Appendix A Supplementary experimental results

The Wikipedia dataset is built by crawling all documents in all subcategories within 3 layers below the *science* category. The Enron dataset is from the Enron email corpus [17]. After usual cleaning steps, the Wikipedia dataset has 114, 274 documents with an average 512 words per document; the Enron dataset has 186, 501 emails with average 91 words per email.

Table 5: Squared residual norm on top 10 recovered eigenvectors of 1000d tensors and running time (excluding I/O and sketch building time) for plan (exact) and sketched robust tensor power methods. Two vectors are considered mismatched (wrong) if $\|\boldsymbol{v} - \hat{\boldsymbol{v}}\|_2^2 > 0.1$.

		No	o. of v	vrong	vecto	ors		Running time (min.)										
	$\log_2(b)$:	12	13	14	15	16	12	13	14	15	16	12	13	14	15	16		
	B = 20	.40	.19	.10	.09	.08	8	6	3	0	0	.85	1.6	3.5	7.4	16.6		
.01	B = 30	.26	.10	.09	.08	.07	7	5	2	0	0	1.3	2.4	5.3	11.3	24.6		
	B = 40	.17	.10	.08	.08	.07	7	4	0	0	0	1.8	3.3	7.3	15.2	33.0		
Ь	Exact	.07					0					293.	293.5					
	B = 20	.52	3.1	.21	.18	.17	8	7	4	0	0	.84	1.6	3.5	7.5	16.8		
	B = 30	4.0	.24	.19	.17	.16	7	5	3	0	0	1.3	2.5	5.4	11.6	26.2		
н Б	B = 40	.30	.22	.18	.17	.16	7	4	0	0	0	1.8	3.3	7.3	15.5	33.5		
	Exact	.16					0					271.	8					

Table 6: Selected negative log-likelihood and running time (min) for fast and exact spectral methods on Wikipedia (top) and Enron (bottom) datasets.

		k	= 50			k	= 100			k = 200				
		Fast RB	RB	ALS		Fast RB	RB	ALS	F	ast RB	RB	ALS		
·H	like.	8.01	7.94	8.16		7.90	7.81	7.93		7.86	7.77	7.89		
Vik	time	2.2	2.2 97.7			6.8	135	29.3		57.3	423	677		
>	$\log_2 b$	10	-	-		12	-	-		14	-	-		
n	like.	8.31	8.28	8.22		8.18	8.09	8.30		8.26	8.18	8.27		
Enron	time	2.4	45.8	5.2		3.7	93.9	40.6		6.4	219	660		
	$\log_2 b$	11	-	-		11	-	-		11	-	-		

Appendix B Fast tensor power method via symmetric sketching

In this section we show how to do fast tensor power method using symmetric tensor sketches. More specifically, we explain how to approximately compute T(u, u, u) and T(I, u, u) when colliding hashes are used.

For symmetric tensors A and B, their inner product can be approximated by

$$\langle \mathbf{A}, \mathbf{B} \rangle \approx \langle \tilde{s}_{\mathbf{A}}, \tilde{s}_{\widetilde{\mathbf{B}}} \rangle,$$
 (10)

where $\widetilde{\mathbf{B}}$ is an "upper-triangular" tensor defined as

$$\widetilde{\mathbf{B}}_{i,j,k} = \begin{cases} \mathbf{B}_{i,j,k}, & \text{if } i \le j \le k; \\ 0, & \text{otherwise.} \end{cases}$$
(11)

Note that in Eq. (10) only the matrix B is "truncated". We show this gives consistent estimates of $\langle A, B \rangle$ in Appendix E.2.

Recall that $\mathbf{T}(u, u, u) = \langle \mathbf{T}, \mathbf{X} \rangle$ where $\mathbf{X} = u \otimes u \otimes u$. The symmetric tensor sketch $\tilde{s}_{\widetilde{\mathbf{X}}}$ can be computed as

$$\tilde{\boldsymbol{s}}_{\widetilde{\mathbf{X}}} = \frac{1}{6} \tilde{\boldsymbol{s}}_{\boldsymbol{u}}^{\otimes 3} + \frac{1}{2} \tilde{\boldsymbol{s}}_{2,\boldsymbol{u}\circ\boldsymbol{u}} * \tilde{\boldsymbol{s}}_{\boldsymbol{u}} + \frac{1}{3} \tilde{\boldsymbol{s}}_{3,\boldsymbol{u}\circ\boldsymbol{u}\circ\boldsymbol{u}},$$
(12)

where
$$\tilde{s}_{2,\boldsymbol{u}\circ\boldsymbol{u}}(t) = \sum_{2h(i)=t} \sigma(i)^2 \boldsymbol{u}_i^2$$
 and $\tilde{s}_{3,\boldsymbol{u}\circ\boldsymbol{u}\circ\boldsymbol{u}}(t) = \sum_{3h(i)=t} \sigma(i)^3 \boldsymbol{u}_i^3$. As a result,
 $\mathbf{T}(\boldsymbol{u},\boldsymbol{u},\boldsymbol{u}) \approx \frac{1}{6} \langle \mathcal{F}(\tilde{\boldsymbol{s}}_{\mathbf{T}}), \mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}})\circ\mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}})\circ\mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}}) \rangle + \frac{1}{2} \langle \mathcal{F}(\tilde{\boldsymbol{s}}_{\mathbf{T}}), \mathcal{F}(\tilde{\boldsymbol{s}}_{2,\boldsymbol{u}\circ\boldsymbol{u}})\circ\mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}}) \rangle + \frac{1}{3} \langle \tilde{\boldsymbol{s}}_{\mathbf{T}}, \tilde{\boldsymbol{s}}_{3,\boldsymbol{u}\circ\boldsymbol{u}\circ\boldsymbol{u}} \rangle$
(13)

Algorithm 2 Fast ALS method

- 1: **Input**: $\mathbf{T} \in \mathbb{R}^{n \times n \times n}$, target rank k, T, B, b.
- 2: Initialize: B independent index hash functions $h^{(1)}, \dots, h^{(B)}$ and $\sigma^{(1)}, \dots, \sigma^{(B)}$; random matrices $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times k}$; $\{\lambda_i\}_{i=1}^k$.
- 3: For $m = 1, \dots, B$ compute $s_{\mathbf{T}}^{(m)} \in \mathbb{C}^{b}$.
- 4: **for** t = 1 to T **do**
- 5: Compute count sketches s_{b_i} , s_{c_i} for $i = 1, \dots, k$. For each $i = 1, \dots, k; m = 1, \dots, b$ compute $v_i^{(m)} \approx \mathbf{T}(\mathbf{I}, b_i, c_i)$.
- 6: $\bar{\boldsymbol{v}}_{ij} \leftarrow \operatorname{med}(\Re(\boldsymbol{v}_{ij}^{(1)}), \Re(\boldsymbol{v}_{ij}^{(2)}), \cdots, \Re(\boldsymbol{v}_{ij}^{(B)})).$
- 7: Set $\widehat{\mathbf{A}} = \{\overline{v}\}_{ij}$ and $\widehat{\lambda}_i = \|\widehat{a}_i\|$; afterwards, normalize each column of \mathbf{A} .
- 8: Update **B** and **C** similarly.
- 9: **Output**: eigenvalues $\{\lambda_i\}_{i=1}^k$; solutions **A**, **B**, **C**.

For $\mathbf{T}(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u})$ recall that $[\mathbf{T}(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u})]_i = \langle \mathbf{T}, \mathbf{Y}_i \rangle$ where $\mathbf{Y}_i = \boldsymbol{e}_i \otimes \boldsymbol{u} \otimes \boldsymbol{u}$. We first symmetrize it by defining $\mathbf{Z}_i = \boldsymbol{e}_i \otimes \boldsymbol{u} \otimes \boldsymbol{u} + \boldsymbol{u} \otimes \boldsymbol{e}_i \otimes \boldsymbol{u} + \boldsymbol{u} \otimes \boldsymbol{u} \otimes \boldsymbol{e}_i$. ⁵ The sketch of $\widetilde{\mathbf{Z}}_i$ can be subsequently computed as

$$\tilde{\boldsymbol{s}}_{\tilde{\boldsymbol{Z}}_{i}} = \frac{1}{2} \tilde{\boldsymbol{s}}_{\boldsymbol{u}} * \tilde{\boldsymbol{s}}_{\boldsymbol{u}} * \tilde{\boldsymbol{s}}_{\boldsymbol{e}_{i}} + \frac{1}{2} \tilde{\boldsymbol{s}}_{2,\boldsymbol{u}\circ\boldsymbol{u}} * \tilde{\boldsymbol{s}}_{\boldsymbol{e}_{i}} + \tilde{\boldsymbol{s}}_{2,\boldsymbol{e}_{i}\circ\boldsymbol{u}} * \tilde{\boldsymbol{s}}_{\boldsymbol{u}} + \tilde{\boldsymbol{s}}_{3,\boldsymbol{e}_{i}\circ\boldsymbol{u}\circ\boldsymbol{u}}.$$
(14)

Consequently,

$$\mathbf{T}(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}) \approx \left\langle \mathcal{F}^{-1}\left(\mathcal{F}(\tilde{\boldsymbol{s}}_{\mathbf{T}}) \circ \overline{\mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}})}\right), \tilde{\boldsymbol{s}}_{2, \boldsymbol{e}_{i} \circ \boldsymbol{u}} \right\rangle + \frac{1}{6} \left\langle \mathcal{F}^{-1}\left(\mathcal{F}(\tilde{\boldsymbol{s}}_{\mathbf{T}}) \circ \overline{\mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}})} \circ \overline{\mathcal{F}(\tilde{\boldsymbol{s}}_{\boldsymbol{u}})}\right), \tilde{\boldsymbol{s}}_{\boldsymbol{e}_{i}} \right\rangle \\ + \frac{1}{6} \left\langle \mathcal{F}^{-1}\left(\mathcal{F}(\tilde{\boldsymbol{s}}_{\mathbf{T}}) \circ \overline{\mathcal{F}(\tilde{\boldsymbol{s}}_{2, \boldsymbol{u} \circ \boldsymbol{u}})}\right), \tilde{\boldsymbol{s}}_{\boldsymbol{e}_{i}} \right\rangle + \left\langle \tilde{\boldsymbol{s}}_{\mathbf{T}}, \tilde{\boldsymbol{s}}_{3, \boldsymbol{e}_{i} \circ \boldsymbol{u} \circ \boldsymbol{u}} \right\rangle.$$
(15)

Note that all of \tilde{s}_{e_i} , $\tilde{s}_{2,e_i \circ u}$ and $\tilde{s}_{3,e_i \circ u \circ u}$ have exactly one nonzero entries. So we can pre-compute all terms on the left sides of inner products in Eq. (15) and then read off the values for each entry in $\mathbf{T}(\mathbf{I}, u, u)$.

Appendix C Fast ALS: method and simulation result

In this section we describe how to use tensor sketching to accelerate the Alternating Least Squares (ALS) method for tensor CP decomposition. We also provide experimental results on synthetic data and compare our fast ALS implementation with the Matlab tensor toolbox [32, 33], which is widely considered to be the state-of-the-art for tensor decomposition.

C.1 Alternating Least Squares

Alternating Least Squares (ALS) is a popular method for tensor CP decompositions [19]. The algorithm maintains $\lambda \in \mathbb{R}^k$, $\mathbf{A}, \mathbf{B}, \mathbf{C} \in \mathbb{R}^{n \times k}$ and iteratively perform the following update steps:

$$\widehat{\mathbf{A}} = \mathbf{T}_{(1)} (\mathbf{C} \odot \mathbf{B}) (\mathbf{C}^{\top} \mathbf{C} \circ \mathbf{B}^{\top} \mathbf{B})^{\dagger}.$$
(16)
$$\widehat{\mathbf{B}} = \mathbf{T}_{(1)} (\widehat{\mathbf{A}} \odot \mathbf{C}) (\widehat{\mathbf{A}}^{\top} \widehat{\mathbf{A}} \circ \mathbf{C}^{\top} \mathbf{C})^{\dagger};$$
$$\widehat{\mathbf{C}} = \mathbf{T}_{(1)} (\widehat{\mathbf{B}} \odot \widehat{\mathbf{A}}) (\widehat{\mathbf{B}}^{\top} \widehat{\mathbf{B}} \circ \widehat{\mathbf{A}}^{\top} \widehat{\mathbf{A}})^{\dagger}.$$

After each update, $\hat{\lambda}_r$ is set to $\|\boldsymbol{a}_r\|_2$ (or $\|\boldsymbol{b}_r\|_2$, $\|\boldsymbol{c}_r\|_2$) for $r = 1, \dots, k$ and the matrix **A** (or **B**, **C**) is normalized so that each column has unit norm. The final low-rank approximation is obtained by $\sum_{i=1}^k \hat{\lambda}_i \hat{\boldsymbol{a}}_i \otimes \hat{\boldsymbol{b}}_i \otimes \hat{\boldsymbol{c}}_i$.

There is no guarantee that ALS converges or gives a good tensor decomposition. Nevertheless, it works reasonably well in most applications [19]. In general ALS requires $O(T(n^3k + k^3))$ computations and $O(n^3)$ storage, where T is the number of iterations.

Table 7: Squared residual norm on top 10 recovered eigenvectors of 1000d tensors and running time (excluding I/O and sketch building time) for plain (exact) and sketched ALS algorithms. Two vectors are considered mismatched (wrong) if $||v - \hat{v}||_2^2 > 0.1$.

	Residual norm								No. of wrong vectors						Running time (min.)					
	$\log_2(b)$:	12	13	14	15	16	-	12	13	14	15	16	-	12	13	14	15	16		
01	B = 20	.71	.41	.25	.17	.12		10	9	7	6	4		.11	.22	.49	1.1	2.4		
•	B = 30	.50	.34	.21	.14	.11		9	8	7	5	3		.17	.33	.75	1.6	3.5		
	B = 40	.46	.28	.17	.10	.07		9	8	6	5	1		.23	.45	1.0	2.2	4.7		
Ь	Exact [†]	.07						1						22.8						
	B = 20	.88	.50	.35	.28	.23		10	8	7	6	6		.13	.32	.78	1.5	3.2		
•	B = 30	.78	.44	.30	.24	.21		9	8	7	5	6		.21	.50	1.1	2.2	4.7		
σ =	B = 40	.56	.38	.28	.19	.16		9	8	6	4	2		.29	.69	1.5	3.5	6.3		
0	Exact [†]	.17						2						32.3						

[†]Calling cp_als in Matlab tensor toolbox. It is run for exactly T = 30 iterations.

C.2 Accelerated ALS via sketching

Similar to robust tensor power method, the ALS algorithm can be significantly accelerated by using the idea of sketching as shown in this work. However, for ALS we cannot use colliding hashes because though the input tensor **T** is symmetric, its CP decomposition is not since we maintain three different solution matrices **A**, **B** and **C**. As a result, we roll back to asymmetric tensor sketches defined in Eq. (1). Recall that given **A**, **B**, **C** $\in \mathbb{R}^{n \times k}$ we want to compute

$$\hat{\mathbf{A}} = \mathbf{T}_{(1)} (\mathbf{C} \odot \mathbf{B}) (\mathbf{C}^{\top} \mathbf{C} \circ \mathbf{B}^{\top} \mathbf{B})^{\dagger}.$$
(17)

When k is much smaller than the ambient tensor dimension n the computational bottleneck of Eq. (17) is $\mathbf{T}_{(1)}(\mathbf{C} \odot \mathbf{B})$, which requires $O(n^3k)$ operations. Below we show how to use sketching to speed up this computation.

Let $x \in \mathbb{R}^{n^2}$ be one row in $\mathbf{T}_{(1)}$ and consider $(\mathbf{C} \odot \mathbf{B})^\top x$. It can be shown that [15]

$$\left[\left(\mathbf{C} \odot \mathbf{B} \right)^{\top} \boldsymbol{x} \right]_{i} = \boldsymbol{b}_{i}^{\top} \mathbf{X} \boldsymbol{c}_{i}, \quad \forall i = 1, \cdots, k,$$
(18)

where $\mathbf{X} \in \mathbb{R}^{n \times n}$ is the reshape of vector \boldsymbol{x} . Subsequently, the product $\mathbf{T}_{(1)}(\mathbf{C} \odot \mathbf{B})$ can be re-written as

$$\mathbf{T}_{(1)}(\mathbf{C} \odot \mathbf{B}) = [\mathbf{T}(\mathbf{I}, \boldsymbol{b}_1, \boldsymbol{c}_1); \cdots; \mathbf{T}(\mathbf{I}, \boldsymbol{b}_k, \boldsymbol{c}_k)].$$
(19)

Using Proposition 1 we can compute each of $\mathbf{T}(\mathbf{I}, \boldsymbol{b}_i, \boldsymbol{c}_i)$ in $O(n + b \log b)$ iterations. Note that in general $\boldsymbol{b}_i \neq \boldsymbol{c}_i$, but Proposition 1 still holds by replacing one of the two \boldsymbol{s}_u sketches. As a result, $\mathbf{T}_{(1)}(\mathbf{C} \odot \mathbf{B})$ can be computed in $O(k(n + b \log b))$ operations once \boldsymbol{s}_T is computed. The pseudocode of fast ALS is listed in Algorithm 2. Its time complexity and space complexity are $O(T(k(n + Bb \log b) + k^3))$ (excluding the time for building \boldsymbol{s}_T) and O(Bb), respectively.

C.3 Simulation results

We compare the performance of fast ALS with a brute-force implementation under various hash length settings on synthetic datasets in Table 7. Settings for generating the synthetic dataset is exactly the same as in Section 5.1. We use the cp_als routine in Matlab tensor toolbox as the reference brute-force implementation of ALS. For fair comparison, exactly T = 30 iterations are performed for both plain and accelerated ALS algorithms. Table 7 shows that when sketch length b is not too small, fast ALS achieves comparable accuracy with exact methods while being much faster in terms of running time.

Appendix D Spectral LDA and fast spectral LDA

Latent Dirichlet Allocation (LDA, [3]) is a powerful tool in topic modeling. In this section we first review the LDA model and introduce the tensor decomposition method for learning LDA models, which was proposed in [1]. We then provide full details of our proposed fast spectral LDA algorithm. Pseudocode for fast spectral LDA is listed in Algorithm 3.

⁵As long as **A** is symmetric, we have $\langle \mathbf{A}, \mathbf{Y}_i \rangle = \langle \mathbf{A}, \mathbf{Z}_i \rangle / 3$.

Algorithm 3 Fast spectral LDA

- 1: **Input**: Unlabeled documents, V, K, α_0, B, b .
- 2: Compute empirical moments $\widehat{\mathbf{M}}_1$ and $\widehat{\mathbf{M}}_2$ defined in Eq. (20,21).
- 3: $[\mathbf{U}, \mathbf{S}, \mathbf{V}] \leftarrow \text{truncatedSVD}(\widehat{\mathbf{M}}_2, k); \mathbf{W}_{ik} \leftarrow \frac{\mathbf{U}_{ik}}{\sqrt{\sigma_k}}.$
- 4: Build *B* tensor sketches of $\widehat{\mathbf{M}}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$.
- 5: Find CP decomposition $\{\lambda_i\}_{i=1}^k$, $\mathbf{A} = \mathbf{B} = \mathbf{C} = \{v_i\}_{i=1}^k$ of $\widehat{\mathbf{M}}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ using either fast tensor power method or fast ALS method.
- 6: **Output:** estimates of prior parameters $\hat{\alpha}_i = \frac{4\alpha_0(\alpha_0+1)}{(\alpha_0+2)^2\lambda_i^2}$ and topic distributions $\hat{\mu}_i = \frac{\alpha_0+2}{2}\lambda_i(\mathbf{W}^{\dagger})^{\top}\boldsymbol{v}_i$.

D.1 LDA and spectral LDA

LDA models a collection of documents by a topic dictionary $\Phi \in \mathbb{R}^{V \times K}$ and a Dirichlet prior $\alpha \in \mathbb{R}^k$, where V is the vocabulary size and k is the number of topics. Each column in Φ is a probability distribution (i.e., non-negative and sum to one) representing the word distribution of a particular topic. For each document d, a topic mixing vector $h_d \in \mathbb{R}^k$ is first sampled from a Dirichlet distribution parameterized by α . Afterwards, words in document d i.i.d. sampled from a categorical distribution parameterized by Φh_d .

A spectral method for LDA based on 3rd-order robust tensor decomposition was proposed in [1] to provably learn LDA model parameters from a polynomial number of training documents. Let $x \in \mathbb{R}^V$ represent a single word; that is, for word w we have $x_w = 1$ and $x_{w'} = 0$ for all $w' \neq w$. Define first, second and third order moments $\mathbf{M}_1, \mathbf{M}_2$ and \mathbf{M}_3 as follows:

$$\mathbf{M}_1 = \mathbb{E}[\boldsymbol{x}_1]; \tag{20}$$

$$\mathbf{M}_2 = \mathbb{E}[\boldsymbol{x}_1 \otimes \boldsymbol{x}_2] - \frac{\alpha_0}{\alpha_0 + 1} \mathbf{M}_1 \otimes \mathbf{M}_1;$$
(21)

$$\mathbf{M}_3 = \mathbb{E}[\boldsymbol{x}_1 \otimes \boldsymbol{x}_2 \otimes \boldsymbol{x}_3] - \frac{\alpha_0}{\alpha_0 + 2} (\mathbb{E}[\boldsymbol{x}_1 \otimes \boldsymbol{x}_2 \otimes \mathbf{M}_1] + \mathbb{E}[\boldsymbol{x}_1 \otimes \mathbf{M}_1 \otimes \boldsymbol{x}_2] + \mathbb{E}[\mathbf{M}_1 \otimes \boldsymbol{x}_1 \otimes \boldsymbol{x}_2])$$

+
$$\frac{2\alpha_0^2}{(\alpha_0+1)(\alpha_0+2)}$$
 $\mathbf{M}_1 \otimes \mathbf{M}_1 \otimes \mathbf{M}_1.$ (22)

Here $\alpha_0 = \sum_k \alpha_k$ is assumed to be a known quantity. Using elementary algebra it can be shown that

$$\mathbf{M}_2 = \frac{1}{\alpha_0(\alpha_0 + 1)} \sum_{i=1}^k \alpha_i \boldsymbol{\mu}_i \boldsymbol{\mu}_i^{\mathsf{T}};$$
(23)

$$\mathbf{M}_{3} = \frac{2}{\alpha_{0}(\alpha_{0}+1)(\alpha_{0}+2)} \sum_{i=1}^{k} \alpha_{i} \boldsymbol{\mu}_{i} \otimes \boldsymbol{\mu}_{i} \otimes \boldsymbol{\mu}_{i}.$$
(24)

To extract topic vectors $\{\boldsymbol{\mu}_i\}_{i=1}^k$ from \mathbf{M}_2 and \mathbf{M}_3 , a simultaneous diagonalization procedure is carried out. More specifically, the algorithm first finds a whitening matrix $\mathbf{W} \in \mathbb{R}^{V \times K}$ with orthonormal columns such that $\mathbf{W}^\top \mathbf{M}_2 \mathbf{W} = \mathbf{I}_{K \times K}$. In practice, this step can be completed by performing a truncated SVD on \mathbf{M}_2 , $\mathbf{M}_2 = \mathbf{U}_K \boldsymbol{\Sigma}_K \mathbf{V}_K$, and set $\mathbf{W}_{ik} = \mathbf{U}_{ik}/\sqrt{\boldsymbol{\Sigma}_{kk}}$. Afterwards, tensor CP decomposition is performed on the whitened third order moment $\mathbf{M}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ ⁶ to obtain a set of eigenvectors $\{\boldsymbol{v}_k\}_{k=1}^K$. The topic vectors $\{\boldsymbol{\mu}_k\}_{k=1}^K$ can be subsequently obtained by multiplying $\{\boldsymbol{v}_k\}_{k=1}^K$ with the pseudoinverse of \mathbf{W} . Note that Eq. (20,21,22) are defined in exact word moments. In practice we use empirical moments (e.g., word frequency vector and cooccurrence matrix) to approximate these exact moments.

⁶For a tensor $\mathbf{T} \in \mathbb{R}^{V \times V \times V}$ and a matrix $\mathbf{W} \in \mathbb{R}^{V \times k}$, the product $\mathbf{Q} = \mathbf{T}(\mathbf{W}, \mathbf{W}, \mathbf{W}) \in \mathbb{R}^{k \times k \times k}$ is defined as $\mathbf{Q}_{i_1, i_2, i_3} = \sum_{j_1, j_2, j_3=1}^{V} \mathbf{T}_{j_1, j_2, j_3} \mathbf{W}_{j_1, i_1} \mathbf{W}_{j_2, i_2} \mathbf{W}_{j_3, i_3}$.

D.2 Fast spectral LDA

To further accelerate the spectral method mentioned in the previous section, it helps to first identify computational bottlenecks of spectral LDA. In general, the computation of $\widehat{\mathbf{M}}_1, \widehat{\mathbf{M}}_2$ and the whitening step are not the computational bottleneck when V is not too large and each document is not too long. The bottleneck comes from the computation of (the sketch of) $\widehat{\mathbf{M}}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ and its tensor decomposition. By Eq. (22), the computation of $\widehat{\mathbf{M}}_3(\mathbf{W}, \mathbf{W}, \mathbf{W})$ reduces to computing $\widehat{\mathbf{M}}_1^{\otimes 3}(\mathbf{W}, \mathbf{W}, \mathbf{W}), \widehat{\mathbb{E}}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \widehat{\mathbf{M}}_1](\mathbf{W}, \mathbf{W}, \mathbf{W}), ^7$ and $\widehat{\mathbb{E}}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3](\mathbf{W}, \mathbf{W}, \mathbf{W})$. The first term $\widehat{\mathbf{M}}_1^{\otimes 3}(\mathbf{W}, \mathbf{W}, \mathbf{W})$ poses no particular challenge as it can be written as $(\mathbf{W}^{\top} \widehat{\mathbf{M}}_1)^{\otimes 3}$. Its sketch can then be efficiently obtained by applying techniques in Section 3.1. In the remainder of this section we focus on efficient computation of the sketch of the other two terms mentioned above.

We first show how to efficiently sketching $\hat{\mathbb{E}}[\boldsymbol{x}_1 \otimes \boldsymbol{x}_2 \otimes \boldsymbol{x}_3](\mathbf{W}, \mathbf{W}, \mathbf{W})$ given the whitening matrix \mathbf{W} and D training documents. Let $\mathbf{T}\hat{\mathbb{E}}[\boldsymbol{x}_1 \otimes \boldsymbol{x}_2 \otimes \boldsymbol{x}_3](\mathbf{W}, \mathbf{W}, \mathbf{W})$ denote the whitened $k \times k \times k$ tensor to be sketched and write $\mathbf{T} = \sum_{d=1}^{D} \mathbf{T}_d$, where \mathbf{T}_d is the contribution of the dth training document to \mathbf{T} . By definition, \mathbf{T}_d can be expressed as $\mathbf{T}_d = \mathbf{N}_d(\mathbf{W}, \mathbf{W}, \mathbf{W})$, where \mathbf{W} is the $V \times k$ whitening matrix and \mathbf{N}_d is the $V \times V \times V$ empirical moment tensor computed on the dth document. More specifically, for $i, j, k \in \{1, \dots, V\}$ we have

$$\mathbf{N}_{d,ijk} = \frac{1}{m_d(m_d - 1)(m_d - 2)} \begin{cases} n_{di}(n_{dj} - 1)(n_{dk} - 2), & i = j = k; \\ n_{di}(n_{di} - 1)n_{dk}, & i = j, j \neq k; \\ n_{di}n_{dj}(n_{dj} - 1) & j = k, i \neq j; \\ n_{di}(n_{di} - 1)n_{dj}, & i = k, i \neq j; \\ n_{di}n_{dj}n_{dk}, & \text{otherwise.} \end{cases}$$

Here m_d is the length (i.e., number of words) of document d and $\mathbf{n}_d \in \mathbb{R}^V$ is the corresponding word count vector. Previous straightforward implementation require at least $O(k^3 + m_d k^2)$ operations per document to build the tensor \mathbf{T} and $O(k^4 L T)$ to decompose it [30, 29], which is prohibitively slow for real-world applications. In section 3 we discussed how to decompose a tensor efficiently once we have its sketch. We now show how to build the sketch of \mathbf{T} efficiently from document word counts $\{\mathbf{n}_d\}_{d=1}^D$.

By definition, \mathbf{T}_d can be decomposed as

$$\mathbf{T}_{d} = \boldsymbol{p}^{\otimes 3} - \sum_{i=1}^{V} n_{i} (\boldsymbol{w}_{i} \otimes \boldsymbol{w}_{i} \otimes \boldsymbol{p} + \boldsymbol{w}_{i} \otimes \boldsymbol{p} \otimes \boldsymbol{w}_{i} + \boldsymbol{p} \otimes \boldsymbol{w}_{i} \otimes \boldsymbol{w}_{i}) + \sum_{i=1}^{V} 2n_{i} \boldsymbol{w}_{i}^{\otimes 3}, \qquad (25)$$

where $p = \mathbf{W}n$ and $w_i \in \mathbb{R}^k$ is the *i*th row of the whitening matrix \mathbf{W} . A direct implementation is to sketch each of the low-rank components in Eq. (25) and compute their sum. Since there are $O(m_d)$ tensors, building the sketch of \mathbf{T}_d requires $O(m_d)$ FFTs, which is unsatisfactory. However, note that $\{w_i\}_{i=1}^V$ are fixed and shared across documents. So when scanning the documents we maintain the sum of n_i and $n_i p$ and add the incremental after all documents are scanned. In this way, we only need O(1) FFT per document with an additional O(V) FFTs. Since the total number of documents D is usually much larger than V, this provides significant speed-ups over the naive method that sketches each term in Eq. (25) independently. As a result, the sketch of \mathbf{T} can be computed in $O(k(\sum_d m_d) + (D + V)b \log b)$ operations, which is much more efficient than the $O(k^2(\sum_d m_d) + Dk^3)$ brute-force computation.

We next turn to the term $\widehat{\mathbb{E}}[x_1 \otimes x_2 \otimes \widehat{\mathbf{M}}_1](\mathbf{W}, \mathbf{W}, \mathbf{W})$. Fix a document d and let $p = \mathbf{W}n_d$. Define $q = \mathbf{W}\widehat{\mathbf{M}}_1$. By definition, the whitened empirical moment can be decomposed as

$$\widehat{\mathbb{E}}[\boldsymbol{x}_1 \otimes \boldsymbol{x}_2 \otimes \widehat{\mathbf{M}}_1](\mathbf{W}, \mathbf{W}, \mathbf{W}) = \sum_{i=1}^V n_i \boldsymbol{p} \otimes \boldsymbol{p} \otimes \boldsymbol{q},$$
(26)

Note that Eq. (26) is very similar to Eq. (25). Consequently, we can apply the same trick (i.e., adding p and $n_i p$ up before doing sketching or FFT) to compute Eq. (26) efficiently.

⁷and also $\hat{\mathbb{E}}[\boldsymbol{x}_1 \otimes \widehat{\mathbf{M}}_1 \otimes \boldsymbol{x}_2](\mathbf{W}, \mathbf{W}, \mathbf{W}), \hat{\mathbb{E}}[\widehat{\mathbf{M}}_1 \otimes \boldsymbol{x}_1 \otimes \boldsymbol{x}_2](\mathbf{W}, \mathbf{W}, \mathbf{W})$ by symmetry.

Appendix E Proofs

E.1 Proofs of some technical propositions

Proof of Proposition 2. We prove the proposition for the case q = 2 (i.e., \tilde{H} is 2-wise independent). This suffices for our purpose in this paper and generalization to q > 2 cases is straightforward. For notational simplicity we omit all modulo operators. Consider two *p*-tuples $\boldsymbol{l} = (l_1, \dots, l_p)$ and $\boldsymbol{l}' = (l'_1, \dots, l'_p)$ such that $\boldsymbol{l} \neq \boldsymbol{l}'$. Since \tilde{H} is permutation invariant, we assume without loss of generality that for some s < p and $1 \le i \le s$ we have $l_i = l'_i$. Fix $t, t' \in [b]$. We then have

$$\Pr[\tilde{H}(\boldsymbol{l}) = t \land \tilde{H}(\boldsymbol{l}') = t'] = \sum_{a} \sum_{h(l_1) + \dots + h(l_s) = a} \Pr[h(l_1) + \dots + h(l_s) = a]$$

$$\cdot \sum_{\substack{r_{s+1} + \dots + r_p = t-a \\ r'_{s+1} + \dots + r'_p = t'-a}} \Pr[h(l_{s+1}) = r_1 \land \dots \land h(l_p) = r_p \land h(l'_{s+1}) = r'_1 \land \dots \land h(l'_p) = r'_p].$$
(27)

Since h is 2p-wise independent, we have

$$\Pr[h(l_1) + \dots + h(l_s) = a] = \sum_{r_1 + \dots + r_s = a} \Pr[h(l_1) = r_1 \wedge \dots \wedge h(l_s) = r_s] = b^{s-1} \cdot \frac{1}{b^s} = \frac{1}{b};$$

$$\sum_{\substack{r_{s+1} + \dots + r_p = t-a \\ r'_{s+1} + \dots + r'_p = t-a}} \Pr[h(l_{s+1}) = r_1 \wedge \dots \wedge h(l_p) = r_p \wedge h(l'_{s+1}) = r'_1 \wedge \dots \wedge h(l'_p) = r'_p]$$

$$= b^{2(p-s-1)} \cdot \frac{1}{b^{2(p-s)}} = \frac{1}{b^2}.$$

Summing everything up we get $\Pr[\tilde{H}(l) = t \wedge \tilde{H}(l') = t'] = 1/b^2$, which is to be demonstrated. \Box

Proof of Proposition 1. Since both FFT and inverse FFT preserve inner products, we have

$$\langle \mathbf{s}_{\mathbf{T}}, \mathbf{s}_{1,\boldsymbol{u}} * \mathbf{s}_{2,\boldsymbol{u}} * \mathbf{s}_{3,\boldsymbol{e}_{i}} \rangle = \langle \mathcal{F}(\mathbf{s}_{\mathbf{T}}), \mathcal{F}(\mathbf{s}_{1,\boldsymbol{u}}) \circ \mathcal{F}(\mathbf{s}_{2,\boldsymbol{u}}) \circ \mathcal{F}(\mathbf{s}_{3,\boldsymbol{e}_{i}}) \rangle$$

$$= \langle \mathcal{F}(\mathbf{s}_{\mathbf{T}}) \circ \overline{\mathcal{F}(\mathbf{s}_{1,\boldsymbol{u}})} \circ \overline{\mathcal{F}(\mathbf{s}_{2,\boldsymbol{u}})}, \mathcal{F}(\mathbf{s}_{3,\boldsymbol{e}_{i}}) \rangle$$

$$= \langle \mathcal{F}^{-1}(\mathcal{F}(\mathbf{s}_{\mathbf{T}}) \circ \overline{\mathcal{F}(\mathbf{s}_{1,\boldsymbol{u}})} \circ \overline{\mathcal{F}(\mathbf{s}_{2,\boldsymbol{u}})}), \mathbf{s}_{3,\boldsymbol{e}_{i}} \rangle.$$

E.2 Analysis of tensor sketch approximation error

Proofs of Theorem 1 is based on the following two key lemmas, which states that $\langle \tilde{s}_A, \tilde{s}_{\widetilde{B}} \rangle$ is a consistent estimator of the true inner product $\langle A, B \rangle$; furthermore, the variance of the estimator decays linearly with the hash length *b*. The lemmas are interesting in their own right, providing useful tools for proving approximation accuracy in a wide range of applications when colliding hash and symmetric sketches are used.

Lemma 1. Suppose $\mathbf{A}, \mathbf{B} \in \bigotimes^{p} \mathbb{R}^{n}$ are two symmetric real tensors and let $\tilde{\mathbf{s}}_{\mathbf{A}}, \tilde{\mathbf{s}}_{\widetilde{\mathbf{B}}} \in \mathbb{C}^{b}$ be the symmetric tensor sketches of \mathbf{A} and $\widetilde{\mathbf{B}}$. That is,

$$\tilde{s}_{\mathbf{A}}(t) = \sum_{\tilde{H}(i_1, \cdots, i_p) = t} \sigma_{i_1} \cdots \sigma_{i_p} \mathbf{A}_{i_1, \cdots, i_p};$$
(28)

$$\tilde{s}_{\tilde{\mathbf{B}}}(t) = \sum_{\substack{\tilde{H}(i_1, \cdots, i_p) = t\\i_1 < \cdots < i_p}} \sigma_{i_1} \cdots \sigma_{i_p} \mathbf{B}_{i_1, \cdots, i_p}.$$
(29)

Assume $\tilde{H}(i_1, \dots, i_p) = (h(i_1) + \dots + h(i_p)) \mod b$ are drawn from a 2-wise independent hash family. Then the following holds:

$$\mathbb{E}_{h,\sigma}\left[\langle \tilde{s}_{\mathbf{A}}, \tilde{s}_{\widetilde{\mathbf{B}}} \rangle\right] = \langle \mathbf{A}, \mathbf{B} \rangle, \tag{30}$$

$$\mathbb{V}_{h,\sigma}\left[\langle \tilde{\boldsymbol{s}}_{\mathbf{A}}, \tilde{\boldsymbol{s}}_{\widetilde{\mathbf{B}}} \rangle\right] \leq \frac{4^{p} \|\mathbf{A}\|_{F}^{2} \|\mathbf{B}\|_{F}^{2}}{b}.$$
(31)

Lemma 2. Following notations and assumptions in Lemma 1. Let $\{\mathbf{A}_i\}_{i=1}^m$ and $\{\mathbf{B}_i\}_{i=1}^m$ be symmetric real $n \times n \times n$ tensors and fix real vector $\boldsymbol{w} \in \mathbb{R}^m$. Then we have

$$\mathbb{E}\left[\sum_{i,j} w_i w_j \langle \tilde{\mathbf{s}}_{\mathbf{A}_i}, \tilde{\mathbf{s}}_{\widetilde{\mathbf{B}}_j} \rangle\right] = \sum_{i,j} w_i w_j \langle \mathbf{A}_i, \mathbf{B}_j \rangle;$$
(32)

$$\mathbb{V}\left[\sum_{i,j} w_i w_j \langle \tilde{\boldsymbol{s}}_{\mathbf{A}_i}, \tilde{\boldsymbol{s}}_{\widetilde{\mathbf{B}}_j} \rangle\right] \leq \frac{4^p \|\boldsymbol{w}\|^4 (\max_i \|\mathbf{A}_i\|_F^2) (\max_i \|\mathbf{B}_i\|_F^2)}{b}.$$
 (33)

Proof of Lemma 1. We first define some notations. Let $\mathbf{l} = (l_1, \dots, l_p) \in [d]^p$ be a *p*-tuple denoting a multi-index. Define $\mathbf{A}_{\mathbf{l}} := \mathbf{A}_{l_1,\dots,l_p}$ and $\sigma(\mathbf{l}) := \sigma_{l_1} \cdots \sigma_{l_p}$. For $\mathbf{l}, \mathbf{l}' \in [n]^p$, define $\delta(\mathbf{l}, \mathbf{l}') = 1$ if $h(l_1) + \dots + h(l_p) \equiv h(l'_1) + \dots + h(l'_p) \pmod{b}$ and $\delta(\mathbf{l}, \mathbf{l}') = 0$ otherwise. For a *p*-tuple $\mathbf{l} \in [n]^p$, let $\mathcal{L}(\mathbf{l}) \in [n]^p$ denote the *p*-tuple obtained by re-ordering indices in \mathbf{l} in ascending order. Let $\mathcal{M}(\mathbf{l}) \in \mathbb{N}^b$ denote the "expanded version" of \mathbf{l} . That is, $[\mathcal{M}(\mathbf{l})]_i$ denote the number of occurrences of the index *i* in \mathbf{l} . By definition, $\|\mathcal{M}(\mathbf{l})\|_1 = p$. Finally, by definition $\widetilde{\mathbf{B}}_{\mathbf{l}'} = \mathbf{B}_{\mathbf{l}'}$ if $\mathbf{l}' = \mathcal{L}(\mathbf{l}')$ and $\widetilde{\mathbf{B}}_{\mathbf{l}'} = 0$ otherwise.

Eq. (30) is easy to prove. By definition and linearity of expectation we have

$$\mathbb{E}[\langle \tilde{\boldsymbol{s}}_{\mathbf{A}}, \tilde{\boldsymbol{s}}_{\widetilde{\mathbf{B}}} \rangle] = \sum_{\boldsymbol{l}, \boldsymbol{l}'} \delta(\boldsymbol{l}, \boldsymbol{l}') \sigma(\boldsymbol{l}) \mathbf{A}_{\boldsymbol{l}} \bar{\sigma}(\boldsymbol{l}') \widetilde{\mathbf{B}}_{\boldsymbol{l}'}.$$
(34)

Note that δ and σ are independent and

$$\mathbb{E}_{\sigma}[\sigma(\boldsymbol{l})\sigma(\boldsymbol{l}')] = \begin{cases} 1, & \text{if } \mathcal{L}(\boldsymbol{l}) = \mathcal{L}(\boldsymbol{l}'); \\ 0, & \text{otherwise.} \end{cases}$$
(35)

Also $\delta(l, l') = 1$ with probability 1 whenever $\mathcal{L}(l) = \mathcal{L}(l')$. Note that $\widetilde{\mathbf{B}}_{l'} = 0$ whenever $l' \neq \mathcal{L}(l')$. Consequently,

$$\mathbb{E}[\langle \tilde{s}_{\mathbf{A}}, \tilde{s}_{\tilde{\mathbf{B}}} \rangle] = \sum_{\boldsymbol{l} \in [n]^{p}} \mathbf{A}_{\boldsymbol{l}} \tilde{\mathbf{B}}_{\mathcal{L}(\boldsymbol{l})} = \langle \mathbf{A}, \mathbf{B} \rangle.$$
(36)

For the variance, we have the following expression for $\mathbb{E}[\langle \tilde{s}_{\mathbf{A}}, \tilde{s}_{\mathbf{B}} \rangle^2]$:

$$\mathbb{E}[\langle \tilde{s}_{\mathbf{A}}, \tilde{s}_{\widetilde{\mathbf{B}}} \rangle^{2}] = \sum_{l,l',r,r'} \mathbb{E}[\delta(l,l')\delta(r,r')] \cdot \mathbb{E}[\sigma(l)\bar{\sigma}(l')\bar{\sigma}(r)\sigma(r')] \cdot \mathbf{A}_{l}\mathbf{A}_{r}\widetilde{\mathbf{B}}_{l'}\widetilde{\mathbf{B}}_{r'} \quad (37)$$
$$=: \sum_{l,l',r,r'} \mathbb{E}[t(l,l',r,r')]. \quad (38)$$

$$=: \sum_{\boldsymbol{l},\boldsymbol{l}',\boldsymbol{r},\boldsymbol{r}'} \mathbb{E}[t(\boldsymbol{l},\boldsymbol{l}',\boldsymbol{r},\boldsymbol{r}')].$$
(38)

We remark that $\mathbb{E}[\sigma(l)\bar{\sigma}(l')\bar{\sigma}(r)\sigma(r')] = 0$ if $\mathcal{M}(l) - \mathcal{M}(l') \neq \mathcal{M}(r) - \mathcal{M}(r')$. In the remainder of the proof we will assume that $\mathcal{M}(l) - \mathcal{M}(l') = \mathcal{M}(r) - \mathcal{M}(r')$. This can be further categorized into two cases:

Case 1: $l' = \mathcal{L}(l)$ and $r' = \mathcal{L}(r)$. By definition $\mathbb{E}[\sigma(l)\bar{\sigma}(l')\sigma(r)\bar{\sigma}(r')] = 1$ and $\mathbb{E}[\delta(l, l')\delta(r, r')] = 1$. Subsequently $\mathbb{E}[t(l, l', r, r')] = \mathbf{A}_l \mathbf{A}_r \widetilde{\mathbf{B}}_{l'} \widetilde{\mathbf{B}}_{r'}$ and hence

$$\sum_{\boldsymbol{l},\boldsymbol{r},\boldsymbol{l}'=\mathcal{L}(\boldsymbol{l}),\boldsymbol{r}'=\mathcal{L}(\boldsymbol{r})} \mathbb{E}[t(\boldsymbol{l},\boldsymbol{l}',\boldsymbol{r},\boldsymbol{r}')] = \sum_{\boldsymbol{l},\boldsymbol{r}} \mathbf{A}_{\boldsymbol{l}} \mathbf{A}_{\boldsymbol{r}} \mathbf{B}_{\boldsymbol{l}} \mathbf{B}_{\boldsymbol{r}} = \langle \mathbf{A}, \mathbf{B} \rangle^{2}.$$
(39)

Case 2: $\mathbf{l}' \neq \mathcal{L}(\mathbf{l})$ or $\mathbf{r}' \neq \mathcal{L}(\mathbf{r})$. Since $\mathcal{M}(\mathbf{l}) - \mathcal{M}(\mathbf{l}') = \mathcal{M}(\mathbf{r}) - \mathcal{M}(\mathbf{r}') \neq 0$ we have $\mathbb{E}[\delta(\mathbf{l}, \mathbf{l}')\delta(\mathbf{r}, \mathbf{r}')] = 1/b$ because h is a 2-wise independent hash function. In addition, $\mathbb{E}[|\sigma(\mathbf{l})\overline{\sigma}(\mathbf{l}')\sigma(\mathbf{r})\overline{\sigma}(\mathbf{r}')|] \leq 1$.

To enumerate all (l, l', r, r') tuples that satisfy the colliding condition $\mathcal{M}(l) - \mathcal{M}(l') = \mathcal{M}(r) - \mathcal{M}(r') \neq 0$, we fix ⁸ $\|\mathcal{M}(l) - \mathcal{M}(l')\|_1 = 2q$ and fix q positions each in l and r (for l' and r' the positions of these indices are automatically fixed because indices in l' and r' must be in ascending

⁸Note that sum($\mathcal{M}(l)$) = sum($\mathcal{M}(l')$) and hence $||\mathcal{M}(l) - \mathcal{M}(l')||_1$ must be even. Furthermore, the sum of positive entries in ($\mathcal{M}(l) - \mathcal{M}(l')$) equals the sum of negative entries.

order). Without loss of generality assume the fixed q positions for both l and r are the first q indices. The 4-tuple (l, r, l', r') with $\|\mathcal{M}(l) - \mathcal{M}(l')\|_1 = 2q$ can then be enumerated as follows:

$$\sum_{\substack{\mathbf{l},\mathbf{r},\mathbf{l}',\mathbf{r}'\\\mathcal{M}(\mathbf{l})-\mathcal{M}(\mathbf{l}')=\mathcal{M}(\mathbf{r})-\mathcal{M}(\mathbf{r}')\\\|\mathcal{M}(\mathbf{l})-\mathcal{M}(\mathbf{l}')\|_{1}=2q}} t(\mathbf{l},\mathbf{l}',\mathbf{r},\mathbf{r}')$$

$$=\sum_{i\in[n]^{q}}\sum_{\substack{\mathbf{j}\in[n]^{q}\\\mathbf{r}\in[n]^{p-q}}} \sum_{\mathbf{l}\in[n]^{p-q}} t(\mathbf{i}\circ\mathbf{l},\mathcal{L}(\mathbf{j}\circ\mathbf{l}),\mathbf{i}\circ\mathbf{r},\mathcal{L}(\mathbf{j}\circ\mathbf{r}))$$

$$\leq \frac{1}{b}\sum_{\substack{\mathbf{i},\mathbf{j}\in[n]^{q}\\\mathbf{l},\mathbf{r}\in[n]^{p-q}}} \mathbf{A}_{i\circ\mathbf{l}}\mathbf{A}_{i\circ\mathbf{r}}\mathbf{B}_{\mathbf{j}\circ\mathbf{l}}\mathbf{B}_{\mathbf{j}\circ\mathbf{r}}$$

$$=\frac{1}{b}\sum_{\substack{\mathbf{i},\mathbf{j}\in[n]^{q}\\\mathbf{l},\mathbf{j}\in[n]^{q}}} \langle \mathbf{A}(\mathbf{e}_{i_{1}},\cdots,\mathbf{e}_{i_{q}},\mathbf{I},\cdots,\mathbf{I}), \mathbf{B}(\mathbf{e}_{j_{1}},\cdots,\mathbf{e}_{j_{q}},\mathbf{I},\cdots,\mathbf{I})\rangle^{2}$$

$$\leq \frac{1}{b}\sum_{\substack{\mathbf{i},\mathbf{j}\in[n]^{q}\\\mathbf{l},\mathbf{j}\in[n]^{q}}} \|\mathbf{A}(\mathbf{e}_{i_{1}},\cdots,\mathbf{e}_{i_{q}},\mathbf{I},\cdots,\mathbf{I})\|_{F}^{2} \|\mathbf{B}(\mathbf{e}_{j_{1}},\cdots,\mathbf{e}_{j_{q}},\mathbf{I},\cdots,\mathbf{I})\|_{F}^{2}$$

$$=\frac{\|\mathbf{A}\|_{F}^{2}\|\mathbf{B}\|_{F}^{2}}{b}.$$
(40)

Here \circ denotes concatenation, that is, $\mathbf{i} \circ \mathbf{l} = (i_1, \cdots, i_q, l_1, \cdots, l_{p-q}) \in [n]^p$. The fourth equation is Cauchy-Schwartz inequality. Finally note that there are no more than 4^p ways of assigning q positions to \mathbf{l} and \mathbf{l}' each. Combining Eq. (39) and (40) we get

$$\mathbb{V}[\langle \tilde{\boldsymbol{s}}_{\mathbf{A}}, \tilde{\boldsymbol{s}}_{\widetilde{\mathbf{B}}} \rangle] = \mathbb{E}[\langle \tilde{\boldsymbol{s}}_{\mathbf{A}}, \tilde{\boldsymbol{s}}_{\widetilde{\mathbf{B}}} \rangle^2] - \langle \mathbf{A}, \mathbf{B} \rangle^2 \le \frac{4^p \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2}{b},$$

he proof.

which completes the proof.

Proof of Lemma 2. Eq. (32) immediately follows Eq. (28) by adding everything together. For the variance bound we cannot use the same argument because in general the m^2 random variables are neither independent nor uncorrelated. Instead, we compute the variance by definition. First we compute the expected square term as follows:

$$\mathbb{E}\left[\left(\sum_{i,j} w_{i}w_{j}\langle \tilde{\mathbf{s}}_{\mathbf{A}_{i}}, \tilde{\mathbf{s}}_{\tilde{\mathbf{B}}_{j}} \rangle\right)^{2}\right] = \sum_{\substack{i,j,i',j'\\ \mathbf{l},\mathbf{l}',\mathbf{r},\mathbf{r}'}} w_{i}w_{j}w_{i'}w_{j'} \cdot \mathbb{E}[\delta(\mathbf{l},\mathbf{l}')\delta(\mathbf{r},\mathbf{r}')] \cdot \mathbb{E}[\sigma(\mathbf{l})\bar{\sigma}(\mathbf{l}')\bar{\sigma}(\mathbf{r})\sigma(\mathbf{r}')] \cdot [\mathbf{A}_{i}]_{\mathbf{l}}[\mathbf{A}_{i'}]_{\mathbf{r}}[\widetilde{\mathbf{B}}_{j}]_{\mathbf{l}'}[\widetilde{\mathbf{B}}_{j'}]_{\mathbf{r}'}.$$
(41)

Define $\mathbf{X} = \sum_{i} w_i \mathbf{A}_i$ and $\mathbf{Y} = \sum_{i} w_i \mathbf{B}_i$. The above equation can then be simplified as

$$\mathbb{E}\left[\left(\sum_{i,j} w_i w_j \langle \tilde{\mathbf{s}}_{\mathbf{A}_i}, \tilde{\mathbf{s}}_{\widetilde{\mathbf{B}}_j} \rangle\right)^2\right] = \sum_{\boldsymbol{l}, \boldsymbol{l}', \boldsymbol{r}, \boldsymbol{r}'} \mathbb{E}[\delta(\boldsymbol{l}, \boldsymbol{l}')\delta(\boldsymbol{r}, \boldsymbol{r}')] \cdot \mathbb{E}[\sigma(\boldsymbol{l})\bar{\sigma}(\boldsymbol{l}')\bar{\sigma}(\boldsymbol{r})\sigma(\boldsymbol{r}')] \cdot \mathbf{X}_{\boldsymbol{l}} \mathbf{X}_{\boldsymbol{r}} \widetilde{\mathbf{Y}}_{\boldsymbol{l}'} \widetilde{\mathbf{Y}}_{\boldsymbol{r}'}$$
(42)

-

Applying Lemma 1 we have

-

$$\mathbb{V}\left[\sum_{i,j} w_i w_j \langle \tilde{\mathbf{s}}_{\mathbf{A}_i}, \tilde{\mathbf{s}}_{\widetilde{\mathbf{B}}_j} \rangle\right] \le \frac{4^p \|\mathbf{X}\|_F^2 \|\mathbf{Y}\|_F^2}{b}.$$
(43)

Finally, note that

$$\|\mathbf{X}\|_{F}^{2} = \sum_{i,j} w_{i} w_{j} \langle \mathbf{A}_{i}, \mathbf{A}_{j} \rangle \leq \sum_{i,j} w_{i} w_{j} \|\mathbf{A}_{i}\|_{F} \|\mathbf{A}_{j}\|_{F} \leq \|\boldsymbol{w}\|^{2} \max_{i} \|\mathbf{A}_{i}\|_{F}^{2}.$$
(44)

With Lemma 1 and 2, we can easily prove Theorem 1.

Proof of Theorem 1. First we prove the $\varepsilon_1(u)$ bound. Let $\mathbf{A} = \mathbf{T}$ and $\mathbf{B} = u^{\otimes 3}$. Note that $\|\mathbf{A}\|_F = \|\mathbf{T}\|_F$ and $\|\mathbf{B}\|_F = \|\mathbf{u}\|^2 = 1$. Note that $[\mathbf{T}(\mathbf{I}, u, u)]_i = \mathbf{T}(e_i, u, u)$. Next we consider $\varepsilon_2(u)$ and let $\mathbf{A} = \mathbf{T}, \mathbf{B} = e_i \otimes u \otimes u$. Again we have $\|\mathbf{A}\|_F = \|\mathbf{T}\|_F$ and $\|\mathbf{B}\|_F = 1$. A union bound over all $i = 1, \dots, n$ yields the result. For the inequality involving w we apply Lemma 2.

E.3 Analysis of fast robust tensor power method

In this section, we prove Theorem 3, a more refined version of Theorem 2 in Section 4.2. We structure the section by first demonstrating the convergence behavior of noisy tensor power method, and then show how error accumulates with deflation. Finally, the overall bound is derived by combining these two parts.

E.3.1 Recovering the principal eigenvector

Define the angle between two vectors v and u to be $\theta(v, u)$. First, in Lemma 3 we show that if the initialization vector u_0 is randomly chosen from the unit sphere, then the angle θ between the iteratively updated vector u_t and the largest eigenvector of tensor \mathbf{T} , v_1 , will decrease to a point that $\tan \theta(v_1, u_t) < 1$. Afterwards, in Lemma 4, we use a similar approach as in [35] to prove that the error between the final estimation and the ground truth is bounded.

Suppose \mathbf{T} is the exact low-rank ground truth tensor and Each noisy tensor update can then be written as

$$\tilde{\boldsymbol{u}}_{t+1} = \mathbf{T}(\mathbf{I}, \boldsymbol{u}_t, \boldsymbol{u}_t) + \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t), \tag{45}$$

where $\tilde{\varepsilon}(u_t) = \mathbf{E}(\mathbf{I}, u_t, u_t) + \varepsilon_{2,T}(u_t)$ is the noise coming from statistical and tensor sketch approximation error.

Before presenting key lemmas, we first define γ -separation, a concept introduced in [1].

Definition 1 (γ -separation, [1]). Fix $i^* \in [k]$, $u \in \mathbb{R}^n$ and $\gamma > 0$. u is γ -separated with respect to v_{i^*} if the following holds:

$$\lambda_{i^*} \langle \boldsymbol{u}, \boldsymbol{v}_{i^*} \rangle - \max_{i \in [k] \setminus \{i^*\}} \lambda_i \langle \boldsymbol{u}, \boldsymbol{v}_i \rangle \ge \gamma \lambda_{i^*} \langle \boldsymbol{u}, \boldsymbol{v}_{i^*} \rangle.$$
(46)

Lemma 3 analyzes the first phase of the noisy tensor power algorithm. It shows that if the initialization vector u_0 is γ -separated with respect to v_1 and the magnitude of noise $\tilde{\varepsilon}(u_t)$ is small at each iteration t, then after a short number of iterations we will have inner product between u_t and v_1 at least a constant.

Lemma 3. Let $\{\boldsymbol{v}_1, \boldsymbol{v}_2, \cdots, \boldsymbol{v}_k\}$ and $\{\lambda_1, \lambda_2, \cdots, \lambda_k\}$ be eigenvectors and eigenvalues of tensor $\mathbf{T} \in \mathbb{R}^{n \times n \times n}$, where $\lambda_1 |\langle \boldsymbol{v}_1, \boldsymbol{u}_0 \rangle| = \max_{i \in [k]} \lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_0 \rangle|$. Denote $\mathbf{V} = (\boldsymbol{v}_1, \cdots, \boldsymbol{v}_k) \in \mathbb{R}^{n \times k}$ as the

matrix for eigenvectors. Suppose that for every iteration t the noise satisfies

$$\left| \langle \boldsymbol{v}_i, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle \right| \le \epsilon_1 \quad \forall i \in [n] \quad and \quad \left\| \mathbf{V}^\top \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \right\| \le \epsilon_2; \tag{47}$$

suppose also the initialization u_0 is γ -separated with respect to v_1 for some $\gamma \in (0.5, 1)$. If $\tan \theta (v_1, u_0) > 1$, and

$$\epsilon_{1} \leq \min\left(\frac{1}{4\frac{\max_{i \in [k]}\lambda_{i}}{\lambda_{1}} + 2}, \frac{1 - (1 + \alpha)/2}{2}\right)\lambda_{1}\left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0}\right\rangle^{2} \text{ and } \epsilon_{2} \leq \frac{1 - (1 + \alpha)/2}{2\sqrt{2}(1 + \alpha)}\lambda_{1}\left|\left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0}\right\rangle\right|$$

$$(48)$$

for some $\alpha > 0$, then for a small constant $\rho > 0$, there exists a $T > \log_{1+\alpha} (1+\rho) \tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_0)$ such that after T iteration, we have $\tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_T) < \frac{1}{1+\rho}$,

18

Proof. Let $\tilde{u}_{t+1} = \mathbf{T}(\mathbf{I}, u_t, u_t) + \tilde{\varepsilon}(u_t)$ and $u_{t+1} = \tilde{u}_{t+1} / \|\tilde{u}_{t+1}\|$. For $\alpha \in (0, 1)$, we try to prove that there exists a T such that for t > T

$$\frac{1}{\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{t+1}\right)} = \frac{|\langle \boldsymbol{v}_{1},\boldsymbol{u}_{t+1}\rangle|}{\left(1-\langle \boldsymbol{v}_{1},\boldsymbol{u}_{t+1}\rangle^{2}\right)^{1/2}} = \frac{|\langle \boldsymbol{v}_{1},\tilde{\boldsymbol{u}}_{t+1}\rangle|}{\left(\sum\limits_{i=2}^{n}\langle \boldsymbol{v}_{i},\tilde{\boldsymbol{u}}_{t+1}\rangle^{2}\right)^{1/2}} \ge 1.$$
(49)

First we examine the numerator. Using the assumption $|\langle \boldsymbol{v}_i, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle| \leq \epsilon_1$ and the fact that $\langle \boldsymbol{v}_i, \tilde{\boldsymbol{u}}_{t+1} \rangle = \lambda_i \langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle^2 + \langle \boldsymbol{v}_i, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle$, we have

$$|\langle \boldsymbol{v}_{i}, \tilde{\boldsymbol{u}}_{t+1} \rangle| \geq \lambda_{i} \langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle^{2} - \epsilon_{1} \geq |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle| \left(\lambda_{i} \left| \langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle \right| - \epsilon_{1} / \left| \langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle \right| \right).$$
(50) nominator by Hölder's inequality we have

For the denominator, by Hölder's inequality we have 1/2

$$\left(\sum_{i=2}^{n} \langle \boldsymbol{v}_i, \tilde{\boldsymbol{u}}_{t+1} \rangle^2 \right)^{1/2} = \left(\sum_{i=2}^{n} \left(\lambda_i \langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle^2 + \langle \boldsymbol{v}_i, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle \right)^{1/2} \right)$$
(51)

$$\leq \left(\sum_{i=2}^{n} \lambda_{i}^{2} \langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle^{4}\right)^{1/2} + \left(\sum_{i=2}^{n} \langle \boldsymbol{v}_{i}, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_{t}) \rangle^{2}\right)^{1/2}$$
(52)

$$\leq \max_{i \neq 1} \lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle| \left(\sum_{i=2}^n \langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle^2 \right)^{1/2} + \epsilon_2$$
(53)

$$\leq \left(1 - \langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle^2\right)^{1/2} \left(\max_{i \neq 1} \lambda_i \left| \langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle \right| + \epsilon_2 / \left(1 - \langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle^2\right)^{1/2}\right)$$
(54)

Equation (50) and (51) yield

the last

$$\frac{1}{\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{t+1}\right)} \geq \frac{\left|\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle\right|}{\left(1-\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle^{2}\right)^{1/2}} \frac{\lambda_{1}\left|\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle\right| - \epsilon_{1}/\left|\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle\right|}{\max_{i\neq1}\lambda_{i}\left|\left\langle\boldsymbol{v}_{i},\boldsymbol{u}_{t}\right\rangle\right| + \epsilon_{2}/\left(1-\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle^{2}\right)^{1/2}} \qquad (55)$$

$$= \frac{1}{\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right)} \frac{\lambda_{1}\left|\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle\right| - \epsilon_{1}/\left|\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle\right|}{\max_{i\neq1}\lambda_{i}\left|\left\langle\boldsymbol{v}_{i},\boldsymbol{u}_{t}\right\rangle\right| + \epsilon_{2}/\left(1-\left\langle\boldsymbol{v}_{1},\boldsymbol{u}_{t}\right\rangle^{2}\right)^{1/2}} \qquad (56)$$

To prove that the second term is larger than $1 + \alpha$, we first show that when t = 0, the inequality holds. Since the initialization vector is a γ -separated vector, we have

$$\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle| - \max_{i \in [k]} \lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \rangle| \geq \gamma \lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle|, \qquad (57)$$
$$\max_{i \in [k]} \lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \rangle| \leq (1 - \gamma) \lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle| \leq 0.5 \lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle|, \qquad (58)$$

the last inequality holds since
$$\gamma > 0.5$$
. Note that we assume $\tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_0) > 1$ and hence $\langle \boldsymbol{v}_1, \boldsymbol{u}_0 \rangle^2 < 0.5$. Therefore,

$$\epsilon_{2} \leq \frac{1 - (1 + \alpha)/2}{2\sqrt{2}(1 + \alpha)} \lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle| \leq \frac{\left(1 - \langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle^{2}\right)^{1/2} (1 - (1 + \alpha)/2)}{2(1 + \alpha)} \lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle|.$$
(59)

Thus, for t = 0, using the condition for ϵ_1 and ϵ_2 we have

$$\frac{\lambda_{1} \left| \left\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \right\rangle \right| - \epsilon_{1} / \left| \left\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \right\rangle \right|}{\max_{i \neq 1} \lambda_{i} \left| \left\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \right\rangle \right| + \epsilon_{2} / \left(1 - \left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \right\rangle^{2} \right)^{1/2}} \geq \frac{\lambda_{1} \left| \left\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \right\rangle \right| - \epsilon_{1} / \left| \left\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \right\rangle \right|}{0.5\lambda_{1} \left| \left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \right\rangle \right| + \epsilon_{2} / \left(1 - \left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \right\rangle^{2} \right)^{1/2}} \geq 1 + \alpha.$$
(60)

The result yields $1/\tan\theta(\boldsymbol{v}_1, \boldsymbol{u}_1) > (1+\alpha)/\tan\theta(\boldsymbol{v}_1, \boldsymbol{u}_0)$. This also indicates that $|\langle \boldsymbol{v}_1, \boldsymbol{u}_1 \rangle| > (1+\alpha)/\tan\theta(\boldsymbol{v}_1, \boldsymbol{u}_1)| > (1+\alpha)/\tan\theta(\boldsymbol{v}_1, \boldsymbol{u}_1)$ $|\langle m{v}_1,m{u}_0
angle|$, which implies that

$$\epsilon_{1} \leq \min\left(\frac{1}{4\frac{\max_{i \in [k]}\lambda_{i}}{\lambda_{1}} + 2}, \frac{1 - (1 + \alpha)/2}{2}\right)\lambda_{1}\left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t}\right\rangle^{2} \text{ and } \epsilon_{2} \leq \frac{1 - (1 + \alpha)/2}{2\sqrt{2}(1 + \alpha)}\lambda_{1}\left|\left\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t}\right\rangle\right|$$

$$(61)$$

also holds for t = 1. Next we need to make sure that for $t \ge 0$

$$\max_{i \neq 1} \lambda_i \left| \langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle \right| \le 0.5 \lambda_1 \left| \langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle \right|.$$
(62)

In other words, we need to show that $\frac{\lambda_1 |\langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle|}{\max_{i \neq 1} \lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle|} \geq 2$. From Equation (58), for t = 0, $\tfrac{\lambda_1|\langle \boldsymbol{v}_1, \boldsymbol{u}_t\rangle|}{\max\limits_{i\neq 1}\lambda_i|\langle \boldsymbol{v}_i, \boldsymbol{u}_t\rangle|} \geq \tfrac{1}{1-\gamma} \geq 2. \text{ For every } i\in[k],$

$$|\langle \boldsymbol{v}_i, \tilde{\boldsymbol{u}}_{t+1} \rangle| \leq \lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle|^2 + \epsilon_1 \leq |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle| \left(\lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle| + \epsilon_1 / |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle|\right).$$
(63)
With equation (50), we have

$$\frac{\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t+1} \rangle|}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t+1} \rangle|} = \frac{\lambda_{1} |\langle \boldsymbol{v}_{1}, \tilde{\boldsymbol{u}}_{t+1} \rangle|}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \tilde{\boldsymbol{u}}_{t+1} \rangle|} \ge \frac{\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle| \left(\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle| - \frac{\epsilon_{1}}{|\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle|}\right)}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle| \left(\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle| - \frac{\epsilon_{1}}{|\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle|}\right)}$$
(64)

$$= \left(\frac{\lambda_1 |\langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle|}{\lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle|}\right)^2 \frac{1 - \frac{\epsilon_1}{\lambda_1 \langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle^2}}{1 + \frac{\lambda_i}{\lambda_1} \frac{\epsilon_1}{\lambda_1 \langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle^2} \left(\frac{\lambda_1 |\langle \boldsymbol{v}_1, \boldsymbol{u}_t \rangle|}{\lambda_i |\langle \boldsymbol{v}_i, \boldsymbol{u}_t \rangle|}\right)^2} \quad (65)$$

$$\geq \left(\frac{\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle|}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle|}\right)^{2} \frac{1 - \frac{\epsilon_{1}}{\lambda_{1} \langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle^{2}}}{1 + \frac{i \in [k]}{\lambda_{1}} \frac{\epsilon_{1}}{\lambda_{1} \langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle^{2}} \left(\frac{\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{t} \rangle|}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{t} \rangle|}\right)^{2}}$$
(66)

$$=\frac{1-\frac{\epsilon_1}{\lambda_1\langle \boldsymbol{v}_1, \boldsymbol{u}_t\rangle^2}}{\left(\frac{\lambda_1|\langle \boldsymbol{v}_1, \boldsymbol{u}_t\rangle|}{\lambda_i|\langle \boldsymbol{v}_i, \boldsymbol{u}_t\rangle|}\right)^2}+\frac{\max_{i\in[k]}\lambda_i}{\lambda_1}\frac{\epsilon_1}{\lambda_1\langle \boldsymbol{v}_1, \boldsymbol{u}_t\rangle^2}}.$$
(67)

Let $\kappa = \frac{\max_{i \in [k]} \lambda_i}{\lambda_1}$. For t = 0, with conditions on ϵ_1 the following holds:

=

$$\frac{\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{1} \rangle|}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{1} \rangle|} \geq \frac{1 - \frac{\epsilon_{1}}{\lambda_{1} \langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle^{2}}}{\left(\frac{1}{\lambda_{1} |\langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle|}{\lambda_{i} |\langle \boldsymbol{v}_{i}, \boldsymbol{u}_{0} \rangle|}\right)^{2} + \frac{\max_{i \in [k]} \lambda_{i}}{\lambda_{1}} \frac{\epsilon_{1}}{\lambda_{1} \langle \boldsymbol{v}_{1}, \boldsymbol{u}_{0} \rangle^{2}}}.$$

$$\geq \frac{1 - \frac{1}{4\kappa + 2}}{\frac{1}{4} + \frac{\kappa}{4\kappa + 2}} = 2$$
(69)

With the two conditions stated in Equation (61), following the same step in (60), we have $\frac{1}{\tan\theta(\boldsymbol{v}_1,u_2)} \ge (1+\alpha)\frac{1}{\tan\theta(\boldsymbol{v}_1,u_1)}$. By induction, $\frac{1}{\tan\theta(\boldsymbol{v}_1,u_{t+1})} \ge (1+\alpha)\frac{1}{\tan\theta(\boldsymbol{v}_1,t)}$. for $t \ge 0$. Subsequently,

$$\frac{1}{\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{T}\right)} \geq (1+\alpha)^{T} \frac{1}{\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{0}\right)}.$$
(70)
of by setting $T > \log_{1+\alpha}\left(1+\rho\right)\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{0}\right).$

Finally, we complete the proof by setting $T > \log_{1+\alpha} (1+\rho) \tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_0)$.

Next, we present Lemma 4, which analyzes the second phase of the noisy tensor power method. The second phase starts with $\tan \theta(v_1, u_0) < 1$, that is, the inner product of v_1 and u_0 is lower bounded by 1/2.

Lemma 4. Let v_1 be the principal eigenvector of a tensor T and let u_0 be an arbitrary vector in \mathbb{R}^d that satisfies $\tan \theta(v_1, u_0) < 1$. Suppose at every iteration t the noise satisfies

$$4\|\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t)\| \le \epsilon \left(\lambda_1 - \lambda_2\right) \text{ and } 4|\langle \boldsymbol{v}_1, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t)\rangle| \le \left(\lambda_1 - \lambda_2\right)\cos^2\theta\left(\boldsymbol{v}_1, \boldsymbol{u}_0\right)$$
(71)

for some $\epsilon < 1$. Then with high probability there exists $T = O\left(\frac{\lambda_1}{\lambda_1 - \lambda_2}\log(1/\epsilon)\right)$ such that after T *iteration we have* $\tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_T) \leq \epsilon$.

Proof. Define $\Delta := \frac{\lambda_1 - \lambda_2}{4}$ and $\mathbf{X} := \mathbf{v}_1^{\perp}$. We have the following chain of inequalities:

$$\tan\theta\left(\boldsymbol{v}_{1},\mathbf{T}\left(\mathbf{I},\boldsymbol{u},\boldsymbol{u}\right)+\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}\right)\right)\leq\frac{\left\|\mathbf{X}^{T}\left(\mathbf{T}\left(\mathbf{I},\boldsymbol{u},\boldsymbol{u}\right)+\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u})\right)\right\|}{\left\|\boldsymbol{v}_{1}^{T}\left(\mathbf{T}\left(\mathbf{I},\boldsymbol{u},\boldsymbol{u}\right)+\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u})\right)\right\|}$$
(72)

$$\leq \frac{\left\|\mathbf{X}^{T}\mathbf{T}\left(\mathbf{I},\boldsymbol{u},\boldsymbol{u}\right)\right\| + \left\|\mathbf{V}^{T}\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u})\right\|}{\left\|\boldsymbol{v}_{1}^{T}\mathbf{T}\left(\mathbf{I},\boldsymbol{u},\boldsymbol{u}\right)\right\| - \left\|\boldsymbol{v}_{1}^{T}\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u})\right\|}$$
(73)

$$\leq \frac{\lambda_2 \left\| \mathbf{X}^T \boldsymbol{u} \right\|^2 + \left\| \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}) \right\|}{\lambda_1 \left| \boldsymbol{v}_1^T \boldsymbol{u} \right|^2 - \left| \boldsymbol{v}_1^\top \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}) \right|}$$
(74)

$$=\frac{\left\|\mathbf{X}^{T}\boldsymbol{u}\right\|^{2}}{\left|\boldsymbol{v}_{1}^{T}\boldsymbol{u}\right|^{2}}\frac{\lambda_{2}}{\lambda_{1}-\frac{\left|\boldsymbol{v}_{1}^{\top}\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u})\right|}{\left|\boldsymbol{v}_{1}^{\top}\boldsymbol{u}\right|^{2}}}+\frac{\frac{\left\|\boldsymbol{\varepsilon}(\boldsymbol{u})\right\|}{\left|\boldsymbol{v}_{1}^{\top}\boldsymbol{u}\right|^{2}}}{\lambda_{1}-\frac{\left|\boldsymbol{v}_{1}^{\top}\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u})\right|}{\left|\boldsymbol{v}_{1}^{\top}\boldsymbol{u}\right|^{2}}}$$
(75)

$$\leq \tan^{2} \theta(\boldsymbol{v}_{1}, \boldsymbol{u}) \frac{\lambda_{2}}{\lambda_{2} + 3\Delta} + \frac{\Delta \epsilon \left(1 + \tan^{2} \theta(\boldsymbol{v}_{1}, \boldsymbol{u})\right)}{\lambda_{2} + 3\Delta}$$
(76)

$$\leq \max\left(\epsilon, \frac{\lambda_2 + \Delta\epsilon}{\lambda_2 + 2\Delta} \tan^2 \theta\left(\boldsymbol{v}_1, \boldsymbol{u}\right)\right) \tag{77}$$

$$\leq \max\left(\epsilon, \frac{\lambda_2 + \Delta\epsilon}{\lambda_2 + 2\Delta} \tan\theta\left(\boldsymbol{v}_1, \boldsymbol{u}\right)\right)$$
(78)

The second step follows by triangle inequality. For $u = u_0$, using the condition $\tan(v_1, u_0) < 1$ we obtain

$$\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{1}\right) \leq \max\left(\epsilon,\frac{\lambda_{2}+\Delta\epsilon}{\lambda_{2}+2\Delta}\tan^{2}\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}\right)\right) \leq \max\left(\epsilon,\frac{\lambda_{2}+\Delta\epsilon}{\lambda_{2}+2\Delta}\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}\right)\right) \quad (79)$$
Since $\frac{\lambda_{2}+\Delta\epsilon}{\lambda_{2}+2\Delta} \leq \max\left(\frac{\lambda_{2}}{\lambda_{2}}+\epsilon\right) \leq (\lambda_{2}/\lambda_{1})^{1/4} \leq 1$ we have

Since
$$\lambda_{2}+2\Delta \leq \max\left(\lambda_{2}+\Delta,\epsilon\right) \leq (\lambda_{2}/\lambda_{1}) \leq 1$$
, we have
 $\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{1}\right) = \tan\theta\left(\boldsymbol{v}_{1},\mathbf{T}\left(\mathbf{I},\boldsymbol{u}_{0},\boldsymbol{u}_{0}\right) + \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_{t})\right) \leq \max\left(\epsilon,(\lambda_{2}/\lambda_{1})^{1/4}\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{0}\right)\right) < 1.$
(80)

By induction,

$$\tan \theta \left(\boldsymbol{v}_1, \boldsymbol{u}_{t+1} \right) = \tan \theta \left(\boldsymbol{v}_1, \mathbf{T} \left(\mathbf{I}, \boldsymbol{u}_t, \boldsymbol{u}_t \right) + \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \right) \le \max \left(\epsilon, (\lambda_2 / \lambda_1)^{1/4} \tan \theta \left(\boldsymbol{v}_1, \boldsymbol{u}_t \right) \right) < 1.$$
for every *t*. Eq. (78) then yields

$$\tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{T}\right) \leq \max\left(\epsilon, \max\epsilon, \left(\lambda_{2}/\lambda_{1}\right)^{L/4} \tan\theta\left(\boldsymbol{v}_{1},\boldsymbol{u}_{0}\right)\right).$$
(81)

Consequently, after
$$T = \log_{(\lambda_2/\lambda_1)^{-1/4}}(1/\epsilon)$$
 iterations we have $\tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_T) \leq \epsilon$.

Lemma 5. Suppose v_1 is the principal eigenvector of a tensor \mathbf{T} and let $u_0 \in \mathbb{R}^n$. For some $\alpha, \rho > 0$ and $\epsilon < 1$, if at every step, the noise satisfies

$$\|\tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t)\| \leq \epsilon \frac{\lambda_1 - \lambda_2}{4} \quad and \quad \left| \langle \boldsymbol{v}_1, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle \right| \leq \min\left(\frac{1}{4\frac{\max_{i \in [k]} \lambda_i}{\lambda_1} + 2} \lambda_1, \frac{1 - (1 + \alpha)/2}{2\sqrt{2}(1 + \alpha)} \lambda_1\right) \frac{1}{\tau^2 n},\tag{82}$$

then with high probability there exists an $T = O\left(\log_{1+\alpha}\left(1+\rho\right)\tau\sqrt{n} + \frac{\lambda_1}{\lambda_1-\lambda_2}\log(1/\epsilon)\right)$ such that after T iterations we have $\left\|\left(\boldsymbol{I} - \boldsymbol{u}_T\boldsymbol{u}_T^T\right)\boldsymbol{v}_1\right\| \leq \epsilon$.

Proof. By Lemma 2.5 in [35], for any fixed orthonormal matrix **V** and a random vector \boldsymbol{u} , we have $\max_{i \in [K]} \tan \theta(\boldsymbol{v}_i, \boldsymbol{u}) \leq \tau \sqrt{n}$ with all but $O(\tau^{-1} + e^{-\Omega(d)})$ probability. Using the fact that $\cos \theta(\boldsymbol{v}_1, \boldsymbol{u}_0) \geq 1/(1 + \tan \theta(\boldsymbol{v}_1, \boldsymbol{u}_0)) \geq \frac{1}{\tau \sqrt{n}}$, the following bounds on the noise level imply the conditions in Lemma 3:

$$\begin{split} \left\| \mathbf{V}^T \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \right\| &\leq \frac{1 - (1 + \alpha)/2}{2\sqrt{2}(1 + \alpha)\tau\sqrt{n}} \text{ and } \left| \langle \boldsymbol{v}_1, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle \right| \\ &\leq \min\left(\frac{1}{4\frac{\max_{i \in [k]} \lambda_i}{\lambda_1} + 2} \lambda_1, \frac{1 - (1 + \alpha)/2}{2} \lambda_1 \right) \frac{1}{\tau^2 n}, \quad \forall t. \end{split}$$

Note that $|\langle \boldsymbol{v}_1, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle| \leq \frac{1-(1+\alpha)/2}{2\sqrt{2}(1+\alpha)} \lambda_1 \frac{1}{\tau^2 n}$ implies the first bound in Eq. (83). In Lemma 4, we assume $\tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_0) < 1$ and prove that for every \boldsymbol{u}_t , $\tan \theta (\boldsymbol{v}_1, \boldsymbol{u}_t) < 1$, which is equivalent to saying that at every step, $\cos \theta (\boldsymbol{v}_1, \boldsymbol{u}_t) > \frac{1}{\sqrt{2}}$. By plugging the inequality into the second condition in Lemma 4, we have $|\langle \boldsymbol{v}_1, \tilde{\boldsymbol{\varepsilon}}(\boldsymbol{u}_t) \rangle| \leq \frac{(\lambda_1 - \lambda_2)}{8}$. The lemma then follows by the fact that $\|(\boldsymbol{I} - \boldsymbol{u}_T \boldsymbol{u}_T^T) \boldsymbol{v}_1\| = \sin \theta (\boldsymbol{u}_T, \boldsymbol{v}_1) \leq \tan \theta (\boldsymbol{u}_T, \boldsymbol{v}_1) \leq \epsilon$.

E.3.2 Deflation

In previous sections we have upper bounded the Euclidean distance between the estimated and the true principal eigenvector of an input tensor **T**. In this section, we show that error introduced from previous tensor power updates can also be bounded. As a result, we obtain error bounds between the entire set of base vectors $\{v_i\}_{i=1}^k$ and their estimation $\{\hat{v}_i\}_{i=1}^k$.

Lemma 6. Let $\{v_1, v_2, \dots, v_k\}$ and $\{\lambda_1, \lambda_2, \dots, \lambda_k\}$ be orthonormal eigenvectors and eigenvalues of an input tensor T. Define $\lambda_{\max} := \max_{i \in [k]} \lambda_i$. Suppose $\{\hat{v}_i\}_{i=1}^k$ and $\{\hat{\lambda}_i\}_{i=1}^k$ are estimated eigenvector/eigenvalue pairs. Fix $\epsilon \ge 0$ and any $t \in [k]$. If

$$|\hat{\lambda}_i - \lambda_i| \le \lambda_i \epsilon/2, \quad and \quad ||\hat{\boldsymbol{u}}_i - \boldsymbol{u}_i|| \le \epsilon$$
(83)

for all $i \in [t]$, then for any unit vector u the following holds:

$$\left\|\sum_{i=1}^{t} \left[\lambda \boldsymbol{v}_{i}^{\otimes 3} - \hat{\lambda}_{i} \hat{\boldsymbol{v}}_{i}^{\otimes 3}\right] (\mathbf{I}, \boldsymbol{u}, \boldsymbol{u})\right\|^{2} \leq 4 \left(2.5\lambda_{\max} + (\lambda_{\max} + 1.5)\epsilon\right)^{2} \epsilon^{2} + 9(1 + \epsilon/2)^{2} \lambda_{\max}^{2} \epsilon^{4}$$

$$\tag{84}$$

$$+8(1+\epsilon/2)^2\lambda_{\max}^2\epsilon^2\tag{85}$$

$$\leq 50\lambda_{\max}^2\epsilon^2.$$
 (86)

Proof. Following similar approaches in [1], Lemma B.5, we define $\hat{\boldsymbol{v}}^{\perp} = \hat{\boldsymbol{v}}_i - (\boldsymbol{v}_i^{\top} \hat{\boldsymbol{v}}_i) \boldsymbol{v}_i$ and $\mathbf{D}_i = \left[\lambda \boldsymbol{v}_i^{\otimes 3} - \hat{\lambda}_i \hat{\boldsymbol{v}}_i^{\otimes 3}\right]$. $\mathbf{D}_i(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u})$ can then be written as the sum of scaled \boldsymbol{v}_i and \boldsymbol{v}_i^{\top} products as follows:

$$\mathbf{D}_{i}\left(\mathbf{I},\boldsymbol{u},\boldsymbol{u}\right) = \lambda_{i}(\boldsymbol{u}^{\top}\boldsymbol{v}_{i})^{2}\boldsymbol{v}_{i} - \hat{\lambda}_{i}(\boldsymbol{u}^{\top}\hat{\boldsymbol{v}}_{i})^{2}\hat{\boldsymbol{v}}_{i}$$
(87)

$$=\lambda_{i}(\boldsymbol{u}^{\top}\boldsymbol{v}_{i})^{2}\boldsymbol{v}_{i}-\hat{\lambda}_{i}(\boldsymbol{u}^{\top}\left(\hat{\boldsymbol{v}}_{i}^{\perp}+(\boldsymbol{v}_{i}^{\top}\hat{\boldsymbol{v}}_{i})\boldsymbol{v}_{i}\right))^{2}\left(\hat{\boldsymbol{v}}^{\perp}+(\boldsymbol{v}_{i}^{\top}\hat{\boldsymbol{v}}_{i})\boldsymbol{v}_{i}\right)$$
(88)

$$= \left(\left(\lambda_i - \hat{\lambda}_i (\boldsymbol{v}_i^{\top} \hat{\boldsymbol{v}}_i)^3 \right) (\boldsymbol{u}^{\top} \boldsymbol{v}_i)^2 - 2\hat{\lambda}_i (\boldsymbol{u}^{\top} \hat{\boldsymbol{v}}_i^{\perp}) (\boldsymbol{v}_i^{\top} \hat{\boldsymbol{v}}_i)^2 (\boldsymbol{u}^{\top} \boldsymbol{v}_i) - \hat{\lambda}_i (\boldsymbol{v}_i^{\top} \hat{\boldsymbol{v}}_i) (\boldsymbol{u}^{\top} \hat{\boldsymbol{v}}^{\perp}) \right) \boldsymbol{v}_i$$

$$-\hat{\lambda}_{i}\left\|\hat{\boldsymbol{v}}_{i}^{\perp}\right\|\left((\boldsymbol{u}^{\top}\boldsymbol{v}_{i})(\boldsymbol{v}_{i}^{\top}\hat{\boldsymbol{v}}_{i})+\boldsymbol{u}^{\top}\hat{\boldsymbol{v}}_{i}^{\perp}\right)\left(\hat{\boldsymbol{v}}_{i}^{\perp}/\left\|\hat{\boldsymbol{v}}_{i}^{\perp}\right\|\right)$$

$$(89)$$

Suppose A_i and B_i are coefficients of v_i and $(\hat{v}_i^{\perp} / \| \hat{v}_i^{\perp} \|)$, respectively. The summation of \mathbf{D}_i can be bounded as

$$\begin{split} \left\| \sum_{i=1}^{t} \mathbf{D}_{i} \left(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u} \right) \right\|^{2} &= \left\| \sum_{i=1}^{t} A_{i} \boldsymbol{v}_{i} - \sum_{i=1}^{t} B_{i} \left(\hat{\boldsymbol{v}}_{i}^{\perp} / \left\| \hat{\boldsymbol{v}}_{i}^{\perp} \right\| \right) \right\|_{2}^{2} \\ &\leq 2 \left\| \sum_{i=1}^{t} A_{i} \boldsymbol{v}_{i} \right\|^{2} + 2 \left\| \sum_{i=1}^{t} B_{i} \left(\hat{\boldsymbol{v}}_{i}^{\perp} / \left\| \hat{\boldsymbol{v}}_{i}^{\perp} \right\| \right) \right\|^{2} \\ &\leq \sum_{i=1}^{t} A_{i}^{2} + 2 \left(\sum_{i=1}^{t} |B_{i}| \right)^{2} \end{split}$$

We then try to upper bound $|A_i|$.

$$|A_{i}| \leq \left| \left(\lambda_{i} - \hat{\lambda}_{i} (\boldsymbol{v}_{i}^{\top} \hat{\boldsymbol{v}}_{i})^{3} \right) (\boldsymbol{u}^{\top} \boldsymbol{v}_{i})^{2} - 2\hat{\lambda}_{i} (\boldsymbol{u}^{\top} \hat{\boldsymbol{v}}_{i}^{\perp}) (\boldsymbol{v}_{i}^{\top} \hat{\boldsymbol{v}}_{i})^{2} (\boldsymbol{u}^{\top} \boldsymbol{v}_{i}) - \hat{\lambda}_{i} (\boldsymbol{v}_{i}^{\top} \hat{\boldsymbol{v}}_{i}) (\boldsymbol{u}^{\top} \hat{\boldsymbol{v}}^{\perp}) \right|$$

$$\leq \left(\lambda_{i} + 1 - (\lambda_{i}^{\top} \hat{\boldsymbol{v}}_{i})^{3} + 1 \right) - \hat{\lambda}_{i} \left(\lambda_{i}^{\top} \hat{\boldsymbol{v}}_{i} \right)^{3} (\lambda_{i}^{\top} \hat{\boldsymbol{v}}_{i})^{2} + 2 \left(\lambda_{i} + 1 \right) - \hat{\lambda}_{i} \left(\lambda_{i}^{\top} \hat{\boldsymbol{v}}_{i} \right) (\boldsymbol{u}^{\top} \hat{\boldsymbol{v}}^{\perp}) \right)$$

$$(90)$$

$$\leq \left(1.5 \left\|\boldsymbol{v}_{i}-\hat{\boldsymbol{v}}_{i}\right\|^{2}+\left|\lambda_{i}-\hat{\lambda}_{i}\right|+2\left(\lambda_{i}+\left|\lambda_{i}-\hat{\lambda}_{i}\right|\right)\left\|\boldsymbol{v}_{i}-\hat{\boldsymbol{v}}_{i}\right\|\right)\left|\boldsymbol{u}^{\top}\boldsymbol{v}_{i}\right| + \left(\lambda_{i}+\left|\lambda_{i}-\hat{\lambda}_{i}\right|\right)\left\|\hat{\boldsymbol{v}}_{i}-\boldsymbol{v}_{i}\right\|^{2}$$

$$(92)$$

$$\leq (2.5\lambda_i + (\lambda_i + 1.5)\epsilon)\epsilon |\boldsymbol{u}^{\top}\boldsymbol{v}_i| + (1 + \epsilon/2)\lambda_i\epsilon^2$$
(93)

Next, we bound $|B_i|$ in a similar manner.

$$|B_i| = \left| \hat{\lambda}_i \left\| \hat{\boldsymbol{v}}_i^{\perp} \right\| \left((\boldsymbol{u}^{\top} \boldsymbol{v}_i) (\boldsymbol{v}_i^{\top} \hat{\boldsymbol{v}}_i) + \boldsymbol{u}^{\top} \hat{\boldsymbol{v}}_i^{\perp} \right) \right|$$
(94)

$$\leq 2\left(\lambda_{i} + \left|\lambda_{i} - \hat{\lambda}_{i}\right|\right) \left\|\hat{\boldsymbol{v}}_{i}^{\perp}\right\| \left((\boldsymbol{u}^{\top}\boldsymbol{v}_{i})^{2} + \left\|\hat{\boldsymbol{v}}_{i}^{\perp}\right\|^{2}\right)$$
(95)

$$\leq 2(1+\epsilon/2)\lambda_i\epsilon(\boldsymbol{u}^{\top}\boldsymbol{v}_i)^2 + 2(1+\epsilon/2)\lambda_i\epsilon^3$$
(96)

Combining everything together we have

$$\begin{aligned} \left\|\sum_{i=1}^{t} \mathbf{D}_{i}\left(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}\right)\right\|^{2} &\leq 2\sum_{i=1}^{t} A_{i}^{2} + 2\left(\sum_{i=1}^{t} |B_{i}|\right)^{2} \end{aligned} \tag{97} \\ &\leq \sum_{i=1}^{t} 4\left(5\lambda_{i} + (\lambda_{i} + 1.5)\right)^{2} \epsilon^{2} \left|\boldsymbol{u}^{\top} \boldsymbol{v}_{i}\right|^{2} + 4(1 + \epsilon/2)^{2} \lambda_{i}^{2} \epsilon^{4} \\ &\quad + 2\left(\sum_{i=1}^{t} 2(1 + \epsilon/2)\lambda_{i}\epsilon(\boldsymbol{u}^{\top} \boldsymbol{v}_{i})^{2} + 2(1 + \epsilon/2)\lambda_{i}\epsilon^{3}\right)^{2} \end{aligned} \tag{98} \\ &\leq 4\left(2.5\lambda_{\max} + (\lambda_{\max} + 1.5)\epsilon\right)^{2} \epsilon^{2} \sum_{i=1}^{t} \left|\boldsymbol{u}^{\top} \boldsymbol{v}_{i}\right|^{2} + 4(1 + \epsilon/2)^{2} \lambda_{\max}^{2} \epsilon^{4} \\ &\quad + 2\left(2(1 + \epsilon/2)\lambda_{\max}\epsilon \sum_{i=1}^{t} (\boldsymbol{u}^{\top} \boldsymbol{v}_{i})^{2} + 2(1 + \epsilon/2)\lambda_{\max}\epsilon^{3}\right)^{2} \end{aligned} \tag{99} \\ &\leq 4\left(2.5\lambda_{\max} + (\lambda_{\max} + 1.5)\epsilon\right)^{2} \epsilon^{2} + 9(1 + \epsilon/2)^{2} \lambda_{\max}^{2} \epsilon^{4} + 8(1 + \epsilon/2)^{2} \lambda_{\max}^{2} \epsilon^{2} . \tag{100} \end{aligned}$$

E.3.3 Main Theorem

In this section we present and prove the main theorem that bounds the reconstruction error of fast robust tensor power method under appropriate settings of the hash length b and number of independent hashes B. The theorem presented below is a more detailed version of Theorem 2 presented in Section 4.2.

Theorem 3. Let $\bar{\mathbf{T}} = \mathbf{T} + \mathbf{E} \in \mathbb{R}^{n \times n \times n}$, where $\mathbf{T} = \sum_{i=1}^{k} \lambda_i \boldsymbol{v}_i^{\otimes 3}$ and $\{\boldsymbol{v}_i\}_{i=1}^{k}$ is an orthonormal basis. Suppose $(\hat{\boldsymbol{v}}_1, \hat{\lambda}_1), (\hat{\boldsymbol{v}}_1, \hat{\lambda}_1), \cdots, (\hat{\boldsymbol{v}}_k, \hat{\lambda}_k)$ is the sequence of estimated eigenvector/eigenvalue pairs obtained using the fast robust tensor power method. Assume $\|\mathbf{E}\| = \epsilon$. There exists constant $C_1, C_2, C_3, \alpha, \rho, \tau \ge 0$ such that the following holds: if

$$\epsilon \le C_1 \frac{1}{n\lambda_{\max}}, \text{ and } T = C_2 \left(\log_{1+\alpha} \left(1+\rho \right) \tau \sqrt{n} + \frac{\lambda_1}{\lambda_1 - \lambda_2} \log(1/\epsilon) \right),$$
 (101)

and

$$\sqrt{\frac{\ln(L/\log_2(k/\eta))}{\ln(k)}} \cdot \left(1 - \frac{\ln\left(\ln L/\log_2(k/\eta)\right) + C_3}{4\ln\left(L/\log_2(k/\eta)\right)} - \sqrt{\frac{\ln(8)}{\ln(L/\log_2(k/\eta))}}\right) \ge 1.02 \left(1 + \sqrt{\frac{\ln(4)}{\ln(k)}}\right) \tag{102}$$

Suppose the tensor sketch randomness is independent among all tensor product evaluations. If $B = \Omega(\log(n/\tau))$ and the hash length b is set to

$$b \ge \left\{ \frac{\|\mathbf{T}\|_{F}^{2} \tau^{4} n^{2}}{\min\left(\frac{1}{4 \max_{i \in [k]} (\lambda_{i}/\lambda_{1})+2} \lambda_{1}, \frac{1-(1+\alpha)/2}{2\sqrt{2}(1+\alpha)} \lambda_{1}\right)^{2}}, \frac{16\epsilon^{-2} \|\mathbf{T}\|_{F}^{2}}{\min_{i \in [k]} (\lambda_{i} - \lambda_{i-1})^{2}}, \epsilon^{-2} \|\mathbf{T}\|_{F}^{2} \right\}$$
(103)

with probability at least $1 - (\eta + \tau^{-1} + e^{-n})$, there exists a permutation π on k such that

$$\left\|\boldsymbol{v}_{\pi(j)} - \hat{\boldsymbol{v}}_{i}\right\| \le \epsilon, \ \left|\lambda_{\pi(j)} - \hat{\lambda}_{j}\right| \le \frac{\lambda_{\pi(j)}\epsilon}{2}, \ and \ \left\|\mathbf{T} - \sum_{j=1}^{k} \hat{\lambda}_{j} \hat{\boldsymbol{v}}_{j}^{\otimes 3}\right\| \le c\epsilon, \tag{104}$$

for some absolute constant c.

Proof. We prove that at the end of each iteration $i \in [k]$, the following conditions hold

- 1. For all $j \leq i$, $|\boldsymbol{v}_{\pi(j)} \hat{\boldsymbol{v}}_j| \leq \epsilon$ and $|\lambda_{\pi(j)} \hat{\lambda}_j| \leq \frac{\lambda_i \epsilon}{2}$
- 2. The tensor error satisfies

$$\left\| \left[\left(\tilde{\mathbf{T}} - \sum_{j \le i} \hat{\lambda}_j \hat{\boldsymbol{v}}_j^{\otimes 3} \right) - \sum_{j \ge i+1} \lambda_{\pi(j)} v_{\pi(j)}^{\otimes 3} \right] (\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}) \right\| \le 56\epsilon$$
(105)

First, we check the case when i = 0. For the tensor error, we have

$$\left\| \left\| \left[\tilde{\mathbf{T}} - \sum_{j=1}^{K} \lambda_{\pi(j)} v_{\pi(j)}^{\otimes 3} \right] (\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}) \right\| = \| \boldsymbol{\varepsilon}(\boldsymbol{u}) \| \le \| \boldsymbol{\varepsilon}_{2,T}(\boldsymbol{u}) \| + \| \mathbf{E}(\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}) \| \le \epsilon + \epsilon = 2\epsilon.$$
(106)

The last inequality follows Theorem 1 with the condition for *b*. Next, Using Lemma 5, we have that $\|\boldsymbol{v}_{\pi(1)} - \hat{\boldsymbol{v}}_1\| \leq \epsilon.$ (107)

In addition, conditions for hash length b and Theorem 1 yield

$$\left|\lambda_{\pi(1)} - \hat{\lambda}_{1}\right| \leq \|\boldsymbol{\varepsilon}_{1,T}(\boldsymbol{v}_{1})\| + \|\mathbf{T}(\hat{\boldsymbol{v}}_{1} - \boldsymbol{v}_{1}, \hat{\boldsymbol{v}}_{1} - \boldsymbol{u}, \hat{\boldsymbol{v}}_{1} - \boldsymbol{v}_{1})\| \leq \epsilon \frac{\lambda_{i} - \lambda_{i-1}}{4} + \epsilon^{3} \|\mathbf{T}\|_{F} \leq \frac{\epsilon \lambda_{i}}{2}$$
(108)

Thus, we have proved that for i = 1 both conditions hold. Assume the conditions hold up to i = t-1 by induction. For the *t*th iteration, the following holds:

$$\left\| \begin{bmatrix} \left(\tilde{\mathbf{T}} - \sum_{j \le t} \hat{\lambda}_j \hat{\boldsymbol{v}}_j^{\otimes 3} \right) - \sum_{j \ge t+1} \lambda_{\pi(j)} v_{\pi(j)}^{\otimes 3} \end{bmatrix} (\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}) \right\|$$

$$\leq \left\| \begin{bmatrix} \tilde{\mathbf{T}} - \sum_{j=1}^K \lambda_{\pi(j)} v_{\pi(j)}^{\otimes 3} \end{bmatrix} (\mathbf{I}, \boldsymbol{u}, \boldsymbol{u}) \right\| + \left\| \sum_{j=1}^t \hat{\lambda}_j \hat{\boldsymbol{v}}_j^{\otimes 3} - \lambda_{\pi(j)} v_{\pi(j)}^{\otimes 3} \right\| \le \epsilon + \sqrt{50} \lambda_{\max} \epsilon.$$

For the last inequality we apply Lemma 6. Since the condition is satisfied, Lemma 5 yields

$$\left\|\boldsymbol{v}_{\pi(t+1)} - \hat{\boldsymbol{v}}_{t+1}\right\| \le \epsilon. \tag{109}$$

Finally, conditions for hash length b and Theorem 1 yield

$$\left|\lambda_{\pi(t+1)} - \hat{\lambda}_{t+1}\right| \leq \|\boldsymbol{\varepsilon}_{1,T}(\boldsymbol{v}_1)\| + \|\mathbf{T}(\hat{\boldsymbol{v}}_t - \boldsymbol{v}_1, \hat{\boldsymbol{v}}_1 - \boldsymbol{u}, \hat{\boldsymbol{v}}_1 - \boldsymbol{v}_1)\| \\ \leq \epsilon \frac{\lambda_i - \lambda_{i-1}}{4} + \epsilon^3 \|\mathbf{T}\|_F \leq \frac{\epsilon \lambda_i}{2} \quad (110)$$

Appendix F Summary of notations for matrix/vector products

We assume vectors $a, b \in \mathbb{C}^n$ are indexed starting from 0; that is, $a = (a_0, a_1, \dots, a_{n-1})$ and $b = (b_0, b_1, \dots, b_{n-1})$. Matrices A, B and tensors T are still indexed starting from 1.

Element-wise product For $a, b \in \mathbb{C}^n$, the element-wise product (Hadamard product) $a \circ b \in \mathbb{R}^n$ is defined as

$$\boldsymbol{a} \circ \boldsymbol{b} = (a_0 b_0, a_1 b_1, \cdots, a_{n-1} b_{n-1}).$$
 (111)

Convolution For $a, b \in \mathbb{C}^n$, their convolution $a * b \in \mathbb{C}^n$ is defined as

$$\boldsymbol{a} \ast \boldsymbol{b} = \left(\sum_{(i+j) \mod n=0} a_i b_j, \sum_{(i+j) \mod n=1} a_i b_j, \cdots, \sum_{(i+j) \mod n=n-1} a_i b_j\right).$$
(112)

Inner product For $a, b \in \mathbb{C}^n$, their inner product is defined as

$$\langle \boldsymbol{a}, \boldsymbol{b} \rangle = \sum_{i=1}^{n} a_i \overline{b_i},\tag{113}$$

`

where $\overline{b_i}$ denotes the complex conjugate of b_i . For tensors $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n \times n}$, their inner product is defined similarly as

$$\langle \mathbf{A}, \mathbf{B} \rangle = \sum_{i,j,k=1}^{n} \mathbf{A}_{i,j,k} \overline{\mathbf{B}}_{i,j,k}.$$
 (114)

Tensor product For $a, b \in \mathbb{C}^n$, the tensor product $a \otimes b$ can be either an $n \times n$ matrix or a vector of length n^2 . For the former case, we have

$$\boldsymbol{a} \otimes \boldsymbol{b} = \begin{bmatrix} a_0 b_0 & a_0 b_1 & \cdots & a_0 b_{n-1} \\ a_1 b_0 & a_1 b_1 & \cdots & a_1 b_{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n-1} b_0 & a_{n-1} b_1 & \cdots & a_{n-1} b_{n-1} \end{bmatrix}.$$
(115)

If $a \otimes b$ is a vector, it is defined as the expansion of the output matrix. That is,

$$\boldsymbol{a} \otimes \boldsymbol{b} = (a_0 b_0, a_0 b_1, \cdots, a_0 b_{n-1}, a_1 b_0, a_1 b_1, \cdots, a_{n-1} b_{n-1}).$$
(116)

Suppose T is an $n \times n \times n$ tensor and matrices $\mathbf{A} \in \mathbb{R}^{n \times m_1}$, $\mathbf{B} \in \mathbb{R}^{n \times m_2}$ and $\mathbf{C} \in \mathbb{R}^{n \times m_3}$. The tensor product $\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})$ is an $m_1 \times m_2 \times m_3$ tensor defined by

$$\left[\mathbf{T}(\mathbf{A}, \mathbf{B}, \mathbf{C})\right]_{i,j,k} = \sum_{i',j',k'=1}^{n} \mathbf{T}_{i',j',k'} \mathbf{A}_{i',i} \mathbf{B}_{j',j} \mathbf{C}_{k',k}.$$
(117)

Khatri-Rao product For $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times m}$, their Khatri-Rao product $\mathbf{A} \odot \mathbf{B} \in \mathbb{C}^{n^2 \times m}$ is defined as

$$\mathbf{A} \odot \mathbf{B} = (\mathbf{A}_{(1)} \otimes \mathbf{B}_{(1)}, \mathbf{A}_{(2)} \otimes \mathbf{B}_{(2)}, \cdots, \mathbf{A}_{(m)} \otimes \mathbf{B}_{(m)}),$$
(118)

where $A_{(i)}$ and $B_{(i)}$ denote the *i*th rows of A and B.

Mode expansion For a tensor **T** of dimension $n \times n \times n$, its first mode expansion $\mathbf{T}_{(1)} \in \mathbb{R}^{n \times n}$ is defined as

$$\mathbf{T}_{(1)} = \begin{bmatrix} \mathbf{T}_{1,1,1} & \mathbf{T}_{1,1,2} & \cdots & \mathbf{T}_{1,1,n} & \mathbf{T}_{1,2,1} & \cdots & \mathbf{T}_{1,n,n} \\ \mathbf{T}_{2,1,1} & \mathbf{T}_{2,1,2} & \cdots & \mathbf{T}_{2,1,n} & \mathbf{T}_{2,2,1} & \cdots & \mathbf{T}_{2,n,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{T}_{n,1,1} & \mathbf{T}_{n,1,2} & \cdots & \mathbf{T}_{n,1,n} & \mathbf{T}_{n,2,1} & \cdots & \mathbf{T}_{n,n,n} \end{bmatrix}.$$
 (119)

The mode expansions $T_{(2)}$ and $T_{(3)}$ can be similarly defined.

References

- [31] A. Anandkumar, R. Ge, D. Hsu, S. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- [32] B. Bader, T. Kolda, et al. MATLAB tensor toolbox version 2.5. Available online, 2012.
- [33] B. W. Bader and T. G. Kolda. Algorithm 862: MATLAB tensor classes for fast algorithm prototyping. *ACM Transactions on Mathematical Software*, 32(4):635–653, 2006.
- [34] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *Journal of machine Learning research*, 3:993–1022, 2003.
- [35] M. Hardt and E. Price. The noisy power method: A meta algorithm with applications. In *NIPS*, 2014.
- [36] T. Kolda and B. Bader. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- [37] C. Wang, X. Liu, Y. Song, and J. Han. Scalable moment-based inference for latent dirichlet allocation. In *ECML/PKDD*, 2014.
- [38] Y. Wang and J. Zhu. Spectral methods for supervised topic models. In NIPS, 2014.