# Accelerating Sparse Tensor Decomposition Using Adaptive Linearized Representation

Jan Laukemann, Ahmed E. Helal, S. Isaac Geronimo Anderson, Fabio Checconi, Yongseok Soh, Jesmin Jahan Tithi, Teresa Ranadive, Brian J Gravelle, Fabrizio Petrini, and Jee Choi

**Abstract**—High-dimensional sparse data emerge in many critical application domains such as healthcare and cybersecurity. To extract meaningful insights from massive volumes of these multi-dimensional data, scientists employ unsupervised analysis tools based on tensor decomposition (TD) methods. However, real-world sparse tensors exhibit highly irregular shapes and data distributions, which pose significant challenges for making efficient use of modern parallel processors. This study breaks the prevailing assumption that compressing sparse tensors into coarse-grained structures (i.e., tensor slices or blocks) or along a particular dimension/mode (i.e., mode-specific) is more efficient than keeping them in a fine-grained, mode-agnostic form. Our novel sparse tensor representation, Adaptive Linearized Tensor Order (ALTO), encodes tensors in a compact format that can be easily streamed from memory and is amenable to both caching and parallel execution. In contrast to existing compressed tensor formats, ALTO constructs one tensor copy that is agnostic to both the mode orientation and the irregular distribution of nonzero elements. To demonstrate the efficacy of ALTO, we accelerate popular TD methods that compute the Canonical Polyadic Decomposition (CPD) model across different types of sparse tensors. We propose a set of parallel TD algorithms that exploit the inherent data reuse of tensor computations to substantially reduce synchronization overhead, decrease memory footprint, and improve parallel performance. Additionally, we characterize the major execution bottlenecks of TD methods on multiple generations of the latest Intel Xeon Scalable processors, including Sapphire Rapids CPUs, and introduce dynamic adaptation heuristics to automatically select the best algorithm based on the sparse tensor characteristics. Across a diverse set of real-world data sets, ALTO outperforms the state-of-the-art approaches, achieving more than an order-of-magnitude speedup over the best mode-agnostic formats. Compared to the best mode-specific formats, which require multiple tensor copies, ALTO achieves $5.1\times$ geometric mean speedup at a fraction ($25\%$) of their storage costs. Moreover, ALTO obtains $8.4\times$ geometric mean speedup over the state-of-the-art memoization approach, which reduces computations by using extra memory, while requiring $14\%$ of its memory consumption.

**Index Terms**—Sparse tensors, tensor decomposition, ALTO, MTTKRP, CP-APR, Poisson tensor decomposition, Alternating Poisson Regression, Multi-core CPU, Sapphire Rapids

✦

## 1 INTRODUCTION

Tensors, the higher-order generalization of matrices, can naturally represent complex interrelations in multi-dimensional sparse data, which emerge in important application domains such as healthcare [1], [2], cybersecurity [3], [4], data mining [5], [6], and machine learning [7], [8]. For example, one mode (or dimension) of a tensor may identify users while another mode details their demographic information, and their (potentially incomplete) ratings for a set of products [9]. To effectively analyze such high-dimensional data, tensor decomposition (TD) is used to reveal their principal components, where each component represents a latent property. One of the most popular TD models is the Canonical Polyadic Decomposition (CPD), which approximates a tensor as a sum of a finite number of rank-one tensors such that each rank-one tensor corresponds to a tensor component, or a latent property [10], [11]. An important class of real-world, high-dimensional data sets is sparse tensors with *non-negative count data* [12], [13], which encodes critical information, such as the number of packets exchanged across a network, or the number of criminal activities in a city. For these tensors, the CP Alternating Poisson Regression (CP-APR) algorithm is a powerful tool for detecting anomalies and group relations. In contrast to other CPD algorithms that assume a Gaussian distribution for randomly distributed data (e.g., CP Alternating Least Squares, or CP-ALS), CP-APR [14] assumes a Poisson distribution which better describes the target count data.

Due to the curse of dimensionality, data become more dispersed as the number of modes increases. Hence, high-dimensional data sets typically suffer from highly irregular shapes and data distributions as well as unstructured and extreme sparsity, which make them challenging to represent efficiently. For instance, Figure 1 illustrates the spatial distribution of nonzero elements in a set of sparse tensors. It shows that the number of nonzero elements in a subspace, or a multi-dimensional block, can vary greatly (note the use of logarithmic scale). Furthermore, as the sparsity of tensors increases (e.g., NELL-1, AMAZON, and REDDIT tensors), the likelihood of finding dense structures in the multi-dimensional space significantly decreases, leading to extremely small numbers of nonzero elements per block.

- Jan Laukemann is with Friedrich-Alexander-Universität Erlangen-Nürnberg and Erlangen National High Performance Computing Center. Email: jan.laukemann@fau.de.
- Ahmed E. Helal, Fabio Checconi, Jesmin Jahan Tithi, and Fabrizio Petrini are with Intel Labs. Email: {ahmed.helal, fabio.checconi, jesmin.jahan.tithi, fabrizio.petrini}@intel.com.
- Teresa Ranadive and Brian J Gravelle are with Laboratory for Physical Sciences. Email: {tranadive, bjgrave}@lps.umd.edu.
- S. Isaac Geronimo Anderson, Yongseok Soh, and Jee Choi are with University of Oregon. Email: {igeroni3, ysoh, jeec}@uoregon.edu.
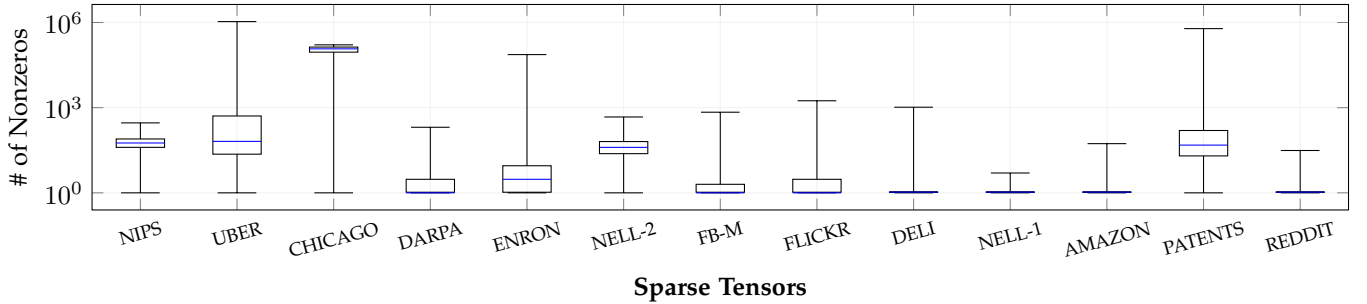
**Sparse Tensors**

Fig. 1: A box plot of the data (nonzero elements) distribution across the multi-dimensional blocks (subspaces) of the hierarchical coordinate storage [15]. The multi-dimensional subspace size is $128^N$, where $N$ is the number of dimensions (modes), as per prior work [16]. The sparse tensors are sorted in an increasing order of their number of nonzero elements. Sparsity is extremely high for many tensors (e.g., NELL-1, AMAZON, and REDDIT) and vary greatly across tensors.

Therefore, efficient execution of TD algorithms on modern parallel processors is challenging because of their low arithmetic intensity [17], random memory access, workload imbalance [18], [19], and data dependencies [17], [19]. Moreover, TD algorithms require computations along *every* mode, and realizing acceptable performance across all modes is difficult without using multiple mode-specific tensor copies.

Prior work on the CP-ALS algorithm improved the parallel performance of the matricized tensor times Khatri-Rao product (MTTKRP) kernel [15], [17], [18], [20], [21], which is the main performance bottleneck of the overall algorithm [20]. The few performance studies conducted on the CP-APR algorithm focused primarily on performance portability [22] and streaming analysis [23], rather than parallel performance optimization. To the best of our knowledge, this paper presents the first in-depth analysis of the key performance bottlenecks of CP-APR, and significantly improves the parallel performance of both CP-ALS and CP-APR over prior state-of-the-art approaches by using a linearized mode-agnostic sparse tensor representation.

Additionally, the previous approaches relied on extending legacy sparse linear algebra formats and algorithms to tensor (multilinear) algebra problems [11], [15], [20], [21], [24], [25], [26], [27], [28], [29]. These techniques can be classified based on their compression of the nonzero elements into raw or compressed formats [30], [31]. Raw formats use simple list-based representations to keep the nonzero elements along with their multi-dimensional coordinates [32]. Hence, they are mode-agnostic and typically require one tensor copy to execute tensor operations along different modes. As a result, the list-based coordinate (COO) format remains the de facto data structure for storing sparse tensors [32] in many tensor libraries (e.g., Tensor Toolbox [11], Tensorflow [33], and Tensorlab [34]). However, due to their unprocessed nature, list-based formats suffer from significant parallel and synchronization overheads [21].

Compressed tensor formats [15], [20], [27], [35] use tree- or block-based structures to organize the nonzero elements, which may decrease the memory footprint of sparse tensors. However, since these approaches rely on finding clusters of nonzero elements in *non-overlapping regions* of the multidimensional space to achieve compression, their efficacy depends on the spatial data distribution, which can be highly irregular and extremely sparse as demonstrated in Figure 1. Therefore, instead of reducing memory storage, compressed formats can introduce substantial memory overhead and degrade the parallel performance of TD algorithms [30].

The most popular compressed format for TD algorithms is compressed sparse fiber (CSF) [20], which extends the classical compressed sparse row (CSR) format to higher-order tensors using tree-like structures. However, CSF-based formats are mode-specific, where the arrangement of nonzero elements depends not only on the mode considered as the root of the index tree but also on the other modes at each subsequent tree level; therefore, they are typically efficient for *only* that specific mode order. This leads to a trade-off between performance and memory, as storing multiple tensor copies, each arranged for a specific mode, yields the best performance at the cost of extra memory [35]. For large-scale tensors, keeping multiple copies may be infeasible, especially for hardware accelerators with limited memory capacity (e.g., GPUs). Although current GPUs (e.g., H100) can have up to 80 GB of device memory [36], encoding large-scale tensors (e.g., AMAZON, PATENTS and REDDIT tensors) using CSF requires hundreds of gigabytes of memory. In contrast, keeping only one tensor copy, arranged for an arbitrarily chosen mode order, yields the smallest memory footprint at the cost of sub-optimal performance [25]. Alternatively, memoization schemes [37], [38] utilize CSF-based formats to reduce computations by keeping and reusing intermediate results across tensor modes, which require substantial storage that largely exceeds the memory consumption of tensor representations [38]. While block-based formats [15], [27] can be mode-agnostic, their storage requirements and parallel performance still depend on the spatial data distributions as well as the parameters of the blocking/tiling schemes [27], [30], which are difficult to determine dynamically.

To overcome these limitations, we present the *Adaptive Linearized Tensor Order (ALTO)* format. ALTO is a mode-agnostic representation that maps a set of $N$-dimensional coordinates onto a single linearized index such that neighboring nonzero elements in the multi-dimensional space are close to each other in memory. This leads to a more cache friendly and memory-scalable tensor storage; that is, ALTO utilizes the inherent data locality of sparse tensors and its storage scales with mode lengths, rather than the number of modes. Additionally, ALTO enables a unified implementation of tensor algorithms that requires a single tensor copy to compute along all modes.

In contrast to prior compressed [15], [20] and linearized [39] TD approaches, we propose a set of parallel algorithms that leverage the ALTO format to address the

performance bottlenecks that have traditionally limited the scalability of sparse tensor computations. Our ALTO-based algorithms generate perfectly balanced tensor partitions in terms of the number of nonzero elements; however, these partitions may divide the multi-dimensional space of a tensor into *overlapping* regions (but each nonzero element belongs to a single partition). Thus, we devise data-aware adaptation heuristics that greatly improve parallel performance by locating the overlapping space between tensor partitions and selecting the best tensor traversal and conflict resolution method according to the inherent data locality of sparse tensors. Moreover, these heuristics choose between recomputing or reusing intermediate results, depending on the properties of the target tensor computations, to maximize the performance while reducing the overall memory footprint. As a result, our ALTO-based TD algorithms deliver substantial performance gains over prior state-of-the-art approaches, while allowing the processing of large-scale tensors. In summary, we make the following contributions:

- We present ALTO, a novel sparse tensor format for high-performance tensor algorithms. Unlike prior compressed formats, ALTO uses a single (mode-agnostic) tensor representation that improves data locality, eliminates workload imbalance, and reduces memory usage, *regardless of the data distribution in the multi-dimensional space* (§3).
- We propose efficient ALTO-based parallel algorithms for the CP-ALS and CP-APR methods as well as input-aware adaptation heuristics to balance data reuse and memory footprint while greatly reducing synchronization cost (§4).
- We conduct an in-depth performance analysis of the main TD kernels and compare against prior state-of-the-art across two generations of the latest Intel Xeon Scalable processors. The results show that on an Intel Sapphire Rapids server, ALTO-based TD algorithms achieve $25.3\times$ and $5.1\times$ geometric mean *speedup* over the best mode-agnostic and mode-specific formats, respectively. Compared to the best memoization method that trades lower computation for higher memory usage, ALTO attains $8.4\times$ geometric mean *speedup*. Furthermore, ALTO requires a fraction ($14\%$ to $25\%$) of the storage used by the state-of-the-art mode-specific and memoization approaches (§5).

## 2 BACKGROUND

This section summarizes popular tensor decomposition methods, sparse tensor formats, and related notations. The survey by Kolda and Bader [10] provides a more detailed discussion of tensor algorithms and their applications.

### 2.1 Tensor Notations

Tensors are $N$-dimensional arrays, where each element has an $N$-tuple index $\mathbf{i} = (i_1, i_2, \ldots, i_N)$. Each index coordinate $i_n$ locates a tensor element along the $n^{\text{th}}$ dimension or mode, with $n \in \{1, 2, \ldots, N\}$ and $i_n \in \{1, 2, \ldots, I_n\}$. Low-dimensional tensors include vectors, where $N = 1$, and matrices, where $N = 2$. In general, a *dense* $N$-dimensional (or a mode-$N$) tensor has $\prod_{n=1}^{N} I_n$ indexed elements. A tensor is said to be sparse if the majority of its elements are zero. The following notations are used in this paper:

1) Scalars are denoted by italicized lowercase letters (e.g. $a$).
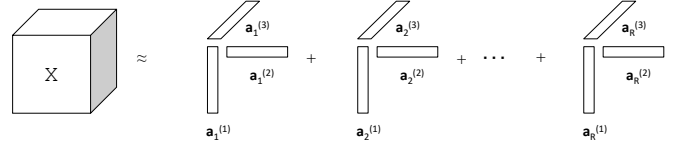2) Vectors are denoted by bold lowercase letters (e.g. $\mathbf{a}$).



Fig. 2: CPD of a mode-3 tensor $\mathcal{X}$. There are $R$ rank-one tensors that are formed by the outer-product between three vectors $\mathbf{a}_r^{(1)}$, $\mathbf{a}_r^{(2)}$, and $\mathbf{a}_r^{(3)}$, where $r \in \{1, 2, \ldots, R\}$. The vectors along the same mode are often grouped together as the columns of a *factor matrix*. For example, the vectors $\mathbf{a}_1^{(1)}, \mathbf{a}_2^{(1)}, \ldots, \mathbf{a}_R^{(1)}$ are the columns of the mode-1 factor matrix $\mathbf{A}^{(1)}$.

3) Matrices are denoted by bold capital letters (e.g. $\mathbf{A}$).
4) Higher-order tensors are denoted by Euler script letters (e.g. $\mathcal{X}$).
5) Fibers are analogous to matrix rows and columns. A mode-$n$ fiber of a tensor $\mathcal{X}$ is any vector formed by fixing all indices of $\mathcal{X}$, *except* the $n^{th}$ index. For example, a matrix column is a mode-1 fiber as it is defined by fixing the second index to a particular value.
6) To indicate every element along a particular mode or dimension, we will use the : symbol. For example, $\mathbf{A}(1,:)$ denotes the first row of the matrix $\mathbf{A}$.
7) Tensor matricization is the process by which a tensor is *unfolded* into a matrix. The mode-$n$ matricization of a tensor is denoted as $\mathbf{X}_{(n)}$, and is obtained by laying out the mode-$n$ fibers of $\mathcal{X}$ as the columns of $\mathbf{X}_{(n)}$.
8) Khatri-Rao product (KRP) [40] is the column-wise Kronecker product between two matrices, and is denoted by the symbol $\odot$. Given matrices $\mathbf{A}^{(1)} \in \mathbb{R}^{I_1 \times R}$ and $\mathbf{A}^{(2)} \in \mathbb{R}^{I_2 \times R}$, their Khatri-Rao product $\mathbf{K}$, denoted $\mathbf{K} = \mathbf{A}^{(1)} \odot \mathbf{A}^{(2)}$, where $\mathbf{K}$ is a $(I_1 \cdot I_2) \times R$ matrix, is defined as: $\mathbf{A}^{(1)} \odot \mathbf{A}^{(2)} = \begin{bmatrix} \mathbf{a}_1^{(1)} \otimes \mathbf{a}_1^{(2)} & \mathbf{a}_2^{(1)} \otimes \mathbf{a}_2^{(2)} & \ldots & \mathbf{a}_R^{(1)} \otimes \mathbf{a}_R^{(2)} \end{bmatrix}$, where $\otimes$ denotes the Kronecker product.
9) Element-wise product and division are denoted by the symbols $*$ and $\oslash$, respectively.

### 2.2 Tensor Decomposition

Tensor decomposition can be considered as a generalization of singular value matrix decomposition and principal component analysis, and it is used to reveal latent information embedded in large multi-dimensional data sets. This work targets algorithms that compute the CPD tensor decomposition model, namely, the Canonical Polyadic Alternating Least Squares (CP-ALS) algorithm for normally distributed data and the Canonical Polyadic Alternating Poisson Regression (CP-APR) algorithm for non-negative count data.

CPD is a popular tensor decomposition model, where a mode-$N$ tensor $\mathcal{X}$ is approximated by the sum of $R$ rank-one tensors. A rank-one tensor is formed by $N$ vectors, each corresponding to a particular mode. The vectors along the same mode can be arranged as the columns of a factor matrix, resulting in $N$ different *factor matrices* so that the decomposition of $\mathcal{X}$ is the outer product of these matrices. An example CPD of a mode-3 tensor is shown in Figure 2.

#### 2.2.1 CP-ALS

Algorithm 1 illustrates the CP-ALS algorithm for iteratively computing the factor matrices of the CPD model. In each CP-ALS iteration, every factor matrix is updated via the

**Algorithm 1** CP-ALS Algorithm

---

**Input:** Tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, initial guess factor matrices $\mathbf{A}^{(1)}, \cdots,$
    $\mathbf{A}^{(N)}$.
**Output:** $\boldsymbol{\lambda}, \mathbf{A}^{(1)}, \cdots, \mathbf{A}^{(N)}$
1: **repeat**
2:     **for** n = 1, $\cdots$, N **do**
3:         $\mathbf{G}^{(1)} \leftarrow \mathbf{A}^{(1)T} \mathbf{A}^{(1)}$
4:         $\cdots$
5:         $\mathbf{G}^{(n-1)} \leftarrow \mathbf{A}^{(n-1)T} \mathbf{A}^{(n-1)}$
6:         $\mathbf{G}^{(n+1)} \leftarrow \mathbf{A}^{(n+1)T} \mathbf{A}^{(n+1)}$
7:         $\cdots$
8:         $\mathbf{G}^{(N)} \leftarrow \mathbf{A}^{(N)T} \mathbf{A}^{(N)}$
9:         $\mathbf{V} \leftarrow \mathbf{G}^{(1)} * \cdots * \mathbf{G}^{(n-1)} \mathbf{G}^{(n+1)} * \cdots * \mathbf{G}^{(N)}$
10:        $\mathbf{K} \leftarrow \mathbf{A}^{(1)} \odot \cdots \odot \mathbf{A}^{(n-1)} \odot \mathbf{A}^{(n+1)} \odot \cdots \odot \mathbf{A}^{(N)}$
11:        $\mathbf{M} \leftarrow \mathbf{X}_{(n)} \mathbf{K}$                    ▷ MTTKRP
12:        $\mathbf{A}^{(n)} \leftarrow \mathbf{M} \mathbf{V}^{\dagger}$              ▷ Pseudoinverse
13:        $\boldsymbol{\lambda} \leftarrow$ column normalize $\mathbf{A}^{(n)}$ and store norms as $\boldsymbol{\lambda}$
14:    **end for**
15: **until** fit ceases to improve or maximum iterations reached

---

**Algorithm 2** CP-APR MU Algorithm

---

**Input:** Tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times \cdots \times I_N}$, initial guess factor matrices
    $\mathbf{A}^{(1)}, \cdots, \mathbf{A}^{(N)}$, and algorithmic parameters:
  • $k_{max}$, maximum number of outer iterations
  • $l_{max}$, maximum number of inner iterations
  • $\tau$, convergence tolerance on KKT conditions
  • $\kappa$, inadmissible zero avoidance adjustment
  • $\kappa_{tol}$, tolerance for potential inadmissible zero
  • $\epsilon$, minimum divisor to prevent divide-by-zero
**Output:** $\mathbf{A}^{(1)}, \cdots, \mathbf{A}^{(N)}$
1: **for** $k = 1, \ldots, k_{max}$ **do**
2:     isConverged $\leftarrow$ true
3:     **for** $n = 1, \ldots, N$ **do**
4:         $\mathbf{S}(i,r) \leftarrow \begin{cases} \kappa, & \text{if } k > 1, \mathbf{A}^{(n)}(i,r) < \kappa_{tol}, \text{and } \mathbf{\Phi}^{(n)}(i,r) > 1, \\ 0, & \text{otherwise} \end{cases}$
5:         $\mathbf{B} \leftarrow \left( \mathbf{A}^{(n)} + \mathbf{S} \right) \mathbf{\Lambda}$
6:         $\mathbf{\Pi} \leftarrow \left( \bigodot_{\forall m \neq n} \mathbf{A}^{(m)} \right)^T$          ▷ Khatri-Rao product (KRP)
7:         **for** $l = 1, \ldots, l_{max}$ **do**
8:             $\mathbf{\Phi}^{(n)} \leftarrow \left( \mathbf{X}_{(n)} \oslash (\max(\mathbf{B}\mathbf{\Pi}, \epsilon)) \right) \mathbf{\Pi}^T$
9:             **if** $|\min(\mathbf{B}, \mathbf{E} - \mathbf{\Phi}^{(n)})| < \tau$ **then**       ▷ Convergence check
10:                break
11:            **end if**
12:            isConverged $\leftarrow$ false
13:            $\mathbf{B} \leftarrow \mathbf{B} * \mathbf{\Phi}^{(n)}$
14:        **end for**
15:        $\boldsymbol{\lambda} \leftarrow e^T\mathbf{B}, \quad \mathbf{A}^{(n)} \leftarrow \mathbf{B}\mathbf{\Lambda}^{-1}$
16:    **end for**
17:    **if** isConverged = true **then**
18:        break
19:    **end if**
20: **end for**

---

alternating least squares (ALS) method whereby every other factor matrix but the one being updated is fixed to yield the best approximation of $\mathcal{X}$. Line 11 shows the matricized tensor times Khatri-Rao product (MTTKRP) operation [20], which involves *tensor matricization* and *Khatri-Rao product*. For a mode-3 tensor $\mathcal{X}$, the mode-1 MTTKRP operation can be expressed as $\mathbf{X}_{(1)} \left( \mathbf{A}^{(2)} \odot \mathbf{A}^{(3)} \right)$. MTTKRP is typically the most expensive tensor kernel of CP-ALS, and it is performed along all modes in every CP-ALS iteration. Since MTTKRP operations are similar across all modes, for brevity, we only discuss mode-1 MTTKRP in this paper.

### 2.2.2 CP-APR

While CP-ALS can decompose sparse count data, CP-APR better describes the random variations in the data by representing it using a Poisson distribution [14], which considers a discrete number of events and assumes zero probability for observing fewer than zero events. As a result, CP-APR is more expensive to compute compared to CP-ALS [14]. There are three common methods for computing CP-APR: (i) Multiplicative update (MU), (ii) Projected damped Newton for row-based sub-problems (PDN-R), and (iii) Projected quasi-Newton for row-based sub-problems (PQN-R).

PDN-R and PQN-R employ second-order information to independently solve *row sub-problems*, whereas MU uses a form of scaled steepest-descent with bound constraints over *all rows* during each iteration [41]. Although MU needs more iterations to converge than PDN-R and PQN-R, it remains the most popular method due to its lower iteration cost and higher parallelism and data reuse, which makes it amenable for efficient execution on modern parallel architectures. As such, we focus our efforts exclusively on optimizing the MU method in this study. For a detailed discussion of the three CP-APR algorithms, we refer our readers to [41] and [14].

Algorithm 2 shows CP-APR with the MU method. Using two nested loops, CP-APR computes the decomposition of a tensor by successively updating each factor matrix while holding the other factor matrices fixed. Lines 6 and 8 show the $\mathbf{\Pi}$ (Khatri-Rao product) and $\mathbf{\Phi}$ (model update) kernels, respectively, which make up majority of the CP-APR execution time. Note that $\mathbf{\Pi}$ is calculated once per *outer* loop *for each mode*, and $\mathbf{\Phi}$ is calculated once per *inner* loop. Since the inner loop is executed at most $l_{max}$ times *per mode per outer loop*, the cost of calculating $\mathbf{\Pi}$ can more easily be amortized as the number of inner iterations goes up.

## 2.3 Sparse Tensor Storage Formats

We present an overview of raw and compressed sparse tensor storage using three popular formats: coordinate (COO), hierarchical coordinate (HiCOO), and compressed sparse fiber (CSF). Figure 3 shows a comparison of the different formats for a $4 \times 8 \times 2$ tensor with six nonzero elements.

### 2.3.1 Coordinate (COO)

COO is the canonical and simplest sparse format, as it *lists* the nonzero elements and their $N$-dimensional coordinates, without any compression. This mode-agnostic form allows tensor algorithms to use one tensor copy across modes. Figure 3(a) shows an example tensor in the COO format.

Decomposing a sparse tensor stored in the COO format typically involves iterating over each nonzero element and updating the corresponding factor matrix row. For example, during mode-1 computation, a nonzero element with coordinates $(i_1, i_2, i_3)$ updates row $i_1$ of the mode-1 factor matrix after reading rows $i_2$ and $i_3$ from mode-2 and mode-3 factor matrices, respectively. Since multiple threads can simultaneously update the same row of a factor matrix, these updates must be done atomically, which can be expensive on parallel processors with a large number of threads.

### 2.3.2 Hierarchical Coordinate (HiCOO)

HiCOO [15], [27] is a block-based sparse tensor format that employs multi-dimensional tiling for data compression. Like COO, HiCOO is mode-agnostic but its compression efficacy depends completely on the properties of the target
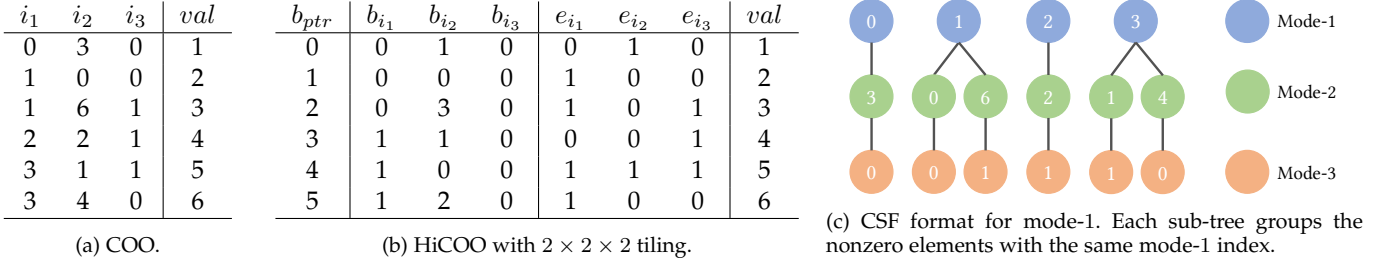
| $i_1$ | $i_2$ | $i_3$ | $val$ |
|---|---|---|---|
| 0 | 3 | 0 | 1 |
| 1 | 0 | 0 | 2 |
| 1 | 6 | 1 | 3 |
| 2 | 2 | 1 | 4 |
| 3 | 1 | 1 | 5 |
| 3 | 4 | 0 | 6 |

(a) COO.

| $b_{ptr}$ | $b_{i_1}$ | $b_{i_2}$ | $b_{i_3}$ | $e_{i_1}$ | $e_{i_2}$ | $e_{i_3}$ | $val$ |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 0 | 1 | 0 | 0 | 2 |
| 2 | 0 | 3 | 0 | 1 | 0 | 1 | 3 |
| 3 | 1 | 1 | 0 | 0 | 0 | 1 | 4 |
| 4 | 1 | 0 | 0 | 1 | 1 | 1 | 5 |
| 5 | 1 | 2 | 0 | 1 | 0 | 0 | 6 |

(b) HiCOO with $2 \times 2 \times 2$ tiling.



(c) CSF format for mode-1. Each sub-tree groups the nonzero elements with the same mode-1 index.

Fig. 3: Different sparse tensor representations of a $4 \times 8 \times 2$ tensors with six nonzero elements.

tensor, such as its shape, density, and data distribution, and determining the optimal parameters for compression (e.g., the tile size) is non-trivial. In some cases, rearranging the nonzero elements to create dense tiles is necessary to achieve any compression [27]. In addition, scheduling the resulting HiCOO blocks can suffer from limited parallelism, due to conflicting updates across blocks, as well as workload imbalance if some blocks have significantly more nonzero elements than others. Figure 3(b) shows the example sparse tensor encoded in the HiCOO format. The memory required to keep the hierarchical indices (i.e., $b_{i_1}$, $e_{i_1}$, etc.) can be lower than storing the actual indices (i.e., $i_1$, $i_2$, and $i_3$) only if each tile has a sufficient number of nonzero elements.

### 2.3.3 Compressed Sparse Fiber (CSF)

CSF stores a tensor as a collection of sub-trees, where each sub-tree represents a group of *all* nonzero elements that update the same factor matrix row. Given a CSF representation with a mode ordering of 1-2-3, where 1 is the root mode and 3 is the leaf mode, the root nodes represent the rows that will be updated, and the leaf nodes represent the nonzero elements that contribute to that update. Thus, iterating over the nonzero elements involves a bottom-up traversal of each sub-tree, such that at each non-leaf node, the partial results from its children are merged and pushed up, and this propagates until results from every node in the tree are merged at the root. Figure 3(c) illustrates the CSF sub-trees created from the example sparse tensor.

One advantage of CSF is that updates to factor matrices can be done asynchronously by assigning one thread to each sub-tree. However, CSF requires $N$ tensor copies to maintain synchronization-free updates across every mode, which can be impractical for large tensors and/or devices with limited memory capacity (e.g., GPUs). Alternatively, the sub-trees can be traversed both bottom-up and top-down, merging partial results at the tree level corresponding to the destination mode. While this approach allows a single tensor copy (with the root mode chosen arbitrarily) to be used across all modes, it requires synchronization to avoid update conflicts and entirely different tree traversal algorithms [25]. Additionally, regardless of the strategy used, CSF suffers from workload imbalance, as some sub-trees can have significantly more nonzero elements than the others.

## 3 ALTO FORMAT

To tackle the highly irregular shapes and distributions of real-world sparse data, our ALTO format maps the coordinates of a nonzero element in the $N$-dimensional space that represents a tensor to a mode-agnostic index in a compact *linear space*. Specifically, ALTO uses a data-aware recursive encoding to partition every mode of the original Cartesian space into multiple regions such that each distinct mode has a variable number of regions to adapt to the unequal cardinalities of different modes and to minimize the storage requirements. This adaptive linearization and recursive partitioning of the multi-dimensional space ensures that neighboring points in space are close to each other on the resulting compact line, thereby maintaining the inherent data locality of tensor algorithms. Moreover, the ALTO format is not only locality-friendly, but also parallelism-friendly as it allows partitioning of the multi-dimensional space into perfectly balanced (in terms of workload) subspaces. Further, it intelligently arranges the modes in the derived subspaces based on their cardinality (dimension length) to further reduce the overhead of resolving the update conflicts that typically occur in parallel sparse tensor computations.

What follows is a detailed description and discussion of the ALTO format generation (§3.1) using a walk-through example. Additionally, we present the ALTO-based sequential algorithm for the MTTKRP operation (§3.2).

### 3.1 ALTO Tensor Generation

Formally, an ALTO tensor $\mathcal{X} = \{x_1, x_2, \ldots, x_M\}$ is an ordered set of nonzero elements, where each element $x_i = \langle v_i, p_i \rangle$ is represented by a value $v_i$ and a position $p_i$. The position $p_i$ corresponds to a compact mode-agnostic encoding of the indexing metadata, which is used to quickly generate the tuple $(i_1, i_2, \ldots, i_N)$ that locates a nonzero element in the multi-dimensional Cartesian space.

The generation of an ALTO tensor is carried out in two stages: linearization and ordering. First, ALTO constructs the indexing metadata using a compressed encoding scheme, based on the cardinalities of tensor modes, to map each nonzero element to a position on a compact line. Second, it arranges the nonzero elements in linearized storage according to their line positions, i.e., the values of their ALTO index. Typically, the ordering stage dominates the format generation time. However, compared to the other compressed sparse tensor formats [15], [20], [25], [27], [28], [29], ALTO requires a minimal generation time because ordering the linearized tensors incurs a fraction of the cost required to sort multi-dimensional index sets (due to the reduction in comparison operations, as detailed in §5).

Figure 4 shows the ALTO format for the sparse tensor from Figure 3(a). The multi-dimensional indices ($i_1$, $i_2$, and $i_3$) are color coded and the $r^{th}$ bit of their binary representation is denoted $b_{i_n,r}$. Specifically, ALTO keeps the value of a nonzero element along with a linearized index, where each bit of this index is selected to partition the multi-dimensional space into two hyperplanes. For example, the
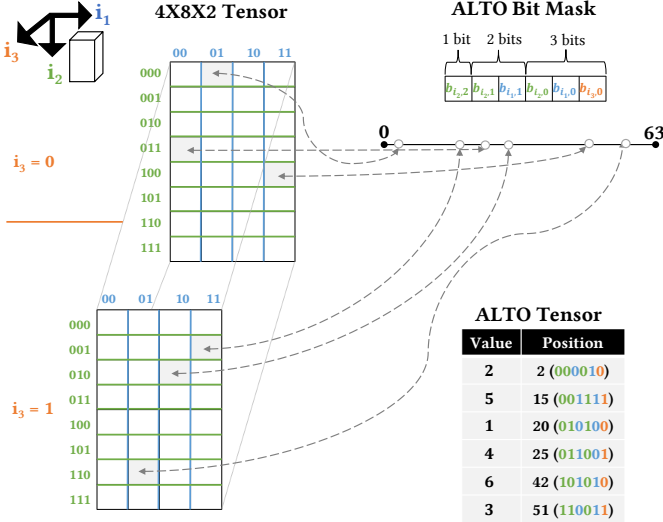
Fig. 4: An example of the ALTO sparse encoding and representation for the three-dimensional tensor in Figure 3(a).

ALTO encoding in Figure 4 uses a compact line of length 64 (i.e., a 6-bit linearized index) to represent the target tensor of size $4 \times 8 \times 2$. This index consists of three groups of bits with variable sizes (resolutions) to efficiently handle high-order data of arbitrary dimensionality. Within each bit group, ALTO arranges the modes in an increasing order of their length (i.e., the shortest mode first), which is equivalent to partitioning the multi-dimensional space along the longest mode first. Such an encoding aims to generate a balanced linearization of irregular Cartesian spaces, where the position resolution of a nonzero element increases with every consecutive bit, starting from the most significant bit. Therefore, the line segments encode subspaces with mode intervals of equivalent lengths, e.g., the line segments $[0-31]$, $[0-15]$, and $[0-7]$ encode subspaces of $4 \times 4 \times 2$, $4 \times 2 \times 2$, and $2 \times 2 \times 2$ dimensions, respectively.

Due to this adaptive encoding, ALTO represents the resulting linearized index using the minimum number of bits, and it improves data locality across all modes of a given sparse tensor. Hence, a mode-$N$ tensor, whose dimensions are $I_1 \times I_2 \times \cdots \times I_N$, can be efficiently represented using a single ALTO format with indexing metadata of size:

$$S_{\text{ALTO}} = M \times \left( \sum_{n=1}^{N} \log_2 I_n \right) \text{ bits,} \qquad (1)$$

where $M$ is the number of nonzero elements.

As a result, compared to the de facto COO format, ALTO reduces the storage requirement regardless of the tensor's characteristics. That is, the metadata compression ratio of the ALTO format relative to COO is always $\geq 1$. On a hardware architecture with a word-level memory addressing mode, this compression ratio is given by:

$$\frac{S_{\text{COO}}}{S_{\text{ALTO}}} = \frac{\sum_{n=1}^{N} \left\lceil \frac{\log_2 I_n}{W_b} \right\rceil}{\left\lceil \frac{\sum_{n=1}^{N} \log_2 I_n}{W_b} \right\rceil}, \qquad (2)$$

where $W_b$ is the word size in bits. For example, on an architecture with byte-level addressing (i.e., $W_b = 8$ bits), the sparse tensor in Figure 4 needs three bytes to store the mode indices (coordinates) for each nonzero element in
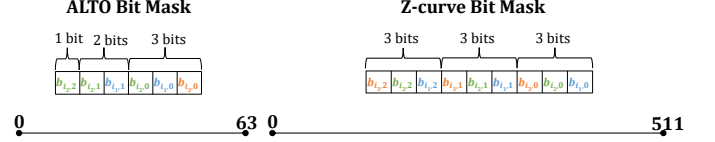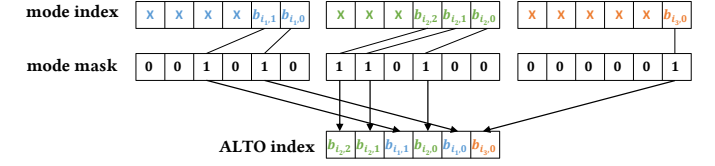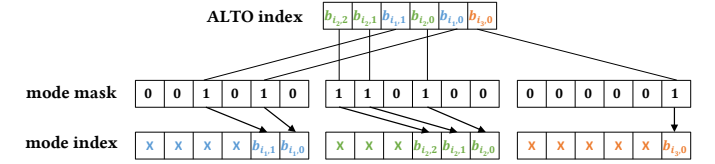


Fig. 5: For the example in Figure 4, ALTO generates a non-fractal, yet more compact encoding compared to traditional space-filling curves, such as the Z-Morton order.



(a) ALTO generates its linearized index using bit-level gather operations.



(b) To recover the multi-dimensional indices, ALTO decodes the linearized indexing metadata using bit-level scatter operations.

Fig. 6: The ALTO-based bit encoding and decoding mechanisms for the example in Figure 4.

the COO format, whereas only a single byte is required to store the linearized index in the ALTO format: the metadata compression ratio of ALTO compared to COO is *three*.

Moreover, the ALTO format not only reduces the memory traffic of sparse tensor computations, but also decreases the number of memory transactions required to access the indexing metadata of a sparse tensor, as reading the linearized index requires fewer accesses compared to reading several multi-dimensional indices. In addition, this natural coalescing of the multi-dimensional indices into a single linearized index increases the memory transaction size to make more efficient use of the main memory bandwidth.

It is important to note that ALTO uses a non-fractal[1] encoding scheme, unlike the traditional space-filling curves (SFCs) [42]. In contrast, SFCs (e.g., Z-Morton order [43]) are based on continuous self-similar (or fractal) functions that target dense data, which can be extremely inefficient when used to encode the irregularly shaped multi-dimensional spaces that emerge in higher-order sparse tensor algorithms as they require indexing metadata of size:

$$S_{\text{SFC}} = M \times \left( N \times \max_{n=1}^{N} (\log_2 I_n) \right) \text{ bits.} \qquad (3)$$

Therefore, in sparse tensor computations, SFCs have been *only used to reorder the nonzero elements to improve their data locality* rather than compressing the indexing metadata [15]. Figure 5 shows the compact encoding generated by ALTO compared to the fractal encoding scheme of the Z-Morton curve. In this example, ALTO reduces the length of the encoding line by a factor of eight, which not only decreases the overall size of the indexing metadata, but also reduces the linearization/de-linearization time required to

---

1. A fractal pattern is a hierarchically self-similar pattern that looks the same at increasingly smaller scales.

map the multi-dimensional space to/from the encoding line.

To allow fast indexing of linearized tensors during sparse tensor computations, the ALTO encoding is implemented using a set of simple $N$ bit masks, where $N$ is the number of modes, on top of common data processing primitives. Figure 6 shows the linearization and de-linearization mechanisms, which are used during the ALTO format generation and the sparse tensor computations, respectively. The linearization process is implemented as a bit-level gather, while the de-linearization is performed as a bit-level scatter. Thus, although the compressed representation of the proposed ALTO format comes at the cost of a de-linearization (decompression) overhead, such a computational overhead is minimal and can be effectively overlapped with the memory accesses of the memory-intensive sparse tensor operations, as shown in §5.

### 3.2 ALTO-based Tensor Operations

Since high-dimensional data analytics is becoming increasingly popular in rapidly evolving areas [1], [3], [6], [8], a fundamental goal of the proposed ALTO format is to deliver superior performance without compromising the productivity of end users to allow fast development of tensor algorithms. To this end, Algorithm 3 illustrates the popular MTTKRP tensor operation using the ALTO format.

First, unlike mode-specific (e.g., CSF-based) formats, ALTO enables end users to perform tensor operations using *a unified code implementation that works on a single tensor copy* regardless of the different mode orientations of such operations. Second, by decoupling the representation of a sparse tensor from the distribution of its nonzero elements, ALTO does not require *manual* tuning to select the optimal *format parameters* for this tensor, in contrast to prior approaches such as HiCOO and CSF. Instead, the ALTO format is automatically generated based on the shape and dimensions of the target sparse tensor (as explained in §3.1).

---

**Algorithm 3** Mode-1 sequential MTTKRP-ALTO algorithm.

---

**Input:** A third-order ALTO sparse tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ with $M$ nonzero elements, dense factor matrices $\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}$
**Output:** Updated dense factor matrix $\tilde{\mathbf{A}} \in \mathbb{R}^{I_1 \times R}$
1: **for** $x = 1, \ldots, M$ **do**
2:    $\mathbf{i} = \text{EXTRACT}(pos(x), MASK)$     ▷ De-linearization.
3:    $\tilde{\mathbf{A}}(i_1,:) += val(x) \times \mathbf{A}^{(2)}(i_2,:) \times \mathbf{A}^{(3)}(i_3,:)$
4: **end for**
5: **return** $\tilde{\mathbf{A}}$

---

## 4 PARALLEL LINEARIZED TENSOR ALGORITHMS

We devise a set of ALTO-based parallel algorithms for accelerating sparse tensor computations and demonstrate how they are employed in popular TD operations.

### 4.1 Workload Partitioning and Scheduling

Prior compressed sparse tensor formats partition the tensor space into *non-overlapping* regions and cluster the data into coarse-grained structures, such as blocks, slices, and/or fibers. Due to the irregular shapes and distributions of higher-order data, such coarse-grained approaches can suffer from severe workload imbalance and limited parallel performance/scalability. In contrast, by employing the
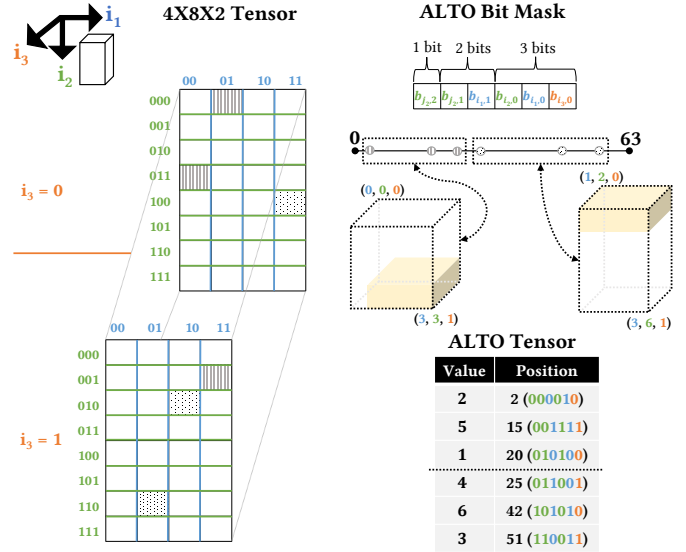


Fig. 7: ALTO partitioning of the example in Figure 4, which generates balanced partitions in terms of workload (nonzero elements) for efficient parallel execution.

ALTO format, linearized tensor algorithms work at the finest granularity level (i.e., a single nonzero element), which allows perfect load-balancing and scalable parallel execution. While a non-overlapping partitioning can be obtained from the ALTO encoding scheme by using a subset of the index bits, the workload balance of such a partitioning still depends on the sparsity patterns of the tensor.

To decouple the performance of sparse tensor computations from the distribution of nonzero elements, we divide the multi-dimensional space into potentially *overlapping* regions and allow workload distribution at the granularity of nonzero elements, which result in perfectly balanced partitions in terms of workload. Figure 7 depicts an example of ALTO's workload decomposition when applied to the sparse tensor in Figure 4. ALTO divides the line segment containing the nonzero elements of the target tensor into two smaller line segments: $[2 - 20]$ and $[25 - 51]$, which have different lengths (i.e., 18 and 26) but the same number of nonzero elements, thus perfectly splitting the workload.

Once the linearized sparse tensor is divided into multiple line segments, ALTO identifies the basis mode intervals (coordinate ranges) of the multi-dimensional subspaces that correspond to these segments. For example, the line segments $[2-20]$ and $[25-51]$ correspond to three-dimensional subspaces bounded by the mode intervals $\{[0-3], [0-3], [0-1]\}$ and $\{[1-3], [2-6], [0-1]\}$, respectively. While the derived multi-dimensional subspaces of the line segments may overlap, as highlighted in yellow in Figure 7, each nonzero element is assigned to exactly one line segment. That is, ALTO imposes a partitioning on a given linearized tensor that generates a disjoint set of *non-overlapping and balanced line segments*, but it does not guarantee that such a partitioning will divide the multi-dimensional space of the tensor into non-overlapping subspaces. In contrast, the prior compressed formats split the multi-dimensional space into non-overlapping (yet highly imbalanced) regions, namely, tensor slices and fibers in CSF-based formats and multi-dimensional blocks in block-based formats (e.g., HiCOO).

**Algorithm 4** Parallel mode-1 MTTKRP-ALTO algorithm. ALTO automatically uses either recursive or output-oriented tensor traversal, based on the reuse of output fibers, to efficiently resolve update conflicts.

---

**Input:** A third-order ALTO sparse tensor $\mathcal{X} \in \mathbb{R}^{I_1 \times I_2 \times I_3}$ with $M$ nonzero elements, dense factor matrices $\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \mathbf{A}^{(3)}$
**Output:** Updated dense factor matrix $\tilde{\mathbf{A}} \in \mathbb{R}^{I_1 \times R}$

1: **for** $l = 1, \ldots, L$ **in parallel do** ▷ ALTO line segments.
2:    **for** $\forall x \in \mathcal{X}_l$ **do**
3:        $p = pos(x) \;\; v = val(x)$
4:        $p = pos\_out(x) \;\; v = val\_out(x)$
5:        $\mathbf{i} = \text{EXTRACT}(p, MASK(1))$ ▷ De-linearization.
6:        $\mathbf{Temp}_l(i_1 - T_{l,1}^s, :) \mathrel{+}= v \times \mathbf{A}^{(2)}(i_2, :) \times \mathbf{A}^{(3)}(i_3, :)$
7:        **if** $i_1$ *is* boundary = true **then**
8:           $\text{ATOMIC}(\tilde{\mathbf{A}}(i_1, :) \mathrel{+}= v \times \mathbf{A}^{(2)}(i_2, :) \times \mathbf{A}^{(3)}(i_3, :))$
9:        **else**
10:          $\tilde{\mathbf{A}}(i_1, :) \mathrel{+}= v \times \mathbf{A}^{(2)}(i_2, :) \times \mathbf{A}^{(3)}(i_3, :)$
11:        **end if**
12:    **end for**
13: **end for**
14: **for** $b = 1, \ldots, I_1$ **in parallel do** ▷ Pull-based reduction.
15:    **for** $\forall l$ where $b \in [T_{l,1}^s, T_{l,1}^e]$ **do**
16:        $\tilde{\mathbf{A}}(b, :) \mathrel{+}= \mathbf{Temp}_l(b - T_{l,1}^s, :)$
17:    **end for**
18: **end for**
19: **return** $\tilde{\mathbf{A}}$

---

Formally, a set of $L$ line segments partitions an ALTO tensor $\mathcal{X}$, which encodes $N$-dimensional sparse data, such that $\mathcal{X} = \mathcal{X}_1 \cup \mathcal{X}_2 \cdots \cup \mathcal{X}_L$ and $\mathcal{X}_i \cap \mathcal{X}_j = \phi \forall i$ and $j$, where $i \neq j$. Each line segment $\mathcal{X}_l$ is an ordered set of nonzero elements that are bounded in an $N$-dimensional space by a set of $N$ closed mode intervals $T_l = \{[T_{l,1}^s, T_{l,1}^e], [T_{l,2}^s, T_{l,2}^e], \cdots [T_{l,N}^s, T_{l,N}^e]\}$, where each mode interval $T_{l,n}$ is delineated by a start $T_{l,n}^s$ and an end $T_{l,n}^e$. The intersection of two mode interval sets represents the subspace overlap between their corresponding line segments (partitions), as highlighted in yellow in Figure 7. This overlap information is used to more efficiently resolve conflicts between partitions, as described in §4.2.

## 4.2 Adaptive Conflict Resolution

Because processing the nonzero elements of a tensor in parallel (e.g., line 1 in Algorithm 3) can result in write conflicts across threads (e.g., line 3 in Algorithm 3), we devise an adaptive parallel algorithm to handle these conflicts by exploiting the inherent data reuse of target tensors. Our adaptive algorithm chooses between 1) recursively traversing the tensor to maximize the reuse of both the input and output fibers, at the cost of global parallel reduction, or 2) traversing the tensor elements in output-oriented ordering, and only synchronizing at partition boundaries.

Algorithm 4 describes the proposed parallel execution scheme using a representative MTTKRP operation that works on a sparse tensor stored in the ALTO format. After ALTO imposes a partitioning on a sparse tensor, as detailed in §4.1, each partition can be assigned to a different thread. To resolve the update/write conflicts that may happen during parallel sparse tensor computations, ALTO uses an *adaptive conflict resolution* approach that *automatically* selects the appropriate tensor traversal and synchronization technique (highlighted by the different gray backgrounds) based
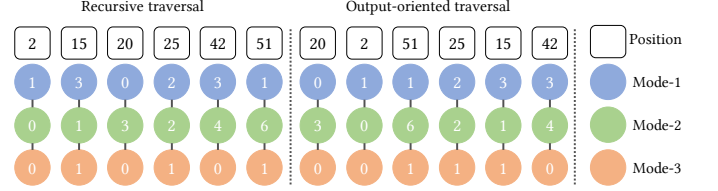


Fig. 8: Recursive vs. output-oriented traversal of the example tensor in Figure 4, where mode-1 is the target/output mode. The coordinates of each nonzero (mode indices) are extracted from its line position (linearized index) as detailed in Figure 6.

on the reuse of the target fibers. This metric represents the average number of nonzero elements per fiber (i.e., the generalization of a matrix row or column) and it is computed in constant time by simply dividing the total number of nonzero elements by the number of fibers along the target mode.

When a sparse tensor operation exhibits high fiber reuse, ALTO recursively traverses the tensor (by accessing the nonzero elements in an increasing order of their linearized index or line position as illustrated in Figure 8) and it uses a limited amount of temporary (local) storage to capture the local updates of different partitions (line 6). Next, ALTO merges the conflicting global updates (lines 14–18) using an efficient pull-based parallel reduction, where the final results are computed by pulling the partial results from the ALTO partitions. When computing the factor matrix of the target mode (e.g., mode-1), such a recursive tensor traversal 1) increases the likelihood that both input (mode-2/3) and output (mode-1) fibers remain in fast memories, and 2) reduces the size of temporary storage (partial copy of mode-1 factor matrix) needed for each partition, which in turn decreases the overhead of the pull-based reduction.

ALTO considers the fiber reuse large enough to use local staging memory for conflict resolution, if the average number of nonzero elements per fiber is *more than* the maximum cost of using this two-stage (buffered) accumulation process, which consists of initialization (omitted for brevity), local accumulation (line 6 in Algorithm 4), and global accumulation (lines 14–18). In the worst (no reuse) case, the buffered accumulation cost is four memory operations (two read and two write operations). As explained in §4.1, each line segment $\mathcal{X}_l$ is bounded in an $N$-dimensional space by a set of $N$ closed mode intervals $T_l$, which is computed during the partitioning of an ALTO tensor; thus, the size of the temporary storage accessed during the accumulation of $\mathcal{X}_l$'s updates along a mode $n$ is directly determined by the mode interval $[T_{l,n}^s, T_{l,n}^e]$ (see lines 6 and 15).

When the target tensor computations suffer from limited fiber reuse, ALTO uses output-oriented tensor traversal, where the nonzero elements are accessed in an increasing order of their target/output mode (e.g. mode-1 as depicted in Figure 8). That way, the data reuse of output fibers is fully captured and synchronization across threads can be avoided. Specifically, ALTO needs to resolve the conflicting updates across its line segments (partitions) using direct atomic operations (line 8) only if the output fiber is at the boundary between different partitions/threads. This output-oriented traversal resembles the CSF-based tree traversal (see Figure 3(c)); however, in contrast to CSF, ALTO

**Algorithm 5** Parallel mode-1 $\Phi$-ALTO kernel. ALTO performs *either* recursive or output-oriented tensor traversal, based on fiber reuse, to efficiently resolve update conflicts. In addition, it determines whether to reuse or recompute intermediate results.

```
 1: for  l = 1, . . . , L in parallel do          ▷ ALTO line segments
 2:     for ∀x ∈ X_l do
 3:         p = pos(x)  v = val(x)
 4:         p = pos_out(x)  v = val_out(x)
 5:         i = EXTRACT(p, MASK)              ▷ Delinearization
 6:         if preCompute Π = true then
 7:             krp ← Π(x, :)
 8:         else
 9:             krp ← (∗_{∀m≠n} A^{(m)}(i_m, :))
10:         end if
11:         Temp_l(i_n − T^s_{l,n}, :) += ( val(x) / max(B(i_n,:)krp^T, ε) ) krp
12:         if i_n is boundary = true then
13:             ATOMIC ( Φ^{(n)}(i_n, :) += ( v / max(B(i_n,:)krp^T, ε) ) krp )
14:         else
15:             Φ^{(n)}(i_n, :) += ( v / max(B(i_n,:)krp^T, ε) ) krp
16:         end if
17:     end for
18:     end for
19: for  b = 1, . . . , I_n in parallel do          ▷ Pull-based reduction
20:     for ∀l where b ∈ [T^s_{l,n}, T^e_{l,n}] do
21:         Φ^{(n)}(b, :) += Temp_l(b − T^s_{l,n}, :)
22:     end for
23: end for
24: return Φ^{(n)}
```

uses a fine-grained compact index (line position) to encode nonzero coordinates instead of a coarse-grained index tree, which requires a single tensor copy (instead of one copy per mode) and allows perfect load balancing during parallel execution. In addition, our output-oriented tensor traversal is only used when fiber reuse is limited; otherwise, the recursive traversal method is employed because of its superior data locality and parallel performance, thanks to reusing both input and output fibers as well as amortizing the overhead of synchronization operations (pull-based reduction).

### 4.3 Adaptive Memory Management

In many tensor decomposition algorithms, such as CP-APR, the intermediate results of tensor kernels are typically stored and then reused during the iterative optimization loop (as shown in Algorithm 2). However, storing these important calculations can substantially increase memory traffic and require prohibitive amount of memory, especially for large tensors and high decomposition ranks. Hence, in contrast to the traditional algorithm that pre-computes and reuses the intermediate values (ALTO-PRE), we introduce an ALTO-based algorithm variant that recomputes these values on-the-fly (ALTO-OTF). Moreover, we propose a heuristic to dynamically decide whether to reuse or recompute the intermediate results of tensor kernels based on the characteristics of the target data sets and tensor computations. It is important to note that such pre-computation is different from prior memoization approaches [37], [38], which use a non-trivial decision making process to reduce computations by *reformulating* tensor operations and reusing intermediate results *across modes*. In contrast, ALTO-PRE uses easily calcu-

lable metrics to decide whether or not to reuse intermediate results *within a mode*, and it performs the same tensor operations as ALTO-OTF.

To demonstrate our adaptive memory management technique, Algorithm 5 shows how the model update ($\Phi$) kernel in CP-APR (Line 8 from Algorithm 2) is parallelized using the ALTO format. The $\Pi$ matrix from Line 6 in Algorithm 2 calculates a dense matrix for a given mode $n$ that is the Khatri-Rao product (KRP) between all factor matrices *excluding* the mode-$n$ factor matrix. However, not every row of $\Pi$ is required for a sparse tensor but only rows that correspond to nonzero elements are necessary and actually calculated, leading to a $\Pi \in \mathbb{R}^{M \times R}$ matrix, where $M$ and $R$ are the number of nonzero elements and the decomposition rank, respectively. In Algorithm 5, if we select to use precomputed $\Pi$, the kernel reads in the $\Pi$ matrix row corresponding to the nonzero element $x$ from memory (Line 7); otherwise, it computes the required KRP from the factor matrices (Line 9) using the delineariezd coordinates. Precomputing the $\Pi$ matrix is simple; in line 6 of Algorithm 2 each nonzero element can calculate its Khatri-Rao product in parallel, using the equation from Line 9 in Algorithm 5.

Next, for each nonzero, the corresponding KRP row is used to update the $\Phi$ matrix, and the conflicting updates are resolved using our adaptive conflict resolution as detailed in §4.2. Specifically, if there is high fiber reuse, recursive tensor traversal is conducted and the update is made to the temporary scratchpad memory **Temp** (Line 11), which is later reduced (Lines 19 to 23) to decrease memory contention; otherwise, output-oriented traversal is used to avoid synchronization and atomic operations are utilized to update the $\Phi$ matrix (Line 13) only at the boundaries between different ALTO partitions/threads.

ALTO employs a simple heuristic for determining which algorithm variant (ALTO-PRE or ALTO-OTF) to use based on the fast memory size of hardware architectures as well as the fiber reuse and size of factor matrices of sparse tensors. Similar to our conflict resolution heuristic (illustrated in §4.2), we use low fiber reuse to infer that on-the-fly computation of KRP is expensive, due to the cost of fetching data from memory. Hence, ALTO decides to use pre-computation (ALTO-PRE) when sparse tensors suffer from low fiber reuse and the size of their factor matrices is substantially larger than the fast memory size. Otherwise, the on-the-fly algorithm variant (ALTO-OTF) is used because of its superior data locality and lower memory consumption.

## 5 EVALUATION

We evaluate ALTO-based tensor algorithms against the state-of-the-art sparse tensor libraries and representations in terms of parallel performance, tensor storage, and format generation cost. We conduct a thorough study of key tensor decomposition operations (§2.2) and demonstrate the performance characteristics across the third and fourth generation of Intel Xeon Scalable processors, codenamed Ice Lake (ICX) and Sapphire Rapids (SPR), respectively.

### 5.1 Experimental Setup

#### 5.1.1 Platform

The experiments were conducted on an Intel Xeon Platinum 8360Y CPU with Ice Lake (ICX) micro-architecture, and

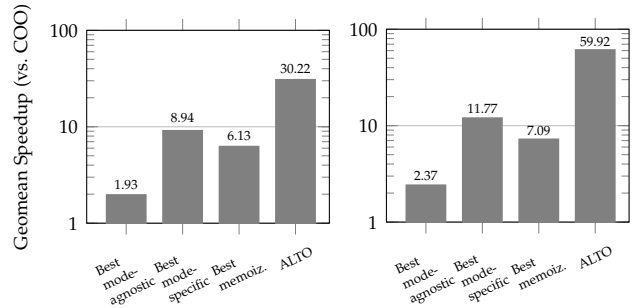TABLE 1: Characteristics of the target sparse tensor data sets. Non-negative count tensors are underlined.

| Tensor | Dimensions | #NNZs | Density | Fib. reuse |
|---|---|---|---|---|
| LBNL | $1.6K \times 4.2K \times 1.6K \times 4.2K \times 868.1K$ | $1.7M$ | $4.2 \times 10^{-14}$ | Limited |
| NIPS | $2.5K \times 2.9K \times 14K \times 17$ | $3.1M$ | $1.8 \times 10^{-06}$ | High |
| UBER | $183 \times 24 \times 1.1K \times 1.7K$ | $3.3M$ | $3.8 \times 10^{-04}$ | High |
| CHICAGO | $6.2K \times 24 \times 77 \times 32$ | $5.3M$ | $1.5 \times 10^{-02}$ | High |
| VAST | $165.4K \times 11.4K \times 2 \times 100 \times 89$ | $26M$ | $7.8 \times 10^{-07}$ | High |
| DARPA | $22.5K \times 22.5K \times 23.8M$ | $28.4M$ | $2.4 \times 10^{-09}$ | Limited |
| ENRON | $6K \times 5.7K \times 244.3K \times 1.2K$ | $54.2M$ | $5.5 \times 10^{-09}$ | High |
| LANL-2 | $3.8K \times 11.2K \times 8.7K \times 75.2K \times 9$ | $69.1M$ | $1.9 \times 10^{-10}$ | High |
| NELL-2 | $12.1K \times 9.2K \times 28.8K$ | $76.9M$ | $2.4 \times 10^{-05}$ | High |
| FB-M | $23.3M \times 23.3M \times 166$ | $99.6M$ | $1.1 \times 10^{-09}$ | Limited |
| FLICKR | $319.7K \times 28.2M \times 1.6M \times 731$ | $112.9M$ | $1.1 \times 10^{-14}$ | Limited |
| DELI | $532.9K \times 17.3M \times 2.5M \times 1.4K$ | $140.1M$ | $4.3 \times 10^{-15}$ | Medium |
| NELL-1 | $2.9M \times 2.1M \times 25.5M$ | $143.6M$ | $9.1 \times 10^{-13}$ | Medium |
| AMAZON | $4.8M \times 1.8M \times 1.8M$ | $1.7B$ | $1.1 \times 10^{-10}$ | High |
| PATENTS | $46 \times 239.2K \times 239.2K$ | $3.6B$ | $1.4 \times 10^{-03}$ | High |
| REDDIT | $8.2M \times 177K \times 8.1M$ | $4.7B$ | $4.0 \times 10^{-10}$ | High |



(a) Performance on a 72-core ICX. (b) Performance on a 104-core SPR.

Fig. 9: The performance of MTTKRP using ALTO in comparison with an oracle selecting the best state-of-the-art variant for each implementation category for all tensors used in this paper.

an Intel Xeon Platinum 8470 CPU with Sapphire Rapids (SPR) micro-architecture. The ICX system has 256 GiB main memory and it consists of two sockets, each with a 54 MiB L3 cache and 36 physical cores running at a fixed frequency of 2.4 GHz for accurate measurements. The SPR system also comprises two sockets with 52 physical cores each, and its frequency was fixed to 2.4 GHz. All cores within a socket share a 105 MiB L3 cache and the overall main memory in the SPR server is 1 TB. The experiments use all hardware threads (72 and 104, respectively) on the target platforms. Both servers run AlmaLinux 8.8 Linux distribution and they are configured to enable "Transparent Huge Pages" and to support two and four NUMA domains per socket on ICX and SPR, respectively.

The code is built using Intel C/C++ compiler (v2021.6.0) with the optimization flags `-O3 -qopt-zmm-usage=high -xHost` to fully utilize vector units. For performance counter measurements and thread pinning, we use the LIKWID tool suite v5.3 [44].

### 5.1.2 Datasets

The experiments consider a gamut of real-world tensor data sets with various characteristics. These tensors are often used in related works and they are publicly available in the FROSTT [9] and HaTen2 [45] repositories. Table 1 shows the detailed features of the target tensors, ordered by size, in terms of dimensions, number of nonzero elements (#NNZs), and density. Additionally, the tensors are classified based on the average reuse of their fibers into high, medium, or limited reuse. We consider a given mode to have high reuse, if its fibers are reused more than eight times on average; when the fibers are reused between five to eight times, they have medium reuse; otherwise, the fibers suffer from limited reuse. Since TD operations access fibers along all modes, a tensor with at least one mode of limited/medium reuse is considered to have an overall limited/medium fiber reuse. In the evaluation, we use all tensors and non-negative count tensors for CP-ALS and CP-APR experiments, respectively.

### 5.1.3 Configurations

We evaluate the proposed ALTO format[2], compared to the mode-agnostic COO and HiCOO formats [15], [22] as well as the mode-specific CSF formats [20], [25], [38]. Specifically, we use the latest code of the state-of-the-art sparse tensor libraries for CPUs, namely, ParTI[3], SPLATT[4], and STeF[5] for normally distributed data and SparTen[6] for non-negative count data. On the ICX and SPR systems, we evaluate the target data sets that can fit in memory for all tensor libraries. We report the best-achieved performance across the different configurations of the COO format; that is, with or without thread privatization (which keeps local copies of the output factor matrix). For the HiCOO format, its performance and storage are highly sensitive to the block and superblock (SB) sizes, which benefit from tuning. Since the current HiCOO implementation does not auto-tune these parameters, we use a block size of 128 ($2^7$) and two superblock sizes of $2^{10}$ ("HiCOO-SB10") and $2^{14}$ ("HiCOO-SB14") according to prior work [16]. We evaluate two variants of the mode-specific formats: CSF and CSF with tensor tilling ("CSF-tile"), both of which use $N$ representations for an order-$N$ sparse tensor to achieve the best performance. For STeF, we use its data movement model to decide which results to memoize and reuse between modes based on the cache size. We report the best performance of STeF across the different cache configurations, which was achieved when setting the cache size to the size of L3 cache on the target ICX and SPR CPUs.

Similar to previous studies [17], [20], [25], the experiments use double-precision arithmetic and 64-bit (native word) integers. To compute the CP-APR model for non-negative count data, we use 32-bit integers to represent the input tensor values. While the target data sets require a linearized index of size between 32 and 80 bits, we configured ALTO to select the size of its linearized index to be multiples of the native word size (i.e., 64 and 128 bits) for simplicity. We use a decomposition rank $R = 16$ for all experiments and set the maximum number of inner iterations ($l_{max}$) in CP-APR to 10, as per prior work [14].

2. Available at: https://github.com/IntelLabs/ALTO
3. Available at: https://github.com/hpcgarage/ParTI
4. Available at: https://github.com/ShadenSmith/splatt
5. Available at: https://github.com/HPCRL/STeF
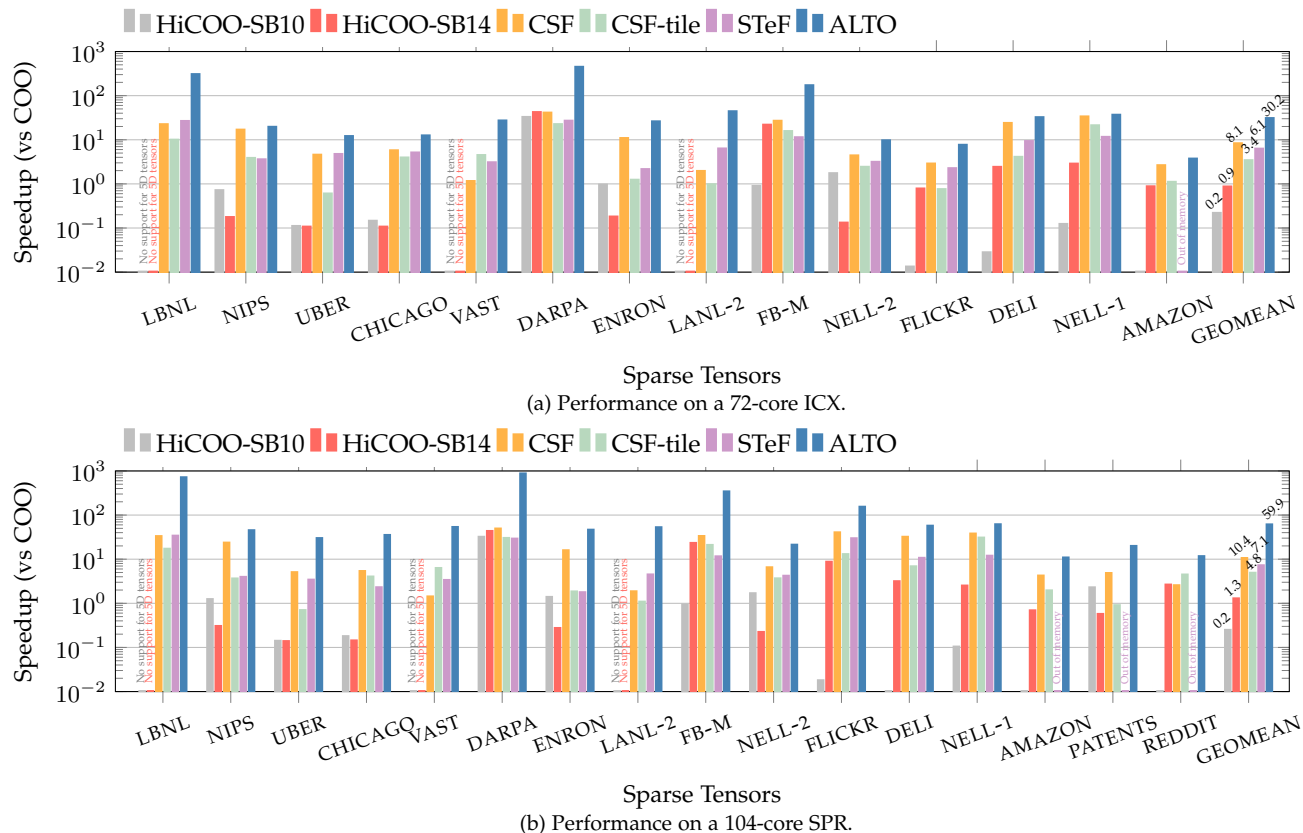6. Available at: https://github.com/sandialabs/sparten

Fig. 10: The overall parallel speedup of MTTKRP (all modes) using the different sparse tensor implementations compared to COO. The sparse tensors are sorted in increasing order of their size (number of nonzero elements).

## 5.2 CP-ALS Performance

We compare our ALTO-based CP-ALS algorithm to a variety of CP-ALS implementations in the state-of-the-art libraries SPLATT, ParTI, and STeF for each tensor dataset. The implementations can be grouped into three categories: 1) mode-agnostic or general formats (COO, HiCOO-SB10, and HiCOO-SB14), which use one tensor copy, 2) mode-specific formats (CSF and CSF-tile), which keep multiple tensor copies (one per mode) for best performance, and 3) memoization techniques (STeF), which retain intermediate results across modes along with the tensor representation to reduce computations.

Figures 9 and 10 show that ALTO outperforms the best mode-agnostic and mode-specific formats as well as memoization schemes in terms of the speedup of tensor operations (MTTKRP on all modes). In addition, the results indicate that ALTO can effectively reduce synchronization and utilize the larger caches on SPR (relative to ICX) to further improve the performance compared to the state-of-the-art libraries. Specifically, ALTO achieves $15.7\times$ and $25.3\times$ geometric mean (GEOMEAN) speedup on the ICX and SPR CPUs, respectively, compared to the best mode-agnostic formats. Although the mode-specific (CSF-based) formats require substantial storage to keep multiple tensor copies, ALTO still delivers $3.4\times$ and $5.1\times$ geometric mean speedup on the ICX and SPR servers, respectively. While memoization methods reduce computations, it comes at the cost of increasing memory traffic and consuming substantial amount of extra memory, which significantly limit their scalability. As a result, ALTO realizes $4.9\times$ and $8.4\times$

geometric mean speedup over STeF on the ICX and SPR CPUs, respectively, while effectively handling all large-scale tensors that cause out-of-memory errors with STeF. Furthermore, ALTO shows scalable performance for the sparse tensors with high data reuse. Compared to its sequential version, ALTO achieves up to $60\times$ and $80\times$ speedup on the 72-core ICX and 104-core SPR CPUs, respectively. For the other tensors (with limited/medium data reuse), ALTO is bounded by memory bandwidth, and as a result it has an average speedup of around $20\times$ and $30\times$ on ICX and SPR, respectively. Conversely, the performance of the previous approaches is highly variable across data sets as it depends on the shape of sparse tensors as well as the spatial distribution of their nonzero elements rather than their inherent data reuse [30]. Specifically, the tree-based (CSF, CSF-tile, and STeF) and block-based (HiCOO-SB10/SB14) methods depend on grouping the nonzero elements into tensor slices and blocks for effective compression. As illustrated in Figure 1, finding *balanced* clusters of nonzero elements in sparse tensors is highly unlikely. Thus, the prior tree- and block-based techniques suffer from workload imbalance and inefficient compression, which in turn lead to limited parallel performance when scaling to a large number of cores.

## 5.3 CP-APR Performance

We compare our ALTO-based CP-APR algorithm to the state-of-the-art SparTen library for non-negative count tensors. SparTen uses a variant of the COO format that keeps indexing as well as scheduling arrays for every tensor mode, which requires more than double the storage of COO [26]. In

(a) Performance on a 72-core ICX.
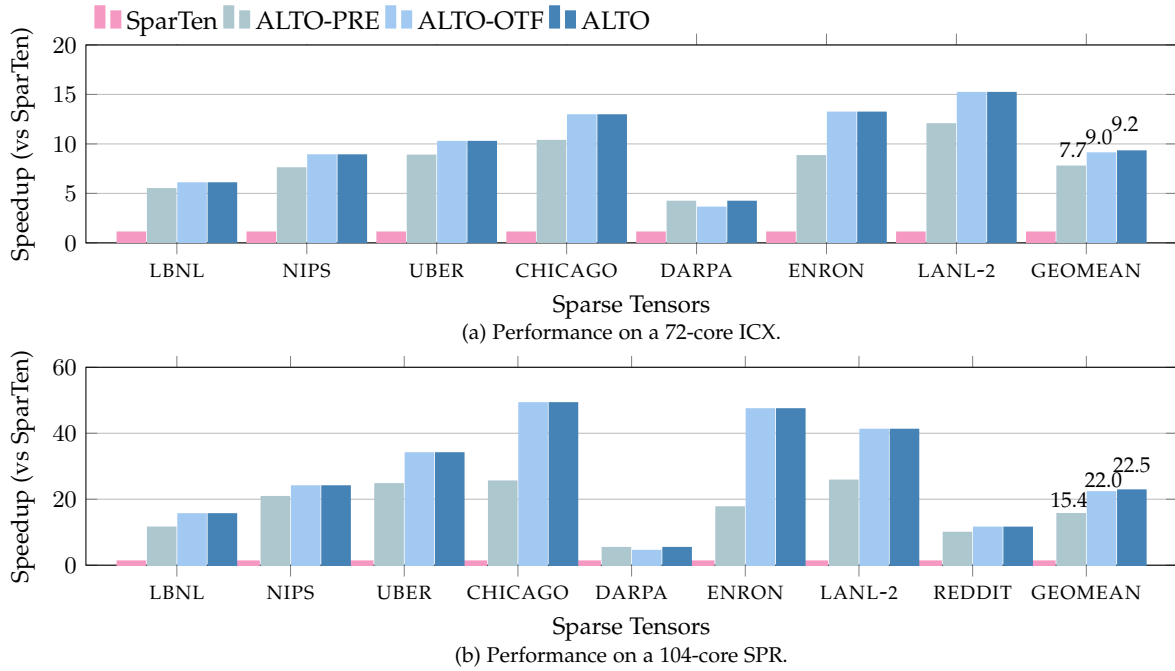


(b) Performance on a 104-core SPR.

Fig. 11: The achieved parallel speedup of the model update in CP-APR. The speedup is reported compared to the state-of-the-art SparTen library. ALTO chooses between ALTO-PRE and ALTO-OTF using our adaptive memory management heuristic (§4.3) to maximize performance. The sparse tensors are sorted in increasing order of their size (number of nonzero elements).

addition, SparTen computes CP-APR using the traditional method that stores and then reuses intermediate results rather than recomputing them. Hence, for large sparse tensors, such as REDDIT, SparTen fails to compute the CP-APR model on the ICX platform, even with 256 GiB of memory. In contrast, ALTO supports both recomputing (ALTO-OTF) or storing and then reusing (ALTO-PRE) intermediate results, which enables our CP-APR implementation to handle large-scale tensors.

Figure 11 shows the parallel performance of ALTO-based CP-APR compared to SparTen on the ICX and SPR CPUs. As the vast majority of time (more than 99 %) is spent in the model update ($\mathbf{\Phi}$) kernel (see Algorithm 5), the performance is evaluated based on the computation time of this tensor kernel. Note that ALTO-PRE and ALTO-OTF in Figure 11 represent the speedup achieved when the respective algorithms are used for all input tensors, and ALTO represents the speedup achieved when our adaptive memory management heuristic (§4.3) is used to choose between the two algorithms to maximize performance. Like CP-ALS, ALTO-based CP-APR realizes more scalable performance for tensors with high fiber reuse, and it further improves the performance on SPR relative to ICX by reducing synchronization and leveraging the larger fast memories on SPR. Therefore, ALTO delivers substantial performance gains compared to the SparTen library, achieving $9.2\times$ and $22.5\times$ speedup on the ICX and SPR CPUs, respectively. Furthermore, as the tensor size and data reuse increase, our on-the-fly (ALTO-OTF) algorithmic variant not only outperforms the traditional pre-computing approach (ALTO-PRE) but also realizes better scalability, even when the intermediate results can fit in memory.

### 5.4 Performance Characterization

Unlike prior compressed tensor formats, the parallel performance of ALTO depends on the inherent data reuse of

sparse tensors rather than the spatial distribution of their nonzero elements. To better understand the performance characteristics of ALTO, we created a Roofline model [46] for the SPR platform and collected performance counters across a set of representative parallel runs. For the Roofline model, an upper performance limit $\mathcal{P}$ is given by $\mathcal{P} = \min(\mathcal{P}_{\text{peak}}, B_{\text{peak}} \times OI)$, where $\mathcal{P}_{\text{peak}}$ is the peak performance, $B_{\text{peak}}$ is the peak memory bandwidth, and $OI$ is the *operational intensity* (i.e., the ratio of floating-point operations per byte). Moreover, we enhance our Roofline model to consider the cache bandwidth. The L2/L3 cache and main memory bandwidth are measured using `likwid-bench` from the Likwid tool suite. Since L1 bandwidth measurements are error-prone, we use the theoretical L1 bandwidth of two cache lines per cycle per core. The peak performance, $\mathcal{P}_{\text{peak}}$, is calculated based on the ability of the cores to execute two fused multiply-add (FMA) instructions on eight-element double precision vector registers (due to the availability of AVX-512) per cycle.

Since prior work detailed the performance analysis of the CP-ALS algorithm and its MTTKRP kernel [30], we focus on characterizing the performance of CP-APR in this study. Figure 12 shows the performance of the parallel $\mathbf{\Phi}$-ALTO kernel (Algorithm 5) for several representative tensors. To quantify the operational intensity, we calculated the required data movement from/to main memory as $l_{\text{avg}}mN(3R + RN + 1)$ for on-the-fly computation (ALTO-OTF) and $l_{\text{avg}}mN(3R+1)$ for pre-computation (ALTO-PRE), where $l_{\text{avg}}$ is the average number of inner iterations, $m$ is the number of nonzero elements, $N$ is the number of modes, and $R$ is the decomposition rank. We obtain the number of FLOPs required for the model update ($\mathbf{\Phi}$) by measuring hardware performance counters using `likwid-perfctr` from the Likwid tool suite. The results indicate that although such memory-intensive computations suffer from low operational intensity, ALTO can still exceed the peak
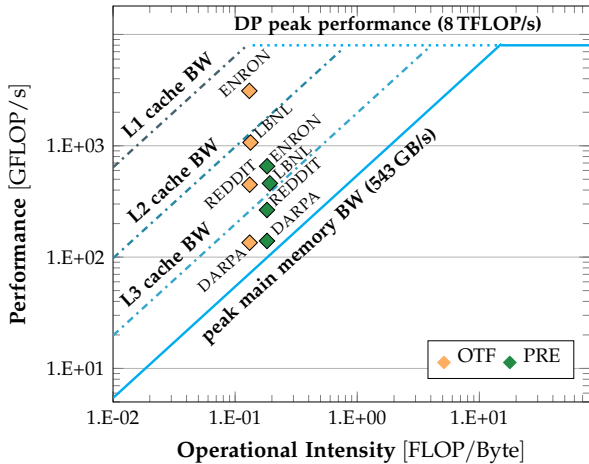
Fig. 12: The parallel performance of the model update ($\Phi$) in the CP-APR algorithm using ALTO on a 104-core SPR system. Orange diamonds indicate on-the-fly computation (ALTO-OTF), while green diamonds represent pre-computation (ALTO-PRE).

main memory bandwidth by exploiting the inherent data reuse and by efficiently resolving update conflicts in caches.

Specifically, the Roofline plot shows the performance of ALTO-OTF and ALTO-PRE for two tensors with high data reuse (ENRON and REDDIT) and two tensor with limited data reuse (LANL and DARPA). For ENRON, we observe that ALTO provides data access in a manner that allows the computation to be handled mainly from L1 and L2 cache. However, as a medium-sized tensor with high data reuse, it does not benefit from pre-computation and shows superior performance for the on-the-fly algorithm. The REDDIT tensor, despite having good fiber reuse, is highly sparse and it is the largest tensor in our set of experiments (with 4.6 billion nonzero elements). This increases the memory pressure and effectively leads to more data accesses from slower memory, which reduces the performance gap between ALTO-OTF and ALTO-PRE relative to the ENRON tensor.

While LBNL is extremely sparse and has limited data reuse, it is also the smallest of all tensors in the experiments. This allows ALTO to handle most of the data from the caches and to benefit from on-the-fly computation; however, the performance of LBNL is lower than denser tensors such as ENRON. Finally, the DARPA tensor, even though being similar in size to ENRON, has very limited fiber reuse (along mode-3). This leads to a significantly lower performance compared to any of the other tensors, yet the hybrid (recursive and output-oriented) tensor traversal of ALTO still captures some data reuse from caches and allows both ALTO-PRE and ALTO-OTF to realize superior performance, exceeding the main memory bandwidth. For DARPA we can observe a slightly better performance when using pre-computation. Hence, the performance analysis indicates that ALTO-PRE is especially relevant for large tensors that additionally show hyper-sparsity and limited data reuse.

### 5.5 Memory Storage

Figure 13 details the relative storage of the different sparse tensor formats compared to COO. Due to its efficient linearization, as detailed in §3.1, ALTO requires less storage than the CSF, CSF-tile and raw (COO) formats for all investigated tensors. The tree-based, mode-specific formats (CSF

and CSF-tile) consume significantly more storage space than COO because they require multiple tensor representations for the different mode orientations. While it can be beneficial for computation, imposing a tilling over the tensors (as done by CSF-tile) increases memory storage. The memory consumption of the block-based formats (HiCOO-SB10/SB14) highly depends not only on the spatial distribution of the nonzero elements, but also on the block and superblock sizes. Compared to the COO format, HiCOO can reduce the memory footprint of tensors when the resulting blocks are relatively dense, i.e., the number of nonzero elements per block is high. However, for hyper-sparse tensors such as DELI, NELL-1, AMAZON, and REDDIT, HiCOO requires more storage by up to a factor of 2.6. While the tensor format of STeF ("STeF (tensor)" in Figure 13) only requires storage on-par or even smaller than ALTO, the additional memory needed for memoization leads to a higher memory footprint ("STeF") compared to ALTO in all cases.

### 5.6 Format Generation Cost

Figure 14 details the generation cost of the different sparse tensor representations from a sparse tensor in the COO format on the SPR platform. Instead of processing nonzero elements in a multi-dimensional form as HiCOO and CSF, ALTO works on a linearized representation that needs substantially lower number of comparison operations to sort nonzero elements. Furthermore, the HiCOO formats require additional clustering of elements based on their multi-dimensional coordinates, as well as scheduling of the blocks and superblocks for avoiding conflicts, while STeF requires additional sorting of the nonzero elements along a specific mode order for best performance. Thus, ALTO achieves substantial geometric mean speedup for format generation compared to HiCOO-SB10 ($50\times$), HiCOO-SB14 ($75\times$), CSF-tile ($10\times$), CSF ($6\times$), and STeF ($44\times$).

## 6 RELATED WORK

Our mode-agnostic ALTO format was motivated by the linearized coordinate (LCO) format [39], which also flattens sparse tensors but in a mode-specific way, i.e., along a given mode orientation. Hence, LCO requires either multiple tensor copies or permuting tensors for efficient computation. Additionally, the authors limit their focus to sequential algorithms, and it is not clear how LCO can be used to efficiently execute parallel sparse tensor computations.

Researchers proposed variants [21], [26] of the COO format to reduce the synchronization overhead using *mode-specific* scheduling arrays. The state-of-the-art SparTen library [22] uses a COO variant [26] to accelerate the decomposition of non-negative count tensors across different hardware architectures. However, these COO variants adversely affect the input data locality and lead to random access of the nonzero elements [26], especially for sparse tensors with high data reuse. Moreover, keeping fine-grained scheduling information for *all tensor modes* can more than double the memory consumption, compared to the COO format [26].

The popular SPLATT library [20] leverages the CSF format [20], [25] to decompose sparse tensors on multicore CPUs. However, this mode-specific compressed format requires multiple tensor copies for best performance.
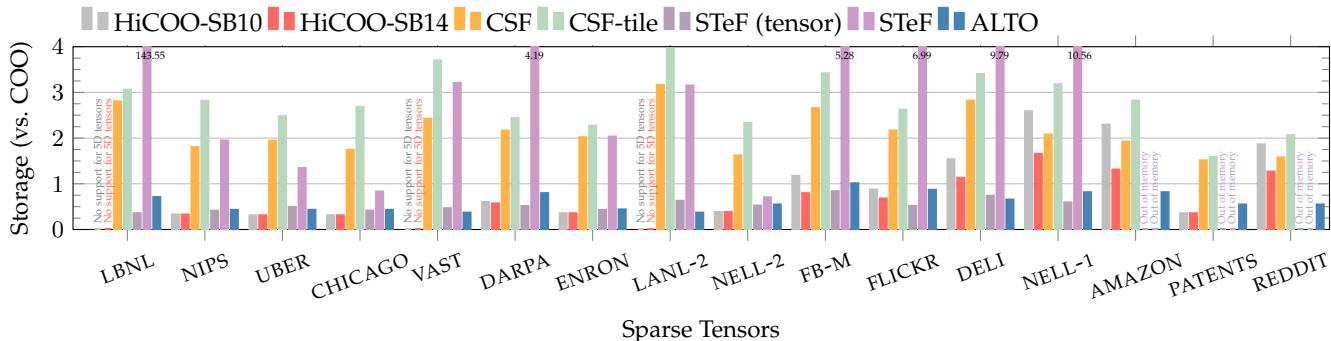
Fig. 13: The tensor storage across the different formats relative to COO. The tensors are sorted in an increasing order of their size.
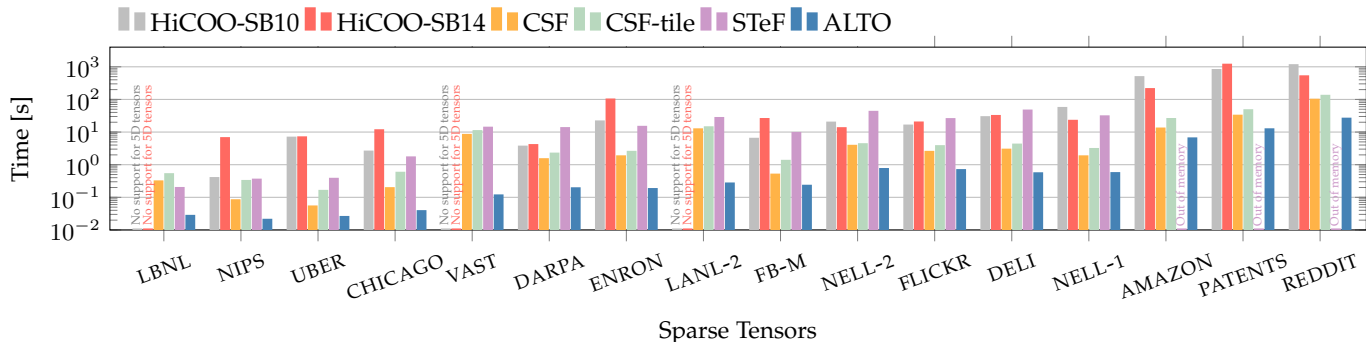


Fig. 14: The format construction cost on SPR in seconds. The sparse tensors are sorted in an increasing order of their size.

In addition, CSF packs the nonzero elements into coarse-grained tensor slices and fibers, which limits its scalability on massively parallel architectures. To improve the performance on GPUs, recent CSF-based formats [28], [29] expose more balanced and fine-grained parallelism but at the expense of substantial synchronization overheads and expensive preprocessing and format generation costs.

Alternatively, the ParTI library uses the mode-agnostic, block-based HiCOO format [15] to decompose sparse tensors using only one tensor copy. Yet, HiCOO is highly sensitive to the characteristics of sparse tensors as well as the block size. Due to the irregular (skewed) data distributions in sparse tensors, the number of nonzero elements per block varies widely across HiCOO blocks, even after expensive mode-specific tensor permutations which in practice further increase workload imbalance [27]. As a result, when tensors are highly sparse, HiCOO can consume more storage than the COO format [15]. Moreover, using small data types for indexing nonzero elements within a block can end up under-utilizing the compute units and memory bandwidth in modern parallel architectures [47], [48], which are optimized for wide memory transactions [49].

STeF [38] leverages the mode-specific CSF format to accelerate all-modes MTTKRP by memoization of partial MTTKRP results. Nevertheless, the additional space required for memoization can be more than double the memory storage of the sparse tensor and factor matrices, which limits STeF applicability to small- and medium-scale tensors [38]. SpTFS [16] utilizes machine learning [50], [51] to predict the best of COO, HiCOO, and CSF formats to compute MTTKRP for a given sparse tensor. FLYCOO [52], [53] extends the COO format to memory-constrained platforms (such as FPGAs) by processing a tensor into small equal-sized shards. However, FLYCOO requires dynamic mode-specific tensor remapping/reordering and it needs more storage than COO to keep sharding information for every mode.

## 7 CONCLUSION

To overcome the limitations of existing sparse tensor formats, this work introduced ALTO, a compact mode-agnostic format to efficiently encode higher-order tensors of irregular shapes and data distributions. Thanks to their adaptive tensor traversal and superior workload balance and data reuse, our ALTO-based parallel algorithms for decomposing normally distributed data (CP-ALS) and non-negative count data (CP-APR) delivered an order-of-magnitude speedup over the best mode-agnostic formats while requiring $\sim 50\%$ of COO storage. Moreover, ALTO achieved $5.1\times$ and $8.4\times$ geometric mean speedup over the best mode-specific and memoization methods, respectively, while needing between $14\%$ to $25\%$ of their overall storage. Our future work will investigate distributed-memory platforms and other common sparse tensor algorithms, besides tensor decomposition.

## REFERENCES

[1] J. C. Ho, J. Ghosh, and J. Sun, "Marble: High-throughput Phenotyping from Electronic Health Records via Sparse Nonnegative Tensor Factorization," in *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '14. New York, NY, USA: Association for Computing Machinery, 2014, p. 115–124, DOI: 10.1145/2623330.2623658.

[2] Y. Wang, R. Chen, J. Ghosh, J. C. Denny, A. Kho, Y. Chen, B. A. Malin, and J. Sun, "Rubik: Knowledge Guided Tensor Factorization and Completion for Health Data Analytics," in *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '15. New York, NY, USA: Association for Computing Machinery, 2015, p. 1265–1274, DOI: 10.1145/2783258.2783395.

[3] T. Kobayashi, A. Sapienza, and E. Ferrara, "Extracting the multi-timescale activity patterns of online financial markets," *Scientific Reports*, vol. 8, no. 1, pp. 1–11, 2018, DOI: 10.1038/s41598-018-29537-w.

[4] H. Fanaee-T and J. Gama, "Tensor-based anomaly detection: An interdisciplinary survey," *Knowledge-Based Systems*, vol. 98, pp. 130–147, 2016. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0950705116000472

[5] T. G. Kolda and J. Sun, "Scalable Tensor Decompositions for Multi-aspect Data Mining," in *2008 Eighth IEEE International Conference on Data Mining*, 2008, pp. 363–372, DOI: 10.1109/ICDM.2008.89.

[6] E. E. Papalexakis, C. Faloutsos, and N. D. Sidiropoulos, "Tensors for Data Mining and Data Fusion: Models, Applications, and Scalable Algorithms," *ACM Trans. Intell. Syst. Technol.*, vol. 8, no. 2, Oct. 2016, DOI: 10.1145/2915921.

[7] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky, "Tensor Decompositions for Learning Latent Variable Models," *Journal of Machine Learning Research*, vol. 15, no. 1, p. 2773–2832, Jan. 2014. [Online]. Available: http://jmlr.org/papers/v15/anandkumar14b.html

[8] N. D. Sidiropoulos, L. De Lathauwer, X. Fu, K. Huang, E. E. Papalexakis, and C. Faloutsos, "Tensor Decomposition for Signal Processing and Machine Learning," *IEEE Transactions on Signal Processing*, vol. 65, no. 13, pp. 3551–3582, 2017, DOI: 10.1109/TSP.2017.2690524.

[9] S. Smith, J. W. Choi, J. Li, R. Vuduc, J. Park, X. Liu, and G. Karypis. (2017) FROSTT: The formidable repository of open sparse tensors and tools. [Online]. Available: http://frostt.io/

[10] T. G. Kolda and B. W. Bader, "Tensor Decompositions and Applications," *SIAM Review*, vol. 51, no. 3, pp. 455–500, 2009, DOI: 10.1137/07070111X.

[11] B. W. Bader and T. G. Kolda, "Efficient MATLAB Computations with Sparse and Factored Tensors," *SIAM Journal on Scientific Computing*, vol. 30, no. 1, pp. 205–231, 2008, DOI: 10.1137/060676489.

[12] D. M. Dunlavy, T. G. Kolda, and E. Acar, "Temporal Link Prediction Using Matrix and Tensor Factorizations," *ACM Trans. Knowl. Discov. Data*, vol. 5, no. 2, Feb. 2011, DOI: 10.1145/1921632.1921636.

[13] J. Sun, D. Tao, and C. Faloutsos, "Beyond Streams and Graphs: Dynamic Tensor Analysis," in *Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '06. New York, NY, USA: Association for Computing Machinery, 2006, p. 374–383, DOI: 10.1145/1150402.1150445.

[14] E. C. Chi and T. G. Kolda, "On Tensors, Sparsity, and Nonnegative Factorizations," *SIAM Journal on Matrix Analysis and Applications*, vol. 33, no. 4, pp. 1272–1299, 2012, DOI: 10.1137/110859063.

[15] J. Li, J. Sun, and R. Vuduc, "HiCOO: Hierarchical Storage of Sparse Tensors," in *SC18: International Conference for High Performance Computing, Networking, Storage and Analysis*, 2018, pp. 238–252, DOI: 10.1109/SC.2018.00022.

[16] Q. Sun, Y. Liu, M. Dun, H. Yang, Z. Luan, L. Gan, G. Yang, and D. Qian, "SpTFS: Sparse Tensor Format Selection for MTTKRP via Deep Learning," in *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, 2020, pp. 1–14, DOI: 10.1109/SC41405.2020.00022.

[17] J. Choi, X. Liu, S. Smith, and T. Simon, "Blocking Optimization Techniques for Sparse Tensor Computation," in *2018 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2018, pp. 568–577, DOI: 10.1109/IPDPS.2018.00066.

[18] O. Kaya and B. Uçar, "Scalable sparse tensor decompositions in distributed memory systems," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '15. New York, NY, USA: Association for Computing Machinery, 2015, DOI: 10.1145/2807591.2807624.

[19] A. Nguyen, A. E. Helal, F. Checconi, J. Laukemann, J. J. Tithi, Y. Soh, T. Ranadive, F. Petrini, and J. W. Choi, "Efficient, Out-of-Memory Sparse MTTKRP on Massively Parallel Architectures," in *Proceedings of the 36th ACM International Conference on Supercomputing*, ser. ICS '22. New York, NY, USA: Association for Computing Machinery, 2022, DOI: 10.1145/3524059.3532363.

[20] S. Smith, N. Ravindran, N. D. Sidiropoulos, and G. Karypis, "SPLATT: Efficient and Parallel Sparse Tensor-Matrix Multiplication," in *2015 IEEE International Parallel and Distributed Processing Symposium*, 2015, pp. 61–70, DOI: 10.1109/IPDPS.2015.27.

[21] B. Liu, C. Wen, A. D. Sarwate, and M. M. Dehnavi, "A Unified Optimization Approach for Sparse Tensor Operations on GPUs," in *2017 IEEE International Conference on Cluster Computing (CLUSTER)*, 2017, pp. 47–57, DOI: 10.1109/CLUSTER.2017.75.

[22] K. Teranishi, D. M. Dunlavy, J. M. Myers, and R. F. Barrett, "SparTen: Leveraging Kokkos for On-node Parallelism in a Second-Order Method for Fitting Canonical Polyadic Tensor Models to Poisson Data," in *2020 IEEE High Performance Extreme Computing Conference (HPEC)*, 2020, pp. 1–7, DOI: 10.1109/HPEC43674.2020.9286251.

[23] P.-D. Letourneau, M. Baskaran, T. Henretty, J. Ezick, and R. Lethin, "Computationally Efficient CP Tensor Decomposition Update Framework for Emerging Component Discovery in Streaming Data," in *2018 IEEE High Performance extreme Computing Conference (HPEC)*, 2018, pp. 1–8, DOI: 10.1109/HPEC.2018.8547700.

[24] M. Baskaran, B. Meister, N. Vasilache, and R. Lethin, "Efficient and scalable computations with sparse tensors," in *2012 IEEE Conference on High Performance Extreme Computing*, 2012, pp. 1–6, DOI: 10.1109/HPEC.2012.6408676.

[25] S. Smith and G. Karypis, "Tensor-Matrix Products with a Compressed Sparse Tensor," in *Proceedings of the 5th Workshop on Irregular Applications: Architectures and Algorithms*, ser. IA³ '15. New York, NY, USA: Association for Computing Machinery, 2015, DOI: 10.1145/2833179.2833183.

[26] E. T. Phipps and T. G. Kolda, "Software for Sparse Tensor Decomposition on Emerging Computing Architectures," *SIAM Journal on Scientific Computing*, vol. 41, no. 3, pp. C269–C290, 2019, DOI: 10.1137/18M1210691.

[27] J. Li, B. Uçar, U. V. Çatalyürek, J. Sun, K. Barker, and R. Vuduc, "Efficient and effective sparse tensor reordering," in *Proceedings of the ACM International Conference on Supercomputing*, ser. ICS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 227–237, DOI: 10.1145/3330345.3330366.

[28] I. Nisa, J. Li, A. Sukumaran-Rajam, R. Vuduc, and P. Sadayappan, "Load-Balanced Sparse MTTKRP on GPUs," in *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2019, pp. 123–133, DOI: 10.1109/IPDPS.2019.00023.

[29] I. Nisa, J. Li, A. Sukumaran-Rajam, P. S. Rawat, S. Krishnamoorthy, and P. Sadayappan, "An Efficient Mixed-Mode Representation of Sparse Tensors," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '19. New York, NY, USA: Association for Computing Machinery, 2019, DOI: 10.1145/3295500.3356216.

[30] A. E. Helal, J. Laukemann, F. Checconi, J. J. Tithi, T. Ranadive, F. Petrini, and J. Choi, "ALTO: adaptive linearized storage of sparse tensors," in *Proceedings of the 35th ACM International Conference on Supercomputing*, ser. ICS '21. New York, NY, USA: Association for Computing Machinery, 2021, p. 404–416, DOI: 10.1145/3447818.3461703.

[31] Y. Soh, A. E. Helal, F. Checconi, J. Laukemann, J. J. Tithi, T. Ranadive, F. Petrini, and J. W. Choi, "Dynamic Tensor Linearization and Time Slicing for Efficient Factorization of Infinite Data Streams," in *2023 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2023, pp. 402–412, DOI: 10.1109/IPDPS54959.2023.00048.

[32] S. Chou, F. Kjolstad, and S. Amarasinghe, "Format abstraction for sparse tensor algebra compilers," *Proc. ACM Program. Lang.*, vol. 2, no. OOPSLA, Oct. 2018, DOI: 10.1145/3276493.

[33] M. Abadi et al., "TensorFlow: A System for Large-Scale Machine Learning," in *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. Savannah, GA: USENIX Association, Nov. 2016, pp. 265–283. [Online]. Available: https://www.usenix.org/conference/osdi16/technical-sessions/presentation/abadi

[34] N. Vervliet, O. Debals, and L. De Lathauwer, "Tensorlab 3.0 — Numerical optimization strategies for large-scale constrained and coupled matrix/tensor factorization," in *2016 50th Asilomar Conference on Signals, Systems and Computers*, 2016, pp. 1733–1738, DOI: 10.1109/ACSSC.2016.7869679.

[35] S. Smith and G. Karypis, "Accelerating the Tucker Decomposition with Compressed Sparse Tensors," in *Euro-Par 2017: Parallel Processing*, F. F. Rivera, T. F. Pena, and J. C. Cabaleiro, Eds. Cham: Springer International Publishing, 2017, pp. 653–668, DOI: 10.1007/978-3-319-64203-1_47.

[36] J. Choquette, "NVIDIA Hopper H100 GPU: Scaling Performance," *IEEE Micro*, vol. 43, no. 3, pp. 9–17, 2023, DOI: 10.1109/MM.2023.3256796.

[37] J. Li, J. Choi, I. Perros, J. Sun, and R. Vuduc, "Model-Driven Sparse CP Decomposition for Higher-Order Tensors," in *2017 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2017, pp. 1048–1057, DOI: 10.1109/IPDPS.2017.80.

[38] S. E. Kurt, S. Raje, A. Sukumaran-Rajam, and P. Sadayappan, "Sparsity-Aware Tensor Decomposition," in *2022 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, 2022, pp. 952–962, DOI: 10.1109/IPDPS53621.2022.00097.

[39] A. P. Harrison and D. Joseph, "High Performance Rearrangement and Multiplication Routines for Sparse Tensor Arithmetic," *SIAM Journal on Scientific Computing*, vol. 40, no. 2, pp. C258–C281, 2018, DOI: 10.1137/17M1115873.

[40] S. Liu and G. Trenkler, "Hadamard, Khatri-Rao, Kronecker, and Other Matrix Products," *International Journal of Information and Systems Sciences*, vol. 4, no. 1, pp. 160–177, 2008.

[41] T. P. Samantha Hansen and T. G. Kolda, "Newton-based optimization for kullback–leibler nonnegative tensor factorizations," *Optimization Methods and Software*, vol. 30, no. 5, pp. 1002–1029, 2015, DOI: 10.1080/10556788.2015.1009977.

[42] G. Peano, "Sur une courbe, qui remplit toute une aire plane," *Mathematische Annalen*, vol. 36, no. 1, pp. 157–160, Mar. 1890, DOI: 10.1007/BF01199438.

[43] G. M. Morton, "A computer Oriented Geodetic Data Base; and a New Technique in File Sequencing," IBM Ltd., 150 Laurier Ave., Ottawa, Ontario, Canada, Tech. Rep., 1966.

[44] T. Gruber, J. Eitzinger, G. Hager, and G. Wellein, "Likwid," Nov. 2023, DOI: 10.5281/zenodo.10105559.

[45] I. Jeon, E. E. Papalexakis, U. Kang, and C. Faloutsos, "HaTen2: Billion-scale tensor decompositions," in *2015 IEEE 31st International Conference on Data Engineering*, 2015, pp. 1047–1058, DOI: 10.1109/ICDE.2015.7113355.

[46] S. Williams, A. Waterman, and D. Patterson, "Roofline: An insightful visual performance model for multicore architectures," *Commun. ACM*, vol. 52, no. 4, pp. 65–76, 2009, DOI: 10.1145/1498765.1498785.

[47] A. Abel and J. Reineke, "uops.info: Characterizing Latency, Throughput, and Port Usage of Instructions on Intel Microarchitectures," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, ser. ASPLOS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 673–686, DOI: 10.1145/3297858.3304062.

[48] S. Mittal, "A Survey of Techniques for Designing and Managing CPU Register File," *Concurrency and Computation: Practice and Experience*, vol. 29, no. 4, p. e3906, 2017, DOI: 10.1002/cpe.3906.

[49] S. Ghose, T. Li, N. Hajinazar, D. S. Cali, and O. Mutlu, "Demystifying Complex Workload-DRAM Interactions: An Experimental Study," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 3, no. 3, Dec. 2019, DOI: 10.1145/3366708.

[50] Y. Zhao, J. Li, C. Liao, and X. Shen, "Bridging the gap between deep learning and sparse matrix format selection," in *Proceedings of the 23rd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPoPP '18. New York, NY, USA: Association for Computing Machinery, 2018, p. 94–108, DOI: 10.1145/3178487.3178495.

[51] Z. Xie, G. Tan, W. Liu, and N. Sun, "IA-SpGEMM: an input-aware auto-tuning framework for parallel sparse matrix-matrix multiplication," in *Proceedings of the ACM International Conference on Supercomputing*, ser. ICS '19. New York, NY, USA: Association for Computing Machinery, 2019, p. 94–105, DOI: 10.1145/3330345.3330354.

[52] S. Wijeratne, T.-Y. Wang, R. Kannan, and V. Prasanna, "Accelerating Sparse MTTKRP for Tensor Decomposition on FPGA," in *Proceedings of the 2023 ACM/SIGDA International Symposium on Field Programmable Gate Arrays*, ser. FPGA '23. New York, NY, USA: Association for Computing Machinery, 2023, p. 259–269, DOI: 10.1145/3543622.3573179.

[53] S. Wijeratne, R. Kannan, and V. Prasanna, "Dynasor: A Dynamic Memory Layout for Accelerating Sparse MTTKRP for Tensor Decomposition on Multi-core CPU," in *2023 IEEE 35th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD)*, 2023, pp. 23–33, DOI: 10.1109/SBAC-PAD59825.2023.00012.

**Jan Laukemann** is a PhD student at the University of Erlangen-Nürnberg Erlangen and works for the Erlangen National High Performance Computing Center. Previously he worked as a Research Scientist at Intel Parallel Computing Lab. He focuses on application optimization and performance engineering for HPC systems and novel algorithms for scalable linear algebra, tensor decomposition and graph computations. His research interests primarily include x86 and non-x86 computer architectures, their performance behavior on the node level, and vectorization techniques.

**Ahmed E. Helal** is a Research Scientist at the Intel Parallel Computing Lab in Santa Clara, CA, USA. He specializes in program analysis and transformation of irregular applications to optimize their performance on heterogeneous parallel architectures. His research interests include scalable algorithms for high-dimensional data analytics and sparse linear algebra. Helal received his Ph.D. degree from the Bradley Department of Electrical and Computer Engineering at Virginia Tech.

**S. Isaac Geronimo Anderson** S. Isaac Geronimo Anderson is a Ph.D. student advised by Jee Choi and Hank Childs in the department of Computer Science at University of Oregon. He has research interests in sparse tensors, performance portability, and HPC application interoperation. Isaac has interned at Sandia National Laboratories, Oak Ridge National Laboratory, and Cray (now Hewlett Packard Enterprise).

**Fabio Checconi** is a Research Scientist at the Intel Parallel Computing Lab in Santa Clara, CA, USA. His research interests include large-scale graph algorithms and sparse linear algebra. Checconi received his Ph.D. degree in computer engineering from Scuola Superiore S. Anna in Pisa, Italy.

**Yongseok Soh** is a Ph.D. student at the University of Oregon, where he is an active member of the High Performance Computing Research Group. His research primarily concentrates on enhancing the efficiency of sparse tensor factorization algorithms on HPC systems involving algorithmic and architectural adjustments and optimizations.

**Dr. Jesmin Jahan Tithi** is an AI Research Scientist and Tech Lead at Intel Corporation, specializing in algorithm engineering, high-performance computing (HPC), and hardware-software codesign. With a Ph.D. from Stony Brook University and a B.Sc. from Bangladesh University of Engineering and Technology, she has completed internships at Google, Intel, and the Pacific Northwest National Laboratory. Dr. Tithi has over 40 peer-reviewed publications and 12 published patents.

**Teresea Ranadive** has been a researcher at the Laboratory for Physical Sciences since 2018. She received a B.S. in Mathematics from St. Vincent College (2011), and an M.S. (2013) and Ph.D. (2016) in Applied Mathematics from the University of Maryland, Baltimore County. Between 2016 and 2018, Dr. Ranadive was an Assistant Research Professor at Johns Hopkins University in the Applied Mathematics and Statistics department. Her research interests include numerical optimization, tensor decomposition, and Ising spin glass models.

**Brian J. Gravelle** Brian J Gravelle is a computer systems researcher at the Laboratory for Physical Sciences in Maryland, USA. He studied Computer Engineering at Gonzaga University before obtaining a PhD in Computer Science at the University of Oregon. In past research, he has evaluated the performance of various scientific and data analytics algorithms running on supercomputers and worked to improve the runtime of key kernels of computation. Currently, he studies new methods of improving the performance and programmability of distributed systems and GPUs with a particular focus on applications which rely on sparse linear algebra computation.

**Fabrizio Petrini** is a Senior Principal Engineer at the Intel Parallel Computing Lab in Santa Clara, CA, USA. His research interests include data-intensive algorithms for graph analytics and sparse linear algebra, Exascale computing, high-performance interconnection networks and novel architectures. He is the principal investigator of the TCStream CS project, Co-Principal investigator of the Intel PIUMA architecture, developed under the DARPA HIVE and SDH programs, and the upcoming Intel TIGRE system, as part of the IARPA AGILE program. Petrini received his Ph.D. degree in computer science from Pisa University, Italy. He is a Senior Member of IEEE.

**Jee Choi** is an Assistant Professor in the Department of Computer Science at the University of Oregon. During his PhD, he worked on designing parallel and scalable algorithms for scientific applications, and modeling their performance and energy efficiency on the latest high-performance computing (HPC) systems. After graduation, he worked as a research staff member at the IBM T. J. Watson Research Center on designing and optimizing tensor decomposition algorithms for Big Data analytics.