Linearized Tensor Format for Performance-Portable Sparse Tensor Computation

Andy Nguyen¹, **Ahmed E. Helal**², Fabio Checconi², Jan Laukemann², Jesmin Jahan Tithi², Yongseok Soh¹, Teresa Ranadive³, Fabrizio Petrini², Tammy Kolda⁴, Jee Whan Choi¹

¹ University of Oregon ² Intel Labs ³ Laboratory for Physical Sciences ⁴ Sandia National Laboratory



- A fundamental problem in *sparse* tensor computation is how to *store, group,* and *organize* the nonzero elements to
- 1. reduce tensor storage





- A fundamental problem in sparse tensor computation is how to *store, group,* and *organize* the nonzero elements to
- 1. reduce tensor storage



2. improve data locality

		••
•••	•••	•••
•••	•••	•••
••••	••••	••••



- A fundamental problem in sparse tensor computation is how to *store, group,* and *organize* the nonzero elements to
- 1. reduce tensor storage



- 2. improve data locality
- 3. increase parallelism





- A fundamental problem in sparse tensor computation is how to *store, group,* and *organize* the nonzero elements to
- 1. reduce tensor storage



- 2. improve data locality
- 3. increase parallelism
- 4. decrease workload imbalance & synchronization overhead





- A fundamental problem in sparse tensor computation is how to *store, group,* and *organize* the nonzero elements to
- 1. reduce tensor storage
- 2. improve data locality
- 3. increase parallelism
- 4. decrease workload imbalance & synchronization overhead





goals

- Can be classified based on their encoding of the indexing metadata into:
- 1. List-based formats
- 2. Tree-based formats
- 3. Block-based formats



• List-based format: COO (coordinate)

i	j	k	V
1	1	1	1
1	1	2	2
1	3	3	3
2	1	2	4
2	1	3	5
3	1	2	6
3	4	4	7
4	2	1	8
4	2	2	9
4	3	3	10
4	3	4	11
4	4	4	12



• Tree-based format: CSF (compressed sparse fiber)



i	j	k	V
1	1	1	1
1	1	2	2
1	3	3	3
2	1	2	4
2	1	3	5
3	1	2	6
3	4	4	7
4	2	1	8
4	2	2	9
4	3	3	10
4	3	4	11
4	4	4	12



• Block-based format: HiCOO (hierarchical COO)

_	<i>b</i> _{ptr}	b_i	b_j	b_k	e_i	e_{j}	e_k	V
Γ	1	0	0	0	1	1	1	1
2x2x2 block	1	0	0	1	1	1	0	2
	1	1	0	1	0	1	0	4
_	2	1	0	1	0	1	1	5
	3	0	1	1	1	1	1	3
	4	1	1	0	2	0	1	8
	4	1	0	1	1	1	0	6
	4	1	1	1	2	0	0	9
	5	1	1	1	2	1	1	10
	5	1	1	1	1	2	2	7
	5	1	1	1	2	1	2	11
	5	1	1	1	2	2	2	12





•••	•••	• • •
•••	•••	• • •
•••	•••	•••

Formats	Granularity	Mode Orientation	Data Locality	Parallelism	Load balance
List-based (COO)	Nonzero element	Mode-agnostic	Poor	Suffer from conflicts	Maximized
Tree-based (CSF)	Compressed tree	Mode-specific	Improved for a specific mode	Improved for a specific mode	Work imbalance (especially in short modes)
Block-based (HiCOO)	Compressed block	Mode-agnostic	Improved	Suffer from conflicts	Work imbalance across blocks

- Tree- and block-based formats are extensions of sparse matrix formats
- However, they are ill-suited for sparse tensors

- Tree- and block-based formats are extensions of sparse matrix formats
- However, they are ill-suited for sparse tensors
- 1. Tensor computation often operate over every dimension
 - This contrasts with sparse matrix-vector multiply (*one dimension*) and sparse matrix-matrix multiply (*two dimensions*)

- Tree- and block-based formats are extensions of sparse matrix formats
- However, they are ill-suited for sparse tensors
- 1. Tensor computation often operate over every dimension
 - This contrasts with sparse matrix-vector multiply (*one dimension*) and sparse matrix-matrix multiply (*two dimensions*)
- 2. Tensor sparsity >> matrix sparsity
 - Due to their dimensionality, tensors are extremely sparse it's difficult to find dense blocks
 - Tensor sparsity ranges from 1.5×10^{-2} to 4.3×10^{-15} [1]

- Tree- and block-based formats are extensions of sparse matrix formats
- However, they are ill-suited for sparse tensors
- 1. Tensor computation often operate over every dimension
 - This contrasts with sparse matrix-vector multiply (*one dimension*) and sparse matrix-matrix multiply (*two dimensions*)
- 2. Tensor sparsity >> matrix sparsity
 - Due to their dimensionality, tensors are extremely sparse it's difficult to find dense blocks
 - Tensor sparsity ranges from 1.5×10^{-2} to 4.3×10^{-15} [1]
- Our hypothesis simple mode-agnostic, list-based formats are the best for sparse tensors
 - How do we improve data locality and reduce conflicts?
 - How do we make it simple and performance-portable?

Linearized Formats

- ALTO (Adaptive Linearized Tensor Order) for CPUs
- BLCO (Block Linearized Coordinate) for GPUs
- Application of linearized formats to
 - *streaming* tensor decomposition [1]
 - *on-the-fly* Khatri-Rao product for CP-APR (WIP)
 - non-negative sparse tensor factorization for GPUs via PLANC (WIP)



- ALTO interleaves the index bits by grouping them by their positions
- Within each group, the bits are arranged from the longest mode (most-significant) to the shortest (least-significant)



ALTO bit mask



0

63

K



ALTO			
Value	Position		
<i>x</i> _{1,0,0}	2 (000010)		
<i>x</i> _{3,1,1}	15 (001111)		
<i>x</i> _{0,3,0}	20 (010100)		
<i>x</i> _{2,2,1}	25 (011001)		
<i>x</i> _{3,4,0}	42 (101010)		
$x_{1,6,1}$	51 (110011)		







K

k = **0**

k = 1



ALTO			
Value	Position		
<i>x</i> _{1,0,0}	2 (000010)		
<i>x</i> _{3,1,1}	15 (001111)		
<i>x</i> _{0,3,0}	20 (010100)		
<i>x</i> _{2,2,1}	25 (011001)		
<i>x</i> _{3,4,0}	42 (101010)		
$x_{1,6,1}$	51 (110011)		

4x8x2 tensor

◀ /100 01/10 11





ALTO			
Value	Position		
<i>x</i> _{1,0,0}	2 (000010)		
<i>x</i> _{3,1,1}	15 (001111)		
<i>x</i> _{0,3,0}	20 (010100)		
<i>x</i> _{2,2,1}	25 (011001)		
<i>x</i> _{3,4,0}	42 (101010)		
<i>x</i> _{1,6,1}	51 (110011)		





4x8x2 tensor



ALTO			
Value	Position		
<i>x</i> _{1,0,0}	2 (000010)		
<i>x</i> _{3,1,1}	15 (001111)		
<i>x</i> _{0,3,0}	20 (010100)		
<i>x</i> _{2,2,1}	25 (011001)		
<i>x</i> _{3,4,0}	42 (101010)		
<i>x</i> _{1,6,1}	51 (110011)		

Partitioning

K

	ALTO				
	Value	Position			
	<i>x</i> _{1,0,0}	2 (000010)			
	<i>x</i> _{3,1,1}	15 (001111)			
	<i>x</i> _{0,3,0}	20 (010100)			
	<i>x</i> _{2,2,1}	25 (011001)			
	<i>x</i> _{3,4,0}	42 (101010)			
	$x_{1,6,1}$	51 (110011)			







Par	rsity of oregon	k
I	ALTO	
Value	Position	
<i>x</i> _{1,0,0}	2 (000010)	
<i>x</i> _{3,1,1}	15 (001111)	
<i>x</i> _{0,3,0}	20 (010100)	
 <i>x</i> _{2,2,1}	25 (011001)	
<i>x</i> _{3,4,0}	42 (101010)	
<i>x</i> _{1,6,1}	51 (110011)	







Index Encoding & Decoding



for l = 1,...,L in parallel do
for all
$$x \in X_1$$
 do
fetch i_x, j_x, k_x, v_x
fetch $v_1 = A^{(2)}(j_x, :)$
fetch $v_2 = A^{(3)}(k_x, :)$
scratch += $v_x \cdot (v_1 * v_2)$
 $A^{(1)}(i_x, :)$ += scratch
endfor

for l = 1,...,L in parallel do for all x $\in X_1$ do fetch i_x, j_x, k_x, v_x fetch $v_1 = A^{(2)}(j_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch += $v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :)$ += scratch endfor

endfor

for l = 1, ..., L in parallel do $\leftarrow ----$ Each partition l could be a HiCOO block or a CSF sub-tree for all $x \in X_1$ do $\leftarrow ----$ Each nonzero element in the partition

for l = 1,...,L in parallel do for all x $\in X_1$ do fetch i_x, j_x, k_x, v_x fetch $v_1 = A^{(2)}(j_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch $+= v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :) +=$ scratch endfor

- Each partition 1 could be a HiCOO block or a CSF sub-tree Each nonzero element in the partition
- Fetch the indices and value associated with nonzero x

for l = 1, ..., L in parallel do \leftarrow ----- Ea for all $x \in X_1$ do \leftarrow ----- Ea fetch i_x, j_x, k_x, v_x \leftarrow ---- Fe fetch $v_1 = A^{(2)}(j_x, :)$ \leftarrow ---- Fe fetch $v_2 = A^{(3)}(k_x, :)$ \leftarrow ---- Fe scratch += $v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :)$ += scratch

- Each partition 1 could be a HiCOO block or a CSF sub-tree Each nonzero element in the partition
- Fetch the indices and value associated with nonzero x
- Fetch row j_x from mode-2 factor matrix
- – Fetch row k_x from mode-3 factor matrix

endfor

for l = 1, ..., L in parallel do $\leftarrow ----$ Each partition l could be a HiCOO block or a CSF sub-tree for all $x \in \mathcal{X}_1$ do fetch i_x, j_x, k_x, v_x fetch $v_1 = A^{(2)}(j_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch += $v_x \cdot (v_1 * v_2) \blacktriangleleft - - - -$ **4**---- $A^{(1)}(i_x,:)$ += scratch endfor

Each nonzero element in the partition Fetch the indices and value associated with nonzero X Fetch row j_x from mode-2 factor matrix Fetch row k_x from mode-3 factor matrix

MTTKRP

Update row i_x from mode-1 factor matrix, either atomically (e.g., COO, or HiCOO), or freely (e.g., CSF)

for l = 1,...,L in parallel do \leftarrow ALTO line segment for all x $\in X_1$ do fetch i_x, j_x, k_x, v_x fetch $v_1 = A^{(2)}(j_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch $+= v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :) +=$ scratch endfor

endfor

Matricized Tensor Times Khatri-Rao Product (MTTKRP)

for l = 1,...,L in parallel do
$$\leftarrow$$
 ALTO line segment
for all x $\in X_1$ do
fetch i_x, j_x, k_x, v_x ----- fetch a_x (ALTO index), v_x
fetch $v_1 = A^{(2)}(j_x, :)$ $i_x, j_x, k_x = decode(a_x)$
fetch $v_2 = A^{(3)}(k_x, :)$
scratch += $v_x \cdot (v_1 * v_2)$
 $A^{(1)}(i_x, :)$ += scratch
endfor



for $l = 1, \dots, L$ in parallel do $\leftarrow -----$ ALTO line segment for all $x \in X_1$ do This can be done efficiently fetch i_x, j_x, k_x, v_x fetch a_x (ALTO index), v_x using parallel bit $i_x, j_x, k_x = decode(a_x)$ fetch $v_1 = A^{(2)}(j_x, :)$ extract/deposit (pext/pdep) instructions on x86 CPUs fetch $v_2 = A^{(3)}(k_x, :)$ scratch += $v_x \cdot (v_1 * v_2)$ Adaptive update mechanism $A^{(1)}(i_x,:)$ += scratch If rows (i_x) have limited reuse \rightarrow update using atomic ops. ٠ endfor If rows (i_x) have large reuse \rightarrow use scratchpad to combine endfor local updates and parallel reduction to merge globally

• reuse = nnz / mode_length



ALTO vs. Z-curve

ALTO Bit Mask

Z-curve Bit Mask



- In contrast to Z-ordering, ALTO uses a non-fractal encoding to
 - 1. adapt to irregularly shaped tensors,
 - 2. further reduce storage, and
 - 3. reduce encoding/decoding time

ALTO vs. State-of-the-art





Formats	Granularity	Mode Orientation	Data Locality	Parallelism	Load balance
List-based (COO)	Nonzero element	Mode-agnostic	Poor	Suffer from conflicts	Maximized
Tree-based (CSF)	Compressed tree	Mode-specific	Improved for a specific mode	Improved for a specific mode	Work imbalance (especially in short modes)
Block-based (HiCOO)	Compressed block	Mode-agnostic	Improved	Suffer from conflicts	Work imbalance across blocks
ALTO	Nonzero element	Mode-agnostic	Improved (via nearest-neighbor traversal)	Improved (via adaptive update)	Maximized

Performance Summary

- Intel Cascade Lake-X
 - 28 x 2 cores @ 1.8 GHz
- Oracle selects the best mode-agnostic and mode-specific format for each of the 15 tensors
- Mode-agnostic formats: COO and HiCOO
- Mode-specific formats: CSF and CSF with tiling with N copies





Speedup

- Speedup against serial MTTKRP ALTO implementation
- HiCOO does not support 5D tensors
- HiCOO and CSF runs out of memory for REDDIT





Speedup

- ALTO achieves ~80% of realizable speedup across tensors with different characteristics
- For high reuse, ALTO achieves 47x geo-mean speedup
- For limited reuse, ALTO achieves 16x geo-mean speedup



Helal et al. ALTO: Adaptive Linearized Storage of Sparse Tensors. ICS'21



Speedup

- CSF shows slightly better performance for NIPS, AMAZON, and NELL-1 (geo-mean speedup of 1.2x over ALTO) by keeping N copies of the tensor, optimized for each mode
- Performance of COO, HiCOO, and CSF are sensitive to irregular tensor shape and data distribution





Storage

• ALTO always requires less storage than COO due to linearization



Helal et al. ALTO: Adaptive Linearized Storage of Sparse Tensors. ICS'21



Storage

- ALTO always requires less storage than COO due to linearization
- HiCOO storage depends on the block/superblock sizes and spatial distribution of nonzero elements
 - For hyper-sparse tensors, HiCOO consumes more storage than COO





ALTO Generation Overhead

- ALTO substantially decreases sorting time by reducing the # of comparison operations
- Block-based formats require expensive clustering (and sometimes reordering) and scheduling of non-zero elements
- CSF is generated from **pre-sorted** tensors, but are still slower on average to construct than ALTO



Performance Portability

- Can ALTO perform well on GPUs?
- If not, what are the challenges?

```
for l = 1,...,L in parallel do

for all x \in \mathcal{X}_1 do

fetch i_x, j_x, k_x, v_x

fetch v_1 = A^{(2)}(j_x, :)

fetch v_2 = A^{(3)}(k_x, :)

scratch += v_x \cdot (v_1 * v_2)

A^{(1)}(i_x, :) += scratch

endfor
```

for l = 1, ..., L in parallel do \leftarrow Each partition l could be a thread block for all $x \in X_1$ do \leftarrow Each thread operates on a 1~k non-zero element fetch i_x, j_x, k_x, v_x fetch $v_1 = A^{(2)}(j_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch $+= v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :) +=$ scratch endfor

for l = 1, ..., L in parallel do \leftarrow Each partition l could be a thread block for all $x \in X_1$ do \leftarrow Each thread operates on a 1~k non-zero element fetch i_x, j_x, k_x, v_x \leftarrow fetch a_x (ALTO index), v_x fetch $v_1 = A^{(2)}(j_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch $+= v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :) +=$ scratch endfor

for l = 1, ..., L in parallel do \leftarrow ----- Each partition l constraints for all $x \in X_1$ do \leftarrow ----- Each thread operation for all $x \in X_1$ do \leftarrow ----- Each thread operation for i_x, j_x, k_x, v_x \leftarrow ----- fetch a_x (ALT fetch $v_1 = A^{(2)}(j_x, :)$ $i_x, j_x, k_x = deconstraints for <math>v_2 = A^{(3)}(k_x, :)$ fetch $v_2 = A^{(3)}(k_x, :)$ scratch $+= v_x \cdot (v_1 * v_2)$ $A^{(1)}(i_x, :) +=$ scratch

 Each partition 1 could be a thread block
 Each thread operates on a 1~k non-zero element
 fetch a_x (ALTO index), v_x i_x, j_x, k_x = decode(a_x) GPUs do NOT support parallel bit extract/deposit instructions

endfor





GPUs have limited on-chip memory



Encoding/Decoding on GPUs

- Blocked Linearized Coordinates (BLCO)
 - Linearize and order using ALTO
 - Re-linearize indices by grouping bits from same index together
 - Bit-wise shift and masking is cheap on GPUs

	l	v	
0	$(00000)_2$	1.0	
4	$(000100)_2$	2.0	
5	$(000101)_2$	4.0	
10	$(001010)_2$	8.0	
12	$(001100)_2$	6.0	
15	$(001111)_2$	9.0	
33	$(100001)_2$	5.0	
48	$(110000)_2$	3.0	
57	$(111001)_2$	10.0	
61	$(111101)_2$	11.0	
62	$(11110)_2$	7.0	
63	$(111111)_2$	12.0	

ALTO

BLCO

	l	v
0	$(00000)_2$	1.0
16	$(10000)_2$	2.0
17	$(10001)_2$	4.0
6	$(00110)_2$	8.0
18	$(10010)_2$	6.0
23	(1 <mark>0111)</mark> 2	9.0
1	$(00001)_2$	5.0
8	$(01000)_2$	3.0
11	$(01011)_2$	10.0
27	(1 1011) ₂	11.0
30	(1 1110) ₂	7.0
31	$(1111)_2$	12.0

Encoding/Decoding on GPUs

- Blocked Linearized Coordinates (BLCO)
 - Linearize and order using ALTO
 - Re-linearize indices by grouping bits from same index together
 - Bit-wise shift and masking is cheap on GPUs
- Adaptive Blocking
 - Use uppermost bits from every mode that exceeds target integer size (e.g., 64 bits) to form the initial blocks, then further divided to meet GPU memory constraints
 - Leverages native integer instructions
 - Reduces overall storage
 - Does not require expensive tuning
 - Leverages GPU schedulers to workload balance



 $(111110)_2$

 $(111111)_2$

62

63

ALTO

BLCO

b		l	v
	0	$(00000)_2$	1.0
	16	$(10000)_2$	2.0
0	17	(1 0001) ₂	4.0
0	6	$(00110)_2$	8.0
	18	(1 0010) ₂	6.0
	23	(1 <mark>0111</mark>) ₂	9.0
	1	$(00001)_2$	5.0
	8	$(01000)_2$	3.0
1	11	$(01011)_2$	10.0
	27	(1 1011) ₂	11.0
	30	(1 1110) ₂	7.0
	31	$(11111)_2$	12.0

Nguyen et al. Efficient, Out-of-Memory Sparse MTTKRP on Massively Parallel Architectures ICS'22

7.0

12.0

Synchronization



Performance

- MM-CSF is the baseline
- Amazon, Patents, and Reddit do not run on any prior frameworks
- F-COO only supports 3D tensors and <code>segfaults</code> on some tensors

BLCO F-COO GenTen MM-CSF



Performance

- Throughput limited by host-to-device data transfer
- In-memory performance on par with other tensors (measured via Nsight System Profiler)

BLCO (No data exchange) BLCO 5 Throughput (TB/s) 4 3 2 0 2 2 3 2 3 3 Reddit Patents Amazon Tensor

Performance

 MM-CSF is sensitive to data distribution, resulting in significant variation in performance across different modes (e.g., Uber, DARPA, Enron, and FB-M)

Intel Device1 V100 A100



Format Construction Cost



Format Construction Cost



Q&A

Backup Slides

ALTO



 The amount of information about the spatial position of a nonzero element decreases with each consecutive bit









K

 The amount of information about the spatial position of a nonzero element decreases with each consecutive bit







ALTO



K

 The amount of information about the spatial position of a nonzero element decreases with each consecutive bit







ALTO



K

 The amount of information about the spatial position of a nonzero element decreases with each consecutive bit







ALTO



K

 The amount of information about the spatial position of a nonzero element decreases with each consecutive bit







ALTO



K

 This is equivalent to partitioning the multi-dimensional space along the longest mode first







Spatial data distribution



A box plot of the data (nonzero elements) distribution across multi-dimensional blocks. The multi-dimensional subspace size is 128^N , where *N* is the number of dimensions (modes).