

NeuKron: Constant-Size Lossy Compression of Sparse Reorderable Matrices and Tensors

Taehyung Kwon*
Kim Jaechul Graduate
School of AI, KAIST
Seoul, South Korea
taehyung.kwon@kaist.ac.kr

Jihoon Ko*
Kim Jaechul Graduate
School of AI, KAIST
Seoul, South Korea
jihoonko@kaist.ac.kr

Jinhong Jung
Dept. of CSE, Jeonbuk
National University
Jeonju, South Korea
jinhongjung@jbnu.ac.kr

Kijung Shin
Kim Jaechul Graduate
School of AI, KAIST
Seoul, South Korea
kijungshin@kaist.ac.kr

ABSTRACT

Many real-world data are naturally represented as a sparse reorderable matrix, whose rows and columns can be arbitrarily ordered (e.g., the adjacency matrix of a bipartite graph). Storing a sparse matrix in conventional ways requires an amount of space linear in the number of non-zeros, and lossy compression of sparse matrices (e.g., Truncated SVD) typically requires an amount of space linear in the number of rows and columns. In this work, we propose NEUKRON for compressing a sparse reorderable matrix into a constant-size space. NEUKRON generalizes Kronecker products using a recurrent neural network with a constant number of parameters. NEUKRON updates the parameters so that a given matrix is approximated by the product and reorders the rows and columns of the matrix to facilitate the approximation. The updates take time linear in the number of non-zeros in the input matrix, and the approximation of each entry can be retrieved in logarithmic time. We also extend NEUKRON to compress sparse reorderable tensors (e.g. multi-layer graphs), which generalize matrices. Through experiments on ten real-world datasets, we show that NEUKRON is **(a) Compact**: requiring up to five orders of magnitude less space than its best competitor with similar approximation errors, **(b) Accurate**: giving up to $10\times$ smaller approximation error than its best competitors with similar size outputs, and **(c) Scalable**: successfully compressing a matrix with over 230 million non-zero entries.

CCS CONCEPTS

• **Information systems** → **Data mining**; **Data compression**.

KEYWORDS

Data Compression, Sparse Matrix, Sparse Tensor

ACM Reference Format:

Taehyung Kwon, Jihoon Ko, Jinhong Jung, and Kijung Shin. 2023. NeuKron: Constant-Size Lossy Compression of Sparse Reorderable Matrices and Tensors. In *Proceedings of the ACM Web Conference 2023 (WWW '23)*, May 1–5, 2023, Austin, TX, USA. ACM, New York, NY, USA, 11 pages. <https://doi.org/10.1145/3543507.3583226>

*Both authors contributed equally to this research.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WWW '23, May 1–5, 2023, Austin, TX, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-9416-1/23/04...\$15.00
<https://doi.org/10.1145/3543507.3583226>

1 INTRODUCTION

We consider a matrix to be *sparse* if the number of non-zero entries is much smaller than that of all entries. Sparse matrices naturally represent many types of data from various domains, as follows:

- **E-commerce**: User-item matrices represent how many times each user purchased each item [13, 28].
- **Search Engines**: Document-keyword matrices represent how many times each document contains each keyword [27]. User-ad matrices indicate how many times each user clicked each ad given by search engines [31].
- **Social Media**: The adjacency matrices of social networks indicate friendship between users [8, 30]. User-group matrices indicate which user belongs to each group [41].
- **Bibliography**: Author-paper matrices represent who authored each paper [32]. The adjacency matrices of collaboration networks represent co-authorships between authors [41].

Despite their sparsity, many real-world matrices require considerable space. Examples include user-ad matrices [37] and the adjacency matrices of web graphs [3] with billions of rows or columns; and keyword-document matrices [27] and the adjacency matrices of online social networks [8, 30] with tens of billions of non-zeros.

Compression of such large sparse matrices becomes important as smartphones and IoT devices become popular. Such memory-limited mobile devices are often required to process a large amount of data without sending them to clouds or servers, due to potential privacy risks [19]. Moreover, as the size of large-scale matrices grows rapidly, storing them is challenging also in desktops and servers [2, 8, 30], and for federate learning, compressing matrices is required to reduce communication costs [15]. As a result, a large number of lossy matrix-compression techniques [2, 9, 36] have been developed over the last few decades.

To the best of our knowledge, existing lossy-compression methods for sparse matrices create outputs whose sizes are at least linear in the numbers of rows and columns of the input matrix. For example, given an N -by- M matrix \mathbf{A} and a positive integer K , truncated singular value decomposition (T-SVD) [11, 35] outputs two matrices of which the numbers of entries are $O(KN)$ and $O(KM)$. Recent methods [2, 9, 36] have the same limitations, while they provide a better trade-off between space and information loss than T-SVD.

Can we compress a matrix into a constant-size space, which can even be smaller than the number of rows and columns? In this paper, we exploit the fact that **many real-world sparse matrices are reorderable**, i.e., the rows and columns of the matrices can be arbitrarily ordered.¹ All of the matrices discussed in the first

¹ A matrix is *non-reorderable* if the orders of rows and columns in it convey information. For example, images and multivariate time series are non-reorderable matrices since the orders of rows and columns in them indicate spatial and temporal adjacency.

paragraph, which are essentially bipartite graphs (nodes of one type correspond to rows, and nodes of the other type correspond to columns), are reorderable. For example, in the case of a user-item matrix built based on e-commerce data, which user (item) comes next to which user (item) does not matter. Our key idea is to **order rows and columns** to facilitate our model to learn and exploit meaningful patterns in the input matrix for compression.

Specifically, we present NEUKRON, a constant-size lossy compression method for sparse reorderable matrices. It consists of a machine-learning model and novel training schemes. The model generalizes the Kronecker power and enhances its expressive power using a recurrent neural network with a constant number of parameters. The training scheme, which is crucial for performance, is to reorder rows and columns in the input matrix to create patterns that the machine-learning model can exploit for better compression. Consider an N -by- M matrix with L non-zeros, where $N \leq M$ without loss of generality. The model and the training schemes are designed carefully so that each training epoch takes $O(M+L \log M)$ time, and after training, the approximation of each entry can be retrieved in $O(\log M)$ time. Note that the time complexity of training depends only on the number of non-zeros instead of all entries.

In addition, we extend NEUKRON for lossy-compression of sparse reorderable tensors while maintaining its strengths. Tensors (i.e., multi-dimensional arrays) generalize matrices to higher dimensions, and in other words, matrices are 2-order tensors. Sparse tensors have been used widely for various purposes, including context-aware recommender systems [17] and knowledge base completion [21]), and for lossy compression of them, tensor decomposition methods (e.g., CP [1, 5] and Tucker [1, 38]) have been developed.

For evaluation, we perform extensive experiments using ten real-world matrices (spec., bipartite graphs) and tensors. The results reveal the following advantages of NEUKRON:

- **Compact:** Its output is up to **5 orders of magnitude smaller** than competitors' with similar approximation error.
- **Accurate:** It achieves up to **10.1× smaller approximation error** than its best competitors that give similar-size outputs.
- **Scalable:** Its running time is **linear** in the number of non-zero entries, and it successfully compresses matrices with **over 230 millions of non-zero entries** on commodity GPUs.

Reproducibility: The code and datasets are available at [20].

Remarks on non-reorderable matrices: While we focus on reorderable matrices in this paper, NEUKRON can also be applied to non-reorderable matrices if the mapping between the original and new orders of rows and columns are stored additionally. We present a related experiment in Appendix C.

2 RELATED WORKS

In this section, we review lossy-compression methods for matrices and tensors. Those for lossy compression of sparse matrices or tensors of any size are compared in Table 1 and also in Section 6.

Factorization-based matrix compression: Given a matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$, singular value decomposition (SVD) [12] decomposes \mathbf{A} into $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ where $\mathbf{U} \in \mathbb{R}^{N \times R}$, $\mathbf{V} \in \mathbb{R}^{M \times R}$, $\mathbf{\Sigma}$ is a diagonal matrix with its singular values, and R is its rank. Truncated SVD (T-SVD) [11, 35] outputs the K ($\leq R$) largest singular values and the corresponding vectors of \mathbf{U} and \mathbf{V} from which the rank- K approximation of \mathbf{A} best in terms of the Frobenius norm can be obtained [35]. Its outputs

Table 1: Comparison of lossy-compression methods for sparse matrices and tensors. For simplicity, we treat the tensor order and all hyperparameters as constants. Comparisons are relative, and we provide details in [20].

Methods	Space & Accuracy Trade-off	Training Complexity (per iteration)	Inference Complexity (per entry)	Number of Hyperparameters	Training Time (total)
NEUKRON	Strong	\propto #non-zeros	$\propto \log(N_{\max})^*$	4**	Long
T-SVD [11, 39]	Weak	\propto #non-zeros	constant	1	Short
CMD [36], CUR [9]	Moderate	\propto #non-zeros	constant	2	Moderate
ACCAMS [2]	Moderate	\propto #all-entries	constant	2	Moderate
bACCAMS [2]	Moderate	\propto #all-entries	constant	4	Long
KronFit [25, 26]	Weak	\propto #non-zeros	$\propto \log(N_{\max})^*$	4	Long
CP [5], Tucker [38]	Weak	\propto #non-zeros	constant	1	Moderate

* Here $N_{\max} = \max(N_1, \dots, N_D)$ is the maximum dimensionality (i.e., mode length).
 ** The learning rate, the optimizer, the weight parameter for the criterion of switching, and the size of hidden dimensions in LSTM.

have $O(K(M+N))$ real values, and typically most of them are non-zero. For further compression, CUR decomposition [9] aims to yield sparse outputs. Specifically, a sparse matrix \mathbf{A} is decomposed into CUR (i.e., $\mathbf{A} \approx \mathbf{C}\mathbf{U}\mathbf{R}$), where $\mathbf{C} \in \mathbb{R}^{N \times K}$ and $\mathbf{R} \in \mathbb{R}^{K \times M}$ are constructed by sampling K columns and rows from \mathbf{A} , respectively. The matrix $\mathbf{U} \in \mathbb{R}^{K \times K}$ is dense but small, and it is determined by \mathbf{C} and \mathbf{R} so that the approximation error is minimized. Compact matrix decomposition (CMD) [36] keeps only unique columns and rows in \mathbf{C} and \mathbf{R} for further efficiency.

Co-clustering-based matrix compression: ACCAMS and bACCAMS [2] use an additive combination of small co-clusters to approximate a given matrix. While the numbers of parameters of them are linear in the numbers of rows and columns, they produce intermediate results whose size is linear in the number of (potentially zero) known entries. Thus, they are computationally and memory inefficient when most entries are known but zero.

Kronecker product-based matrix compression: The adjacency matrix of a Kronecker graph [24] is a Kronecker power of a fixed seed matrix (e.g., 2-by-2 matrix). KronFit [25, 26] searches for a seed matrix whose Kronecker power approximates the adjacency matrix of a given graph. While KronFit is designed for adjacency matrices, it can be easily extended to matrices of any size, and the output seed matrix can be considered as a constant-size lossy compression of a given matrix. However, the approximation error is considerable, even when the seed matrix is large, due to the inflexibility of the Kronecker product, as shown in Section 6.2.

Tensor compression: CP decomposition (CP) [5] and Tucker decomposition (Tucker) [38] generalize the aforementioned T-SVD to higher-order tensors. They approximate a given tensor using the sums and products (e.g., outer product and n -mode product) of much smaller low-rank tensors and matrices, which can be considered as a lossy compression of the given tensor. Efficient CP and Tucker methods for sparse tensors have been developed [1]. For lossless compression of sparse tensors, compressed sparse fiber (CSF) [33, 34] is available.

Other related works: Unipartite-graph summarization algorithms [22, 23, 29] can be used for compressing adjacency matrices of unipartite graphs, while they cannot be directly applied to weighted and/or non-symmetric matrices, which we aim to compress. Matrix sketching methods replace a given large matrix with a more compact matrix that follows the properties of the input matrix, for

example, by leaving only important columns (rows) of the input matrix [9, 10]. These methods, however, cannot be applied to our problem because the entries of the input matrix cannot be estimated directly from their outputs.

3 NOTATIONS AND PROBLEM DEFINITION

In this section, we introduce basic concepts and give a formal problem definition. See Table 2 for common notations.

3.1 Notations and Concepts

Sparse reorderable matrix and tensor: A matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$ is a 2-dimensional array with N rows and M columns, and real entries. A D -order tensor $\mathcal{X} \in \mathbb{R}^{N_1 \times \dots \times N_D}$ is a D -dimensional array of size $N_1 \times \dots \times N_D$ with real entries. We use a_{ij} or $A(i, j)$ to denote the (i, j) -th entry of \mathbf{A} , and we use x_{i_1, \dots, i_D} to denote the (i_1, \dots, i_D) -th entry of \mathcal{X} . We consider a matrix or a tensor to be *sparse* if the number of non-zero entries is much smaller than that of all entries.² We call a matrix *reorderable* if its rows and columns can be arbitrarily ordered. We provide some examples of reorderable matrices where the orders of rows and columns do not convey any information and some examples of non-reorderable ones (see Footnote 1) in Section 1. Similarly, we call a tensor reorderable if the indices in each mode can be arbitrarily ordered.

Approximation error: The *Frobenius norm* is a function $\|\cdot\|_F : \mathbb{R}^{N \times M} \rightarrow \mathbb{R}$ defined as the square root of the square sum of all entries in the given matrix. Similarly, the Frobenius norm of a tensor is defined as the square root of the square sum of all entries in the given tensor. The *approximation error* of a matrix $\tilde{\mathbf{A}}_\Theta$ that approximates \mathbf{A} is defined as $\|\mathbf{A} - \tilde{\mathbf{A}}_\Theta\|_F^2$. Similarly, the approximation error of $\tilde{\mathcal{X}}_\Theta$ that approximates \mathcal{X} is defined as $\|\mathcal{X} - \tilde{\mathcal{X}}_\Theta\|_F^2$.

Kronecker product and power: Given two matrices $\mathbf{A} \in \mathbb{R}^{N \times M}$ and $\mathbf{B} \in \mathbb{R}^{P \times Q}$, the *Kronecker product* $\mathbf{A} \otimes \mathbf{B} \in \mathbb{R}^{NP \times MQ}$ is a large matrix formed by multiplying \mathbf{B} by each element of \mathbf{A} , i.e.,

$$\mathbf{A} \otimes \mathbf{B} := \begin{bmatrix} a_{11}\mathbf{B} & \dots & a_{1M}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{N1}\mathbf{B} & \dots & a_{NM}\mathbf{B} \end{bmatrix}.$$

We denote the l -th *Kronecker power* of \mathbf{A} as $\mathbf{A}^{\otimes l}$, where $\mathbf{A}^{\otimes l} = \mathbf{A}^{\otimes(l-1)} \otimes \mathbf{A}$ and $\mathbf{A}^{\otimes 1} = \mathbf{A}$.

3.2 Problem Definition

The constant-size lossy matrix compression problem that we address in this paper is defined in Problem 1. It should be noted that the given constant k can be even smaller than N and M . The problem of *constant-size lossy compression of a sparse reorderable tensor* can be defined by simply replacing the matrix \mathbf{A} with a tensor \mathcal{X} and $\|\mathbf{A} - \tilde{\mathbf{A}}_\Theta\|_F^2$ with $\|\mathcal{X} - \tilde{\mathcal{X}}_\Theta\|_F^2$.

PROBLEM 1. (Constant-size Lossy Compression of a Sparse Reorderable Matrix)

- **Given:** (1) a sparse and reorderable matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$,
(2) a constant $k = O(1)$,
- **Find:** a model Θ
- **to Minimize:** the approximation error $\|\mathbf{A} - \tilde{\mathbf{A}}_\Theta\|_F^2$, where $\tilde{\mathbf{A}}_\Theta$ is the matrix approximated from Θ .
- **Subject to:** the number of parameters in Θ is at most k .

²The ratio is at most 0.0046 in the datasets considered in the paper.

Table 2: Frequently-used notations

Symbol	Definition
$\mathbf{A} \in \mathbb{R}^{N \times M}$	an N -by- M sparse matrix
a_{ij} or $\mathbf{A}(i, j)$	(i, j) -th entry of \mathbf{A}
$\mathbf{A}_{i,:}, \mathbf{A}_{:,j}$	i -th row of \mathbf{A} , j -th column of \mathbf{A}
$\mathcal{X} \in \mathbb{R}^{N_1 \times \dots \times N_D}$	tensor
D	order of \mathcal{X}
x_{i_1, \dots, i_D}	(i_1, \dots, i_D) -th entry of \mathcal{X}
$\text{nnz}(\mathbf{A}), \text{nnz}(\mathcal{X})$	number of non-zero entries in \mathbf{A} and \mathcal{X}
$\ \mathbf{A}\ _F, \ \mathcal{X}\ _F$	Frobenius norm of \mathbf{A} and \mathcal{X}
\otimes	Kronecker product
$\mathbf{A}^{\otimes l}, \mathcal{X}^{\otimes l}$	l -th Kronecker power of \mathbf{A} and \mathcal{X}
Θ	a NEUKRON model which compresses \mathbf{A} and \mathcal{X}
$\tilde{\mathbf{A}}_\Theta, \tilde{\mathcal{X}}_\Theta$	approximated matrix and tensor of \mathbf{A} and \mathcal{X} by Θ
q	a parameter for the scale of model outputs
h	hidden dimension in LSTM
$[n]$	a set of integers from 1 to n (i.e., $\{1, 2, \dots, n\}$)

4 PROPOSED METHOD

In this section, we present NEUKRON, a constant-space lossy compression method for sparse reorderable matrices and tensors. We first describe its neural network model and then the training strategies for it. After that, we analyze the computational complexity of NEUKRON. For ease of explanation, we assume that the input is a matrix through the section, and then we describe the extensions for tensors in Section 5.

4.1 Model

4.1.1 Overview. When designing a neural network model Θ for NEUKRON, we aim to achieve the following goals:

- **G1. Constant Size:** The number of parameters of the model should be constant, regardless of the size of the input matrix.
- **G2. Exploitation of Sparsity:** It should be possible to fit the model to the input by accessing only non-zero entries.
- **G3. Fast Approximation:** From the trained model, it should be possible to approximate each entry of the input matrix in sublinear time (preferably, in constant or logarithmic time).

For **G1**, given a matrix \mathbf{A} to be compressed, we encode the position (i, j) of each entry a_{ij} as a sequence and use an auto-regressive sequence model, specifically LSTM [14], which has a constant number of parameters, to process the sequence. For our purpose, LSTM performs similarly with GRU [7] and outperforms the decoder layer of Transformer [40], as shown empirically in [20]. For an entry a_{ij} , the sequence encoding the position (i, j) is fed into LSTM, and the outputs of LSTM are combined for its approximation \tilde{a}_{ij} in logarithmic time, achieving **G3** (see Theorem 1 in Section 4.3). Moreover, regarding **G2**, the outputs of LSTM are combined so that the sparsity can be exploited for efficient computation of the objective and its gradient (see Section 4.2.2). The details of encoding inputs and combining outputs are described in the following subsections. Regarding **G3**, it should be noticed that many factorization-based methods approximate each entry even in constant time (see Table 1).

4.1.2 Encoding inputs (lines 1-3 of Algorithm 1). For simplicity, we assume an input matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$ where $N = M = 2^l$ (see Section 4.1.4 for generalization to matrices of any size). Algorithm 1 depicts how NEUKRON approximates such \mathbf{A} .

For each entry a_{ij} of \mathbf{A} , NEUKRON encodes its position (i, j) in a sequence of length $l = \log_2 M$ by recursively subdividing \mathbf{A} in

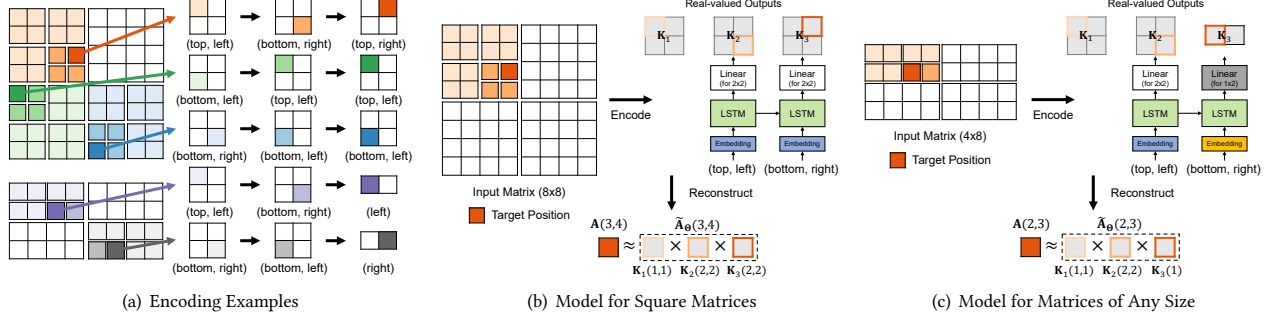


Figure 1: The overall approximation process of NEUKRON. It encodes the input position into a sequence by recursively dividing the input matrix. The sequence is fed into LSTM, and the outputs of LSTM are aggregated based on the Kronecker product.

a top-down manner. NEUKRON first chooses the partition where a_{ij} lies when A is divided into 2×2 partitions of the same size (i.e., $2^{l-1} \times 2^{l-1}$). Each division gives four partitions at top left (TL), top right (TR), bottom left (BL), and bottom right (BR). Then, NEUKRON repeats the process on the chosen partition until only the target entry a_{ij} is left. The sequence of the positions of the chosen partition is used to encode a_{ij} . In our implementation, each entry of the sequence, which is a position, is converted into a tuple in $\{1, 2\} \times \{1, 2\}$. Specifically, the k -th entry of the sequence that encodes the position (i, j) is $(t(i, k), t(j, k))$ where

$$t(i, k) := \left(\left\lfloor \frac{(i-1)}{2^{l-k}} \right\rfloor \bmod 2 \right) + 1. \quad (1)$$

EXAMPLE 1 (ENCODING IN SQUARE MATRICES). Suppose we encode the position $(3, 4)$ of the square matrix in Figure 1(a), where $l = 3$. The position $(3, 4)$ is located in the *top-left* partition of the input matrix, and it is located in the *bottom-right* part of the chosen partition. Lastly, the position $(3, 4)$ is located at the *top-right* one of the lastly chosen partition. Thus, the position $(3, 4)$ is encoded in the sequence $TL \rightarrow BR \rightarrow TR$, which becomes $(1, 1) \rightarrow (2, 2) \rightarrow (1, 2)$ based on t (Eq. (1)).

Each tuple in the sequence, except for the last one, goes through an embedding layer (line 2) to be converted into a corresponding embedded vector of size h , where h is a hyperparameter. Then, the vector is fed into LSTM (line 3).

4.1.3 Handling outputs (lines 4-6 of Algorithm 1). Below, we present how NEUKRON produces an approximation. See Figure 1(b) for a pictorial description. We again assume an input matrix $A \in \mathbb{R}^{N \times N}$ where $N = M = 2^l$ for ease of explanation. Given the position (i, j) of a target entry a_{ij} , NEUKRON creates $K_1 \in \mathbb{R}^{2 \times 2}, \dots, K_l \in \mathbb{R}^{2 \times 2}$. Specifically, given the sequence of tuples that encode (i, j) (see Section 4.1.2 for encoding), for each $k \in [l-1]$, the k -th LSTM cell receives the embedding of the k -th tuple, and then the hidden state of the cell goes through the linear layer and the Softplus activation to produce K_{k+1} (line 5). The entries of K_1 are separate learnable parameters. The approximation \tilde{a}_{ij} is computed from the (i, j) -th entry of their Kronecker product $K_1 \otimes \dots \otimes K_l$ as follows (line 6):

$$\tilde{a}_{ij} := \sqrt{q} \cdot \prod_{k=1}^l K_k(t(i, k), t(j, k)) / \|K_k\|_F, \quad (2)$$

where $\prod_{k=1}^l K_k(t(i, k), t(j, k))$ is the (i, j) -th entry of the Kronecker product, and q is a learnable parameter. It should be noticed that the entire Kronecker product does not have to be computed. By combining the outputs of LSTM using Eq.(2), G2 in Section 4.1.1

Algorithm 1: Approximation process of NEUKRON for an N -by- $N (= 2^l)$ matrix A

- Input:** (a) a position: $(i, j) \in [N] \times [N]$ where $N = 2^l$
 (b) parameters of Embedding, LSTM, and the linear layer (W, b)
 (c) scale parameter q and the first matrix of Kronecker products K_1
Output: an approximation \tilde{a}_{ij} of a_{ij} , which is the (i, j) -th entry of the input matrix $A \in \mathbb{R}^{N \times N}$
- 1 **for** $k \leftarrow 1$ to l **do**
 - 2 $x_k \leftarrow \text{Embedding}(t(i, k), t(j, k))$ ▶ Sect. 4.1.2
 - 3 $y_2, \dots, y_l \leftarrow \text{LSTM}(x_1, x_2, \dots, x_{l-1})$
 - 4 **for** $k \leftarrow 2$ to l **do**
 - 5 $K_k \leftarrow \text{Softplus}(W y_k + b)$ ▶ Sect. 4.1.3
 - 6 **return** $\tilde{a}_{ij} \leftarrow \sqrt{q} \cdot \prod_{k=1}^l K_k(t(i, k), t(j, k)) / \|K_k\|_F$

can be achieved. Specifically, using Eq.(2) enables the exploitation of the sparsity of the input matrix A for linear-time training, as described in detail in Section 4.2.2 (see Lemma 1).

4.1.4 Handling matrices of any size. Below, we describe how the above processes of NEUKRON are generalized to compress a matrix of any size. For a given matrix $A \in \mathbb{R}^{N \times M}$, we consider integers l_{row} and l_{col} such that $2^{l_{\text{row}}} \geq N$ and $2^{l_{\text{col}}} \geq M$. Then, $A \in \mathbb{R}^{N \times M}$ is extended to the $2^{l_{\text{row}}}$ -by- $2^{l_{\text{col}}}$ matrix with additional rows and columns filled with zeros. Specifically, NEUKRON sets $l_{\text{row}} \leftarrow \lceil \log_2 N \rceil$ and set $l_{\text{col}} \leftarrow \lceil \log_2 M \rceil$ so that the number of new entries is minimized.

Without loss of generality, we assume $N \leq M$ and thus $l_{\text{row}} \leq l_{\text{col}}$. If $l_{\text{row}} = l_{\text{col}}$, the extended square matrix is considered as the input and processed as described in Sections 4.1.2 and 4.1.3. Otherwise (i.e., if $l_{\text{row}} < l_{\text{col}}$), to encode the position (i, j) of a target entry a_{ij} , NEUKRON first recursively divides A into 2×2 partitions, l_{row} times, to obtain a partition has a size of $1 \times 2^{l_{\text{col}} - l_{\text{row}}}$, and then it recursively divides the partition into two partitions of the same size (i.e., 1×2), $l_{\text{col}} - l_{\text{row}}$ times. Each division gives two partitions at left (L) and right (R). Specifically, the k -th entry of the sequence that encodes the position (i, j) is $(t_{\text{row}}(i, k), t_{\text{col}}(j, k))$, where $\forall d \in \{\text{row}, \text{col}\}$,

$$t_d(i, k) = \begin{cases} \left(\left\lfloor \frac{(i-1)/2^{l_d-k}}{2} \right\rfloor \bmod 2 \right) + 1, & \text{if } k \leq l_d, \\ 0, & \text{otherwise.} \end{cases} \quad (3)$$

EXAMPLE 2 (ENCODING IN RECTANGULAR MATRICES). Suppose we encode the position $(2, 3)$ of the non-square matrix in Figure 1(a), where $(l_{\text{row}}, l_{\text{col}}) = (2, 3)$. The position $(2, 3)$ is located in the *top-left* partition and in the *bottom-right* partition, respectively, in the first two divisions. In the last division, the position $(2, 3)$ is located in the

left one. Thus, the position (2, 3) is encoded in the sequence TL → BR → L, which becomes (1, 1) → (2, 2) → (0, 1) based on t_d (Eq. (3)).

As in Section 4.1.3, NEUKRON produces an approximation of a_{ij} using the (i, j) -th entry of the modified Kronecker product in Eq. (2) $\mathbf{K}_1 \otimes \cdots \otimes \mathbf{K}_{l_{\text{col}}}$. The only difference is that $\mathbf{K}_{l_{\text{row}}+1}, \dots, \mathbf{K}_{l_{\text{col}}}$ are matrices of size 1×2 , and for them, a separate embedding and linear layers are used, as described in Figure 1(c).

4.1.5 Comparison with Kronecker Graphs. Our model Θ generalizes the Kronecker graph model [25, 26] in two ways:

- While the Kronecker graph model uses the power of a single seed matrix, Θ uses the Kronecker product of potentially different matrices (i.e., $\mathbf{K}_1, \dots, \mathbf{K}_l$) for approximation.
- In Θ , the matrices $\mathbf{K}_1, \dots, \mathbf{K}_l$ may vary depending on the position of the target entry to be approximated. Specifically, \tilde{a}_{ij} is computed using the (i, j) -th entry of $\mathbf{K}_1^{(f_1(i), f_1(j))} \otimes \mathbf{K}_2^{(f_2(i), f_2(j))} \otimes \cdots \otimes \mathbf{K}_l^{(f_l(i), f_l(j))}$, where $f_k(i) = \lfloor (i-1)/2^{l-k} \rfloor$.

This generalization leads to a significantly better trade-off between parameter size and approximation error in practice, as shown in Section 6.2. Notably, there are also two differences:

- While the Kronecker graph model is trained under a log-likelihood objective, Θ uses the squared Frobenius norm and normalizes the matrices to apply the tricks in Eq. (5) and Eq. (6).
- As specified in Eq. (2), each matrix (i.e., $\mathbf{K}_1, \dots, \mathbf{K}_l$) is normalized and mapped onto the unit hypersphere.

4.2 Training Strategies

In this subsection, we propose novel training schemes for NEUKRON's model Θ . We first present how to fit Θ to a given sparse reorderable matrix while exploiting its sparsity. Then, we present how to reorder the rows and columns of the input matrix so that Θ can be better fit to it. These two steps are alternated until convergence, as described in Algorithm 2. Below, we assume a matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$ where $(N, M) = (2^{l_{\text{row}}}, 2^{l_{\text{col}}})$. As described in Section 4.1.4, a matrix of any size can be extended by zero-padding to satisfy this condition. We also assume $N \leq M$, without loss of generality.

4.2.1 Update of row/column orders. It is crucial to properly order the rows and columns of a given reorderable matrix for NEUKRON's model Θ better fit the matrix. This is because proper ordering reveals patterns (e.g., self-similarity and co-clusters), which Θ can exploit for accurate compression.

Overall process: For initialization, any co-clustering algorithms can be used. In our implementation, the matrix reordering scheme in [16] is used (see Section 6.3 for the effect of initialization). After initialization, NEUKRON repeats (a) sampling two rows (or columns), (b) measuring the change in the approximation error (i.e., $\|\mathbf{A} - \tilde{\mathbf{A}}_{\Theta}\|_F^2$), and (c) determining whether to swap the sampled rows (or columns) or not probabilistically using the following criterion:

$$u < \exp(-\gamma \cdot \Delta), \quad (4)$$

where $u \sim U(0, 1)$, Δ is the change in the approximation error, and $\gamma > 0$ is a hyperparameter that controls the probability of accepting swaps that increase the approximation error.

Similarity-aware sampling: Below, we describe how NEUKRON samples candidate pairs of rows (or columns) to be potentially

Algorithm 2: Overall training process of NEUKRON

Input: (a) a sparse reorderable matrix \mathbf{A}
 (b) a number T_p of permutation updates
Output: a NEUKRON model Θ

- 1 Initialize Θ
- 2 **while not converged do**
- 3 **for** $k \leftarrow 1$ to T_p **do**
- 4 $\mathbf{A} \leftarrow \text{UPDATEROWORDER}(\mathbf{A})$ ▶ Sect. 4.2.1
- 5 $\mathbf{A} \leftarrow \text{UPDATECOLORDER}(\mathbf{A})$ ▶ Sect. 4.2.1
- 6 $\Theta \leftarrow \text{UPDATEMODEL}(\mathbf{A}, \Theta)$ ▶ Sect. 4.2.2
- 7 **return** Θ

swapped. Compared to a naive uniform sampling, the proposed sampling method has two advantages: **(a) effective:** it samples pairs based on the similarity of rows (or columns) so that swapping the pairs is likely to reduce the approximation error, and **(b) easy-to-parallelize:** it samples disjoint pairs, which can be processed in parallel. The main idea is to select candidate pairs so that swapping pairs is likely to make similar rows (or columns) close to each other and thus to make them encoded in similar sequences in Section 4.1.2. Below, we describe the sampling method step by step for sampling row pairs. Column pairs are sampled similarly.

- **Estimating similarity:** In order to quickly estimate the similarity, min-hashing [4] is used. Specifically, for a uniform random bijective function $h_{\text{col}} : [M] \rightarrow [M]$ for the columns, the shingle $\min_{a_{ij} \neq 0} (h_{\text{col}}(j))$ of each i -th row is computed. It can be shown that two rows have the same shingle with probability proportional to the Jaccard similarity of the column indices of their non-zeros [4].
- **Locating similar rows/cols nearby:** We match rows with the same shingle disjointly, and for each matched rows, we sample pairs of rows to be swapped so that they are located in *nearby positions*, which we define as positions whose binary representations differ in only 1 bit. Let $p(i, k)$ be the position whose binary representation differs with that of i only in the k -th bit. Specifically, if two rows in the i_1 -th and i_2 -th positions are matched, we sample $(i_1, p(i_2, k))$ and $(i_2, p(i_1, k))$ so that i_1 and i_2 become nearby after swaps. The position $k \in [l_{\text{col}}]$ is sampled probabilistically (see Appendix B for details).
- **Pairing unmatched rows:** The rows remaining unmatched are randomly matched, and for each matched rows, we sample pairs as described above.

We describe the entire process of reordering for rows in Algorithm 3.

4.2.2 Update of model parameters. The objective function of optimization is $\|\mathbf{A} - \tilde{\mathbf{A}}_{\Theta}\|_F^2$, as in Problem 1. Naively computing it takes $\Omega(NM \log M)$ time since all NM entries are approximated and approximating each entry takes $\Theta(\log M)$ time (see Theorem 1 in Section 4.3).

For its efficient computation, we reformulate the error as

$$\begin{aligned} \|\mathbf{A} - \tilde{\mathbf{A}}_{\Theta}\|_F^2 &= \sum_{i=1}^N \sum_{j=1}^M (a_{ij} - \tilde{a}_{ij})^2 = \sum_{a_{ij} \neq 0} (a_{ij} - \tilde{a}_{ij})^2 \quad (5) \\ &+ \sum_{a_{ij}=0} \tilde{a}_{ij}^2 = \sum_{a_{ij} \neq 0} ((a_{ij} - \tilde{a}_{ij})^2 - \tilde{a}_{ij}^2) + \sum_{i=1}^N \sum_{j=1}^M \tilde{a}_{ij}^2. \end{aligned}$$

In our model Θ , the last term, (i.e., the sum of squares) can be immediately computed from a learnable parameter $q \in \mathbb{R}^+$ (which is used in Eq. (2)), as formalized in Lemma 1.

LEMMA 1. *For approximation by Eq. (2), Eq. (6) always holds.*

$$\sum_{i=1}^{2^{l_{\text{row}}}} \sum_{j=1}^{2^{l_{\text{col}}}} \tilde{a}_{ij}^2 = q^{l_{\text{col}}} \quad (6)$$

PROOF. To prove this lemma, we use an induction. For $(l_{\text{row}}, l_{\text{col}}) = (1, 1)$ and $(l_{\text{row}}, l_{\text{col}}) = (0, 1)$, the statement holds trivially. Suppose the statement holds when $(l_{\text{row}}, l_{\text{col}}) = (0, l_2)$. For $(l_{\text{row}}, l_{\text{col}}) = (0, l_2 + 1)$, the statement also holds since

$$\begin{aligned} \sum_{i=1}^{2^{l_{\text{row}}}} \sum_{j=1}^{2^{l_2+1}} \tilde{a}_{ij}^2 &= \frac{q\mathbf{K}_1(1, 1)^2}{\|\mathbf{K}_1\|_F^2} \sum_{i=1}^{2^{l_{\text{row}}}} \sum_{j=1}^{2^{l_2}} \frac{\tilde{a}_{ij}^2}{q\mathbf{K}_1(1, 1)^2 / \|\mathbf{K}_1\|_F^2} \\ &\quad + \frac{q\mathbf{K}_1(1, 2)^2}{\|\mathbf{K}_1\|_F^2} \sum_{i=1}^{2^{l_{\text{row}}}} \sum_{j=2^{l_2}+1}^{2^{l_2+1}} \frac{\tilde{a}_{ij}^2}{q\mathbf{K}_1(1, 2)^2 / \|\mathbf{K}_1\|_F^2} \\ &= q^{l_2} \left(\frac{q\mathbf{K}_1(1, 1)^2}{\|\mathbf{K}_1\|_F^2} + \frac{q\mathbf{K}_1(1, 2)^2}{\|\mathbf{K}_1\|_F^2} \right) = q^{l_2} \cdot q = q^{l_2+1} \end{aligned}$$

Similarly, if the statement holds for $(l_{\text{row}}, l_{\text{col}}) = (l_1, l_2)$ and $l_1 \leq l_2$, the statement also holds for $(l_{\text{row}}, l_{\text{col}}) = (l_1 + 1, l_2 + 1)$. By induction, the statement holds for all $0 \leq l_{\text{row}} \leq l_{\text{col}}$. \square

This property follows from our careful design of Eq. (2), which is based on the Kronecker product. While q can be set so that the square sum of entries of $\tilde{\mathbf{A}}_\Theta$ is equal to that of \mathbf{A} , making it learnable leads to better compression since this gives more degrees of freedom to the model (see Section 6.3). As a result, the error becomes $\sum_{a_{ij} \neq 0} ((a_{ij} - \tilde{a}_{ij})^2 - \tilde{a}_{ij}^2) + q^{l_{\text{col}}}$, and thus the error and its gradient can be computed in time proportional to the number of non-zeros, without having to approximate zero entries in \mathbf{A} explicitly (see Theorem 2 in Section 4.3). It should be noticed that we do use the loss function that encourages the model to fit all entries including zeros, and we speed up its computation without changing it. Gradient descent is used for updating the model parameters.

Implementation in practice: Since candidate pairs are disjoint, processing them, including computing Eq. (4), is performed in parallel in our implementation. Shingles are also computed in parallel.

4.3 Theoretical Analysis

We analyze the time and space complexity of NEUKRON. We assume that (a) $N \leq M$ for the input matrix $\mathbf{A} \in \mathbb{R}^{N \times M}$ and (b) the dimension h of LSTM is a constant (i.e., $O(1)$), which is a user-defined hyperparameter. NEUKRON requires logarithmic time for approximation (Theorem 1), as confirmed empirically in Section 2 of [20]. For training, it requires time proportional to the number of non-zero entries of \mathbf{A} , denoted by $\text{nnz}(\mathbf{A})$ (Theorem 2).

THEOREM 1 (APPROXIMATION TIME FOR EACH ENTRY). *The approximation of each entry by NEUKRON takes $\Theta(\log M)$ time.*

PROOF. First, we need to encode the position of the given entry. Since we need the subdivision $\Theta(\log M)$ times, the time complexity of the encoding step is $\Theta(\log M)$. The computational cost to approximate an entry only depends on the length of the input of the LSTM, so the time complexity for inference is $\Theta(\log M)$. \square

THEOREM 2 (TRAINING TIME). *Each training epoch in NEUKRON takes $O(\text{nnz}(\mathbf{A}) \cdot \log M)$ time.*

PROOF. The time complexity for inference is $O(\log M)$ for each input. Thus, computing the approximation error takes $O(\text{nnz}(\mathbf{A}) \cdot \log M)$ with Eq. (5) (see Lemma 1). The time complexity for computing the gradients is also $O(\text{nnz}(\mathbf{A}) \cdot \log M)$, since the gradient of each component in the model, such as matrix multiplication and taking a non-linearity, does not require a greater time complexity. For optimizing the orders of rows and columns, computing the shingle values for rows and columns takes $O(\text{nnz}(\mathbf{A}))$ time since we need to look up all non-zero entries. Matching the rows and the columns as pairs requires $O(N + M)$ time. Only the entries of the output that correspond to non-zero entries are changed due to swaps and inference of a single element takes $O(\log M)$ time. Thus, checking the criterion in Eq. (4) takes $O(\text{nnz}(\mathbf{A}) \cdot \log M)$ time. Therefore, the overall training time per epoch is $O(\text{nnz}(\mathbf{A}) \cdot \log M)$. \square

While NEUKRON requires space proportional to the number of non-zero entries in the input matrix during training (Theorem 4), it gives a constant-size compression. (Theorem 3). Refer to Appendix D for the proofs of Theorems 3 and 4.

THEOREM 3 (SPACE COMPLEXITY OF OUTPUTS). *The number of model parameters of NEUKRON is $\Theta(1)$.*

THEOREM 4 (SPACE COMPLEXITY DURING TRAINING). *NEUKRON requires $O(\text{nnz}(\mathbf{A}) + M)$ space during training.*

5 EXTENSION TO TENSORS

We extend NEUKRON to sparse reorderable tensors. Theoretical analyses are available at Section 3 of [20].

5.1 Model

For a given D -order tensor $\mathcal{X} \in \mathbb{R}^{N_1 \times \dots \times N_D}$ (we assume $N_1 \leq \dots \leq N_D$ without loss of generality), we first compute $l_i = \lceil \log_2 N_i \rceil$ for each $i \in [D]$ and extend \mathcal{X} to the tensor of size $2^{l_1} \times \dots \times 2^{l_D}$ with additional entries filled with zeros. As in Section 4.1.4, for encoding, NEUKRON first recursively divides the extended tensor into 2^D partitions l_1 times to obtain a partition has a size of $1 \times 2^{l_2 - l_1} \times \dots \times 2^{l_D - l_1}$. Then, it recursively divides the partition as it handles a $(D - 1)$ -order tensor. As a result, the k -th entry of the encoded sequence for the position (i_1, \dots, i_D) is $(t_1(i_1, k), \dots, t_D(i_D, k))$, where t_d is identical to Eq. (3). We provide an example of NEUKRON on a 3-order tensor in Figure 2. After encoding, NEUKRON produces an approximation using the Kronecker product $\mathcal{K}_1 \otimes \dots \otimes \mathcal{K}_{l_D}$ from D linear layers for handling tensors of D different sizes.

5.2 Training Strategies

The main difference in training strategies lies in computing shingles. For a D -order tensor, D random bijective functions are used, thus each mode index has $D - 1$ shingles from those functions except for the function of the same mode. In our extension, we match positions i_1 and i_2 as a pair only if the $D - 1$ shingles of i_1 and those of index j are all the same, and the orders of indices are randomly initialized. All other procedures are identical to the original NEUKRON.

6 EXPERIMENTS

We conducted experiments to answer the following questions:

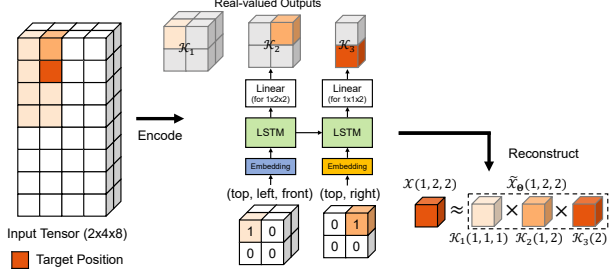


Figure 2: Example of NEUKRON on a 3-order tensor \mathcal{X} .

- Q1. **Compression Performance:** Does NEUKRON perform more compact and accurate compression than its best competitors?
- Q2. **Ablation Study:** How effective are NEUKRON’s training strategies for compression performance?
- Q3. **Scalability and Speed:** Does NEUKRON scale linearly with the number of non-zero entries of input data?
- Q4. **Approximation Analysis:** How does the approximation error of NEUKRON vary depending on entry values?
- Q5. **Effects of Data Properties:** How do the skewness, order, and dimension of the input affect the approximation error?

The answers for Q3, Q4, and Q5 are provided in Appendix A.

6.1 Experiment Specifications

Machine: We ran experiments for NEUKRON on a machine with 4 RTX 2080Ti GPUs and 128GB RAM. For competitors, which do not require GPUs, we ran experiments on a desktop with a 3.8GHz AMD Ryzen 3900X CPU and 128GB RAM. **Note that outputs and compression ratios do not depend on machine specifications.**

Datasets: We used six real-world matrices and four real-world tensors listed in Table 3. All the datasets are weighted (i.e., non-binary matrices and tensors) except for the email and threads datasets. Detailed semantics and structural properties of the datasets are provided in Table 3 and Table 7 of [20], respectively.

Competitors: For matrices, we compared NEUKRON with KronFit [25], T-SVD (truncated SVD), CMD [36], ACCAMS [2], CUR [9], and bCCAMS [2]. In order to compress matrices of any size, we extended KronFit so that it (a) fits a non-square seed matrix, (b) permutes rows and columns separately, and (c) aims to minimize the approximation error in Problem 1. We did not consider methods designed for unipartite and/or unweighted graphs (e.g., [22, 23, 29]) as competitors since they are not applicable to most of the datasets. For tensors, we compared NEUKRON with CP [1] and Tucker [18] decompositions and CSF [33], which is lossless. The competitors are described in Section 2, and see [20] for implementation details.

Experimental Setup: We trained NEUKRON and its competitors under the following stopping condition with the patience of 100 epochs: $\frac{\mathcal{E}_{\min} - \mathcal{E}_{\text{curr}}}{\mathcal{E}_{\min}} < 10^{-5}$, where \mathcal{E}_{\min} is the lowest approximation error so far, and $\mathcal{E}_{\text{curr}}$ is the current approximation error. For all experiments, we set T_p in Algorithm 2 to 2, and set γ in Eq. (4) to 10, after a preliminary study (see Section 6 of [20]). NEUKRON was trained by Adam optimizer whose learning rate was set to 10^{-3} for the email and threads datasets, and 10^{-2} for the others. Unless otherwise stated, we set the hidden dimension h to 30 in the email, nyc, and tky datasets and to 60 in the kasandr, nips, and threads datasets. For the other datasets, we set h to 90. We ran all

Table 3: Real-world datasets used in the paper. All datasets are publicly available, and links to them are available in [20].

Type	Name	Size	# of non-zeros
Matrix	email	$1,005 \times 25,919$	92,159
	nyc	$1,083 \times 38,333$	91,024
	tky	$2,293 \times 61,858$	211,955
	kasandr	$414,520 \times 503,702$	903,366
	threads	$176,445 \times 595,778$	1,457,727
	twitch	$790,100 \times 15,524,309$	234,422,289
Tensor	nips	$2,482 \times 2,862 \times 14,036$	3,101,609
	4-gram	$48K \times 54K \times 55K \times 58K$	7,495,550
	3-gram	$88K \times 100K \times 110K$	9,778,281
	enron	$5,699 \times 6,066 \times 244K$	31,312,375

experiments 5 times with different random seeds and reported the average error. The setups for the competitors are depicted in [20].

6.2 Q1. Compression Performance

We compared the (a) size in bytes³ and (b) approximation error of the compressed output obtained by the considered algorithms. We varied the hidden dimension h of NEUKRON from 5 to 30 for the email, nyc, and tky datasets and from 10 to 60 for the kasandr, nips and threads datasets. For the others, we varied h from 15 to 90. Similarly, we varied the hyperparameters of each competitor as to reveal its trade-off between the size and error (refer to [20]).

For all datasets, NEUKRON achieved the best trade-off between the approximation error and the compressed size. As seen in Figure 3, the size was up to **five orders of magnitude smaller** in NEUKRON than in the competitors when their errors were similar. The error was also up to **10.1× smaller** in NEUKRON than in the competitors when the outputs were of similar size. Note that the errors of KronFit do not always decrease as the number of parameters increases, as previously reported in [25].

Performance on non-reorderable data: NEUKRON can also be applied to non-reorderable data if the mapping between the original and new orders of rows and columns are stored additionally. Even when we assume that the datasets are non-reorderable and consider the extra cost, NEUKRON gives by far the best trade-off between size and approximation error, as shown in Figure 9.

6.3 Q2. Ablation Study

On the four smallest matrices and the two smallest tensors, we demonstrate the effectiveness of the components of NEUKRON illustrated in Section 4 by comparing it with the following variants:

- (a) NEUKRON: the proposed method with all components.
- (b) NEUKRON-H (N-H): a variant that uniformly samples pairs of rows and columns without using min-hashing.
- (c) NEUKRON-I (N-I): a variant that randomly initializes the orders of rows and columns without using the scheme in [16].
- (d) NEUKRON-F (N-F): a variant that fixes q to the sum of the squares of all entries in the input.
- (e) NEUKRON-A (N-A): a variant without any auto-regressive architecture. It only uses two learnable matrices $\mathbf{K}_{\text{square}} \in \mathbb{R}^{2 \times 2}$ and $\mathbf{K}_{\text{rect}} \in \mathbb{R}^{1 \times 2}$, as in KronFit, to compute $\mathbf{K}_{\text{square}}^{\otimes l_{\text{row}}} \otimes \mathbf{K}_{\text{rect}}^{\otimes (l_{\text{col}} - l_{\text{row}})}$ for approximation. Similarly, it uses D learnable tensors to approximate N -order tensors.

³In our implementation, each floating-point number took 4 bytes.

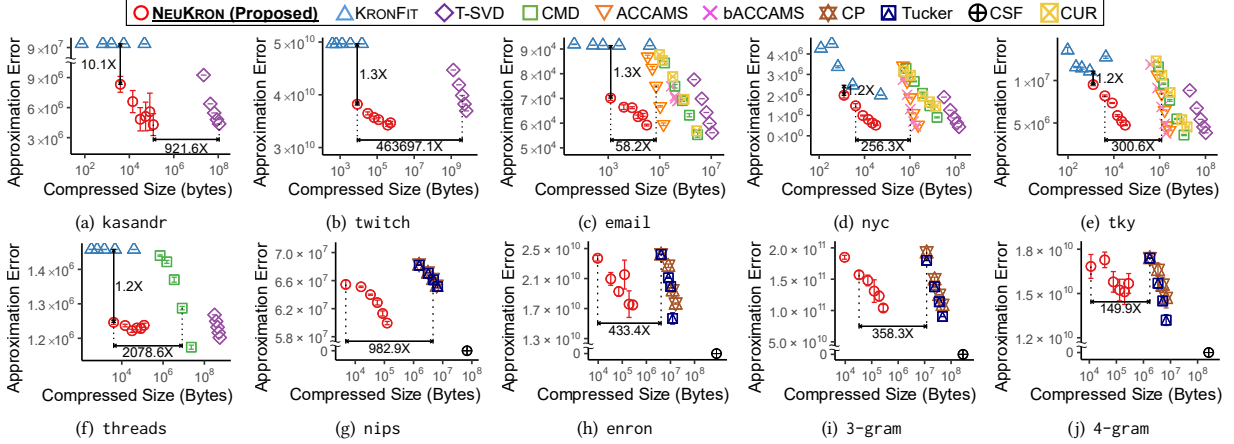


Figure 3: NEUKRON provides concise and accurate compressions. The outputs of NEUKRON are up to five orders of magnitude smaller than those of the competitors when the approximation errors in them are similar. When the sizes of the outputs are similar, the approximation error was up to 10.1 \times smaller in the outputs of NEUKRON than those in the competitors. ACCAMS, bACCAMS, CUR, and CMD ran out of memory in some datasets, and their results do not appear in the corresponding plots. Note that the errors of KronFit do not always decrease as the number of parameters increases, as previously reported in [25].

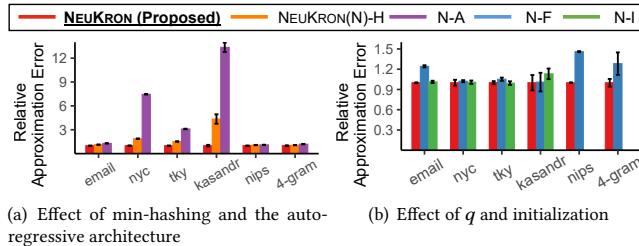


Figure 4: Effectiveness of the components of NEUKRON. We report the approximation errors of variants relative to that of NEUKRON. Results of NEUKRON-I on tensors are omitted since, for tensors, NEUKRON also randomly initializes orders.

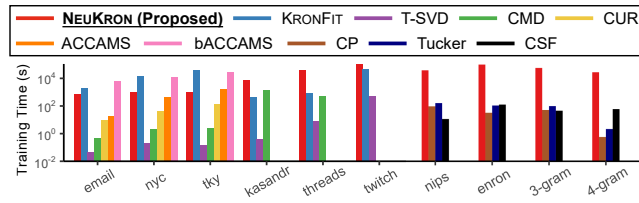


Figure 5: Training time of NEUKRON and the competitors. Note that NEUKRON requires much longer training time than many competitors, while it provides the best trade-off between space and accuracy (Figure 3). See Appendix A.1 for detailed hyperparameter settings for each method.

As seen in Figure 4, NEUKRON outperformed NEUKRON-A and NEUKRON-H, which indicates that the auto-regressive architecture (i.e., LSTM) and the min-hashing technique are crucial to enhance the performance of NEUKRON. Moreover, making q learnable was effective especially on the email, nips, and 4-gram datasets. For the order initialization, NEUKRON-I showed comparable or slightly poor performance than NEUKRON, implying that how the rows and columns are initialized can affect the compression quality.

Extra Results: For details results regarding Q3-Q5, refer to Appendix A. A training time comparison is available in Figure 5.

7 CONCLUSION

We focus on compressing sparse reorderable matrices and tensors into a constant-size space. Our contributions are three-fold:

- **Compact yet Accurate Method:** We proposed NEUKRON, which lossily compresses matrices and fixed-order tensors of any size with a constant number of parameters. NEUKRON provided an output that is up to five orders of magnitude smaller than the outputs of the best competitors when the approximation errors in them are similar (Figure 3).
- **Theoretical Analysis:** We carefully designed NEUKRON so that, for sparse reorderable matrices and fixed-order tensors of any size (a) the number of parameters is constant, (b) each entry is approximated in a logarithmic time, and (c) the model is fitted to an input in time proportional to the number of non-zero entries in it. We proved these desirable properties (Theorems 1-3).
- **Extensive Experiments:** Through extensive experiments on 10 real-world datasets, we demonstrated the effectiveness and scalability of NEUKRON (Figures 3 and 6). Especially, we showed that NEUKRON successfully compressed a matrix with up to 230 millions of non-zero entries.

Reproducibility: The code and datasets used are available at [20].

Acknowledgements

This work was supported by Samsung Electronics Co., Ltd., National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No.2021R1C1C1008526), and Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (No. 2022-0-00157, Robust, Fair, Extensible Data-Centric Continual Learning) (No. 2019-0-00075, Artificial Intelligence Graduate School Program (KAIST)) (No.2021-0-02068, Artificial Intelligence Innovation Hub).

REFERENCES

- [1] Brett W Bader and Tamara G Kolda. 2008. Efficient MATLAB computations with sparse and factored tensors. *SIAM Journal on Scientific Computing* 30, 1 (2008), 205–231.

- [2] Alex Beutel, Amr Ahmed, and Alexander J Smola. 2015. Accams: Additive clustering to approximate matrices succinctly. In *WWW*.
- [3] Paolo Boldi and Sebastiano Vigna. 2004. The webgraph framework I: compression techniques. In *WWW*.
- [4] Andrei Z Broder, Moses Charikar, Alan M Frieze, and Michael Mitzenmacher. 2000. Min-wise independent permutations. *JCSS* 60, 3 (2000), 630–659.
- [5] J Douglas Carroll and Jih-Jie Chang. 1970. Analysis of individual differences in multidimensional scaling via an N-way generalization of “Eckart-Young” decomposition. *Psychometrika* 35, 3 (1970), 283–319.
- [6] Deepayan Chakrabarti, Yiping Zhan, and Christos Faloutsos. 2004. R-MAT: A recursive model for graph mining. In *SDM*.
- [7] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *EMNLP* (2014).
- [8] Laxman Dhulipala, Igor Kabiljo, Brian Karrer, Giuseppe Ottaviano, Sergey Pupyrev, and Alon Shalita. 2016. Compressing graphs and indexes with recursive graph bisection. In *KDD*.
- [9] Petros Drineas, Ravi Kannan, and Michael W Mahoney. 2006. Fast Monte Carlo algorithms for matrices II: Computing a low-rank approximation to a matrix. *SIAM Journal on computing* 36, 1 (2006), 158–183.
- [10] Petros Drineas, Michael W Mahoney, and Shan Muthukrishnan. 2008. Relative-error CUR matrix decompositions. *SIAM J. Matrix Anal. Appl.* 30, 2 (2008), 844–881.
- [11] Carl Eckart and Gale Young. 1936. The approximation of one matrix by another of lower rank. *Psychometrika* 1, 3 (1936), 211–218.
- [12] Gene H Golub and Christian Reinsch. 1971. Singular value decomposition and least squares solutions. In *Linear algebra*. Springer, 134–151.
- [13] Ruining He and Julian McAuley. 2016. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *WWW*.
- [14] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [15] Nam Hyeon-Woo, Moon Ye-Bin, and Tae-Hyun Oh. 2022. FedPara: Low-rank Hadamard Product for Communication-Efficient Federated Learning. In *ICLR*.
- [16] Jinhong Jung and Lee Sael. 2020. Fast and accurate pseudoinverse with sparse matrix reordering and incremental approach. *Machine Learning* 109, 12 (2020), 2333–2347.
- [17] Alexandros Karatzoglou, Xavier Amatriain, Linas Baltrunas, and Nuria Oliver. 2010. Multiverse recommendation: n-dimensional tensor factorization for context-aware collaborative filtering. In *RecSys*.
- [18] Tamara G Kolda and Jimeng Sun. 2008. Scalable tensor decompositions for multi-aspect data mining. In *ICDM*.
- [19] Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. 2016. Federated optimization: Distributed machine learning for on-device intelligence. *arXiv preprint arXiv:1610.02527* (2016).
- [20] Taehyung Kwon, Jihoon Ko, Jinhong Jung, and Kijung Shin. 2023. *NeuKron: Constant-Size Lossy Compression of Sparse Reorderable Matrices and Tensors (Code, Datasets, and Appendix)*. <https://github.com/kbrother/NeuKron>
- [21] Timothée Lacroix, Nicolas Usunier, and Guillaume Obozinski. 2018. Canonical tensor decomposition for knowledge base completion. In *ICML*.
- [22] Kyuhan Lee, Hyeonsoo Jo, Jihoon Ko, Sungsu Lim, and Kijung Shin. 2020. Ssumm: Sparse summarization of massive graphs. In *KDD*.
- [23] Kristen LeFevre and Evimaria Terzi. 2010. GraSS: Graph structure summarization. In *SDM*.
- [24] Jurij Leskovec, Deepayan Chakrabarti, Jon Kleinberg, and Christos Faloutsos. 2005. Realistic, mathematically tractable graph generation and evolution, using kronecker multiplication. In *ECML/PKDD*.
- [25] Jure Leskovec, Deepayan Chakrabarti, Jon Kleinberg, Christos Faloutsos, and Zoubin Ghahramani. 2010. Kronecker graphs: an approach to modeling networks. *JMLR* 11, 2 (2010).
- [26] Jure Leskovec and Christos Faloutsos. 2007. Scalable modeling of real graphs using kronecker multiplication. In *ICML*.
- [27] Chao Liu, Fan Guo, and Christos Faloutsos. 2009. Bbm: bayesian browsing model from petabyte-scale data. In *KDD*.
- [28] Julian McAuley, Jure Leskovec, and Dan Jurafsky. 2012. Learning attitudes and attributes from multi-aspect reviews. In *ICDM*.
- [29] Matteo Riondato, David García-Soriano, and Francesco Bonchi. 2017. Graph summarization with quality guarantees. *DMKD* 31 (2017), 314–349.
- [30] Kijung Shin, Amol Ghoting, Myunghwan Kim, and Hema Raghavan. 2019. Sweg: Lossless and lossy summarization of web-scale graphs. In *WWW*.
- [31] Sumit Sidana, Charlotte Laclau, Massih R Amini, Gilles Vandelle, and André Bois-Crettez. 2017. KASANDR: a large-scale dataset with implicit feedback for recommendation. In *SIGIR*.
- [32] Arnab Sinha, Zhihong Shen, Yang Song, Hao Ma, Darrin Eide, Bo-June (Paul) Hsu, and Kuansan Wang. 2015. An Overview of Microsoft Academic Service (MAS) and Applications. In *WWW*.
- [33] Shaden Smith and George Karypis. 2015. Tensor-matrix products with a compressed sparse tensor. In *Proceedings of the 5th Workshop on Irregular Applications: Architectures and Algorithms*. 1–7.
- [34] Shaden Smith, Niranjay Ravindran, Nicholas D Sidiropoulos, and George Karypis. 2015. SPLATT: Efficient and parallel sparse tensor-matrix multiplication. In *IPDPS*.
- [35] Gilbert W Stewart. 1993. On the early history of the singular value decomposition. *SIAM review* 35, 4 (1993), 551–566.
- [36] Jimeng Sun, Yinglian Xie, Hui Zhang, and Christos Faloutsos. 2007. Less is more: Compact matrix decomposition for large sparse graphs. In *SDM*.
- [37] Daniel Ting. 2018. Count-min: Optimal estimation and tight error bounds using empirical error distributions. In *KDD*.
- [38] Ledyard R Tucker. 1966. Some mathematical notes on three-mode factor analysis. *Psychometrika* 31, 3 (1966), 279–311.
- [39] Nick Vannieuwenhoven, Raf Vandebril, and Karl Meerbergen. 2012. A new truncation strategy for the higher-order singular value decomposition. *SIAM Journal on Scientific Computing* 34, 2 (2012), A1027–A1052.
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.
- [41] Jaewon Yang and Jure Leskovec. 2012. Defining and Evaluating Network Communities Based on Ground-Truth. In *ICDM*.

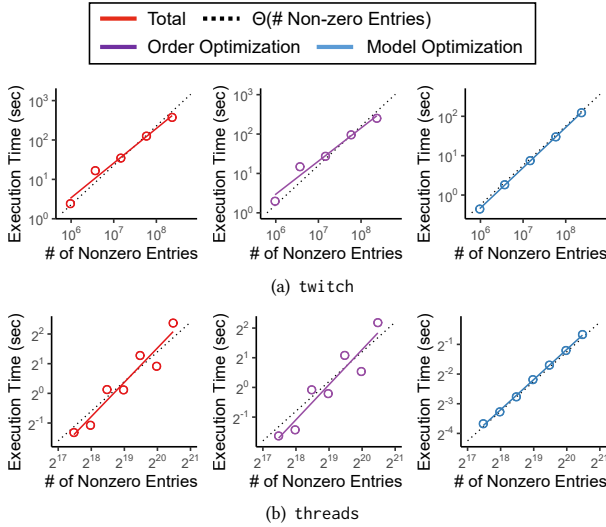


Figure 6: The training process of NEUKRON is scalable. Both model and order optimizations scale near-linearly with the number of non-zeros in the input.

A ADDITIONAL EXPERIMENTAL RESULTS

A.1 Q3. Scalability and Speed

In order to evaluate the scalability of NEUKRON, we generated multiple matrices of various sizes from the threads and twitch datasets by tracking their evolutions over time. In them, we measured the training time per epoch for model and order optimizations in addition to the total training time per epoch. The hidden dimension h was fixed to 60. As shown in Figure 6, the individual and overall training processes of NEUKRON scaled **linearly with the number of non-zeros**, which is consistent with the theoretical results in Section 4.3. We further confirmed the linear scalability of NEUKRON on tensor datasets and in hidden dimensions in Section 8 of [20].

We compared the training time of NEUKRON and the competitors in Figure 5. We followed the hyperparameter settings in Section 6.2. For NEUKRON, we reported the result with the smallest hidden dimensions that we considered. For all competitors except for KronFit, we reported their results when their approximation errors are closest to that of NEUKRON. For KronFit, we reported its result when its output size is closest to that of NEUKRON. Since our optimization problem is a mixed discrete-continuous optimization problem, which is notoriously difficult, the convergence of NEUKRON takes much longer than that of factorization-based methods. While the convergence took long, the approximation error dropped rapidly in early iterations in most cases. The detailed training curves are given in Figure 2 of [20].

A.2 Q4. Approximation Analysis

We analyzed how the approximation error by NEUKRON varies depending on the ground-truth value of approximated entries. In each dataset, we grouped the approximated entries by log-binning of their ground-truth values, and for each group, we computed the root mean squared error (RMSE) of the approximation errors. As seen in Figure 7, RMSE tended to increase with respect to ground-truth entry values. We also checked at most 1,000 largest singular values of matrices obtained by NEUKRON and the two strongest

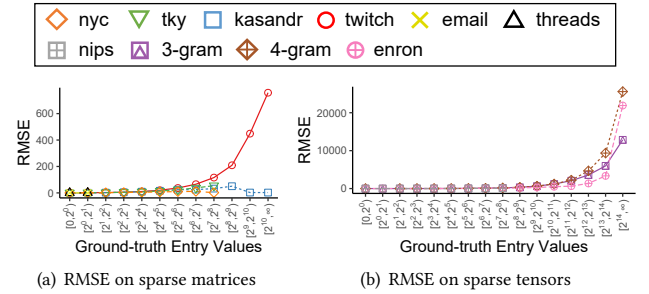


Figure 7: Analysis of approximation errors of NEUKRON. The errors tend to increase with respect to the ground truth values of approximated entries. The x-axis is in the log scale.

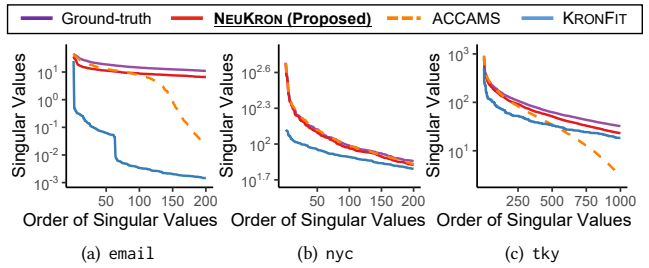


Figure 8: NEUKRON preserves singular values well. The singular values of the matrix obtained by NEUKRON are closest to the ground-truth ones. We used the smallest datasets for this experiment since computing singular values requires approximating all entries, including zeros.

competitors. For each method, we used the hyperparameter settings that led to the least approximation error in Figures 3 and 9. As seen in Figure 8, the singular values obtained by NEUKRON were closest to the singular values of the input matrices.

A.3 Q5. Effects of Data Properties

We investigated the effects of properties of an input tensor \mathcal{X} on the performance of NEUKRON. For this experiment, we synthetically generated tensors using the multi-dimensional extension R-MAT [6]. Specifically, we first split each mode of a tensor into two partitions and then chose either the first partition with probability p or the second one with probability $1 - p$. This process is repeated until the target position is determined. As a default setting, we set (a) p to 0.8, (b) the order D to 3, (c) the sum of all tensor entries to 10^6 , and (d) the number of entries to 2^{30} . We measured *fitness*, which is defined as $1 - \|\mathcal{X} - \tilde{\mathcal{X}}\|_F / \|\mathcal{X}\|_F$ (the higher, the better). The fitness is widely used to compare the errors of approximations to different tensors. We varied the skewness p from 0.65 to 0.85. Note that increasing p makes the distribution of non-zero entries more skewed with distinct patterns, and decreasing p makes the distribution more uniform without patterns. As seen in Figure 10(a), the fitness increased as p increased, implying that NEUKRON provides better performance on skewed tensors with distinct patterns. Next, we changed the order D from 2 to 6, but no significant effect of D was observed, as shown in Figure 10(b). Lastly, we analyzed the effect of dimension (i.e., the number of indices in each mode) by changing the dimension of tensors while fixing the number of non-zeros. As seen in Figure 10(c), the fitness of NEUKRON decreased as

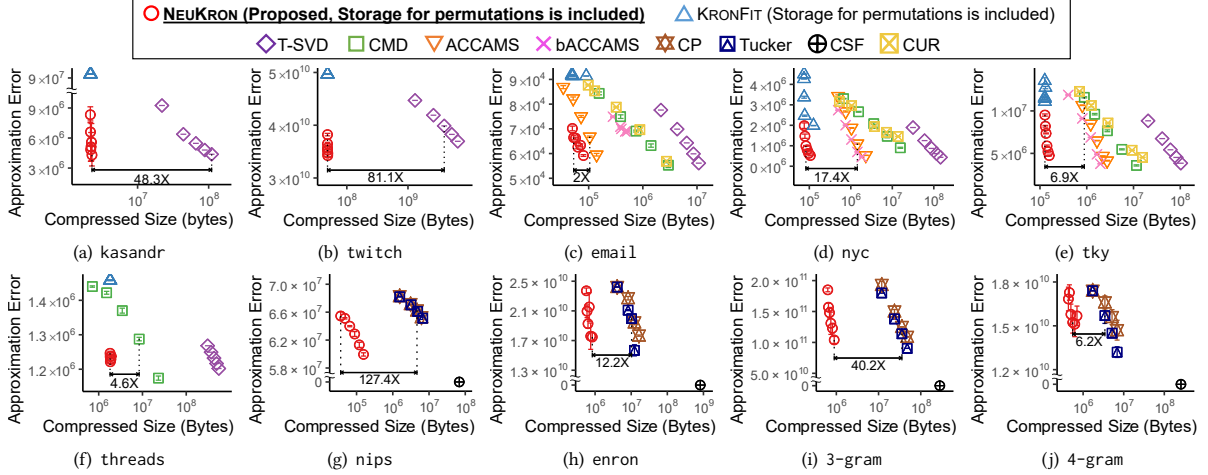


Figure 9: NEUKRON significantly outperforms the competitors even if we assume that matrices and tensors are non-reorderable and separately store the permutations of indices for all modes. Note that the outputs of NEUKRON require up to two orders of magnitude smaller space than those of the competitors with similar approximation error.

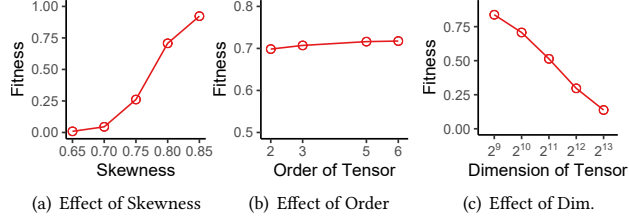


Figure 10: Effects of data properties on NEUKRON. (a) The fitness increases as the skewness p increases. (b) The order of tensors does not significantly affect the fitness. (c) As tensors become bigger, the fitness decreases.

tensors became bigger and thus more entries were approximated by a fixed number of parameters.

B PSEUDOCODE FOR ORDERING ROWS

The pseudocode of UPDATEROWORDER described in Section 4.2.1, is given in Algorithm 3, where binary representations start from 0, while row indices start from 1.

C EFFECTIVENESS OF NEUKRON ON NON-REORDERABLE DATA

NEUKRON can also be applied to non-reorderable matrices and tensors if the mapping between the original and new orders of mode indices are stored additionally. Even with this additional space requirement, NEUKRON still yielded the best trade-off between the approximation error and the compressed size, as seen in Figure 9, where we assume that the input matrices and tensors are not reorderable. Remark that KronFit also needs space for storing orders when it is applied to non-reorderable matrices.

D PROOF OF THEOREMS

Proof of Theorem 3: In NEUKRON, the number of the parameters for LSTM is $\Theta(h^2)$. The embedding layer before the LSTM and the linear layers after LSTM require $\Theta(h)$ of parameters. The number of parameters for \mathbf{K}_1 in Algorithm 1 is 4, which is the number

Algorithm 3: UPDATEROWORDER

Input: (a) a reorderable matrix \mathbf{A} , (b) a hyperparameter γ in Eq. (4)
Output: updated matrix \mathbf{A}

- 1 Sample $k \in \mathbb{N} \cup \{0\}$ so that $P(k=i) = 1/2^{i+1}$
- 2 $k \leftarrow \min(k, l_{\text{row}} - 1)$; $R \leftarrow \emptyset$; $P \leftarrow \emptyset$
- 3 Generate a random hash bijective functions h_{col}
- 4 **foreach** $i \in \{i \in [n] : (i-1) \text{ AND } 2^k = 0\}$ **do**
- 5 $u \sim U(0, 1)$
- 6 **if** $u < 1/2$ **then** $R \leftarrow R \cup \{i\}$
- 7 **else** $R \leftarrow R \cup \{i + 2^k\}$
- 8 **foreach** $i \in R$ **do**
- 9 $f_{\text{row}}(i) \leftarrow \min_{a_{ij} \neq 0} (h_{\text{col}}(j))$
- 10 **while** $\exists (i_1, i_2)$ s.t. $f_{\text{row}}(i_1) = f_{\text{row}}(i_2)$ **do**
- 11 $P \leftarrow P \cup \{(i_1, (i_2 - 1) \text{ XOR } 2^k + 1), (i_2, (i_1 - 1) \text{ XOR } 2^k + 1)\}$
- 12 $R \leftarrow R \setminus \{i_1, i_2\}$
- 13 $R \leftarrow R \cup \{(r-1) \text{ XOR } 2^k + 1 : r \in R\}$
- 14 **while** $R \neq \emptyset$ **do**
- 15 Randomly sample (i_1, i_2) from R
- 16 $P \leftarrow P \cup \{(i_1, i_2)\}$; $R \leftarrow R \setminus \{i_1, i_2\}$
- 17 $P_{\text{accept}} \leftarrow \emptyset$
- 18 **foreach** $(i_1, i_2) \in P$ **do**
- 19 $u \sim U(0, 1)$
- 20 $\Delta \leftarrow$ change in the approximation error
- 21 **if** $u \geq \exp(-\gamma \cdot \Delta)$ **then** $P_{\text{accept}} \leftarrow P_{\text{accept}} \cup \{(i_1, i_2)\}$
- 22 **foreach** $(i_1, i_2) \in P_{\text{accept}}$ **do**
- 23 $\mathbf{A}_{i_1,:}, \mathbf{A}_{i_2,:} \leftarrow \mathbf{A}_{i_2,:}, \mathbf{A}_{i_1,:}$

of entries. We consider $\Theta(h) = \Theta(1)$ as h is a constant; thus, the number of parameters is $\Theta(1)$.

Proof of Theorem 4: Storing the input matrix in a sparse format requires $O(\text{nnz}(\mathbf{A}))$ space. Changing the orders of rows and columns requires $O(M+N)$ space for saving random hash functions, shingles, and changes of losses. If we assume that the batch size and the number of parameters of LSTM are constants, its memory usage during inference and backpropagation is $O(\log M)$ since the input length is $O(\log M)$. Thus, the overall space complexity during training is $O(\text{nnz}(\mathbf{A}) + M)$.

NeuKron: Constant-Size Lossy Compression of Sparse Reorderable Matrices and Tensors (Supplementary Document)

1 ANALYSIS FOR THE TYPE OF THE SEQUENTIAL MODEL (RELATED TO SECTION 4.1)

We compared the performances of auto-regressive sequence models, when they are equipped with NEUKRON. We varied the hidden dimension h of NEUKRON from 5 to 30 for LSTM and GRU, and the model dimension d_{model} from 8 to 32 for the decoder layer of Transformer. As seen in Figure 11, when equipped with NEUKRON, LSTM and GRU performed similarly, outperforming the decoder layer of Transformer.

2 ANALYSIS ON INFERENCE TIME (RELATED TO SECTION 4.3)

We measure the inference time for 10^6 elements randomly chosen from square matrices of which numbers of rows and cols vary from 2^7 to 2^{16} . We ran 5 experiments for each size and report the average of them. As expected from Theorem 1 of the main paper, the approximation of each entry by NEUKRON is almost in $\Omega(\log M)$ (see Figure 13).

3 ANALYSIS OF THE TENSOR EXTENSION (RELATED TO SECTION 5)

Below, we analyze the time and the space complexities of NEUKRON extended to sparse reorderable tensors. For all proofs, we assume a D -order tensor $\mathcal{X} \in \mathbb{R}^{N_1 \times \dots \times N_D}$ where $N_1 \leq N_2 \leq \dots \leq N_D$, without loss of generality. The complexities are the same with those in Section IV of the main paper if we assume a fixed-order tensor (i.e., if $D = O(1)$).

THEOREM 5 (APPROXIMATION TIME FOR EACH ENTRY). *The approximation of each entry by NEUKRON takes $\Theta(D \log N_D)$ time.*

PROOF. For encoding, NEUKRON subdivides the input tensor $O(\log N_D)$ times and each subdivision takes $O(D)$. For approximation, the length of the input of the LSTM equals to the number of the subdivisions, so the time complexity for retrieving each entry is $O(D \log N_D)$. \square

THEOREM 6 (TRAINING TIME). *Each training epoch in NEUKRON takes $O(\text{nnz}(\mathcal{X}) \cdot D \log N_D + D^2 N_D)$.*

PROOF. Since the time complexity for inference is $O(D \log N_D)$ for each input, model optimization takes $O(\text{nnz}(\mathcal{X}) \cdot D \log N_D)$. For reordering, the time complexity of matching the indices as pairs for all the dimensions is bounded above to $O(D \cdot (DN_D)) = O(D^2 N_D)$. For checking the criterion for all pairs, we need to retrieve all the non-zero entries, and it takes $O(\text{nnz}(\mathcal{X}) \cdot D \log N_D)$. Therefore, the overall training time per epoch is $O(\text{nnz}(\mathcal{X}) \cdot D \log N_D + D^2 N_D)$. \square

THEOREM 7 (SPACE COMPLEXITY DURING TRAINING). *NEUKRON requires $O(D \cdot \text{nnz}(\mathcal{X}) + D^2 N_D)$ space during training.*

PROOF. The bottleneck is storing the input tensor in a sparse format, the random hash functions and the shingle values, which require $O(D \cdot \text{nnz}(\mathcal{X}))$, $O(\sum_{i=1}^D N_i)$, and $O((D-1) \cdot \sum_{i=1}^D N_i)$, respectively. Thus, the overall complexity during training is $O(D \cdot \text{nnz}(\mathcal{X}) + (D-1) \cdot \sum_{i=1}^D N_i) = O(D \cdot \text{nnz}(\mathcal{X}) + D^2 N_D)$. \square

THEOREM 8 (SPACE COMPLEXITY OF OUTPUTS). *The number of model parameters of NEUKRON is $\Theta(2^D)$.*

PROOF. In NEUKRON, the number of parameters for LSTM does not depend on the order of the input tensor; thus, it is still in $\Theta(1)$. The embedding layer and the linear layers connected to the LSTM require $\Theta(2 + 2^2 + \dots + 2^D) = \Theta(2^D)$ parameters. \square

4 SEMANTICS AND PROPERTIES OF DATASETS (RELATED TO SECTION 6.1)

We provide the semantics of the datasets in Table 3 and the distributions of degrees, entry values, and connected-component sizes in Table 7. For degrees, we computed the sums of the rows and those of the columns for matrices. For connected-component sizes, we treated sparse matrices as bipartite graphs and used the number of nodes in each connected component as its size. Note that these properties are naturally extended to the tensors.

5 IMPLEMENTATION DETAILS (RELATED TO SECTION 6.1)

We implemented NEUKRON in PyTorch. We implemented the extended version of KronFit in C++. For ACCAMS, bACCAMS, and CMD, we used the open-source implementations provided by the authors. We used the svds function of SciPy for T-SVD. We used the implementations of CP and Tucker decompositions in Tensor Toolbox [?] in MATLAB. Below, we provide the detailed hyperparameter setups of each competitor.

- **KronFit:** The maximum size of the seed matrix was set as follows - email: 32×161 , nyc: 33×196 , tky: 14×40 , kasandr: 75×80 , threads: 57×85 , twit ch: 30×63 . We tested the performance of KronFit when γ is 1 and 10, and fixed γ to 10 because it performs better when γ is set to 10. We performed a grid search for the learning rate in $\{10^{-1}, 10^{-2}, \dots, 10^{-8}\}$.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

TheWebConf'23, May 1–5, 2023, Austin, TX, USA
 © 2023 Association for Computing Machinery.
 ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

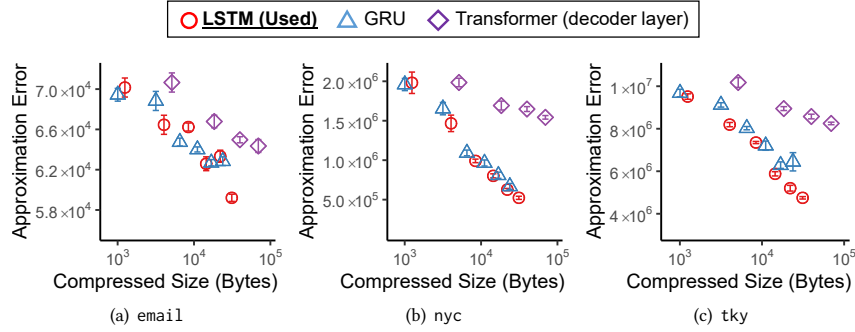


Figure 11: When equipped with NEUKRON, LSTM leads to concise and accurate compression. NEUKRON with LSTM and that with GRU show perform similarly, but NEUKRON with the decoder layer of Transformer requires significantly more space for the same level of approximation error.

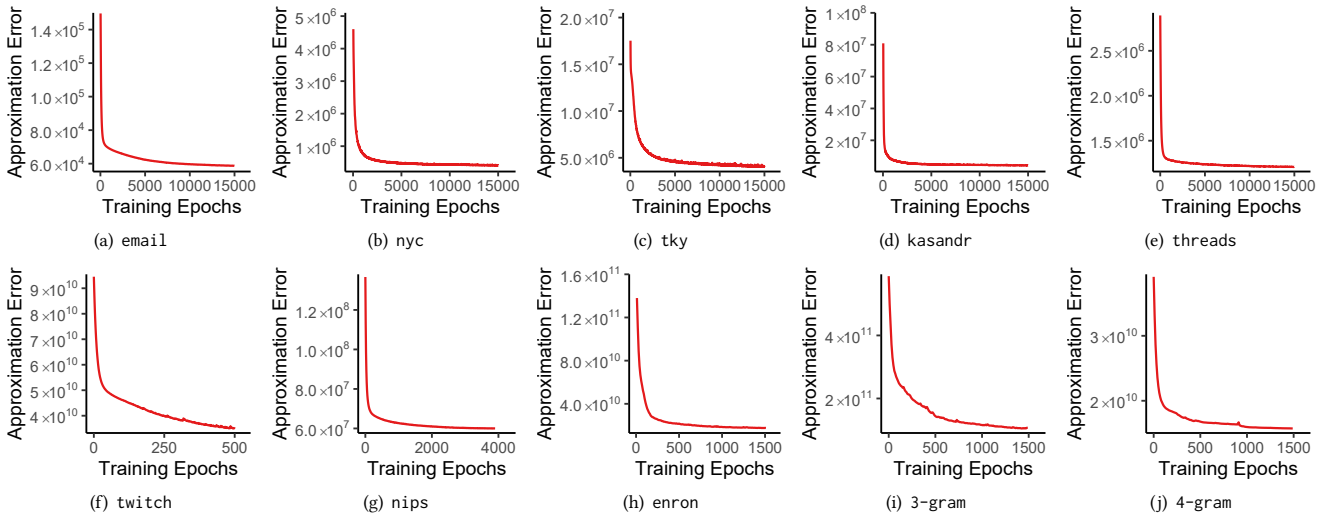


Figure 12: Approximation error of NEUKRON after each epoch. In most cases, the approximation error drops rapidly in early iterations, especially within one third of the total epochs that are determined by the termination condition in Section 6.1.

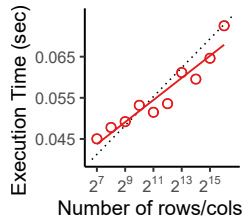


Figure 13: Inference time of NEUKRON is nearly proportional to the logarithm of the number of rows and columns.

- **T-SVD:** The ranks were up to 50 for email, 460 for nyc, 200 for tky, 15 for kasandr, 90 for threads, and 50 for twitch.
- **CUR:** We selected ranks for CUR from {10, 100, 1000}. We sampled {1%, 1.25%, 2.5%, 5%, 10%} of rows and columns in email, {3.3%, 5%, 10%, 14.3%, 20%} of rows and columns in nyc, and {1%, 2%, 4%, 8.3%, 11.1%} of rows and columns in tky
- **CMD:** We sampled (# rows, # columns) as much as {(30, 150), (60, 350), (90, 700), (100, 1400), (150, 2500)} for email, {(65, 2125), (125, 4250), (250, 8500), (500, 17000), (1000, 34000)} for nyc, {(45, 1315), (90, 2625), (175, 5250), (350, 10500), (700, 21000)} for tky, and {(55, 184), (109, 368), (218, 736), (436, 1471), (871, 2941)} for threads.

- **ACCAMS:** We used 5, 50, and 50 stencils for email, nyc, and tky, respectively. We used up to 48, 64, and 40 clusters of rows and columns for the aforementioned datasets, respectively.
- **bACCAMS:** We set the maximum number of clusters of rows and columns to 48, 48, and 24 for email, nyc, and tky, respectively. We used 50 stencils for the datasets.
- **CP:** The ranks were set up to 40 for nips, 8 for enron, 20 for 3-gram, and 4 for 4-gram.
- **Tucker:** We used hypercubes as core tensors. The maximum dimension of a hypercube for each dataset is as follows - nips: 40, enron: 6, 3-gram: 20, and 4-gram: 4.

We followed the default setting in the official code from the authors for the other hyperparameters of ACCAMS and bACCAMS. The implementations of KronFit, T-SVD, CP, and Tucker used 8 bytes for real numbers. The implementations of ACCAMS and bACCAMS used 4 bytes for real numbers and assumed the Huffman coding for clustering results.

6 HYPERPARAMETER ANALYSIS (RELATED TO SECTION 6.1)

We investigate how the approximation error of NEUKRON varies depending on γ values. We considered three γ values and four

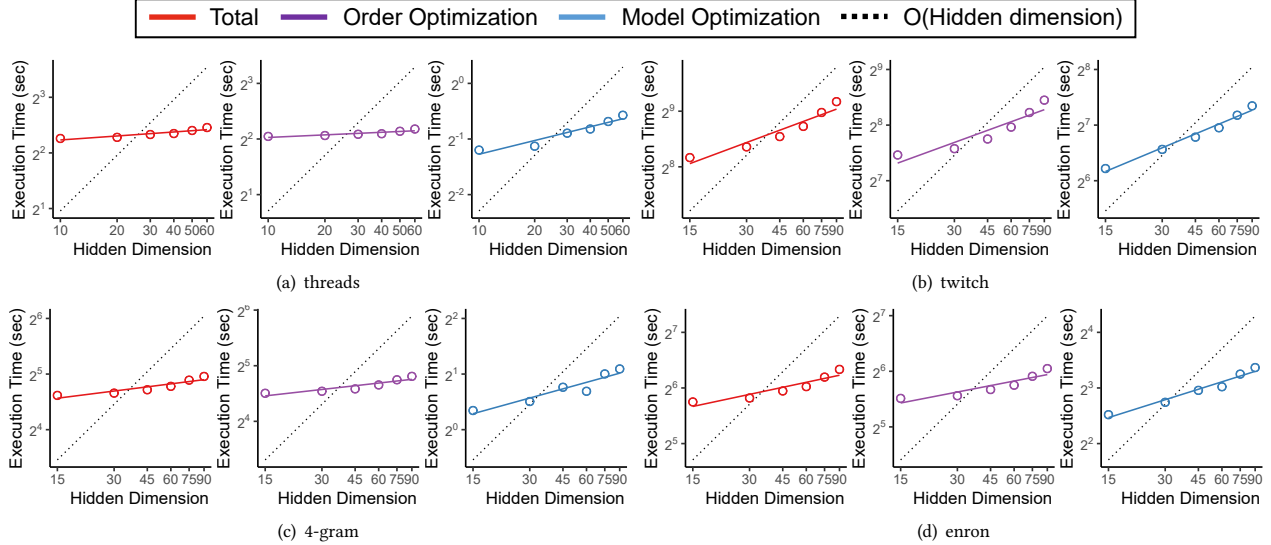


Figure 14: The training time of NEUKRON is empirically sub-linear in the hidden dimension h of NEUKRON. We measure the total elapsed time, the elapsed time for order optimization, and the elapsed time for model optimization.

Table 3: Semantics of Real-world Datasets

Name	Description
email	e-mail addresses \times e-mails [binary]
nyc	venues \times users [check-in counts]
tky	venues \times users [check-in counts]
kasandr	offers \times users [clicks]
threads	users \times threads [participation]
twitch	streamers \times users [watching time]
nips	papers \times authors \times words [counts]
4-gram	words \times words \times words \times words [counts]
3-gram	words \times words \times words [counts]
enron	receivers \times senders \times words [counts]

datasets (email, nyc, tky, and kasandr) and reported the approximation error in Table 4. Note that setting γ to ∞ results in a hill climbing algorithm that switches rows/column in pairs only if the approximation error decreases. The results show that, empirically, the approximation error was smallest when γ was set to 10 on all datasets except for the nyc dataset.

7 SPEED AND SCALABILITY ON HIDDEN DIMENSION (RELATED TO APPENDIX A.1)

We report the average the training time per epoch of NEUKRON in Table 6. The training time per epoch varied from less than 1 second to more than 9 minutes depending on the dataset. As seen in Figure 12, the training plots of all datasets dropped dramatically within one third of total epochs that were determined by the termination condition in Section 6.1. Thus, a model that worked well enough could be obtained before convergence.

We also analyzed the effect of the hidden dimension h on the training time per epoch of NEUKRON. As seen in Figure 14, both the elapsed time for order optimization and the elapsed for model optimization were empirically sublinear in the hidden dimension.

Table 4: The effect of γ on approximation error. We report the means and standard errors of approximation errors on the email, nyc, tky, and kasandr datasets.

Dataset	γ	Approximation error
email	1	90561.25 \pm 467.996
	10	58691.88 \pm 335.143
	∞	59113.75 \pm 891.544
nyc	1	421451.2 \pm 4842.068
	10	402673.6 \pm 17291.959
	∞	397947.5 \pm 2393.016
tky	1	4166292.3 \pm 143013.605
	10	3981669.6 \pm 91907.201
	∞	4034389.1 \pm 48117.964
kasandr	1	6315784.36 \pm 140974.6535
	10	4300280.71 \pm 488804.599
	∞	4385800.32 \pm 496004.629

8 SCALABILITY ON TENSOR DATASETS (RELATED TO APPENDIX A.1)

For the 4-gram and enron datasets, we generated multiple smaller tensors by sampling non-zero entries uniformly at random. The hidden dimension was fixed to 60. Consistently with the results on matrices, the overall training process of NEUKRON is also linearly scalable on sparse tensors, as seen in Figure ??.

9 COMPARISON OF LOSSY COMPRESSION METHODS (RELATED TO SECTION 2)

In Table 5, we provide a comparison of lossy compression methods for sparse matrices and tensors, which supplement Table 1 in the main paper.

Table 5: Comparison of lossy-compression methods for sparse matrices and tensors. $nnz(\mathcal{X})$: the number of non-zeros in a matrix/tensor \mathcal{X} . D : the order of the input tensor. N & M : the numbers of rows and columns of the input matrix. N_{\max} : the maximum dimensionality (i.e., $N_{\max} = \max(N_1, \dots, N_D)$). R : rank. S_c & S_r : the numbers of sampled rows and columns. h : the hidden dimension of the model of NEUKRON. T : the number of iterations of an inner loop. k : the number of clusters of rows and columns. w : the weight parameter for the criterion of switching. α, β : parameters for the probability distributions of clusters. E_r, E_c : the number of rows and columns of a seed matrix.

Methods	Training Complexity	Inference Complexity	Hyperparameters
NEUKRON	$O(h^2 nnz(\mathcal{X}) \log(M))$	$O(h^2 \log(N_{\max}))$	$h, w, \text{optimizer, learning rate}$
T-SVD [? ?]	$O(nnz(\mathcal{X})R + R^3)$	$O(R)$	R
CMD [?]	$O(nnz(\mathcal{X}) + S_c^3 + S_c S_r)$	$O(S_r)$	S_c, S_r
CUR [?]	$O(nnz(\mathcal{X}) + S_c^3 + S_c^2 S_r)$	$O(S_r)$	S_c, S_r
ACCAMS [?]	$O(NM + nnz(\mathcal{X})Tk)$	$O(R)$	k, R
bACCAMS[?]	$O(T\{k(N + M) + nnz(\mathcal{X}) + NM + k^2\})$	$O(R)$	k, R, α, β
KronFit [? ?]	$O(nnz(\mathcal{X}) \log(M))$	$O(\log(N_{\max}))$	$E_r, E_c, \text{optimizer, learning rate}$
CP[?]	$O(nnz(\mathcal{X})DR)$	$O(DR)$	R
Tucker[?]	$O(nnz(\mathcal{X})DR)$	$O(DR^D)$	R

Table 6: Training time per epoch on all datasets. We report the means and standard errors.

Dataset (Hidden Dimension)	Training time
email (30)	0.19 ± 0.010
nyc (30)	0.21 ± 0.004
tky (30)	0.32 ± 0.005
kasandr (60)	1.93 ± 0.005
threads (60)	5.49 ± 0.012
twitch (90)	566.82 ± 3.308
nips (50)	6.31 ± 0.081
enron (90)	80.69 ± 0.266
3-gram (90)	27.19 ± 0.089
4-gram (90)	41.09 ± 0.785

Table 7: Structural properties of real-world datasets.

Dataset	Degrees			Entry Values	Connected Components	
email						
nyc						
tky						
kasandr						
threads						
twitch						
nips						
enron						
3-gram						
4-gram						