbf{A}_n : If δ_n is the sparsity
424 (proportion of non-zero elements) of \mathbf{A}_n , then the sparsity of the true signal

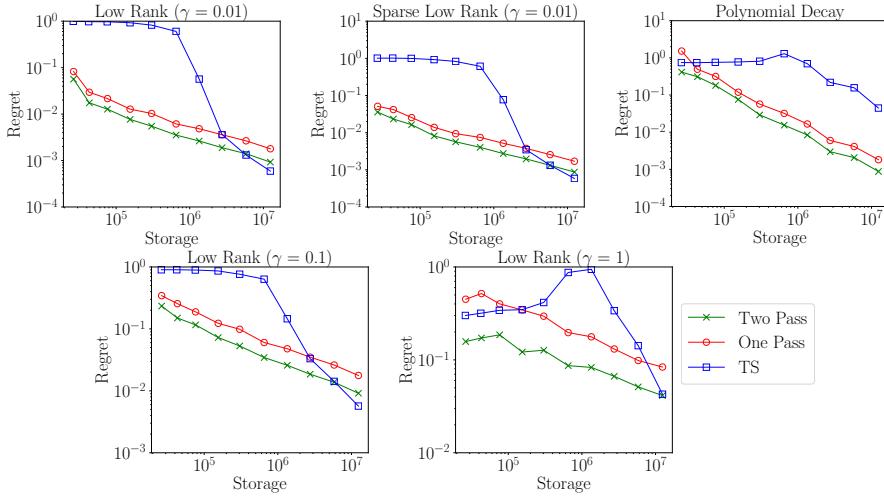


Fig. 4: *Approximations improve with more memory: synthetic data.* We approximate 3D synthetic tensors (see 6.1) with $I = 300$, using T.-TS and our one-pass and two-pass algorithms with the Gaussian TRP to produce approximations with equal ranks $r = 10$. Notice every marker on the plot corresponds to a $2700 \times$ compression!

425 \mathcal{X}^\natural is $\prod_{n=1}^N \delta_n$. We use $\delta_n = 0.2$ unless otherwise specified.
 426 • *Polynomial decay.* We construct the input tensor \mathcal{X} as

427
$$\mathcal{X} = \text{superdiag}(1, \dots, 1, 2^{-t}, 3^{-t}, \dots, (I - r)^{-t}).$$

428 The first r entries are 1. Recall **superdiag** converts a vector to N dimensional
 429 superdiagonal tensor. Our experiments use $t = 1$ (geometric decay).

430 **6.1.1. Different dimension reduction maps perform similarly.** Our first
 431 experiment investigates the performance of our one-pass fixed-rank algorithm as
 432 the sketch size (and hence, required storage) varies, for several types of dimension
 433 reductions maps, including Gaussian, SSRFT, Gaussian TRP, and Sparse TRP. We
 434 generate synthetic data as described above with $\mathbf{r} = (5, 5, 5)$, $I = 600$. 1 shows the
 435 rank- \mathbf{r} approximation error as a function of the compression factor k/I . (Results for
 436 other input tensors are presented as 9 and 10 in I.) We see that the log relative error
 437 for our one-pass algorithm converges to that of HOOI as k increases for all input
 438 tensors. In the low rank case, the convergence rate is lower for higher noise levels. In
 439 general, the performance for different maps are approximately the same, although our
 440 theory only pertains to the Gaussian map.

441 We evaluate the run time for HOOI and our two algorithms with several different
 442 DRMs in 3. We can see that the one-pass algorithm is always slightly faster than the
 443 two-pass algorithm. The TRP generally provides a modest speedup in addition to the
 444 memory advantage. Both our one-pass and two-pass algorithms achieve nearly the
 445 accuracy of HOOI, and are usually much faster.

446 **6.1.2. A second pass reduces error.** The second experiment compares our
 447 two-pass and one-pass algorithm. The design is similar to the first experiment. 2 shows
 448 that the two-pass algorithm typically outperforms the one-pass algorithm, especially

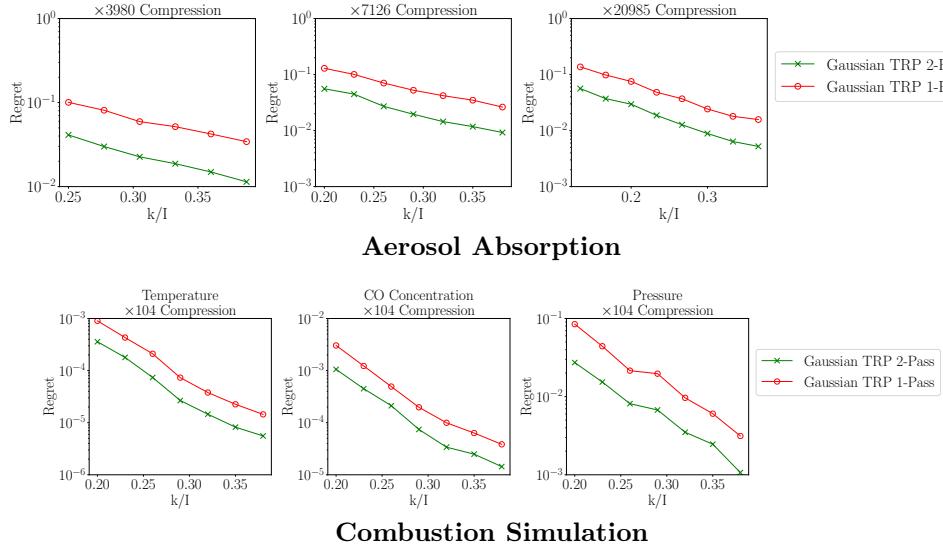


Fig. 5: *Approximations improves with more memory: real data.* We approximate aerosol absorption and combustion data using our one-pass and two-pass algorithms with the Gaussian TRP. We compare three target ranks ($r/I = 0.125, 0.1, 0.067$) for the former, and use the same target rank ($r/I = 0.1$) for each measured quantity in the combustion dataset. Notice $r/I = 0.1$ gives a hundred-fold compression!

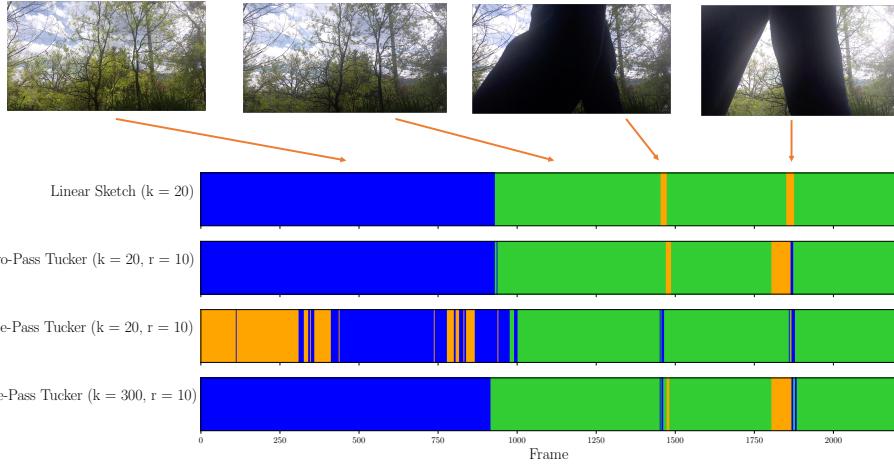
449 in the high-noise, sparse, or rank-decay case. Both converge at the same asymptotic
 450 rate. (Results for other input tensors are available in I.)

451 **6.1.3. Improvement on state-of-the-art.** The third experiment compares the
 452 performance of our two-pass and one-pass algorithms and Tucker TensorSketch (T.-
 453 TS), as described in [25], the only extant one-pass algorithm. For a fair comparison,
 454 we allocate the same storage budget to each algorithm and compare the relative error
 455 of the resulting fixed-rank approximations. We approximate synthetic 3D tensors with
 456 side length $I = 300$ with Tucker rank $r = 10$. We use the suggested parameter settings
 457 for each algorithm: $k = 2r$ and $s = 2k + 1$ for our methods; $K = 10$ for T.-TS. Our
 458 one-pass algorithm (with the Gaussian TRP) uses $((2k + 1)^N + kIN)$ storage, whereas
 459 T.-TS uses $(Kr^{2N} + Kr^{2N-2})$ storage (see 3 in H).

460 Figure 4 shows that our algorithms generally perform as well as T.-TS, and
 461 dramatically outperforms for small storage budgets. For example, our method achieves
 462 $1/50$, $1/50$, $1/7$, and $1/4$ the relative error of T.-TS for low rank and sparse low rank
 463 ($\gamma = 0.01$), low rank ($\gamma = 0.1$), and polynomial-decay input tensors, respectively. For
 464 the low rank ($\gamma = 1$) tensor, the performance of T.-TS is not even monotone as the
 465 storage budget increases! The performance of T.-TS is comparable with that of the
 466 algorithms presented in this paper only when the storage budget is large.

467 *Remark 6.1.* The paper [25] proposes a multi-pass method, Tucker Tensor-Times-
 468 Matrix-TensorSketch (TTMTS) that is dominated by the one-pass method Tucker
 469 TensorSketch(TS) in all numerical experiments; hence we compare only with T.-TS.

470 **6.2. Applications.** We also apply our method to datasets drawn from three
 471 application domains: climate, combustion, and video.



Video Scene Classification

Fig. 6: *Video Scene Classification* ($2200 \times 1080 \times 1980$): We classify frames from the video data from [25] (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

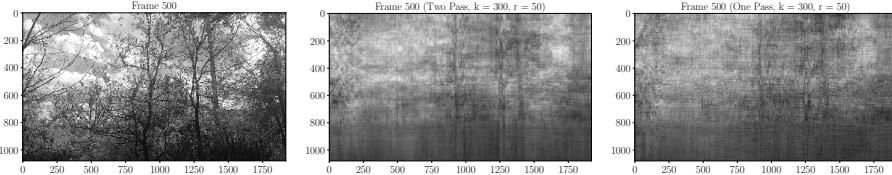


Fig. 7: *Visualizing Video Recovery*: Original frame (left); approximation by two-pass sketch (middle); approximation by one-pass sketch (right).

- 472 • *Climate data*. We consider global climate simulation datasets from the Com-
 473 munity Earth System Model (CESM) Community Atmosphere Model (CAM)
 474 5.0 [17, 18]. The dataset on aerosol absorption has four dimensions: times,
 475 altitudes, longitudes, and latitudes ($240 \times 30 \times 192 \times 288$). The data on net
 476 radiative flux at surface and dust aerosol burden have three dimensions: times,
 477 longitudes, and latitudes ($1200 \times 192 \times 288$). Each of these quantities has a
 478 strong impact on the absorption of solar radiation and on cloud formation.
 479 • *Combustion data*. We consider combustion simulation data from [22]. The
 480 data consists of three measured quantities — pressure, CO concentration, and
 481 temperature — each observed on a $1408 \times 128 \times 128$ spatial grid.
 482 • *Video data*. We use our streaming method to cluster frames of a video, as in

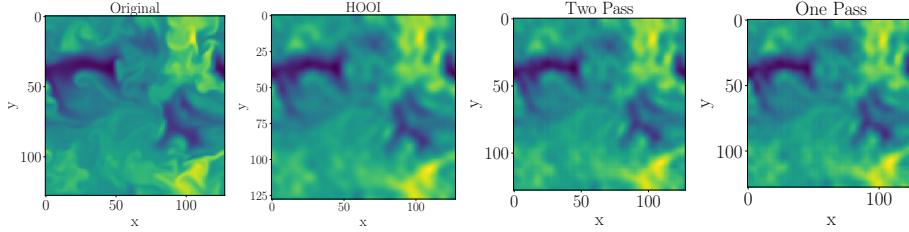


Fig. 8: *Visualizing Combustion Simulation*: All four figures show a slice of the temperature data along the first dimension. The approximation uses $\mathbf{r} = (281, 25, 25)$, $\mathbf{k} = (562, 50, 50)$, $\mathbf{s} = (1125, 101, 101)$, with the Gaussian TRP in the Tucker sketch.

[25]. Here, a low frame rate camera is mounted in a fixed position as people walk by. A 3D tensor is constructed with each video frames as a slice. The video consists of 2493 frames, each of size 1080 by 1980. As a tensor, stored as a `numpy.array`, the video data is 41.4 GB in total.

6.2.1. Data compression. We show that our proposed algorithms are able to successfully compress climate and combustion data even when the full data does not fit in memory. Since the Tucker rank of the original tensor is unknown, we perform experiments for three different target ranks. In this experiment, we hope to understand the effect of different choices of storage budget k to achieve the same compression ratio. We define the compression ratio as the ratio in size between the original input tensor and the output Tucker factors, i.e. $\frac{\prod_{i=1}^N I_i}{\sum_{i=1}^N r_i I_i + \prod_{i=1}^N r_i}$. As in our experiments on simulated data, 5 shows that the two-pass algorithm outperforms the one-pass algorithm as expected. However, as the storage budget k increases, both methods converge to the performance of HOOI. The rate of convergence is faster for smaller target ranks. Performance of our algorithms on the combustion simulation is qualitatively similar, but converges faster to the performance of HOOI. 8 visualizes the recovery of the temperature data in combustion simulation for a slice along the first dimension. We could observe that the recovery for both two-pass and one-pass algorithm approximate the recovery from HOOI. 13 in I shows similar results on another dataset.

6.2.2. Video scene classification. We show how to use our single pass method to classify scenes in the video data described above. The goal is to identify frames in which people appear. We remove the first 100 frames and last 193 frames where the camera setup happened, as in [25]. We stream over the tensor and sketch it using parameters $k = 300, s = 601$. Finally, we compute a fixed-rank approximation with $\mathbf{r} = (10, 10, 10)$ and $(20, 20, 20)$. We apply K-means clustering to the resulting 10 or 20 dimensional vectors corresponding to each of the remaining 2200 frames.

We experimented with clustering vectors found in three ways: from the two-pass or one-pass Tucker approximation, or directly from the factor sketch.

When matching the video frames with the classification result, we can see that the background light is relatively dark at the beginning, thus classified into **Class 0**. After a change in the background light, most other frames of the video are classified into **Class 1**. When a person passes by the camera, the frames are classified into **Class 2**. Right after the person passed by, the frames are classified into **Class 0**, the brighter background scene, due to the light adjustment.

518 Our classification results (using the linear sketch or approximation) are similar
519 to those in [25] while using only 1/500 as much storage; the one pass approximation
520 requires more storage (but still less than [25]) to achieve similar performance. In
521 particular, using the sketch itself, rather than the Tucker approximation, to summarize
522 the data enables very efficient video scene classification.

523 On the other hand, to reconstruct the original video frames we require much larger
524 \mathbf{k} and \mathbf{r} : the video is not very low rank along the spatial dimensions. 7 shows that
525 even with $\mathbf{s} = 601, 601, 601$, $\mathbf{k} = (300, 300, 300)$, $\mathbf{r} = (50, 50, 50)$, the recovered frame is
526 very noisy.

527 **Acknowledgments.** MU, YS, and YG were supported in part by DARPA Award
528 FA8750-17-2-0101. JAT gratefully acknowledges support from ONR Awards N00014-
529 11-10025, N00014-17-12146, and N00014-18-12363. The authors wish to thank Osman
530 Asif Malik and Stephen Becker for their help in understanding and implementing
531 Tucker TensorSketch, and Tamara Kolda for insightful comments on an early draft.

- [1] D. Achlioptas, *Database-friendly random projections: Johnson-Lindenstrauss with binary coins*, Journal of Computer and System Sciences, 66 (2003), pp. 671–687.
- [2] N. Ailon and B. Chazelle, *The fast Johnson-Lindenstrauss transform and approximate nearest neighbors*, SIAM Journal on Computing, 39 (2009), pp. 302–322.
- [3] S. Arora and B. Barak, *Computational complexity: a modern approach*, Cambridge University Press, 2009.
- [4] W. Austin, G. Ballard, and T. G. Kolda, *Parallel tensor compression for large-scale scientific data*, in Parallel and Distributed Processing Symposium, 2016 IEEE International, IEEE, 2016, pp. 912–922.
- [5] R. Ballester-Ripoll, P. Lindstrom, and R. Pajarola, *Tthresh: Tensor compression for multidimensional visual data*, IEEE transactions on visualization and computer graphics, (2019).
- [6] M. Baskaran, B. Meister, N. Vasilache, and R. Lethin, *Efficient and scalable computations with sparse tensors*, in High Performance Extreme Computing (HPEC), 2012 IEEE Conference on, IEEE, 2012, pp. 1–6.
- [7] C. Battaglino, G. Ballard, and T. G. Kolda, *A practical randomized cp tensor decomposition*, SIAM Journal on Matrix Analysis and Applications, 39 (2018), pp. 876–901.
- [8] C. Battaglino, G. Ballard, and T. G. Kolda, *Faster parallel tucker tensor decomposition using randomization*, (2019).
- [9] C. Boutsidis and A. Gittens, *Improved matrix algorithms via the subsampled randomized hadamard transform*, SIAM Journal on Matrix Analysis and Applications, 34 (2013), pp. 1301–1340.
- [10] A. Cichocki, *Tensor decompositions: a new concept in brain data analysis?*, arXiv preprint arXiv:1305.0395, (2013).
- [11] K. L. Clarkson and D. P. Woodruff, *Low-rank approximation and regression in input sparsity time*, Journal of the ACM (JACM), 63 (2017), p. 54.
- [12] G. Cormode and M. Hadjieleftheriou, *Finding frequent items in data streams*, Proceedings of the VLDB Endowment, 1 (2008), pp. 1530–1541.
- [13] L. De Lathauwer, B. De Moor, and J. VandeWalle, *A multilinear singular value decomposition*, SIAM journal on Matrix Analysis and Applications, 21 (2000), pp. 1253–1278.
- [14] V. De Silva and L.-H. Lim, *Tensor rank and the ill-posedness of the best low-rank approximation problem*, SIAM Journal on Matrix Analysis and Applications, 30 (2008), pp. 1084–1127.
- [15] H. Diao, Z. Song, W. Sun, and D. P. Woodruff, *Sketching for Kronecker Product Regression and P-splines*, arXiv e-prints, (2017), arXiv:1712.09473, p. arXiv:1712.09473, <https://arxiv.org/abs/1712.09473>.
- [16] N. Halko, P.-G. Martinsson, and J. A. Tropp, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM review, 53 (2011), pp. 217–288.
- [17] J. W. Hurrell, M. M. Holland, P. R. Gent, S. Ghosh, J. E. Kay, P. J. Kushner, J.-F. Lamarque, W. G. Large, D. Lawrence, K. Lindsay, et al., *The community earth system model: a framework for collaborative research*, Bulletin of the American Meteorological Society, 94 (2013), pp. 1339–1360.
- [18] J. Kay, C. Deser, A. Phillips, A. Mai, C. Hannay, G. Strand, J. Arblaster, S. Bates, G. Danabasoglu, J. Edwards, et al., *The community earth system model (cesm) large ensemble project: A community resource for studying climate change in the presence of internal climate variability*, Bulletin of the American Meteorological Society, 96 (2015), pp. 1333–1349.
- [19] O. Kaya and B. Uçar, *High performance parallel algorithms for the tucker decomposition of sparse tensors*, in Parallel Processing (ICPP), 2016 45th International Conference on, IEEE, 2016, pp. 103–112.
- [20] T. G. Kolda and B. W. Bader, *Tensor decompositions and applications*, SIAM review, 51 (2009), pp. 455–500.
- [21] T. G. Kolda and J. Sun, *Scalable tensor decompositions for multi-aspect data mining*, in 2008 Eighth IEEE International Conference on Data Mining, IEEE, 2008, pp. 363–372.
- [22] S. Lapointe, B. Savard, and G. Blanquart, *Differential diffusion effects, distributed burning, and local extinctions in high karlovitz premixed flames*, Combustion and flame, 162 (2015), pp. 3341–3355.
- [23] J. Li, C. Battaglino, I. Perros, J. Sun, and R. Vuduc, *An input-adaptive and in-place approach to dense tensor-times-matrix multiply*, in High Performance Computing, Networking, Storage and Analysis, 2015 SC-International Conference for, IEEE, 2015,

- 593 pp. 1–12.
- 594 [24] P. LI, T. J. HASTIE, AND K. W. CHURCH, *Very sparse random projections*, in Proceedings of
595 the 12th ACM SIGKDD international conference on Knowledge discovery and data mining,
596 ACM, 2006, pp. 287–296.
- 597 [25] O. A. MALIK AND S. BECKER, *Low-rank tucker decomposition of large tensors using tensors-*
598 *ketch*, in Advances in Neural Information Processing Systems, 2018, pp. 10116–10126.
- 599 [26] S. MUTHUKRISHNAN ET AL., *Data streams: Algorithms and applications*, Foundations and
600 Trends® in Theoretical Computer Science, 1 (2005), pp. 117–236.
- 601 [27] S. OYMAK AND J. A. TROPP, *Universality laws for randomized dimension reduction, with*
602 *applications*, Information and Inference: A Journal of the IMA, (2015).
- 603 [28] M. RUDELSON, *Row products of random matrices*, Advances in Mathematics, 231 (2012),
604 pp. 3199–3231.
- 605 [29] J. SUN, D. TAO, S. PAPADIMITRIOU, P. S. YU, AND C. FALOUTSOS, *Incremental tensor*
606 *analysis: Theory and applications*, ACM Transactions on Knowledge Discovery from Data
607 (TKDD), 2 (2008), p. 11.
- 608 [30] Y. SUN, Y. GUO, J. A. TROPP, AND M. UDELL, *Tensor random projection for low memory*
609 *dimension reduction*, in NeurIPS Workshop on Relational Representation Learning, 2018,
610 <https://r2learning.github.io/assets/papers/CameraReadySubmission%2041.pdf>.
- 611 [31] J. A. TROPP, *Improved analysis of the subsampled randomized hadamard transform*, Advances
612 in Adaptive Data Analysis, 3 (2011), pp. 115–126.
- 613 [32] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Practical sketching algorithms for*
614 *low-rank matrix approximation*, SIAM Journal on Matrix Analysis and Applications, 38
615 (2017), pp. 1454–1485.
- 616 [33] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *More practical sketching algorithms*
617 *for low-rank matrix approximation*, Tech. Report 2018-01, California Institute of Technology,
618 Pasadena, California, 2018.
- 619 [34] J. A. TROPP, A. YURTSEVER, M. UDELL, AND V. CEVHER, *Streaming low-rank matrix*
620 *approximation with an application to scientific simulation*, Submitted to SISC, (2019),
621 <https://arxiv.org/abs/1902.08651>.
- 622 [35] C. E. TSOURAKAKIS, *Mach: Fast randomized tensor decompositions*, in Proceedings of the 2010
623 SIAM International Conference on Data Mining, SIAM, 2010, pp. 689–700.
- 624 [36] L. R. TUCKER, *Some mathematical notes on three-mode factor analysis*, Psychometrika, 31
625 (1966), pp. 279–311.
- 626 [37] N. VANNIEUWENHOVEN, R. VANDEBRIL, AND K. MEERBERGEN, *A new truncation strategy for*
627 *the higher-order singular value decomposition*, SIAM Journal on Scientific Computing, 34
628 (2012), pp. A1027–A1052.
- 629 [38] M. A. O. VASILESCU AND D. TERZOPoulos, *Multilinear analysis of image ensembles: Tensor-*
630 *faces*, in European Conference on Computer Vision, Springer, 2002, pp. 447–460.
- 631 [39] Y. WANG, H.-Y. TUNG, A. J. SMOLA, AND A. ANANDKUMAR, *Fast and guaranteed tensor*
632 *decomposition via sketching*, in Advances in Neural Information Processing Systems, 2015,
633 pp. 991–999.
- 634 [40] D. P. WOODRUFF ET AL., *Sketching as a tool for numerical linear algebra*, Foundations and
635 Trends® in Theoretical Computer Science, 10 (2014), pp. 1–157.
- 636 [41] F. WOOLFE, E. LIBERTY, V. ROKHLIN, AND M. TYGERT, *A fast randomized algorithm for*
637 *the approximation of matrices*, Applied and Computational Harmonic Analysis, 25 (2008),
638 pp. 335–366.
- 639 [42] G. ZHOU, A. CICHOCKI, AND S. XIE, *Decomposition of big tensors with low multilinear rank*,
640 arXiv preprint arXiv:1412.1885, (2014).

641 **Appendix A. Proof of Main Results.**642 **A.1. Error bound for the two pass approximation Algorithm 4.2.**

643 *Proof of Theorem 5.1.* Suppose $\hat{\mathcal{X}}_2$ is the low-rank approximation from 4.2. Use
644 the definition of the mode- n product to see

$$645 \quad \begin{aligned} \hat{\mathcal{X}}_2 &= [\mathcal{X} \times_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top] \times_1 \mathbf{Q}_1 \times_1 \cdots \times_N \mathbf{Q}_N \\ &= \mathcal{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top. \end{aligned}$$

646 Although it seems that we sequentially project tensor \mathcal{X} to column space spanned
647 with \mathbf{Q}_n , but since mode product is exchangeable, in fact $\hat{\mathcal{X}}_2$ is the projection to space
648 $\{\mathcal{X} : \mathcal{X}^{(n)} \in \text{col}(\mathbf{Q}_n)\}$. This is a generalization of projection matrix where [14] has
649 a very detailed explanation and it is referred as multi-linear orthogonal projection.
650 Following exact techniques in Theorem 5.1 in [37] by sequentially applying Pythagorean
651 theory sequentially we can show that

$$652 \quad (\text{A.1}) \quad \|\hat{\mathcal{X}}_2 - \mathcal{X}\|_F^2 \leq \sum_{n=1}^N \|(\mathbf{I} - \mathbf{Q}_n \mathbf{Q}_n^\top) \mathcal{X}^{(n)}\|_F^2. \quad \square$$

653 Then taking expectation on \mathbf{Q}_n , and applying Lemma D.2 we complete the proof.

654 **A.2. Error bound for the one pass approximation Algorithm 4.3.**

655 *Proof of Theorem 5.2.* We show the approximation error can be decomposed as
656 the error due to the factor matrix approximations and the error due to the core
657 approximation. Let $\hat{\mathcal{X}}_1$ be the one pass approximation from 4.3, and let

$$658 \quad (\text{A.2}) \quad \hat{\mathcal{X}}_2 = \mathcal{X} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top,$$

659 be the two pass approximation from 4.2. The difference in one-pass and two-pass
660 approximation is in the core:

$$661 \quad \hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2 = (\mathcal{W} - \mathcal{X} \times_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_n^\top) \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N.$$

662 Thus $\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2$ is in the space defined above: $\{\mathcal{X} : \mathcal{X}^{(n)} \in \text{col}(\mathbf{Q}_n)\}$ while as pointed
663 before $\hat{\mathcal{X}}_2 - \mathcal{X}$ is perpendicular to that space. Therefore,

$$664 \quad (\text{A.3}) \quad \langle \hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2, \hat{\mathcal{X}}_2 - \mathcal{X} \rangle = 0.$$

665 Now we use the (expectation of) the Pythagorean theorem to bound the expected
666 error of the one pass approximation:

$$667 \quad (\text{A.4}) \quad \mathbb{E} \|\hat{\mathcal{X}}_1 - \mathcal{X}\|_F^2 = \mathbb{E} \|\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2\|_F^2 + \mathbb{E} \|\hat{\mathcal{X}}_2 - \mathcal{X}\|_F^2.$$

668 Consider the first term which is due to core approximation. Based in the definition
669 of $\hat{\mathcal{X}}_1$ and $\hat{\mathcal{X}}_2$ we can see that

$$670 \quad \begin{aligned} \|\hat{\mathcal{X}}_1 - \hat{\mathcal{X}}_2\|_F^2 &= \|(\mathcal{W}_1 - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top) \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N\|_F^2 \\ &= \|(\mathcal{W}_1 - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top)\|_F^2, \end{aligned}$$

673 where we use the invariance of the Frobenius norm under orthonormal transformations to get the second line. Now using B.2 to bound for the error due to the core
 674 approximations as

$$676 \quad \mathbb{E}\|\hat{\mathbf{X}}_1 - \hat{\mathbf{X}}_2\|_F^2 \leq \Delta \left[\sum_{n=1}^N \left(1 + \frac{\rho_n}{k_n - \rho_n - 1} \right) (\tau_{\rho_n}^{(n)})^2 \right].$$

677 Finally, as shown in proof for 5.1 to bound the error due to the factor matrix
 678 approximations (the second term in (A.4)) as

$$679 \quad \mathbb{E}\|\hat{\mathbf{X}}_2 - \mathbf{X}\|_F^2 \leq \left[\sum_{n=1}^N \left(1 + \frac{\rho_n}{k_n - \rho_n - 1} \right) (\tau_{\rho_n}^{(n)})^2 \right].$$

680 Summing these two bounds finishes the proof. \square

681 A.3. Error bound for the fixed rank approximation Algorithm 4.4.

682 *Proof of Theorem 5.3.* Our argument follows the proof of [32, Proposition 6.1]:

$$\begin{aligned} & \|\mathbf{X} - [\hat{\mathbf{X}}]_{\mathbf{r}}\|_F \leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\hat{\mathbf{X}} - [\hat{\mathbf{X}}]_{\mathbf{r}}\|_F \\ & \leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\hat{\mathbf{X}} - [\mathbf{X}]_{\mathbf{r}}\|_F \\ & \leq \|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\hat{\mathbf{X}} - \mathbf{X} + \mathbf{X} - [\mathbf{X}]_{\mathbf{r}}\|_F \\ & \leq 2\|\mathbf{X} - \hat{\mathbf{X}}\|_F + \|\mathbf{X} - [\mathbf{X}]_{\mathbf{r}}\|_F. \end{aligned}$$

684 The first and the third line are the triangle inequality, and the second line follows from
 685 the definition of the best rank- r approximation. Take the expectation of $\|\mathbf{X} - \hat{\mathbf{X}}\|_F$
 686 and use Jensen's inequality $\mathbb{E}\|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq \sqrt{\mathbb{E}\|\mathbf{X} - \hat{\mathbf{X}}\|_F^2}$ to finish the proof. \square

687 **Appendix B. Probabilistic Analysis of Core Sketch Error.** This section
 688 contains the most technical part of our proof. We provide a probabilistic error bound
 689 for the difference between the two pass core approximation \mathcal{W}_2 from 4.2 and the one
 690 pass core approximation \mathcal{W}_1 from 4.3.

691 Introduce for each $n \in [N]$ the orthonormal matrix \mathbf{Q}_n^\perp that forms a basis for the
 692 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

$$693 \quad (B.1) \quad \Phi_n^Q = \Phi_n^\top \mathbf{Q}_n, \quad \Phi_n^{Q^\perp} = \Phi_n^\top \mathbf{Q}_n^\perp.$$

694 Recall that the DRMs Φ_n are i.i.d. Gaussian. Hence conditional on \mathbf{Q}_n , Φ_n^Q and $\Phi_n^{Q^\perp}$
 695 are independent.

696 **B.1. Decomposition of Core Approximation Error.** In this section, we
 697 characterize the difference between the one and two pass core approximations $\mathcal{W}_1 -$
 698 $\mathcal{W}_2 = \mathcal{W}_1 - \mathbf{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top$.

699 **LEMMA B.1.** *Suppose that Φ_n has full column rank for each $n \in [N]$. Then*

$$700 \quad \mathcal{W}_1 - \mathcal{W}_2 = \mathcal{W}_1 - \mathbf{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top = \sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum_{j=1}^N i_j \geq 1} \mathcal{Y}_{i_1 \dots i_N},$$

701 where

$$\begin{aligned} 702 \quad (B.2) \quad \mathcal{Y}_{i_1 \dots i_N} &= \mathbf{X} \times_1 \left(\mathbb{1}_{i_1=0} \mathbf{Q}_1^\top + \mathbb{1}_{i_1=1} (\Phi_1^Q)^\dagger \Phi_1^{Q^\perp} (\mathbf{Q}_1^\perp)^\top \right) \\ &\quad \times_2 \cdots \times_N \left(\mathbb{1}_{i_N=0} \mathbf{Q}_N^\top + \mathbb{1}_{i_N=1} (\Phi_N^Q)^\dagger \Phi_N^{Q^\perp} (\mathbf{Q}_N^\perp)^\top \right). \end{aligned}$$

703 *Proof.* Let \mathcal{H} be the core sketch from 4.1. Write \mathcal{W}_1 as

$$\begin{aligned} \mathcal{W}_1 &= \mathcal{H} \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \\ 704 &= (\mathbf{X} - \hat{\mathbf{X}}_2) \times_1 \Phi_1^\top \times_2 \cdots \times_N \Phi_N^\top \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \\ &\quad + \hat{\mathbf{X}}_2 \times_1 \Phi_1^\top \times_2 \cdots \times_N \Phi_N^\top \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger. \end{aligned}$$

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

$$\begin{aligned} \hat{\mathbf{X}} \times_1 \Phi_1^\top \times_2 \cdots \times_N \Phi_N^\top \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \\ 706 &= \mathbf{X} \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \Phi_1^\top \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \Phi_N^\top \mathbf{Q}_N \mathbf{Q}_N^\top \\ &= \mathbf{X} \times_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N^\top, \end{aligned}$$

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

$$708 \quad \mathcal{W}_1 - \mathcal{W}_2 = (\mathbf{X} - \hat{\mathbf{X}}) \times_1 \Phi_1^\top \times_2 \cdots \times_N \Phi_N^\top \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger.$$

709 We continue to simplify this difference:

$$\begin{aligned} 710 \quad (B.3) \quad &(\mathbf{X} - \hat{\mathbf{X}}) \times_1 \Phi_1^\top \times_2 \cdots \times_N \Phi_N^\top \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \\ &= (\mathbf{X} - \hat{\mathbf{X}}) \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \Phi_1^\top \times_2 \cdots \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \Phi_N^\top \\ &= (\mathbf{X} - \hat{\mathbf{X}}) \times_1 (\Phi_1^\top \mathbf{Q}_1)^\dagger \Phi_1^\top (\mathbf{Q}_1 \mathbf{Q}_1^\top + \mathbf{Q}_1^\perp (\mathbf{Q}_1^\perp)^\top) \cdots \\ &\quad \times_N (\Phi_N^\top \mathbf{Q}_N)^\dagger \Phi_N^\top (\mathbf{Q}_N \mathbf{Q}_N^\top + \mathbf{Q}_N^\perp (\mathbf{Q}_N^\perp)^\top) \\ &= (\mathbf{X} - \hat{\mathbf{X}}) \times_1 (\mathbf{Q}_1^\top + (\Phi_1^Q)^\dagger \Phi_1^{Q^\perp} (\mathbf{Q}_1^\perp)^\top) \times_2 \cdots \\ &\quad \times_N (\mathbf{Q}_N^\top + (\Phi_N^Q)^\dagger \Phi_N^{Q^\perp} (\mathbf{Q}_N^\perp)^\top). \end{aligned}$$

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

$$712 \quad 1. (\mathbf{X} - \hat{\mathbf{X}}) \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top = 0.$$

$$713 \quad 2. \text{ For each } n \in [N], \tilde{\mathbf{X}} \times_n (\Phi_n^Q)^\dagger \Phi_n^{Q^\perp} (\mathbf{Q}_n^\perp)^\top = 0.$$

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

$$718 \quad \sum_{(i_1, \dots, i_N) \in \{0,1\}^N, \sum n=1^N i_n \neq 0} \mathbf{Y}_{i_1 \dots i_N}.$$

719 **B.2. Probabilistic Core Error Bound.** In this section, we derive a probabilistic
720 error bound based on the core error decomposition from Lemma B.1.

721 LEMMA B.2. *Sketch the tensor \mathbf{X} using a Tucker sketch with parameters \mathbf{k} and*
722 *$\mathbf{s} > 2\mathbf{k}$ with i.i.d. Gaussian $\mathcal{N}(0, 1)$ DRMs. Define $\Delta = \max_{n=1}^N \frac{k_n}{s_n - k_n - 1}$. Then for*
723 *any natural numbers $1 \leq \rho < \mathbf{k} - 1$,*

$$724 \quad \mathbb{E} \|\mathcal{W}_1 - \mathbf{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top\|_F^2 \leq \Delta \left[\sum_{n=1}^N \left(1 + \frac{\rho_n}{k_n - \rho_n - 1} \right) (\tau_{\rho_n}^{(n)})^2 \right].$$

725 *Proof.* It suffices to show

$$726 \quad (B.4) \quad \mathbb{E} [\|\mathcal{W}_1 - \mathbf{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top\|_F^2 | \Omega_1, \dots, \Omega_N] \leq \Delta \|\mathbf{X} - \hat{\mathbf{X}}_2\|_F^2.$$

727 Then take the expectation with respect to $\Omega_1, \dots, \Omega_N$ and apply results in 5.1 to
 728 bound $\|\mathcal{X} - \hat{\mathcal{X}}_2\|_F^2$ to finish the proof. To show B.4, we will use the fact that the core
 729 DRMs $\{\Omega_n\}_{n \in [N]}$ are independent of the factor matrix DRMs $\{\Phi_n\}_{n \in [N]}$, and that
 730 the randomness in each factor matrix approximation \mathbf{Q}_n comes solely from Ω_n .
 731 For $i \in \{0, 1\}^N$, define $\mathcal{B}_{i_1 \dots i_N}$ =

$$732 \quad \mathcal{X} \times_1 (\mathbb{1}_{i_1=0} \mathbf{Q}_1 \mathbf{Q}_1^\top + \mathbb{1}_{i_1=1} \mathbf{Q}_1^\perp (\mathbf{Q}_1^\perp)^\top) \cdots \times_N (\mathbb{1}_{i_N=0} \mathbf{Q}_N \mathbf{Q}_N^\top + \mathbb{1}_{i_N=1} \mathbf{Q}_N^\perp (\mathbf{Q}_N^\perp)^\top).$$

733 B.1 decomposes the core error as the sum of $\mathcal{Y}_{i_1 \dots i_N}$ where $\sum_{n=1}^N i_n \geq 1$. Applying D.1
 734 and using the orthogonal invariance of the Frobenius norm, we observe

$$735 \quad \mathbb{E} [\|\mathcal{Y}_{i_1 \dots i_N}\|_F^2 \mid \Omega_1 \dots \Omega_N] = \left(\prod_{n=1}^N \Delta_n^{i_n} \right) \|\mathcal{B}_{i_1 \dots i_N}\|_F^2 \leq \Delta \|\mathcal{B}_{i_1 \dots i_N}\|_F^2$$

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

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

$$740 \quad (\text{B.5}) \quad \langle \mathcal{B}_{\mathbf{q}_1}, \mathcal{B}_{\mathbf{q}_2} \rangle = \langle \dots \mathbf{Q}_r^\top \mathbf{Q}_r^\perp \dots \rangle = 0.$$

741 Similarly we can show that the inner product between $\mathcal{Y}_{\mathbf{q}_1}$ and $\mathcal{Y}_{\mathbf{q}_2}$ is zero with different
 742 $\mathbf{q}_1, \mathbf{q}_2$. Noticing that $\mathcal{B}_{0, \dots, 0} = \hat{\mathcal{X}}_2$, we have

$$743 \quad \|\mathcal{X} - \hat{\mathcal{X}}_2\|_F^2 = \left\| \sum_{(i_1, \dots, i_N) \in \{0, 1\}^N, \sum_{n=1}^N i_n \geq 1} \mathcal{B}_{i_1 \dots i_N} \right\|_F^2 = \sum_{\substack{(i_1, \dots, i_N) \in \{0, 1\}^N, \\ \sum_{n=1}^N i_n \geq 1}} \|\mathcal{B}_{i_1 \dots i_N}\|_F^2.$$

744 Putting all these together and using the Pythagorean theorem, to show B.4:

$$\begin{aligned} 745 \quad & \mathbb{E} [\|\mathcal{W} - \mathcal{X} \times_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N^\top\|_F^2 \mid \Omega_1, \dots, \Omega_N] \\ &= \sum_{(i_1, \dots, i_N) \in \{0, 1\}^N, \sum_{n=1}^N i_n \geq 1} \mathbb{E} [\|\mathcal{Y}_{i_1 \dots i_N}\|_F^2 \mid \Omega_1, \dots, \Omega_N] \\ 746 \quad &\leq \Delta \left(\sum_{(i_1, \dots, i_N) \in \{0, 1\}^N, \sum_{n=1}^N i_n \geq 1} \|\mathcal{B}_{i_1 \dots i_N}\|_F^2 \right) = \Delta \|\mathcal{X} - \hat{\mathcal{X}}_2\|_F^2. \end{aligned}$$

747 Appendix C. Proof of fixed rank approximation lemma.

748 *Proof of 4.1.* The target tensor to be approximated is $\mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N$ is
 749 apparently in the space $\{\mathcal{X} : \mathcal{X}^{(n)} \in \text{col}(\mathbf{Q}_n)\}$. For any approximation $\hat{\mathcal{X}}$, we can
 750 project it into this space as

$$751 \quad \hat{\mathcal{X}} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \times_2 \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top$$

752 and by Pythagorean theory,

$$753 \quad (\text{C.1}) \quad \begin{aligned} & \|\hat{\mathcal{X}} - \mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N\|_F \leq \|\hat{\mathcal{X}} - \hat{\mathcal{X}} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top\|_F^2 \\ &+ \|\hat{\mathcal{X}} \times_1 \mathbf{Q}_1 \mathbf{Q}_1^\top \cdots \times_N \mathbf{Q}_N \mathbf{Q}_N^\top - \mathcal{W} \times_1 \mathbf{Q}_1 \cdots \times_N \mathbf{Q}_N\|_F^2, \end{aligned}$$

754 which indicates that the optimal Tucker decomposition resides in the space $\{\mathcal{X} : \mathcal{X}^{(n)} \in$
 755 $\text{col}(\mathbf{Q}_n)\}$. Suppose $[\mathcal{W}; \mathbf{V}_1, \dots, \mathbf{V}_N]$ is the optimal solution to the problem, since its
 756 unfolding is in the space spanned by \mathbf{Q}_n , each \mathbf{V}_n can be written as $\mathbf{Q}_n \mathbf{U}_n$ for some
 757 orthogonal matrix $\mathbf{U}_n \in \mathbb{R}^{k_n \times r_n}$. Then, noticing orthogonal transformation does not
 758 change Frobenius norm,

$$\begin{aligned} & \|\mathcal{W} \times_1 \mathbf{Q}_1 \times \cdots \times_N \mathbf{Q}_N - \mathcal{G} \times_1 \mathbf{Q}_1 \mathbf{U}_1 \times \cdots \times_N \mathbf{Q}_N \mathbf{U}_N\|_F \\ (C.2) \quad & = \|\mathcal{W} - \mathcal{G}\|_F \geq \|\mathcal{W} - [\mathcal{W}]_{\mathbf{r}}\|_F \\ & = \|\mathcal{W} \times_1 \mathbf{Q}_1 \times \cdots \times_N \mathbf{Q}_N - [\mathcal{W}]_{\mathbf{r}} \times_1 \mathbf{Q}_1 \times \cdots \times_N \mathbf{Q}_N\|_F. \end{aligned}$$

760 This finishes the proof. \square

761 Appendix D. Technical Lemmas.

762 **D.1. Random projections of matrices.** Proofs for lemmas in this section
 763 appear in [16, chapters 9 and 10].

764 LEMMA D.1. *Assume that $t > q$. Let $\mathbf{G}_1 \in \mathbb{R}^{t \times q}$ and $\mathbf{G}_2 \in \mathbb{R}^{t \times p}$ be independent
 765 standard normal matrices. For any matrix \mathbf{B} with conforming dimensions,*

$$766 \quad \mathbb{E}\|\mathbf{G}_1^\dagger \mathbf{G}_2 \mathbf{B}\|_F^2 = \frac{q}{t-q-1} \|\mathbf{B}\|_F^2.$$

767 LEMMA D.2. *Suppose that \mathbf{A} is a real $m \times n$ matrix with singular value $\sigma_1 \geq \sigma_2 \geq \dots$, choose a target rank $k \geq 2$ and an oversampling parameter $p \geq 2$, where
 768 $k+p \leq \min\{m, n\}$. Draw an $n \times (k+p)$ standard Gaussian matrix Ω , and construct
 769 the sample matrix $\mathbf{Y} = \mathbf{A}\Omega$, then the expectation of approximation error is*

$$771 \quad \mathbb{E}\|(\mathbf{I} - \mathbf{P}_{\mathbf{Y}})\mathbf{A}\|_F^2 \leq \left(1 + \frac{k}{p-1}\right) \left(\sum_{j>k} \sigma_j^2\right).$$

772 **Appendix E. More Algorithms.** This section provides detailed implemen-
 773 tations for a linear sketch appropriate to a streaming setting (Algorithm E.1) or a
 distributed setting (Algorithm E.2).

Algorithm E.1 Linear Update to Sketches

```

1: function SKETCHLINEARUPDATE( $\mathcal{F}, \mathbf{V}_1, \dots, \mathbf{V}_N, \mathcal{H}; \theta_1, \theta_2$ )
2:   for  $n = 1, \dots, N$  do
3:      $\mathbf{V}_n \leftarrow \theta_1 \mathbf{V}_n + \theta_2 \mathbf{F}^{(n)} \Omega_n$ 
4:   end for
5:    $\mathcal{H} \leftarrow \theta_1 \mathcal{H} + \theta_2 \mathcal{F} \times_1 \Phi_1 \times \cdots \times_N \Phi_N$ 
6:   return  $(\mathbf{V}_1, \dots, \mathbf{V}_N, \mathcal{H})$ 
7: end function

```

774

775 Appendix F. Scrambled Subsampled Randomized Fourier Transform.

776 In order to reduce the cost of storing the test matrices, in particular, $\Omega_1, \dots, \Omega_N$,
 777 we can use the Scrambled Subsampled Randomized Fourier Transform (SSRFT). To
 778 reduce the dimension of a matrix, $\mathbf{X} \in \mathbb{R}^{m \times n}$, along either the row or the column to
 779 size k , we define the SSRFT map Ξ as:

$$780 \quad \Xi = \begin{cases} \mathbf{R}\mathbf{F}^\top \mathbf{I}\mathbf{F}\mathbf{H}^\top \in \mathbb{F}^{k \times m} & (\text{Row linear transform}) \\ (\bar{\mathbf{R}}\bar{\mathbf{F}}^\top \bar{\mathbf{I}}\bar{\mathbf{F}}\bar{\mathbf{H}}^\top)^\top \in \mathbb{F}^{n \times k} & (\text{Column linear transform}), \end{cases}$$

Algorithm E.2 Sketching in Distributed Setting

Require: \mathcal{X}_i is the part of the tensor \mathcal{X} at local machine i and $\mathcal{X} = \sum_{i=1}^m \mathcal{X}_i$.

- 1: **function** COMPUTESKETCHDISTRIBUTED($\mathcal{X}_1, \dots, \mathcal{X}_m$)
- 2: Send the same random generating environment to every local machine.
- 3: Generate the same DRM at each local machine.
- 4: **for** $i = 1 \dots m$ **do**
- 5: $(\mathbf{V}_1^{(i)}, \dots, \mathbf{V}_n^{(i)}, \mathcal{H}^{(i)}) \leftarrow \text{ComputeSketch}(\mathcal{X}_i)$
- 6: **end for**
- 7: **for** $j = 1 \dots n$ **do**
- 8: $\mathbf{V}_j \leftarrow \sum_{i=1}^m \mathbf{V}_j^{(i)}$
- 9: **end for**
- 10: $\mathcal{H} \leftarrow \sum_{i=1}^m \mathcal{H}^{(i)}$
- 11: **return** $(\mathbf{V}_1, \dots, \mathbf{V}_n, \mathcal{H})$
- 12: **end function**

781 where $\Pi, \Pi' \in \mathbb{R}^{m \times m}, \bar{\Pi}, \bar{\Pi}' \in \mathbb{R}^{n \times n}$ are signed permutation matrices. That is, the
 782 matrix has exactly one non-zero entry, 1 or -1 with equal probability, in each row and
 783 column. $\mathbf{F} \in \mathbb{F}^{m \times m}, \mathbf{F} \in \mathbb{F}^{n \times n}$ denote the discrete cosine transform ($\mathbb{F} = \mathbb{R}$) or the
 784 discrete fourier transform ($\mathbb{F} = \mathbb{C}$). The matrix $\mathbf{R}, \bar{\mathbf{R}}$ is the restriction to k coordinates
 785 chosen uniformly at random.

786 In practice, we implement the SSRFT as in [Algorithm F.1](#). It takes only $\mathcal{O}(m)$ or
 787 $\mathcal{O}(n)$ bits to store Ξ , compared to $\mathcal{O}(km)$ or $\mathcal{O}(kn)$ for Gaussian or uniform random
 788 map. The cost of applying Ξ to a vector is $\mathcal{O}(n \log n)$ or $\mathcal{O}(m \log m)$ arithmetic
 789 operations for fast Fourier transform and $\mathcal{O}(n \log k)$ or $\mathcal{O}(m \log k)$ for fast cosine
 790 transform. Though in practice, SSRFT behaves similarly to the Gaussian random
 791 map, its analysis is less comprehensive [9, 31, 2] than the Gaussian case.

Algorithm F.1 Scrambled Subsampled Randomized Fourier Transform (Row Linear Transform)

Require: $\mathbf{X} \in \mathbb{R}^{m \times n}, \mathcal{F} = \mathbb{R}$, **randperm** creates a random permutaion vector, and
randsign creates a random sign vector. **dct** denotes the discrete cosine transform.

- 1: **function** SSRFT(\mathbf{X})
- 2: $\mathbf{coords} \leftarrow \text{randperm}(m, k)$
- 3: $\mathbf{perm}_j \leftarrow \text{randperm}(m)$ for $j = 1, 2$
- 4: $\mathbf{sgn}_j \leftarrow \text{randsign}(m)$ for $j = 1, 2$
- 5: $\mathbf{X} \leftarrow \text{dct}(\mathbf{sgn}_1 \cdot \mathbf{X}[\mathbf{perm}_1, :])$ \triangleright elementwise product
- 6: $\mathbf{X} \leftarrow \text{dct}(\mathbf{sgn}_2 \cdot \mathbf{X}[\mathbf{perm}_2, :])$
- 7: **return** $\mathbf{X}[\mathbf{coords}, :]$
- 8: **end function**

792 **Appendix G. TensorSketch.** Many authors have developed methods to
 793 perform dimension reduction efficiently. In particular [15] proposed a method called
 794 tensor sketching aiming to solve least square problem with design matrix has kroneck
 795 product structure. [25] applied this technique to their one pass Tucker decomposition.
 796 Here we review the definition of tensor sketch and how it be applied in [25].

797 *CountSketch.* [12] proposed the CountSketch method. A comprehensive theoretical
 798 analysis in the context of low-rank approximation problems appears in [11]. To

799 compute the sketch $\mathbf{X}\Omega \in \mathbb{R}^{d \times k}$ for $\mathbf{X} \in \mathbb{R}^{m \times d}$, CountSketch defines $\Omega = \mathbf{D}\Phi$, where
 800 1. $\mathbf{D} \in \mathbb{R}^{d \times d}$ is a diagonal matrix with each diagonal entry equal to $(-1, 1)$ with
 801 probability $(1/2, 1/2)$.
 802 2. $\Phi \in \mathbb{R}^{d \times k}$ is the matrix form of a Hashing function.

803 In total, these two matrices have $2d$ non-zero entries in total, thus requiring much
 804 less storage than the standard kd entries. Furthermore, these two matrices can act as
 805 an operator on each column of \mathbf{X} and require only $\mathcal{O}(kd)$ operations.

806 *TensorSketch*. [25] proposes to use the countsketch inside the HOOI method
 807 for Tucker decomposition. They apply sketching method solve least square problem
 808 appearing in (2.4) and (2.5) in Algorithm 2.2. They use J_1, J_2 to denote the reduced
 809 dimension. Using a standard random map, it will need J_1 -by- $I_{(-n)}$ random matrix for
 810 (2.4) and a J_2 -by- $\prod_{n=1}^N I_n$ random matrix to compute (2.5).

811 But as shown in [25], these two stages can be expressed as

812 (G.1) For $n = 1, \dots, N$, update $\mathbf{U}^{(n)} = \arg \min_{\mathbf{U} \in \mathbb{R}^{I_n \times R_n}} \left\| \left(\bigotimes_{\substack{i=N \\ i \neq n}}^1 \mathbf{U}^{(i)} \right) \mathbf{G}_{(n)}^\top \mathbf{U}^\top - \mathbf{Y}_{(n)}^\top \right\|_F^2.$

813 (G.2) Update $\mathcal{G} = \arg \min_{\mathbf{Z} \in \mathbb{R}^{R_1 \times \dots \times R_N}} \left\| \left(\bigotimes_{i=N}^1 \mathbf{U}^{(i)} \right) \text{vec } \mathbf{Z} - \text{vec } \mathbf{Y} \right\|_2^2,$

814 where \mathbf{Y} is the original data. $\forall i \in [n]$, \mathbf{U}_i is the factor matrix, and \mathcal{G} is the core tensor.
 815 R_1, \dots, R_N denote the rank of the data.

816 As what shown in [12], [25] proposes to apply tensorSketch to the Kronecker
 817 product structure of the input matrix in the sketch construction, i.e. $\bigotimes_{\substack{i=1 \\ i \neq n}}^N \mathbf{U}_i$ in (G.1)
 818 and $\bigotimes_{i=1}^N \mathbf{U}_i$ in (G.2). TensorSketch method combines the CountSketch of each factor
 819 matrix via the Khatri-Rao product and Fast Fourier Transform. Consider sketching
 820 $\bigotimes_{i=1}^N \mathbf{U}_i$ in (G.2). TensorSketch is defined as

821 (G.3) $\Omega \mathbf{X} = \text{FFT}^{-1} \left(\odot_{n=1}^N \left(\text{FFT} \left(\text{CountSketch}^{(n)}(\mathbf{U}^{(n)}) \right)^\top \right)^\top \right)$

822 By only storing $\text{CountSketch}^{(1)}, \dots, \text{CountSketch}^{(N)}$, TensorSketch only requires $2 \sum_{i=1}^N I_n$
 823 storage. Therefore, the storage cost of the sketch is dominated by the sketch size,
 824 $NR^{n-1}J_1 + J_2R^n \approx NKR^{2n-2} + KR^{2n}$, when $J_1 = KR^{n-1}, J_2 = KR^n$.

825 Appendix H. Time and Storage Complexity.

826 **H.1. Comparison Between Algorithm 4.4 and T.-TS [25].** Here we com-
 827 pare the time and storage complexity of the two extant methods for streaming Tucker
 828 approximation: our one-pass method, and T.-TS [25].

829 To compare the storage and time costs of both T.-TS and the one-pass algorithm, we
 830 separate the cost into two parts: one for forming the sketch, the other for each iteration
 831 of ALS. Assume the tensor to approximate has equal side lengths $I_1 = \dots = I_N = I$
 832 and that the target rank for each mode is R .

833 The suggested default parameters for the sketch in [25] are $J_1 = 10R^{N-1}$ and
 834 $J_2 = 10R^N$. Our suggested default parameters are $k = 2r, s = 2k + 1$. Under the
 835 choice of the default parameter, we compare the the cost of storage and time in Table 3
 836 and Table 4. In most problems with data not perfectly low Tucker rank, i.e. $R > 4$, the

suggested default setting of T.-TS typically leads to a higher storage cost. Moreover, our algorithm uses less storage and is faster to compute, particularly for tensors with many modes N .

However, the evaluation of the two algorithms should not be solely based on their default setups. If the memory constraint is set to be the same, our one-pass algorithm performs much better in the low-memory case, but slightly worse in the case with very high-memory as in Figure 4. The memory of our suggested setting typically implies a much smaller memory usage than their suggested setting.

H.2. Computational Complexity of Algorithm 4.4. Here, we will derive a fine-grained computation complexity for our one pass fixed-rank approximation algorithm.

In the sketching stage of the streaming algorithm, we need to first compute the factor sketches, $\mathbf{G}_n = \mathbf{X}\Omega_n, n \in [N]$ with $kN\bar{I}$ flops in total. Then, we need to compute the core tensor sketch \mathcal{Z} by recursively multiplying \mathcal{X} by $\Phi_n, n \in [N]$. We can find the upper bound for the number of flops to be $\frac{s(1-\delta_1^N)}{1-\delta_1}\bar{I}$. Then, in the approximation stage, we first perform "economy size" QR factorizations on $\mathbf{G}_1, \dots, \mathbf{G}_N$ with $\mathcal{O}(k^2(\sum_{n=1}^N I_n))$ to find the orthonormal bases $\mathbf{Q}_1, \dots, \mathbf{Q}_N$. To find the linkage tensor \mathcal{W} , we need to recursively solve linear square problems, with $\frac{k^2s^N(1-(k/s)^N)}{1-k/s}$ flops. Overall, the sketch computation dominates the total time complexity.

The higher order SVD directly acts on \mathcal{X} by first computing the SVD for each unfolding in $\mathcal{O}(kN\bar{I})$, and then multiplying \mathcal{X} by $\mathbf{U}_1^\top, \dots, \mathbf{U}_N^\top$ in $\mathcal{O}(\frac{k(1-\delta_1^N)\bar{I}}{1-\delta_1})$. The total time cost is less than the streaming algorithm with a constant factor. Note: we can use the randomized SVD in the first step to improve the computational cost to $\bar{I}N \log k + \sum_{n=1}^N (I_n + I_{(-n)})k^2$ [16].

Algorithm	Storage Cost ($I = o(r^{2N})$)	
T.-TS	Sketching	$\mathcal{O}(r^{2N})$
	Recovery	$\mathcal{O}(r^{2N})$
Algorithm 4.3 (One Pass)	Sketching	$\mathcal{O}(4^N r^N)$
	Recovery	$\mathcal{O}(4^N r^N)$

Table 3: Storage complexity of Algorithm 4.3 and T.-TS on tensor $\mathcal{X} \in \mathbb{R}^{I \times \dots \times I}$. Algorithm 4.3 uses parameters $(k, s) = (2r, 4r + 1)$ and uses a TRP composed of Gaussian DRMs inside the Tucker sketch. T.-TS uses default values for hyperparameters: $J_1 = 10r^{N-1}, J_2 = 10r^N$.

861 Appendix I. More Numerics.

862 This section provides more numerical results on simulated datasets in Figure 9,
863 Figure 10, Figure 11, and Figure 12.

864 We also provide more numerical results on real datasets in Figure 13.

Algorithm	Time Cost ($I = o(r^{2N})$)	
T.-TS	Sketching	$\mathcal{O}(N \text{nnz}(\mathcal{X}))$
	Recovery	$\mathcal{O}(NIr^N + Nr^{2N-1} + r^{2N})$
Algorithm 4.3 (One Pass)	Sketching	$\mathcal{O}(Nr \text{nnz}(\mathcal{X}))$
	Recovery	$\mathcal{O}(Nr^{N+1})$

Table 4: Time complexity of Algorithm 4.3 and T.-TS on tensor $\mathcal{X} \in \mathbb{R}^{I \times \dots \times I}$. Algorithm 4.3 uses parameters $(k, s) = (2r, 4r + 1)$ and uses a TRP composed of Gaussian DRMs inside the Tucker sketch. T.-TS uses default values for hyper-parameters: $J_1 = 10r^{N-1}$, $J_2 = 10r^N$.

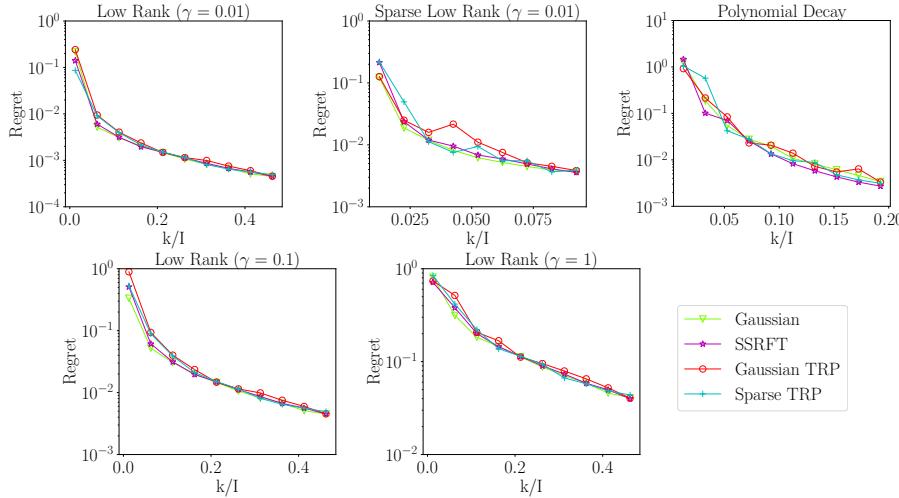


Fig. 9: We approximate 3D synthetic tensors (see subsection 6.1) with $I = 400$, using our one-pass algorithm with $r = 5$ and varying k ($s = 2k + 1$), using a variety of DRMs in the Tucker sketch: Gaussian, SSRFT, Gaussian TRP, or Sparse TRP.

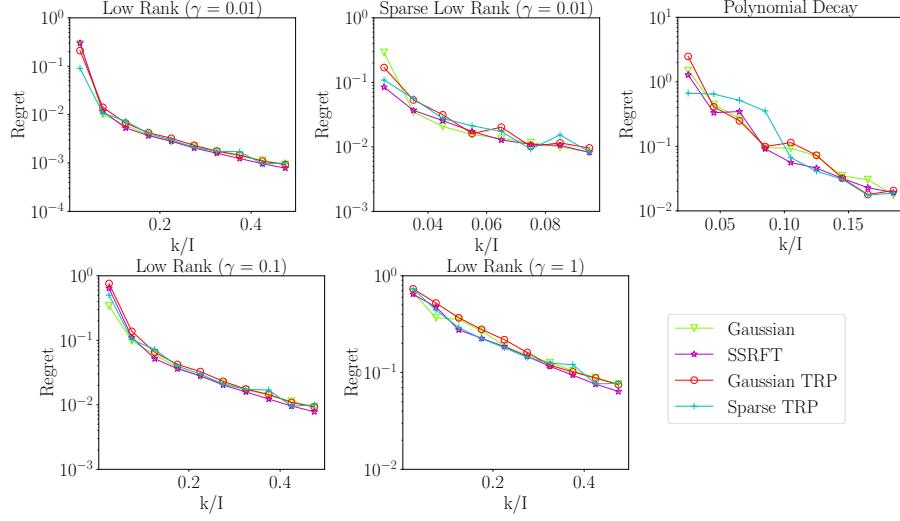


Fig. 10: We approximate 3D synthetic tensors (see subsection 6.1) with $I = 200$, using our one-pass algorithm with $r = 5$ and varying k ($s = 2k + 1$), using a variety of DRMs in the Tucker sketch: Gaussian, SSRFT, Gaussian TRP, or Sparse TRP.

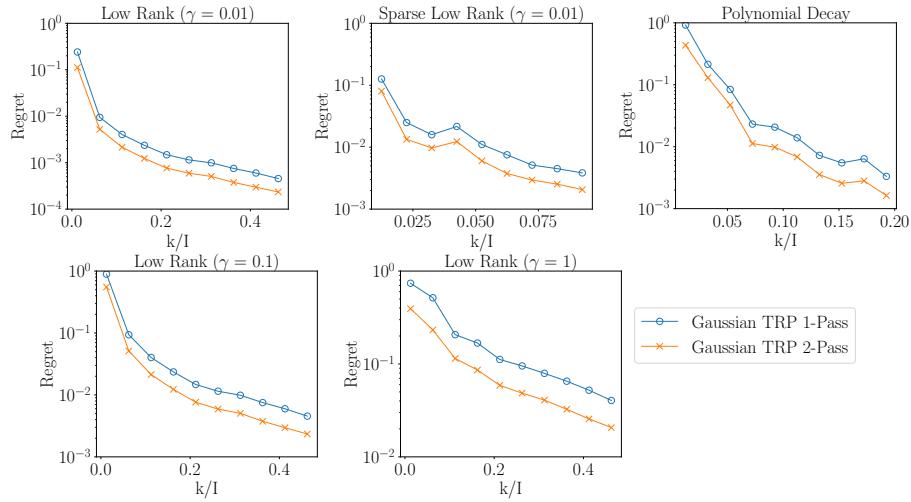


Fig. 11: We approximate 3D synthetic tensors (see subsection 6.1) with $I = 400$, using our one-pass and two-pass algorithms with $r = 5$ and varying k ($s = 2k + 1$), using the Gaussian TRP in the Tucker sketch.

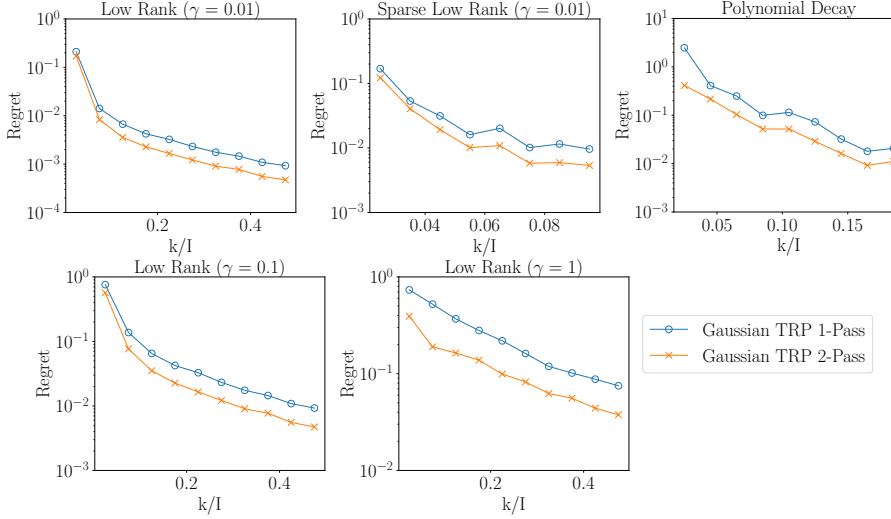


Fig. 12: We approximate 3D synthetic tensors (see subsection 6.1) with $I = 200$, using our one-pass and two-pass algorithms with $r = 5$ and varying k ($s = 2k + 1$), using the Gaussian TRP in the Tucker sketch.

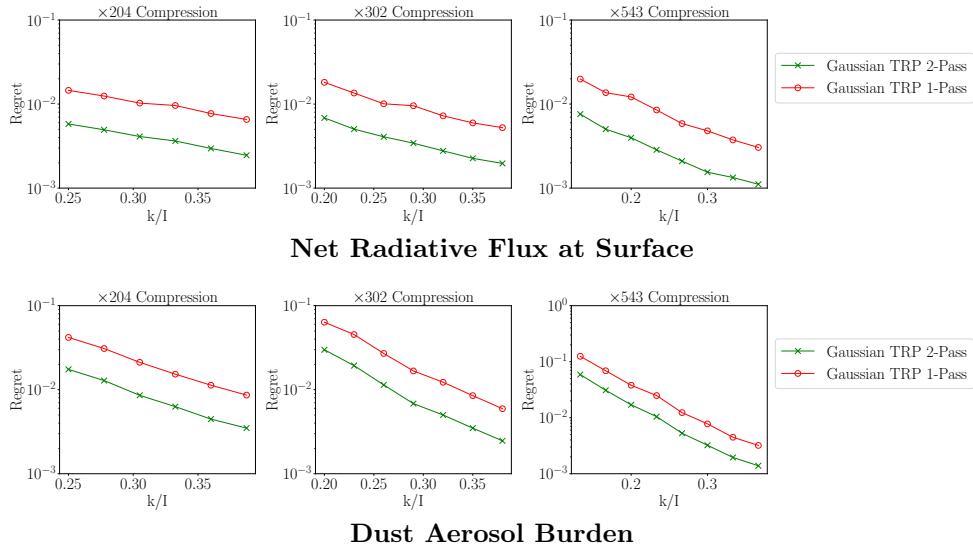


Fig. 13: We approximate the net radiative flux and dust aerosol burden data using our one-pass and two-pass algorithms using Gaussian TRP. We compare the performance under different ranks ($r/I = 0.125, 0.2, 0.067$). The dataset comes from the CESM CAM. The dust aerosol burden measures the amount of aerosol contributed by the dust. The net radiative flux determines the energy received by the earth surface through radiation.