

Supplementary Material of SGD_Tucker: A Novel Stochastic Optimization Strategy for Parallel Sparse Tucker Decomposition

Hao Li, *Student Member, IEEE*, Zixuan Li, Kenli Li, *Senior Member, IEEE*, Jan S. Rellermeyer, Lydia Chen, *Senior Member, IEEE*, Keqin Li, *Fellow, IEEE*.

1 BASIC OPTIMIZATION MODELS

The properties of $\mathbf{A}^{(n)}, n \in \{N\}$ accord to the task which determines the decomposition algorithm, i.e., in HOOI [1], each $a_{:,j_n}^{(n)}, j_n \in \{J_n\}$ where $n \in \{N\}$ is orthogonal and the L_2 norm constraints $\|a_{:,j_n}^{(n)}\|_2 = 1, j_n \in \{J_n\}, n \in \{N\}$. The orthogonality and unity of HOOI can track the optimal low-rank orthogonal subspace, as shown in Algorithm 1.

Algorithm 1 A Training Epoch of HOOI.

Input: Sparse tensor \mathcal{X} , Factor Matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times J_n}, n \in \{N\}$.

Output: Factor matrices $\mathbf{A}^{(n)} \in \mathbb{R}^{I_n \times J_n}, n \in \{N\}$;
Core Tensor \mathcal{G} N th order tensor $\in \mathbb{R}^{J_1 \times J_2 \times \dots \times J_N}$.

- 1: **for** n from 1 to N **do**
 - 2: $\mathcal{Y} \leftarrow \mathcal{X} \times_{(1)} \mathbf{A}^{(1)T} \times_{(2)} \dots \times_{(n-1)} \mathbf{A}^{(n-1)T} \times_{(n+1)} \mathbf{A}^{(n+1)T} \dots \times_{(N)} \mathbf{A}^{(N)T}$;
 - 3: $\mathbf{A}^{(n)} \leftarrow J_n$ leading left singular vectors of $\mathbf{Y}_{(n)}$;
 - 4: **end for**
 - 5: $\mathcal{G} \leftarrow \mathcal{X} \times_{(1)} \mathbf{A}^{(1)T} \times_{(2)} \dots \times_{(n)} \mathbf{A}^{(n)T} \times_{(n+1)} \dots \times_{(N)} \mathbf{A}^{(N)T}$.
-

The other methods are to infer the low-rank factor matrices $\mathbf{A}^{(n)}, n \in \{N\}$ and core tensor \mathcal{G} with L_2 norm regularization [1–5]. L_2 norm regularization can keep the promise of the smoothness [6–9] for the low-rank factor matrices $\mathbf{A}^{(n)}, n \in \{N\}$ and the core tensor \mathcal{G} . The optimization algorithms are divided into two classes, i.e., dense condition

and sparse condition. The optimization problem for dense tensor is presented as:

$$\begin{aligned} \arg \min_{\mathbf{A}^{(n)}, n \in \{N\}, \mathcal{G}} f(\mathcal{X}, \{\mathbf{A}^{(n)}\}, \mathcal{G}) \\ = \left\| \mathcal{X} - \widehat{\mathcal{X}} \right\|_2^2 + \lambda_{\mathcal{G}} \left\| \mathcal{G} \right\|_2^2 + \lambda_{\mathbf{A}} \left\| \mathbf{A}^{(n)} \right\|_2^2, \end{aligned} \quad (1)$$

where $\widehat{\mathcal{X}} = \mathcal{G} \times_1 \mathbf{A}^{(1)} \times_2 \dots \times_n \mathbf{A}^{(n)} \times_{n+1} \dots \times_N \mathbf{A}^{(N)}$ and $\lambda_{\mathcal{G}}$ and $\lambda_{\mathbf{A}}$ are the regularization parameters for core tensor and low-rank factor matrices, respectively. The optimization objective (1) involves variables multiplication, which is non-convex. The non-convex problem can be tackled by convex solution via updating a variable and fixing the others. The optimization problem (1) can be split into updating core tensor \mathcal{G} (n th vectorization form $g^{(n)}$) and updating low-rank factor matrices $\mathbf{A}^{(n)}, n \in \{N\}$ as following [10–12]:

$$\left\{ \begin{array}{l} \arg \min_{g^{(n)}} f\left(g^{(n)} \middle| x^{(n)}, \{\mathbf{A}^{(n)}\}, g^{(n)}\right) \\ \quad = \left\| x^{(n)} - \widehat{x}^{(n)} \right\|_2^2 + \lambda_{g^{(n)}} \left\| g^{(n)} \right\|_2^2; \\ \arg \min_{\mathbf{A}^{(n)}, n \in \{N\}} f\left(\mathbf{A}^{(n)} \middle| \mathbf{X}^{(n)}, \{\mathbf{A}^{(n)}\}, \mathbf{G}^{(n)}\right) \\ \quad = \left\| \mathbf{X}^{(n)} - \widehat{\mathbf{X}}^{(n)} \right\|_2^2 + \lambda_{\mathbf{A}} \left\| \mathbf{A}^{(n)} \right\|_2^2, \end{array} \right. \quad (2)$$

where $\widehat{x}^{(n)} = \mathbf{H}^{(n)} g^{(n)}$, $\widehat{\mathbf{X}}^{(n)} = \mathbf{A}^{(n)} \mathbf{G}^{(n)} \mathbf{S}^{(n)T}$ and $\mathbf{S}^{(n)} = \mathbf{A}^{(N)} \otimes \dots \otimes \mathbf{A}^{(n+1)} \otimes \mathbf{A}^{(n-1)} \otimes \dots \otimes \mathbf{A}^{(1)}$. Equ. (2) discusses the basic optimization model for the two key parts, i.e., core tensor \mathcal{G} and factor matrices $\mathbf{A}^{(n)}, n \in \{N\}$. The convex optimization objective (2) for core tensor $\mathcal{G} \in \mathbb{R}^{J_1 \times \dots \times J_n \times \dots \times J_N}$ cooperates with parameters $\{\mathcal{X}, \mathbf{H}^{(n)}\}$ and the gradient of $g^{(n)}$ on (2) is:

$$\begin{aligned} \frac{\partial f\left(g^{(n)} \middle| x^{(n)}, \{\mathbf{A}^{(n)}\}, g^{(n)}\right)}{\partial g^{(n)}} \\ = -\mathbf{H}^{(n)T} x^{(n)} + \mathbf{H}^{(n)T} \mathbf{H}^{(n)} g^{(n)} + \lambda_{g^{(n)}} g^{(n)}, \end{aligned} \quad (3)$$

-
- Hao Li, Zixuan Li, Kenli Li and Keqin Li are with the College of Computer Science and Electronic Engineering, Hunan University, and National Supercomputing Center in Changsha, Hunan, China, 410082.
 - Corresponding author: Kenli Li, Keqin Li.
E-mail: lihao123@hnu.edu.cn (H.Li-9@tudelft.nl), zixuanli@hnu.edu.cn, kll@hnu.edu.cn, J.S.Rellermeyer@tudelft.nl, Y.Chen-10@tudelft.nl, lik@newpaltz.edu.
 - Hao Li is also with the TU Delft, Netherlands.
 - J.S.Rellermeyer, Lydia Chen are with the TU Delft, Netherlands.
 - Keqin Li is also with the Department of Computer Science, State University of New York, New Paltz, New York 12561, USA.

where:

$$\mathbf{H}^{(n)} = \mathbf{A}^{(N)} \otimes \dots \otimes \mathbf{A}^{(n+1)} \otimes \mathbf{A}^{(n-1)} \otimes \dots \otimes \mathbf{A}^{(1)} \otimes \mathbf{A}^{(n)}, \mathbf{H}^{(n)} \in \mathbb{R}^{\prod_{n=1}^N I_n \times \prod_{n=1}^N J_n}, \quad (4)$$

and the element form is $h_{i,:}^{(n)} = a_{i_N,:}^{(N)} \odot \dots \odot a_{i_{n+1,:}}^{(n+1)} \odot a_{i_{n-1,:}}^{(n-1)} \odot \dots \odot a_{i_1,:}^{(1)} \odot a_{i_n,:}^{(n)}$, $h_{i,:}^{(n)} \in \mathbb{R}^{\prod_{n=1}^N J_n}$, $i \in \Omega_V^{(n)}$, $(i_1, \dots, i_N) \in \Omega$, $\mathbf{H}^{(n)T} \mathbf{H}^{(n)} = \mathbf{P}^{(N)} \otimes \dots \otimes \mathbf{P}^{(n+1)} \otimes \mathbf{P}^{(n-1)} \otimes \dots \otimes \mathbf{P}^{(1)} \otimes \mathbf{P}^{(n)}$, $\mathbf{P}^{(n)} = \mathbf{A}^{(n)T} \mathbf{A}^{(n)}$, and $g^{(n)} \in \mathbb{R}^{\prod_{n=1}^N J_n}$. The overall space overheads for intermediate matrices $\{\mathbf{H}^{(n)}, \mathbf{H}^{(n)T} \mathbf{H}^{(n)}\}$ are $\left\{O\left(\prod_{n=1}^N I_n J_n\right), O\left(\prod_{n=1}^N J_n^2\right)\right\}$, respectively, and the overall computational complexity for the gradient of core tensor (3) is $O\left(\prod_{n=1}^N (I_n J_n) + \prod_{n=1}^N J_n^2 + 3 \prod_{n=1}^N J_n\right)$.

The gradient for the factor matrices $\mathbf{A}^{(n)}$, $n \in \{N\}$ on (2) is:

$$\frac{\partial f\left(\mathbf{A}^{(n)} \mid \mathbf{X}^{(n)}, \{\mathbf{A}^{(n)}\}, \mathbf{G}^{(n)}\right)}{\partial \mathbf{A}^{(n)}} = -\mathbf{X}^{(n)} \mathbf{E}^{(n)T} + \mathbf{A}^{(n)} \mathbf{E}^{(n)} \mathbf{E}^{(n)T} + \lambda_{\mathbf{A}} \mathbf{A}^{(n)}, \quad (5)$$

where:

$$\begin{cases} \mathbf{S}^{(n)} = \mathbf{A}^{(N)} \otimes \dots \otimes \mathbf{A}^{(n+1)} \otimes \mathbf{A}^{(n-1)} \otimes \dots \otimes \mathbf{A}^{(1)}; \\ \mathbf{S}^{(n)} \in \mathbb{R}^{k=1, k \neq n} \prod_{k=1, k \neq n}^N I_k \times \prod_{k=1, k \neq n}^N J_k; \\ \mathbf{E}^{(n)} = \mathbf{G}^{(n)} \mathbf{S}^{(n)T} \in \mathbb{R}^{J_n \times \prod_{k=1, k \neq n}^N I_k}, \end{cases} \quad (6)$$

and the element form is $s_{j,:}^{(n)} = a_{i_N,:}^{(N)} \odot \dots \odot a_{i_{n+1,:}}^{(n+1)} \odot a_{i_{n-1,:}}^{(n-1)} \odot \dots \odot a_{i_1,:}^{(1)}$, $s_{j,:}^{(n)} \in \mathbb{R}^{\prod_{k=1, k \neq n}^N J_k}$, $(i_n, j) \in \Omega_M^{(n)}$, and $(i_1, \dots, i_N) \in \Omega$.

In common conditions, the original tensor is a HOHDST; thus, the space and time overheads for the intermediate matrices are huge which are not practical. There are some works that can construct the optimization objectives following the sparsity model, e.g., ALS [3, 4] and CD [5], etc. ALS should construct the Hessian matrices $\mathbf{C}^{(n)} \in \mathbb{R}^{J_n \times J_n}$, $n \in \{N\}$ and the element-wise form is presented as:

$$\begin{cases} g^{(n)} \leftarrow \left(\mathbf{C} + \lambda_{g^{(n)}} \mathbf{I}\right)^{-1} d, \mathbf{I} \in \mathbb{R}^{\prod_{n=1}^N J_n \times \prod_{n=1}^N J_n}; \\ a_{i_n,:}^{(n)} \leftarrow d_{i_n,:}^{(n)T} \left(\mathbf{C}^{(n)} + \lambda_{\mathbf{A}} \mathbf{I}_n\right)^{-1}, \mathbf{I}_n \in \mathbb{R}^{J_n \times J_n}, \\ n \in \{N\}, \end{cases} \quad (7)$$

where

$$\begin{cases} d = h_{\Omega_V^{(n)},:}^{(n)T} x^{(n)}; \\ c_{j_1, j_2} = h_{\Omega_V^{(n)}, j_1}^{(n)T} h_{\Omega_V^{(n)}, j_2}^{(n)}, j_1, j_2 \in \left\{\prod_{n=1}^N J_n\right\}; \\ d_{i_n,:}^{(n)} = x_{i_n, (\Omega_M^{(n)})_{i_n}}^{(n)} \mathbf{E}_{i_n, (\Omega_M^{(n)})_{i_n}}^{(n)T}; \\ c_{j_1, j_2}^{(n)} = e_{j_1, (\Omega_M^{(n)})_{i_n}}^{(n)} e_{j_2, (\Omega_M^{(n)})_{i_n}}^{(n)T}, j_1, j_2 \in \{J_n\}. \end{cases} \quad (8)$$

The computational complexity for $\left(\mathbf{C} + \lambda_{g^{(n)}} \mathbf{I}\right)^{-1}$ is $O\left(\left(\prod_{n=1}^N J_n\right)^3\right)$ and the computational complexity for $\left(\mathbf{C}^{(n)} + \lambda_{\mathbf{A}} \mathbf{I}_n\right)^{-1}$ is $O\left(J_n^3\right)$. Thus, the total computational complexity is $O\left(\sum_{n=1}^N I_n J_n^3 + \left(\prod_{n=1}^N J_n\right)^3 + \left(\prod_{n=1}^N J_n\right) |\Omega| + \left(\prod_{n=1}^N J_n\right)^2 |\Omega| + \sum_{n=1}^N J_n^2 |\Omega| + \sum_{n=1}^N J_n |\Omega|\right)$ and the space overhead is $O\left(N |\Omega| + \sum_{n=1}^N I_n J_n + \prod_{n=1}^N J_n + \sum_{n=1}^N J_n^2 + \left(\prod_{n=1}^N J_n\right)^2\right)$. CD is a special version of ALS and CD updates each feature element in a feature vector discretely. In the sequel, CD neglects the successive reading and writing for a feature vector which will increase the data addressing overheads. The optimization objective for $\{g_{\alpha}^{(n)}, a_{i_n, j_n}^{(n)}\}$ is presented as:

$$\left\{ \begin{aligned} & \arg \min_{g_{\alpha}^{(n)}, \alpha \in \left\{\prod_{n=1}^N J_n\right\}} f\left(g_{\alpha}^{(n)} \mid x^{(n)}, \{\mathbf{A}^{(n)}\}, g^{(n)}\right) \\ & = \left\| x^{(n)} - \sum_{\beta \neq \alpha}^{\prod_{n=1}^N J_n} h_{\Omega_V^{(n)}, \beta}^{(n)} g_{\beta}^{(n)} - h_{\Omega_V^{(n)}, \alpha}^{(n)} g_{\alpha}^{(n)} \right\|_2^2 \\ & \quad + \lambda_{g^{(n)}} \left\| g_{\alpha}^{(n)} \right\|_2^2; \\ & \arg \min_{a_{i_n, j_n}^{(n)}, n \in \{N\}} f\left(a_{i_n, j_n}^{(n)} \mid \mathbf{X}^{(n)}, \{\mathbf{A}^{(n)}\}, \mathbf{G}^{(n)}\right) \\ & = \left\| x_{i_n, (\Omega_M^{(n)})_{i_n}}^{(n)} - \sum_{j \neq j_n}^{J_n} a_{i_n, j}^{(n)} e_{j, (\Omega_M^{(n)})_{i_n}}^{(n)} - a_{i_n, j_n}^{(n)} e_{j_n, (\Omega_M^{(n)})_{i_n}}^{(n)} \right\|_2^2 \\ & \quad + \lambda_{\mathbf{A}} \left\| a_{i_n, j_n}^{(n)} \right\|_2^2. \end{aligned} \right. \quad (9)$$

The element-wise form for $\{g_{\alpha}^{(n)}, a_{i_n, j_n}^{(n)}\}$ is presented as:

$$\left\{ \begin{aligned} g_{\alpha}^{(n)} & \leftarrow \frac{h_{\Omega_V^{(n)}, \alpha}^{(n)T} \left(x^{(n)} - \sum_{\beta \neq \alpha}^{\prod_{n=1}^N J_n} h_{\Omega_V^{(n)}, \beta}^{(n)} g_{\beta}^{(n)}\right)}{\lambda_{g^{(n)}} + h_{\Omega_V^{(n)}, \alpha}^{(n)T} h_{\Omega_V^{(n)}, \alpha}^{(n)}}; \\ a_{i_n, j_n}^{(n)} & \leftarrow \frac{\left(x_{i_n, (\Omega_M^{(n)})_{i_n}}^{(n)} - \sum_{j \neq j_n}^{J_n} a_{i_n, j}^{(n)} e_{j, (\Omega_M^{(n)})_{i_n}}^{(n)}\right) e_{j_n, (\Omega_M^{(n)})_{i_n}}^{(n)T}}{\lambda_{\mathbf{A}} + e_{j_n, (\Omega_M^{(n)})_{i_n}}^{(n)T} e_{j_n, (\Omega_M^{(n)})_{i_n}}^{(n)}}. \end{aligned} \right. \quad (10)$$

HOOI, ALS and CD rely on the whole training set which results in high computation overhead and memory bottleneck especially in the situation of HOHDST.

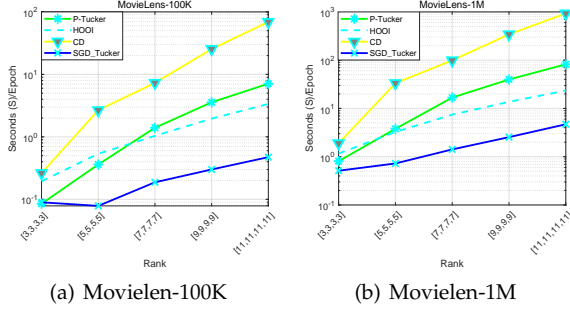


Fig. 1: Rank scalability for time overhead on full threads. The computational scalability for P-Tucker, HOOI, CD, and SGD_Tucker on the 2 datasets with successively increased number of total elements, i.e., MovieLens-100K, MovieLens-1M.

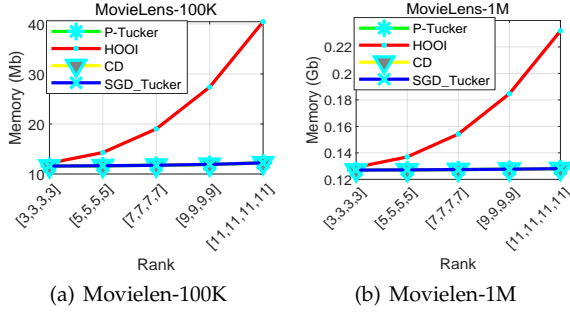


Fig. 2: Rank scalability for memory overhead on a thread running. (On this work, GB refers to GigaBytes and MB refers to Megabytes). The space scalability for P-Tucker, CD, HOOI, and SGD_Tucker on the 2 small datasets with successively increased total elements, i.e., MovieLens-100K, MovieLens-1M.

2 EXPERIMENTS RESULTS

Fig. 5 illustrates the computational time of the key parts, i.e., $TTMc$, $top-N$ of $svds$, and the total computational time. The key and most time-consuming part of $TTMc$ is vectors Kronecker product which is computed by $kron$ of *Armadillo* library and the index access code is borrowed from SPLATT of [1]. The SVD for the intermediate matrix of $TTMc$ is computed by $top-N$ of $svds$ on *Armadillo* library. As the Fig. 5 show, there are two

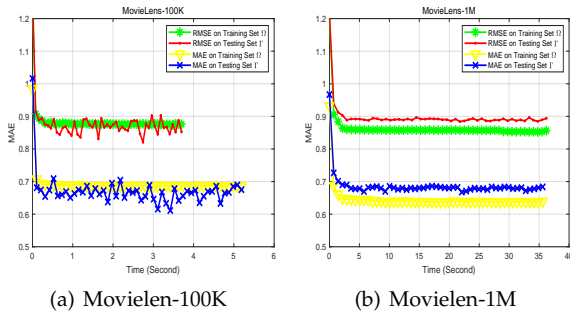


Fig. 3: RMSE and MAE vs time for SGD_Tucker on training set Ω and testing set Γ

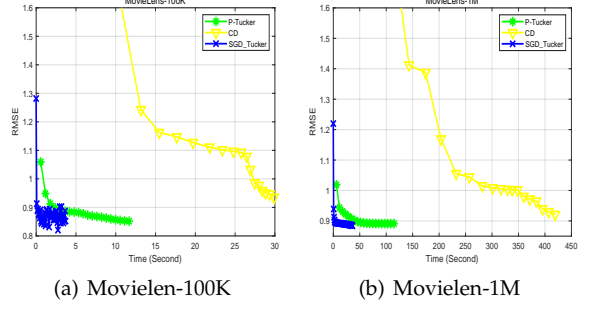


Fig. 4: RMSE comparison of SGD_Tucker, P-Tucker, and CD on 2 small datasets.

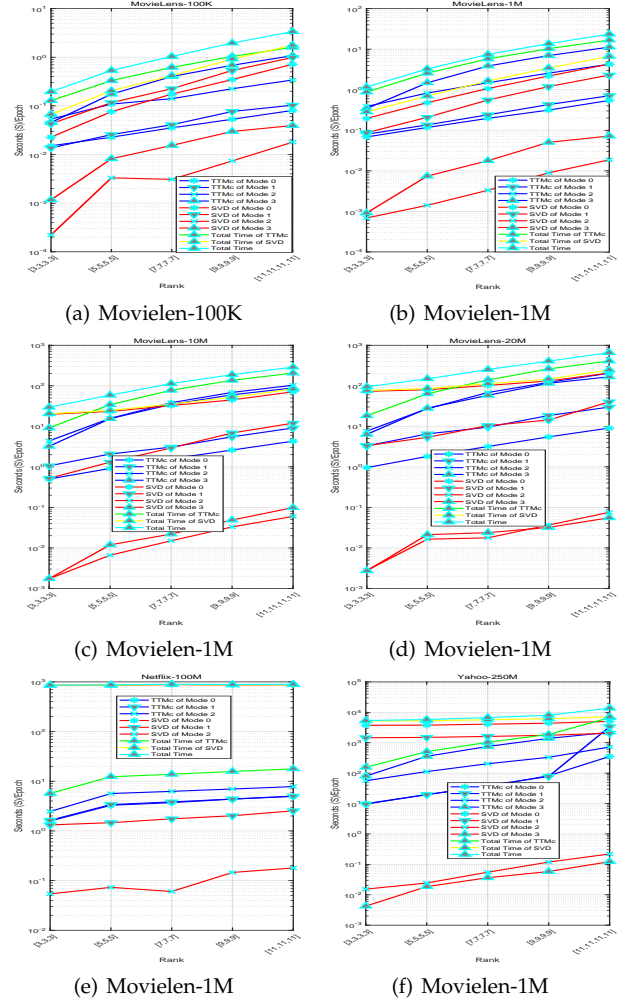


Fig. 5: Computational time of $TTMc$ and $top-N$ $svds$ based on the *Armadillo* library on the 6 datasets. The experiment runs on full threads

regulations as: 1) the computational overhead for $TTMc$ is controlled by the degree of thread balance which is balance degree of $\sum_{i_n \text{ belong to } l\text{th thread}} |(\Omega_M^{(n)})_{i_n}|, n \in \{N\}$ for thread l . 2) the computational overhead for $svds$ is regulated by $\left\{ I_n \times \prod_{k \neq n}^N J_k \mid n \in \{N\} \right\}$. From the results of Fig. 5 and Figs. 5 and 6 in main paper, the time scalability of HOOI is the same with SGD_Tucker. However, SGD_Tucker has lower computational overhead and less space overhead than HOOI.

REFERENCES

- [1] S. Smith and G. Karypis, "Accelerating the tucker decomposition with compressed sparse tensors," in *European Conference on Parallel Processing*. Springer, 2017, pp. 653–668.
- [2] H. Ge, K. Zhang, M. Alfifi, X. Hu, and J. Caverlee, "Distenc: A distributed algorithm for scalable tensor completion on spark," in *2018 IEEE International Conference on Data Engineering*. IEEE, 2018, pp. 137–148.
- [3] S. Oh, N. Park, S. Lee, and U. Kang, "Scalable tucker factorization for sparse tensors-algorithms and discoveries," in *IEEE International Conference on Data Engineering (ICDE)*. IEEE, 2018, pp. 1120–1131.
- [4] S. Oh, N. Park, J.-G. Jang, L. Sael, and U. Kang, "High-performance tucker factorization on heterogeneous platforms," *IEEE Transactions on Parallel and Distributed Systems*, vol. 30, no. 10, pp. 2237–2248, 2019.
- [5] M. Park, J.-G. Jang, and S. Lee, "Vest: Very sparse tucker factorization of large-scale tensors," *arXiv preprint arXiv:1904.02603*, 2019.
- [6] R. Johnson and T. Zhang, "Accelerating stochastic gradient descent using predictive variance reduction," in *Advances in neural information processing systems*, 2013, pp. 315–323.
- [7] H. Yu, R. Jin, and S. Yang, "On the linear speedup analysis of communication efficient momentum sgd for distributed non-convex optimization," in *International Conference on Machine Learning*, 2019, pp. 7184–7193.
- [8] M. Schmidt, N. Le Roux, and F. Bach, "Minimizing finite sums with the stochastic average gradient," *Mathematical Programming*, vol. 162, no. 1-2, pp. 83–112, 2017.
- [9] L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takáč, "Sarah: A novel method for machine learning problems using stochastic recursive gradient," in *Proceedings of the International Conference on Machine Learning*. JMLR. org, 2017, pp. 2613–2621.
- [10] G. T. Minton and E. Price, "Improved concentration bounds for count-sketch," in *Proceedings of the annual ACM-SIAM symposium on Discrete algorithms*. SIAM, 2014, pp. 669–686.
- [11] A. Traoré, M. Berar, and A. Rakotomamonjy, "Singleshot: a scalable tucker tensor decomposition," in *Advances in Neural Information Processing Systems*, 2019, pp. 6301–6312.
- [12] O. A. Malik and S. Becker, "Low-rank tucker decomposition of large tensors using tensorsketch," in *Advances in Neural Information Processing Systems*, 2018, pp. 10 096–10 106.