



Efficient, Out-of-Memory Sparse MTTKRP on Massively Parallel Architectures

Andy Nguyen¹, Ahmed E. Helal², Fabio Checconi², Jan Laukemann²,
Jesmin Jahan Tithi², Yongseok Soh¹, Teresa Ranadive³, Fabrizio Petrini²,
Jee Whan Choi¹

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

February 25th @ SIAM PP '22



Outline

- Introduction to sparse tensors
- Sparse tensor kernel: MTTKRP in CPD-ALS
- Prior state of the art approaches
- Our approach
 - Tensor data format
 - Parallel algorithm
 - Evaluations
- Closing thoughts



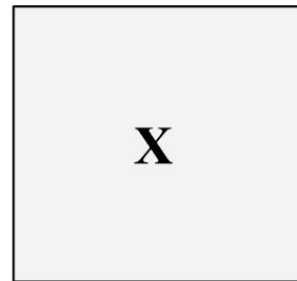
Tensors

Vector



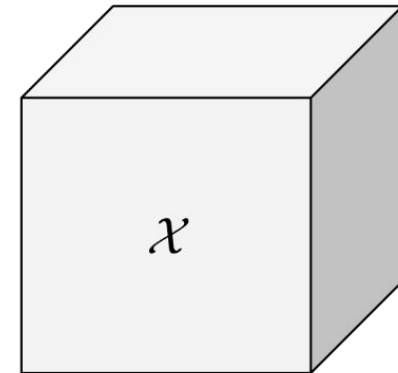
$$\mathbf{x} \in \mathbb{R}^I$$

Matrix



$$\mathbf{X} \in \mathbb{R}^{I \times J}$$

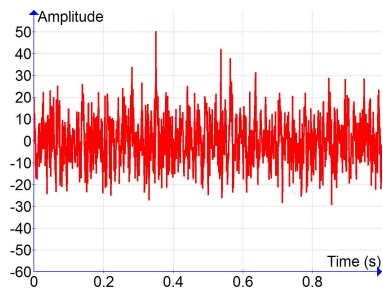
Tensor



$$\mathcal{X} \in \mathbb{R}^{I \times J \times K}$$



Real-World Tensors



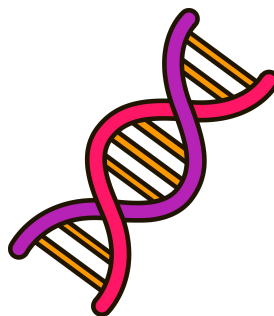
Signal processing

NETFLIX

Recommendation systems



Product reviews



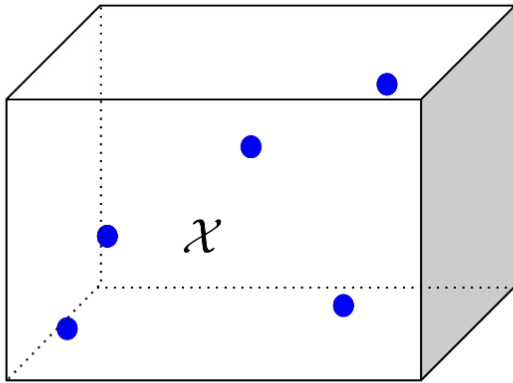
Electronic health records



Law enforcement data

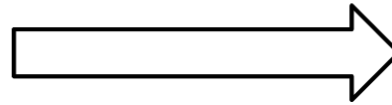
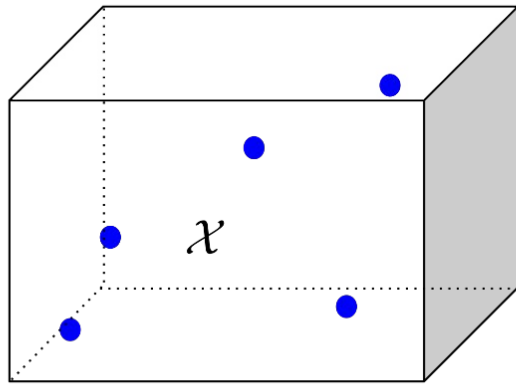


Sparse Tensors





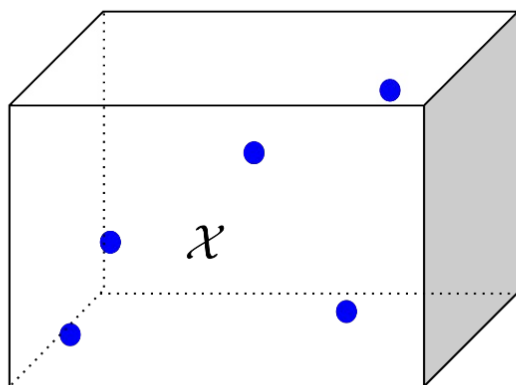
Sparse Tensors



i	j	k	v
1	4	3	5.2
4	5	2	8.7
1	1	3	0.3
5	3	2	4.6
4	1	2	1.9



Sparse Tensors

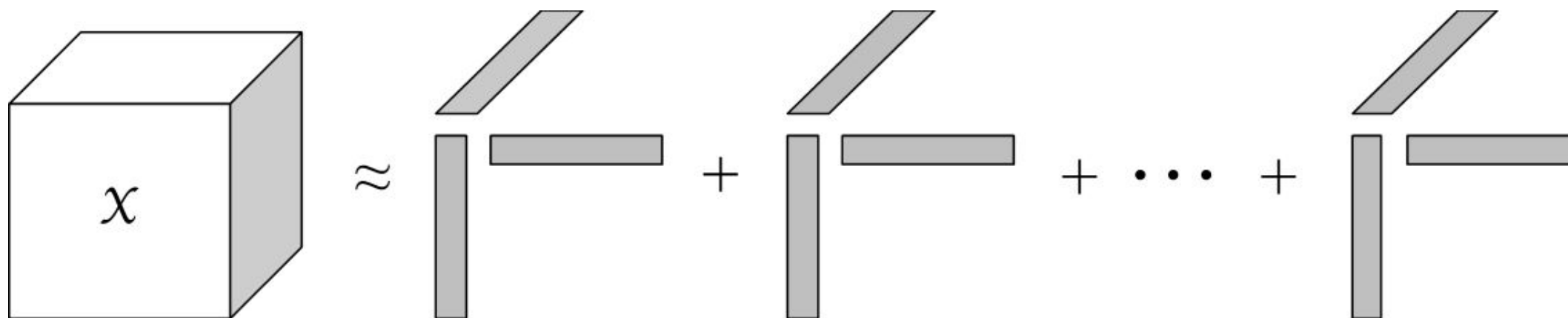


Even harder to optimize than sparse matrix kernels:

- Curse of dimensionality
- Algorithm scalability
- Complex data compression



Tensor Decomposition

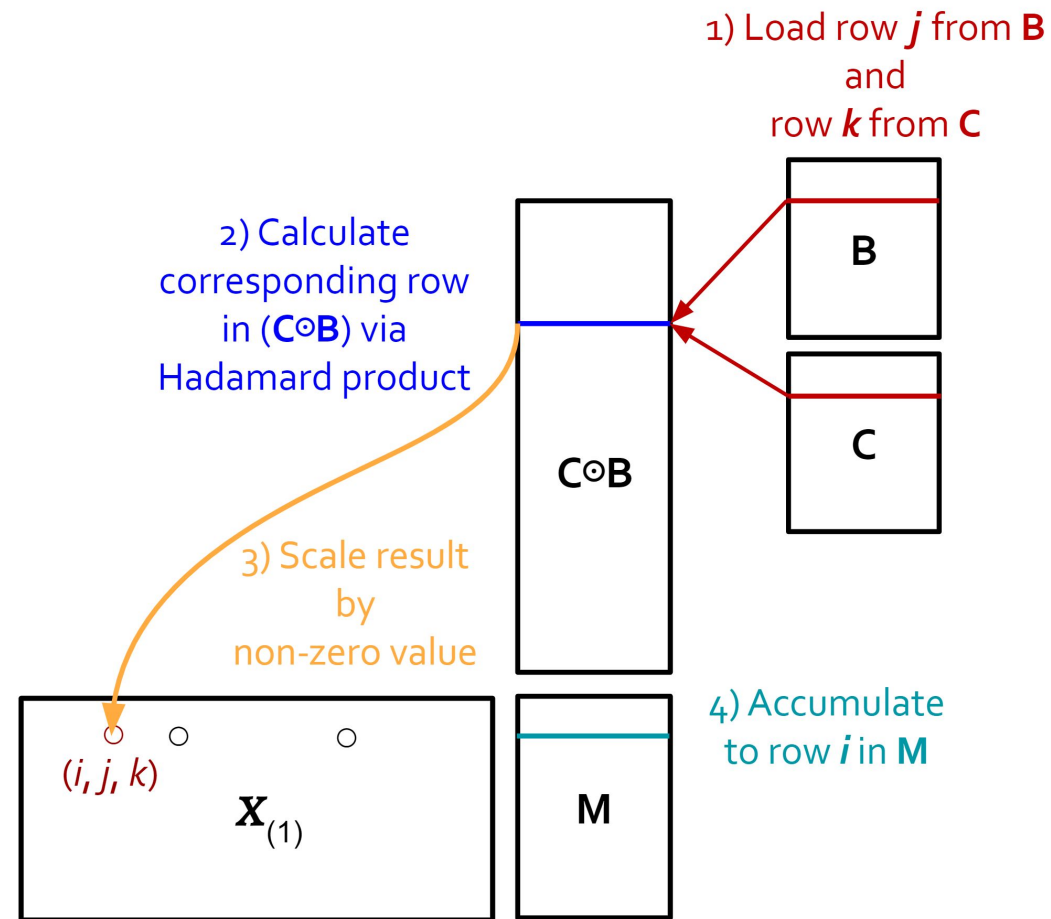


Canonical Polyadic Decomposition (CPD)



MTTKRP

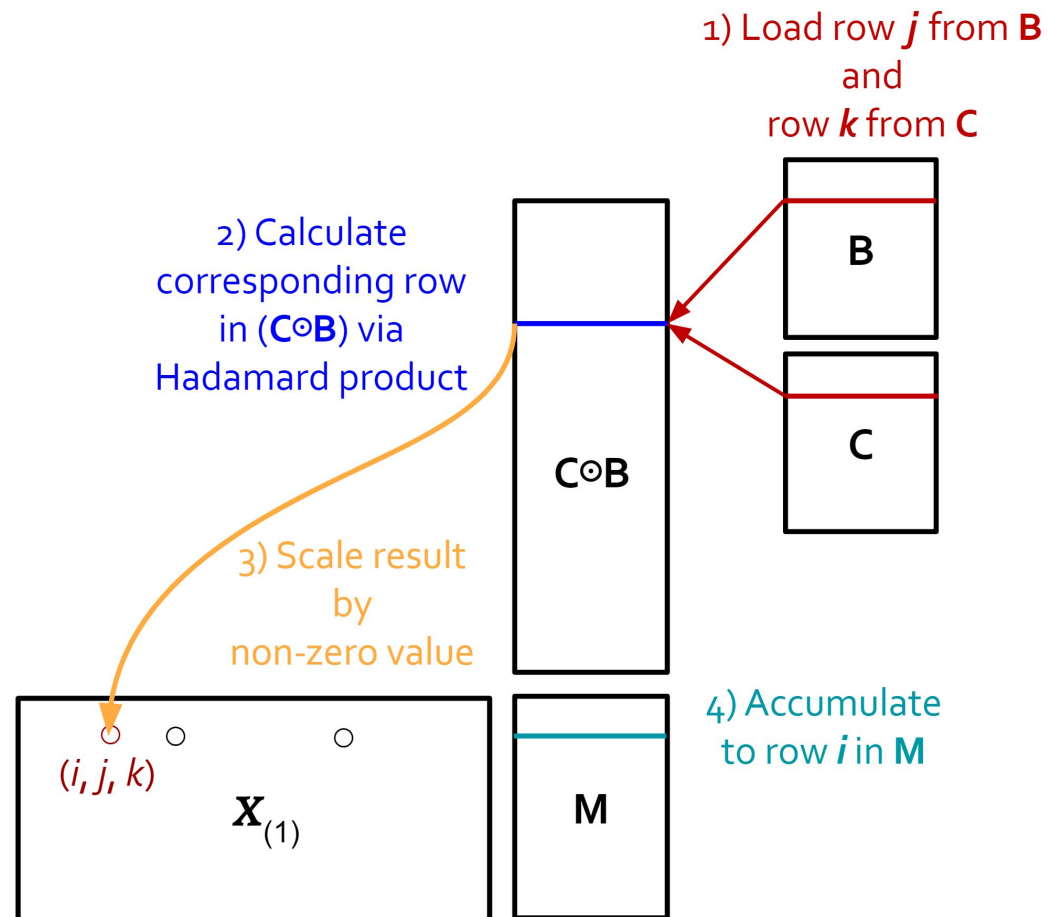
- Bottleneck in well known CPD-ALS algorithm
- $M \leftarrow X_{(1)} C \odot B$





MTTKRP

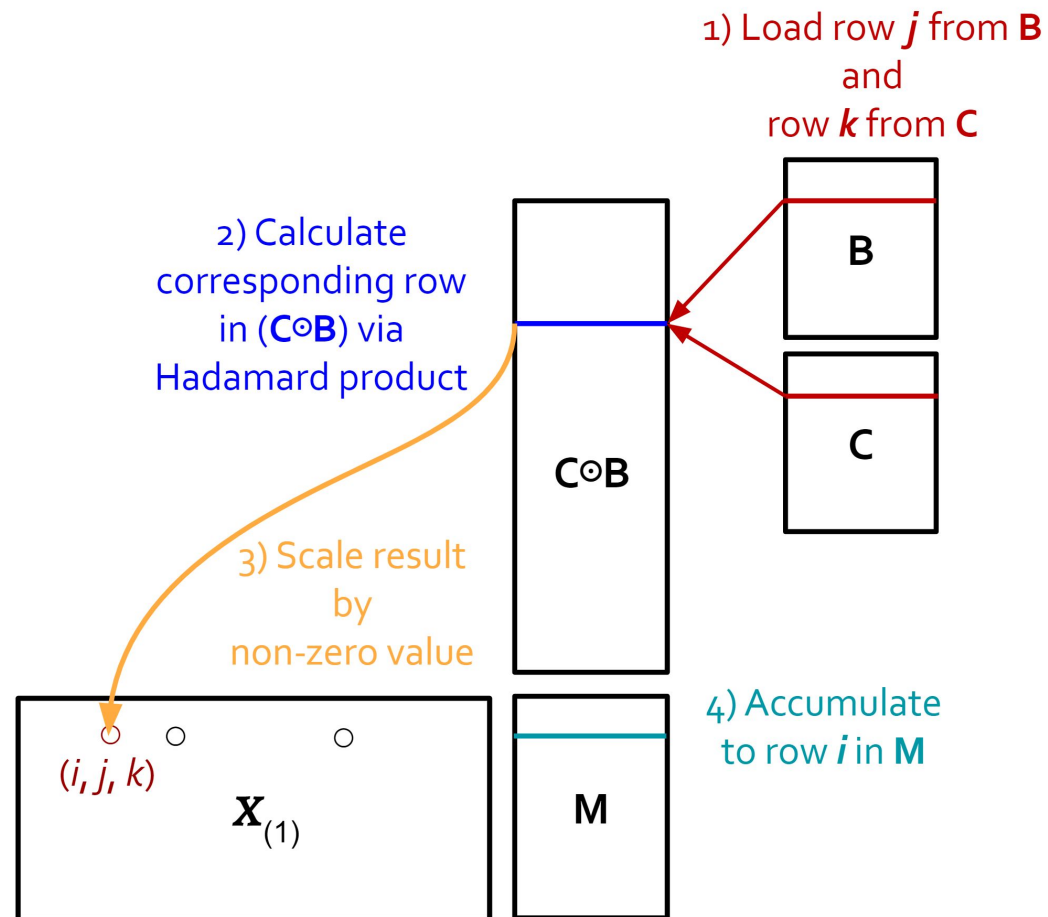
- Challenges in optimizing MTTKRP:





MTTKRP

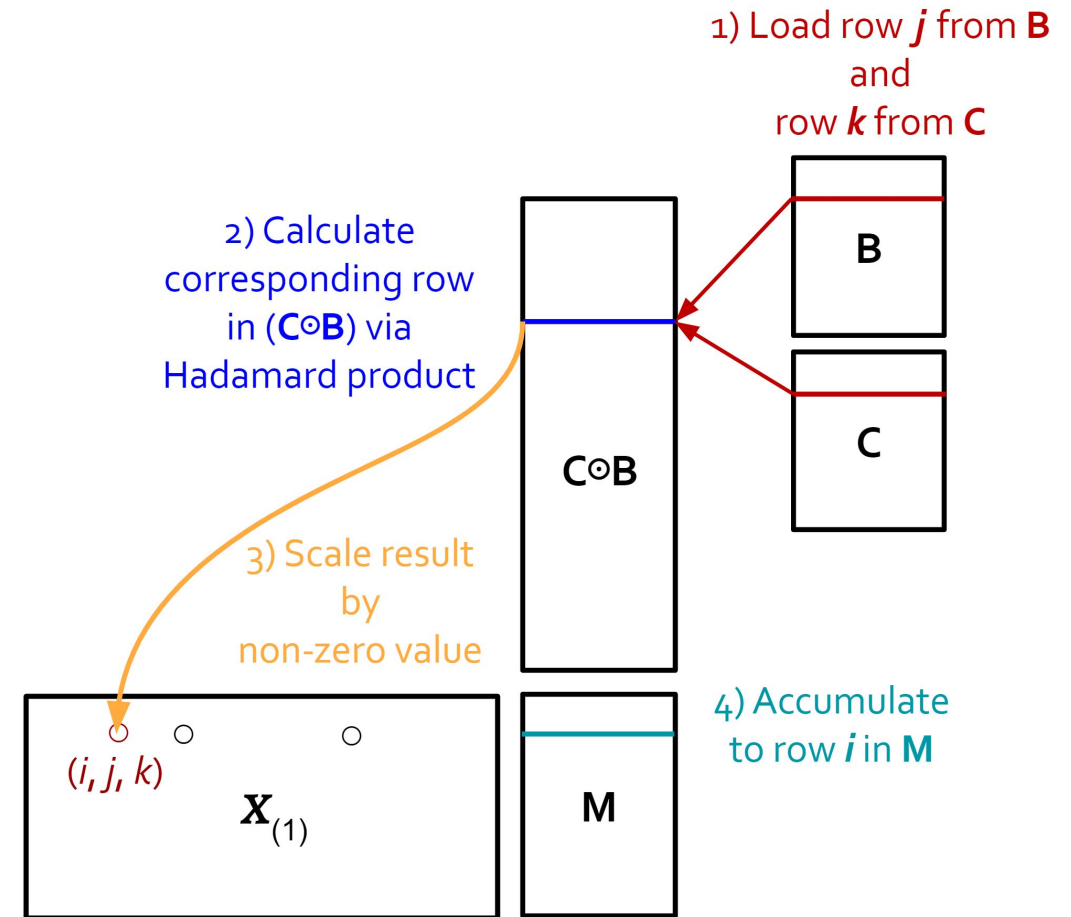
- Challenges in optimizing MTTKRP:
 - **Memory-intensive** kernel





MTTKRP

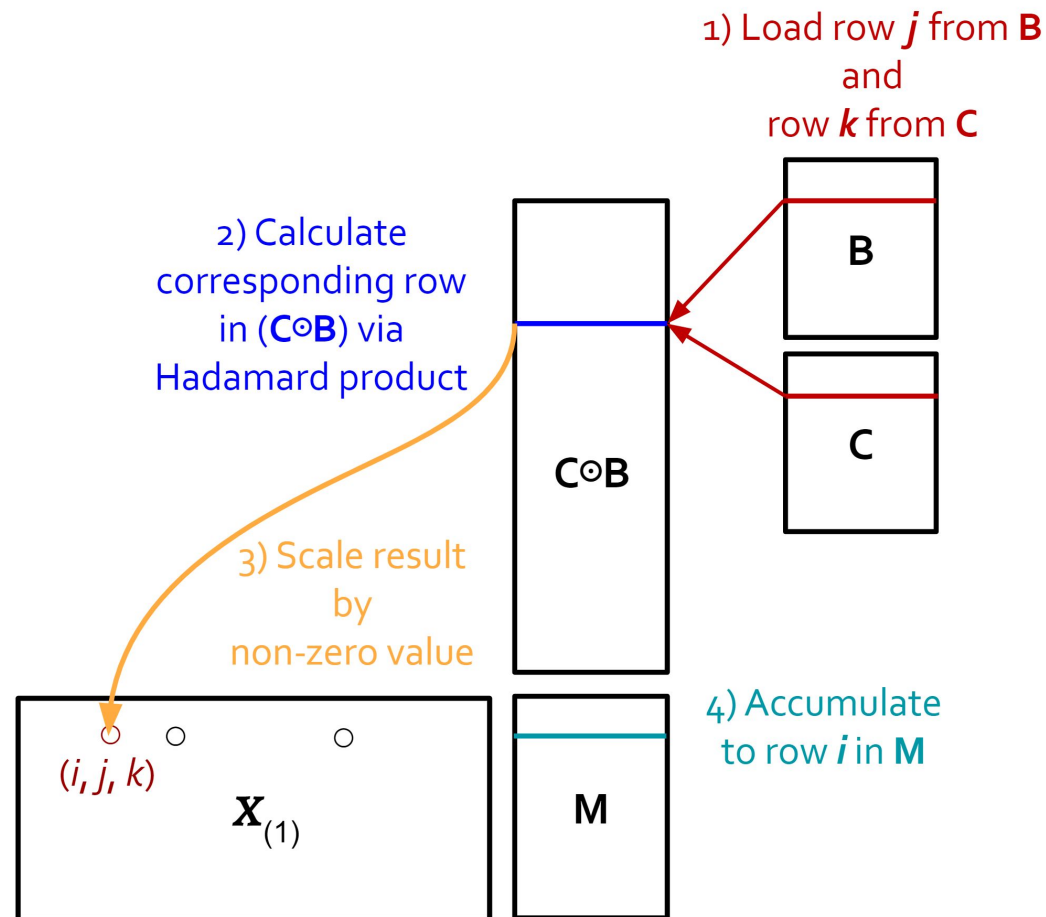
- Challenges in optimizing MTTKRP:
 - **Memory-intensive** kernel
 - **Complex data locality** nature





MTTKRP

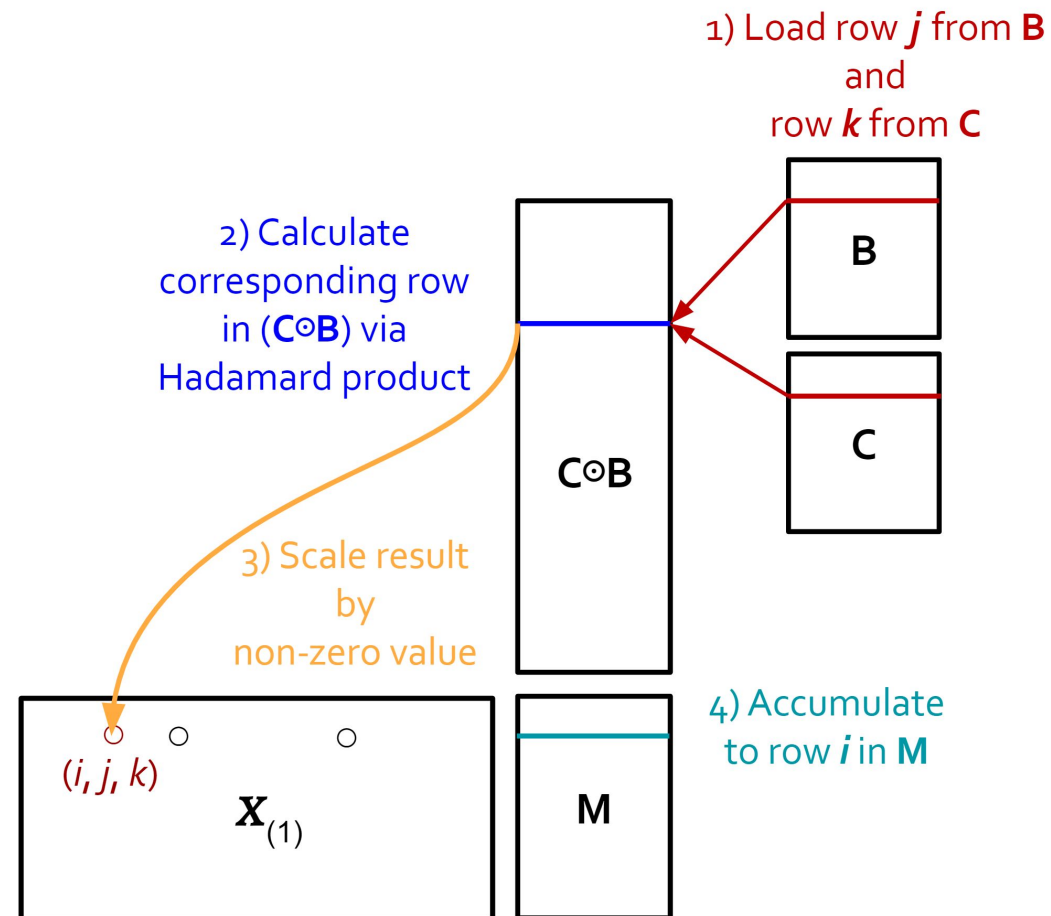
- Challenges in optimizing MTTKRP:
 - **Memory-intensive** kernel
 - **Complex data locality** nature
 - **Synchronization overhead** when updating factor matrices





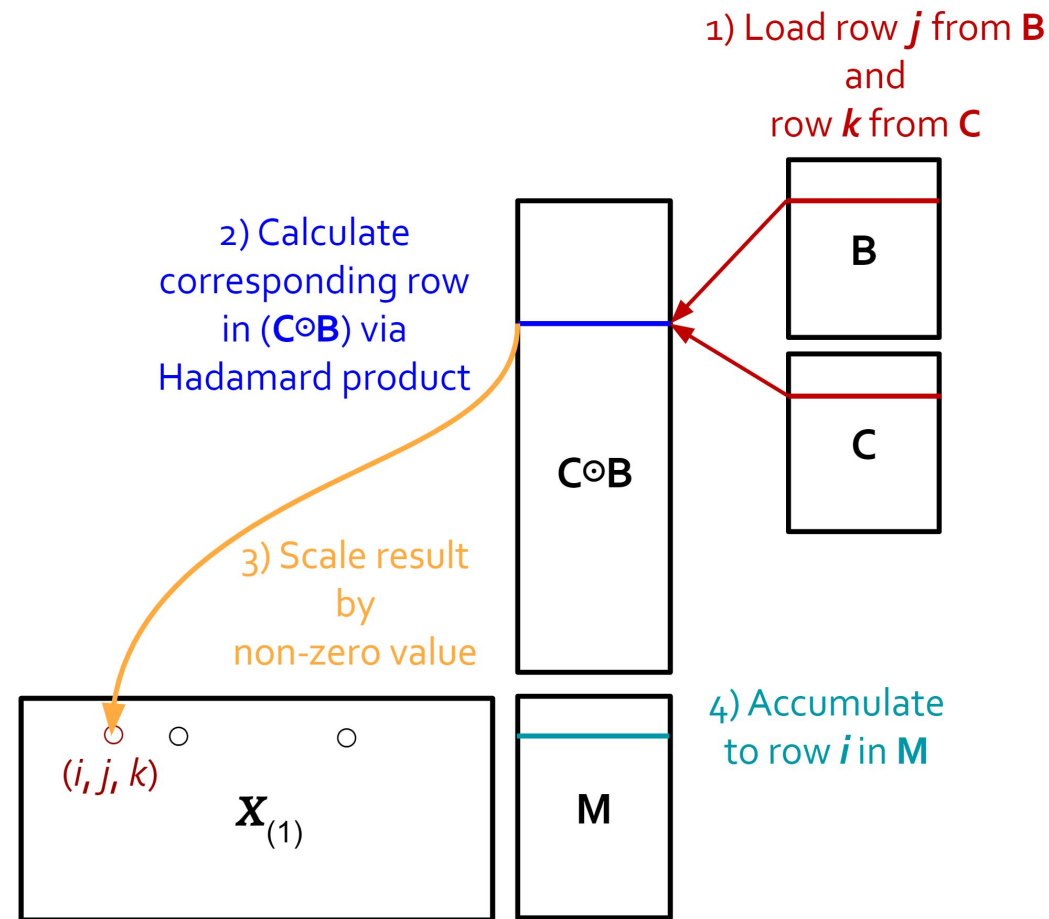
MTTKRP

- Challenges in optimizing MTTKRP:
 - **Memory-intensive** kernel
 - **Complex data locality** nature
 - **Synchronization overhead** when updating factor matrices
 - **Load imbalance** across threads



MTTKRP

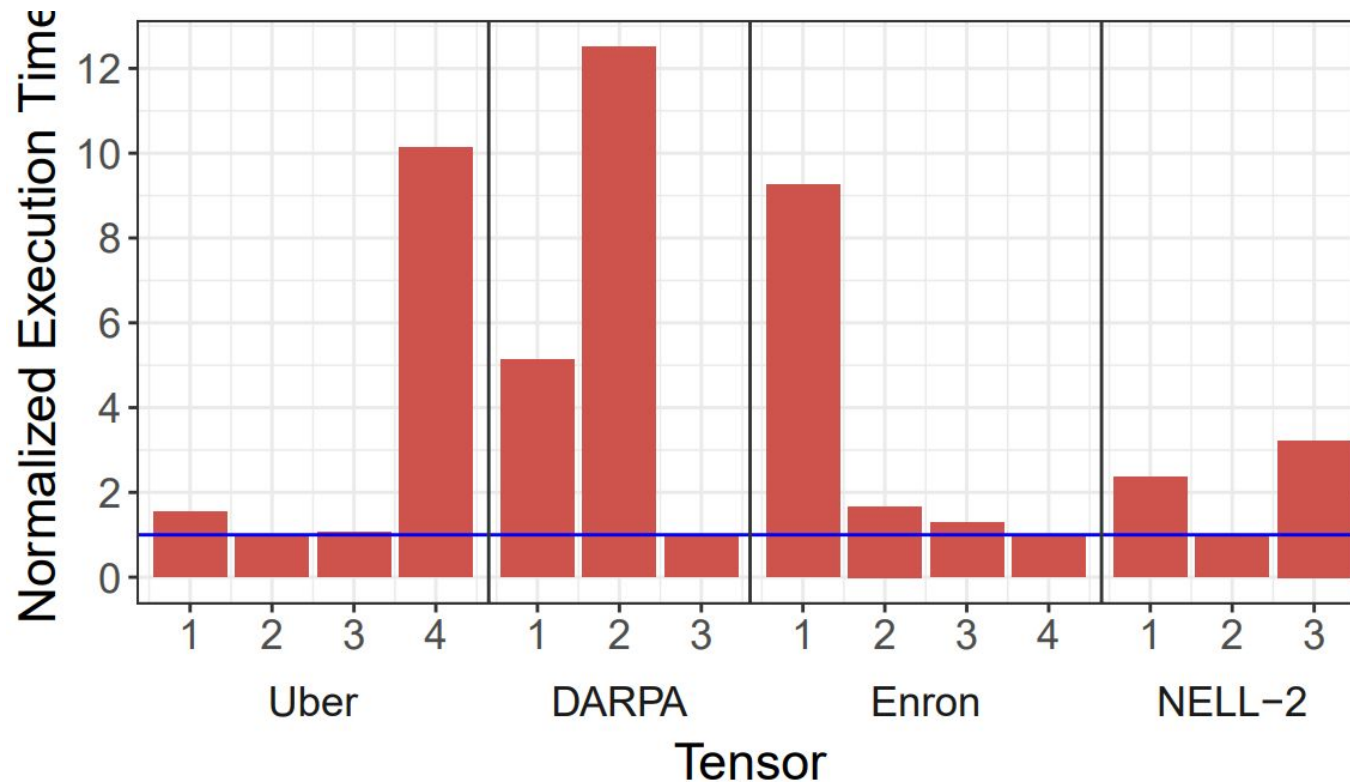
- Challenges in optimizing MTTKRP:
 - **Memory-intensive** kernel
 - **Complex data locality** nature
 - **Synchronization overhead** when updating factor matrices
 - **Load imbalance** across threads
 - **Limited memory** on GPU accelerators





State-of-the-Art Performance Issues

- Setup:
 - Individual mode MTTKRP on MM-CSF framework
 - Normalized against best performing mode
 - Blue line is baseline
 - NVIDIA A100
- Higher is worse

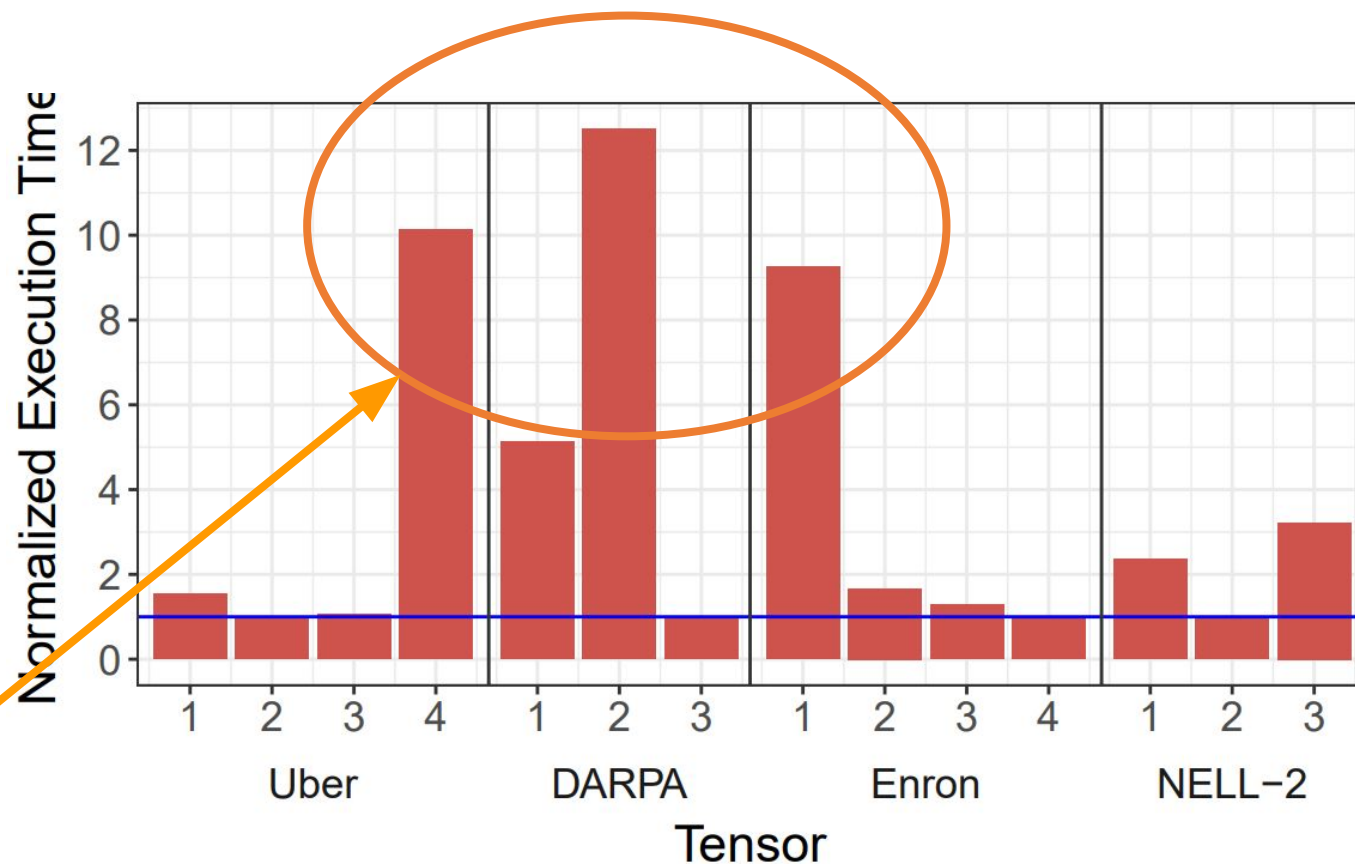




State-of-the-Art Performance Issues

- Setup:
 - Individual mode MTTKRP on MM-CSF framework
 - Normalized against best performing mode
 - Blue line is baseline
 - NVIDIA A100
- Higher is worse

Over an *order of magnitude* worse depending on the mode. All modes require the same number of FLOPs





State-of-the-Art Approaches

Name	Data Format	Load Balance	Synchronization	Mode Orientation	Data Compression	Mode Algorithms
F-COO [Liu et al. 2017]	List-based	Optimal	Reductions in shared memory	Mode-specific	Multiple copies required	1
B-CSF [Nisa et. al, 2019]	Tree-based	Improved	Minimized by data structure	Mode-specific	Multiple copies required	N
MM-CSF [Nisa et. al, 2019]	Tree-based	Improved	Minimized by data structure	Mode-specific	Single copy	N



State-of-the-Art Approaches

Name	Data Format	Load Balance	Synchronization	Mode Orientation	Data Compression	Mode Algorithms
F-COO [Liu et al. 2017]	List-based	Optimal	Reductions in shared memory	Mode-specific	Multiple copies required	1
B-CSF [Nisa et. al, 2019]	Tree-based	Improved	Minimized by data structure	Mode-specific	Multiple copies required	N
MM-CSF [Nisa et. al, 2019]	Tree-based	Improved	Minimized by data structure	Mode-specific	Single copy	N

→ State of the art approaches have tradeoffs that severely impact performance



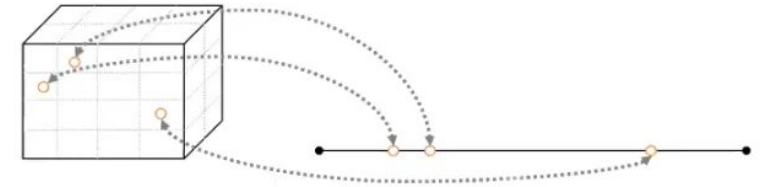
Blocked Linearized CoOrdinates (BLCO)

- A list-based yet highly compressed and massively parallel tensor format



Blocked Linearized CoOrdinates (BLCO)

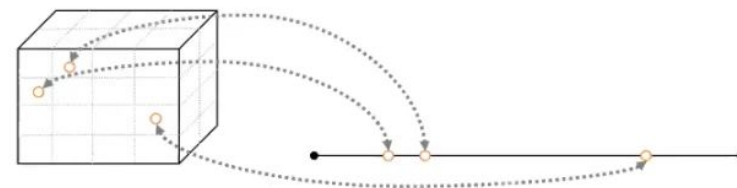
- A list-based yet highly compressed and massively parallel tensor format
- **Index linearization** maps multidimensional coordinates into linear space





Blocked Linearized CoOrdinates (BLCO)

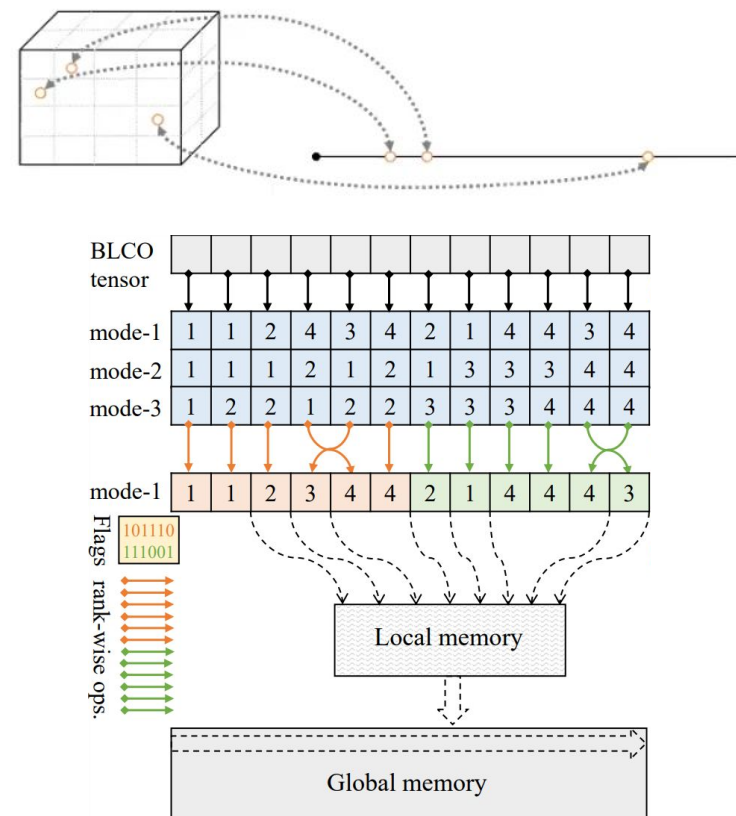
- A list-based yet highly compressed and massively parallel tensor format
- **Index linearization** maps multidimensional coordinates into linear space
- **Adaptive blocking** blocks the tensor to meet target GPU resource constraints





Blocked Linearized CoOrdinates (BLCO)

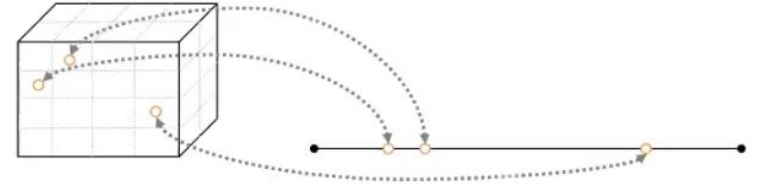
- A list-based yet highly compressed and massively parallel tensor format
- **Index linearization** maps multidimensional coordinates into linear space
- **Adaptive blocking** blocks the tensor to meet target GPU resource constraints
- **Hierarchical conflict resolution** resolves synchronization conflicts at multiple levels of the memory hierarchy





Index Linearization

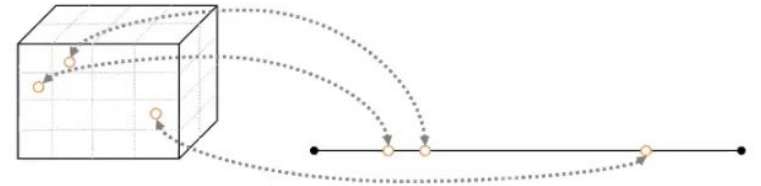
- Element indices are linearized in one-dimensional space





Index Linearization

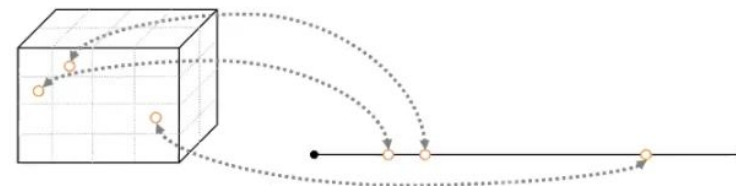
- Element indices are linearized in one-dimensional space
 - Highly compressed representation





Index Linearization

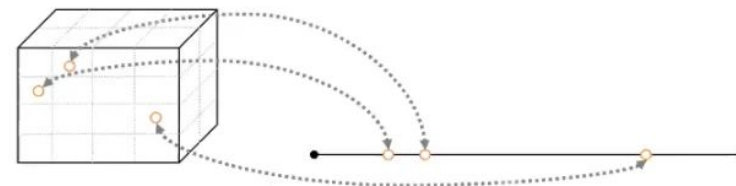
- Element indices are linearized in one-dimensional space
 - Highly compressed representation
 - More effective use of memory bandwidth





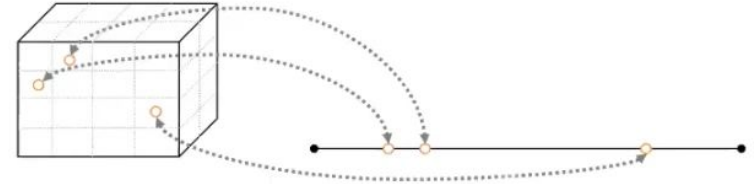
Index Linearization

- Element indices are linearized in one-dimensional space
 - Highly compressed representation
 - More effective use of memory bandwidth
 - Allows for fine-grained parallelism (load balance)



Index Linearization

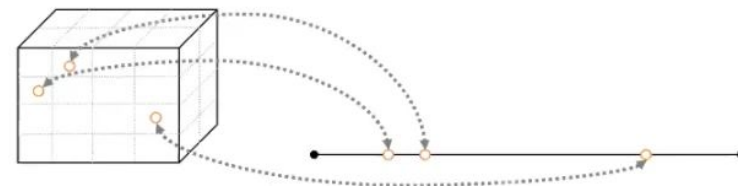
- Element indices are linearized in one-dimensional space
 - Highly compressed representation
 - More effective use of memory bandwidth
 - Allows for fine-grained parallelism (load balance)
- Elements are linearized and ordered based on ALTO [Helal et. al, 2021]



	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	(111110) ₂	7.0
63	(111111) ₂	12.0

Index Linearization

- Element indices are linearized in one-dimensional space
 - Highly compressed representation
 - More effective use of memory bandwidth
 - Allows for fine-grained parallelism (load balance)
- Elements are linearized and ordered based on ALTO [Helal et. al, 2021]
- Indices are then re-linearized
 - Allows for fast delinearization on GPUs



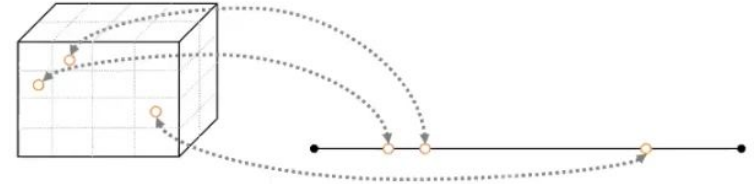
	l	v
0	(00000) ₂	1.0
4	(000100) ₂	2.0
5	(000101) ₂	4.0
10	(001010) ₂	8.0
12	(001100) ₂	6.0
15	(001111) ₂	9.0
33	(100001) ₂	5.0
48	(110000) ₂	3.0
57	(111001) ₂	10.0
61	(111101) ₂	11.0
62	(111110) ₂	7.0
63	(111111) ₂	12.0

	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	(10111) ₂	9.0
1	(00001) ₂	5.0
8	(01000) ₂	3.0
11	(01011) ₂	10.0
27	(11011) ₂	11.0
30	(11110) ₂	7.0
31	(11111) ₂	12.0



Adaptive Blocking

- Linearization may exceed integer width (i.e. 64 bits)

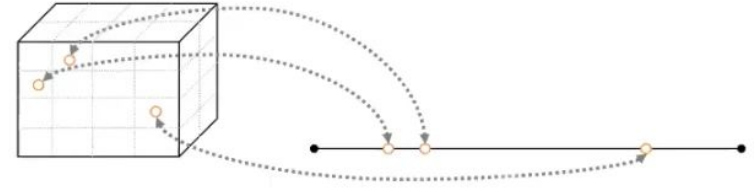


l	v
0 (00000) ₂	1.0
4 (000100) ₂	2.0
5 (000101) ₂	4.0
10 (001010) ₂	8.0
12 (001100) ₂	6.0
15 (001111) ₂	9.0
33 (100001) ₂	5.0
48 (110000) ₂	3.0
57 (111001) ₂	10.0
61 (111101) ₂	11.0
62 (111110) ₂	7.0
63 (111111) ₂	12.0

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 (10111) ₂	9.0
1 (00001) ₂	5.0
8 (01000) ₂	3.0
11 (01011) ₂	10.0
27 (11011) ₂	11.0
30 (11110) ₂	7.0
31 (11111) ₂	12.0

Adaptive Blocking

- Linearization may exceed integer width (i.e. 64 bits)
- Block the tensor based on the excessive bits

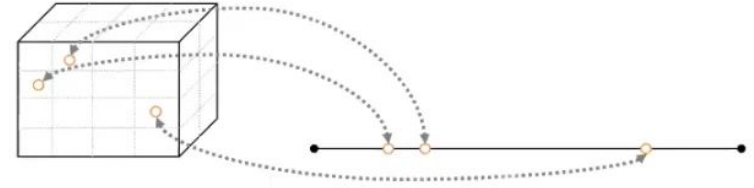


	l	v
0	(00000) ₂	1.0
4	(000100) ₂	2.0
5	(000101) ₂	4.0
10	(001010) ₂	8.0
12	(001100) ₂	6.0
15	(001111) ₂	9.0
33	(100001) ₂	5.0
48	(110000) ₂	3.0
57	(111001) ₂	10.0
61	(111101) ₂	11.0
62	(111110) ₂	7.0
63	(111111) ₂	12.0

b	l	v
0	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 (10111) ₂	9.0
1	1 (00001) ₂	5.0
	8 (01000) ₂	3.0
	11 (01011) ₂	10.0
	27 (11011) ₂	11.0
	30 (11110) ₂	7.0
	31 (11111) ₂	12.0

Adaptive Blocking

- Linearization may exceed integer width (i.e. 64 bits)
- Block the tensor based on the excessive bits
 - Mode-agnostic approach

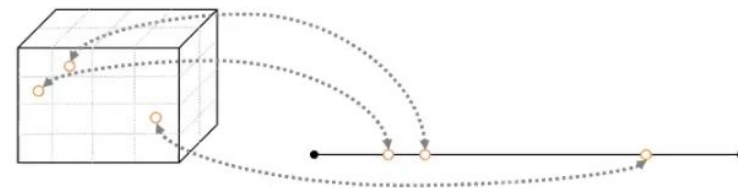


	l	v
0	(00000) ₂	1.0
4	(000100) ₂	2.0
5	(000101) ₂	4.0
10	(001010) ₂	8.0
12	(001100) ₂	6.0
15	(001111) ₂	9.0
33	(100001) ₂	5.0
48	(110000) ₂	3.0
57	(111001) ₂	10.0
61	(111101) ₂	11.0
62	(111110) ₂	7.0
63	(111111) ₂	12.0

b	l	v
0	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 (10111) ₂	9.0
1	1 (00001) ₂	5.0
	8 (01000) ₂	3.0
	11 (01011) ₂	10.0
	27 (11011) ₂	11.0
	30 (11110) ₂	7.0
	31 (11111) ₂	12.0

Adaptive Blocking

- Linearization may exceed integer width (i.e. 64 bits)
- Block the tensor based on the excessive bits
 - Mode-agnostic approach
 - Adds minimal memory overhead

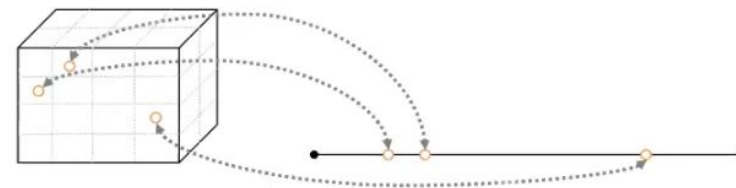


	l	v
0	(00000) ₂	1.0
4	(000100) ₂	2.0
5	(000101) ₂	4.0
10	(001010) ₂	8.0
12	(001100) ₂	6.0
15	(001111) ₂	9.0
33	(100001) ₂	5.0
48	(110000) ₂	3.0
57	(111001) ₂	10.0
61	(111101) ₂	11.0
62	(111110) ₂	7.0
63	(111111) ₂	12.0

b	l	v
0	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 (10111) ₂	9.0
1	1 (00001) ₂	5.0
	8 (01000) ₂	3.0
	11 (01011) ₂	10.0
	27 (11011) ₂	11.0
	30 (11110) ₂	7.0
	31 (11111) ₂	12.0

Adaptive Blocking

- Linearization may exceed integer width (i.e. 64 bits)
- Block the tensor based on the excessive bits
 - Mode-agnostic approach
 - Adds minimal memory overhead
 - Blocks can be streamed for *out-of-memory computation*

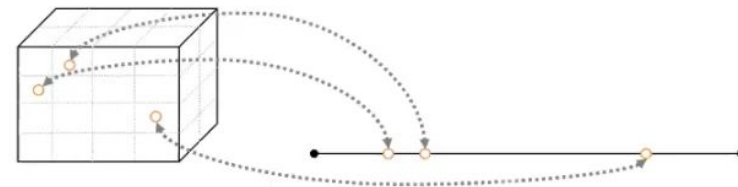


	l	v
0	(00000) ₂	1.0
4	(000100) ₂	2.0
5	(000101) ₂	4.0
10	(001010) ₂	8.0
12	(001100) ₂	6.0
15	(001111) ₂	9.0
33	(100001) ₂	5.0
48	(110000) ₂	3.0
57	(111001) ₂	10.0
61	(111101) ₂	11.0
62	(111110) ₂	7.0
63	(111111) ₂	12.0

b	l	v
0	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 (10111) ₂	9.0
1	1 (00001) ₂	5.0
	8 (01000) ₂	3.0
	11 (01011) ₂	10.0
	27 (11011) ₂	11.0
	30 (11110) ₂	7.0
	31 (11111) ₂	12.0

Adaptive Blocking

- Linearization may exceed integer width (i.e. 64 bits)
- Block the tensor based on the excessive bits
 - Mode-agnostic approach
 - Adds minimal memory overhead
 - Blocks can be streamed for *out-of-memory computation*
 - Leverage list-based storage to expose fine grained parallel processing...



	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	(111110) ₂	7.0
63	(111111) ₂	12.0

b	l	v
0	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 (10111) ₂	9.0
1	1 (00001) ₂	5.0
	8 (01000) ₂	3.0
	11 (01011) ₂	10.0
	27 (11011) ₂	11.0
	30 (11110) ₂	7.0
	31 (11111) ₂	12.0



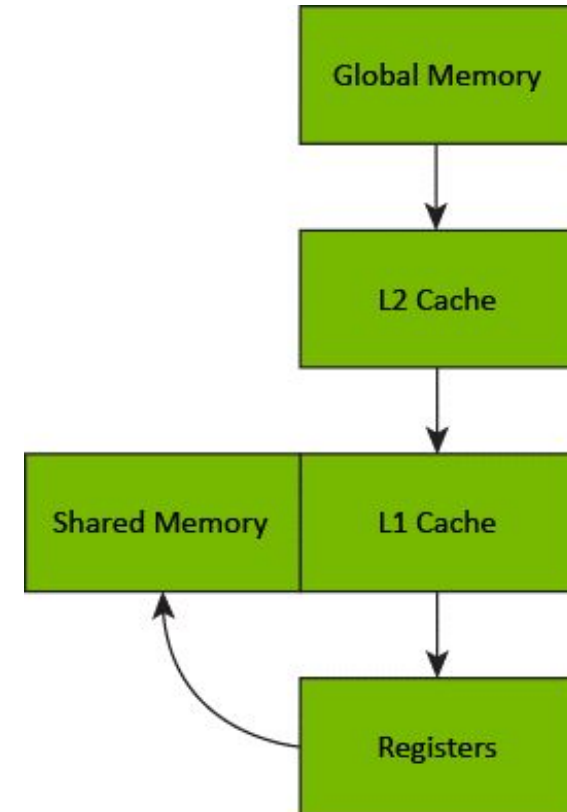
Hierarchical Conflict Resolution

- Synchronization overhead is expensive for list-based formats
 - Especially true on GPUs



Hierarchical Conflict Resolution

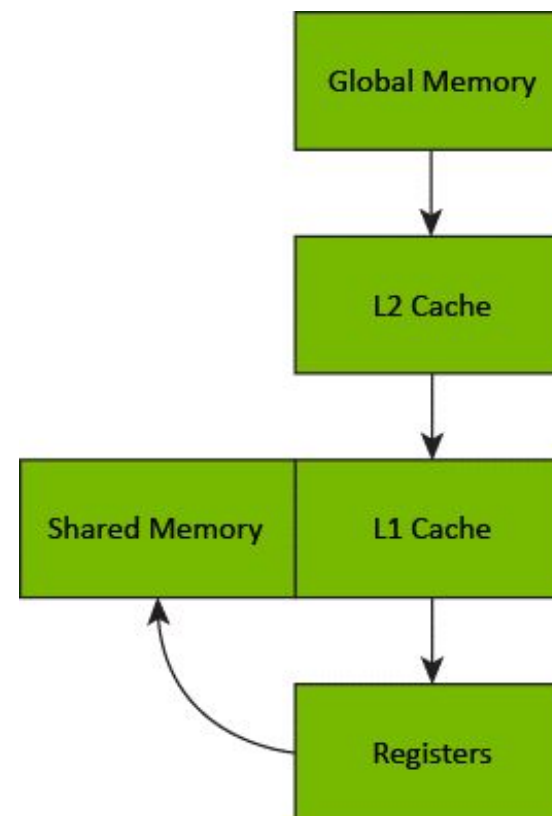
- Synchronization overhead is expensive for list-based formats
 - Especially true on GPUs
- Leverage multiple memory levels hierarchically to minimize atomics





Hierarchical Conflict Resolution

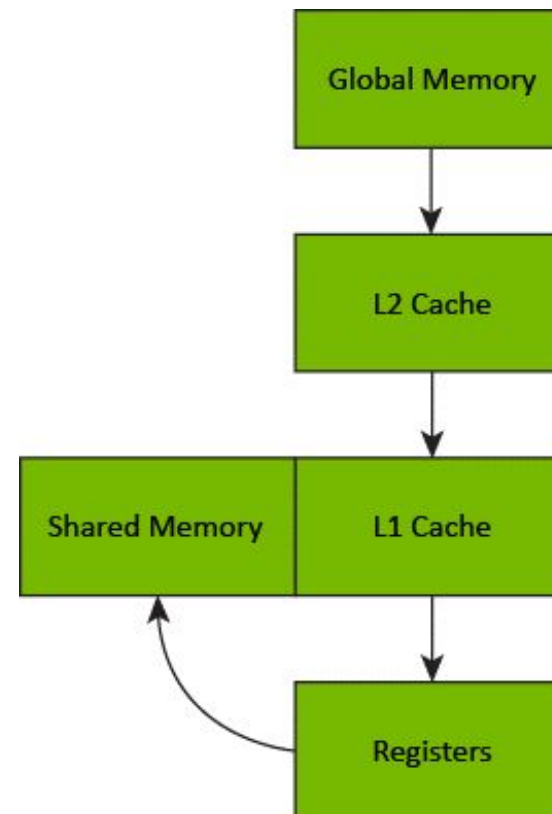
- Synchronization overhead is expensive for list-based formats
 - Especially true on GPUs
- Leverage multiple memory levels hierarchically to minimize atomics
 - Register memory: segmented scan





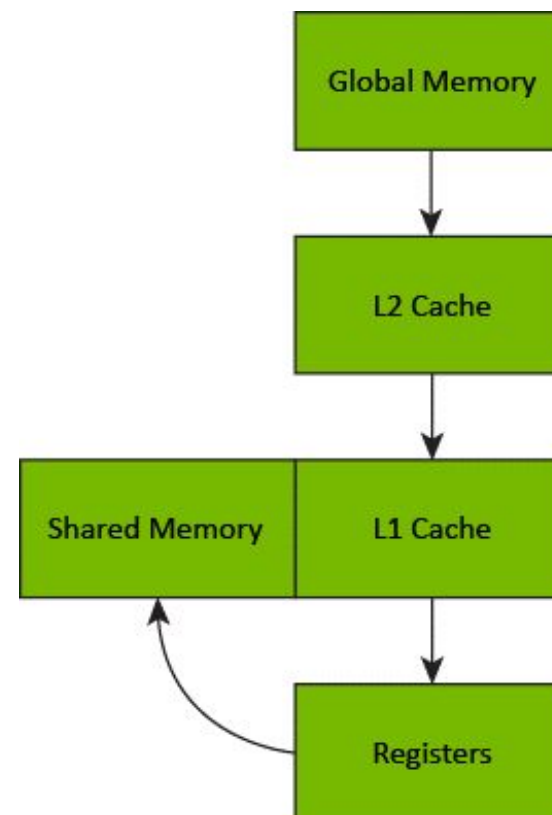
Hierarchical Conflict Resolution

- Synchronization overhead is expensive for list-based formats
 - Especially true on GPUs
- Leverage multiple memory levels hierarchically to minimize atomics
 - Register memory: segmented scan
 - Shared memory: software cache



Hierarchical Conflict Resolution

- Synchronization overhead is expensive for list-based formats
 - Especially true on GPUs
- Leverage multiple memory levels hierarchically to minimize atomics
 - Register memory: segmented scan
 - Shared memory: software cache
 - Global memory: pull-based reduction





State-of-the-Art Approaches

Name	Data Format	Load Balance	Synchronization	Mode Orientation	Data Compression	Mode Algorithms
F-COO [Liu et al. 2017]	List-based	Optimal	Reductions in shared memory	Mode-specific	Multiple copies required	1
B-CSF [Nisa et. al, 2019]	Tree-based	Improved	Minimized by data structure	Mode-specific	Multiple copies required	N
MM-CSF [Nisa et. al, 2019]	Tree-based	Improved	Minimized by data structure	Mode-specific	Single copy	N
BLCO [this work]	List-based	Optimal	Reductions at multiple memory levels	Mode-agnostic	Single copy, streamable	1



Evaluation Platform

GPUS:

- NVIDIA A100
- NVIDIA V100

Considered frameworks:

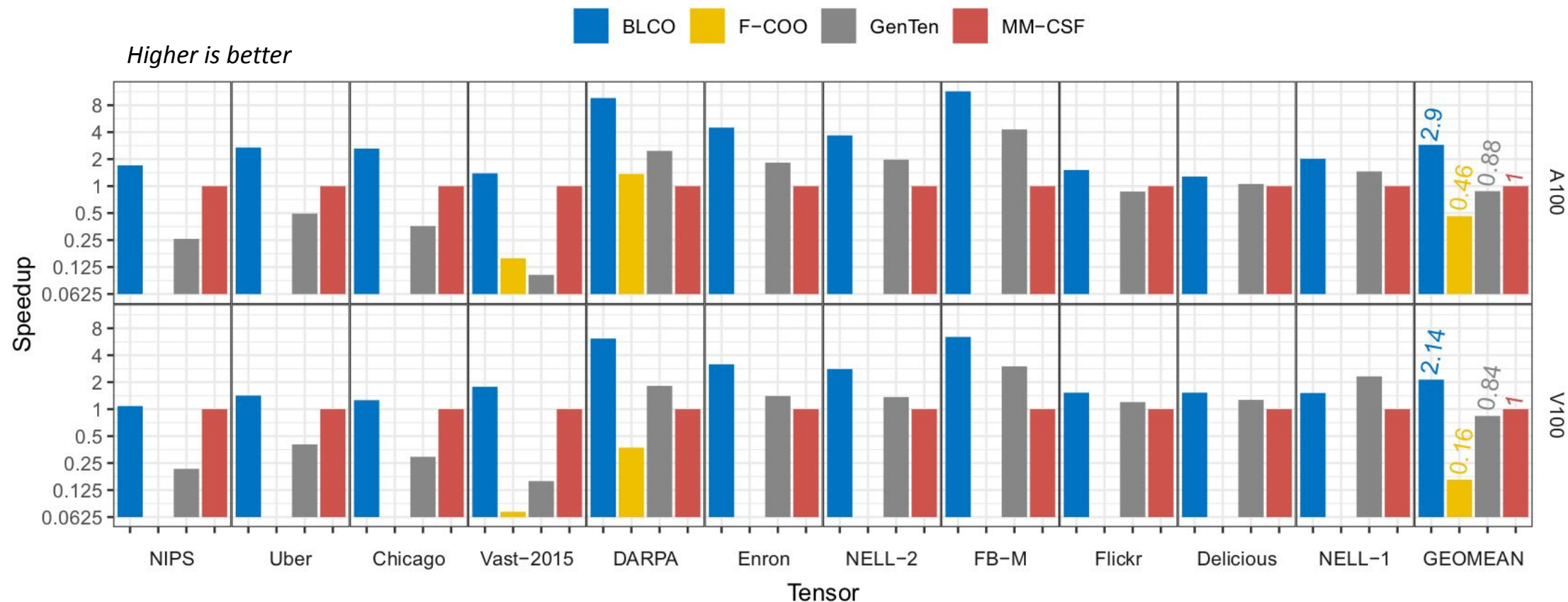
- MM-CSF
- F-COO
- Genten
- BLCO

Benchmark datasets: 11 real-world tensors from FROSTT

Compilation: double-precision / 64-bit data types



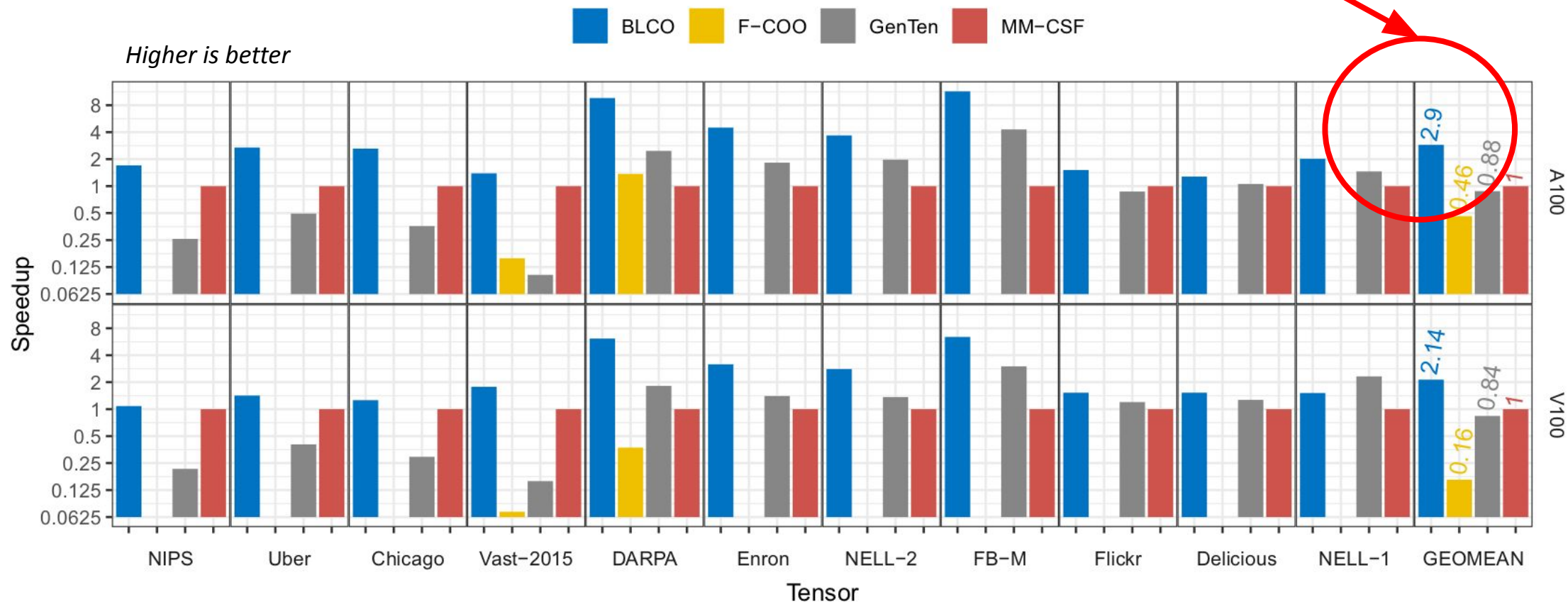
Framework Comparison





Framework Comparison

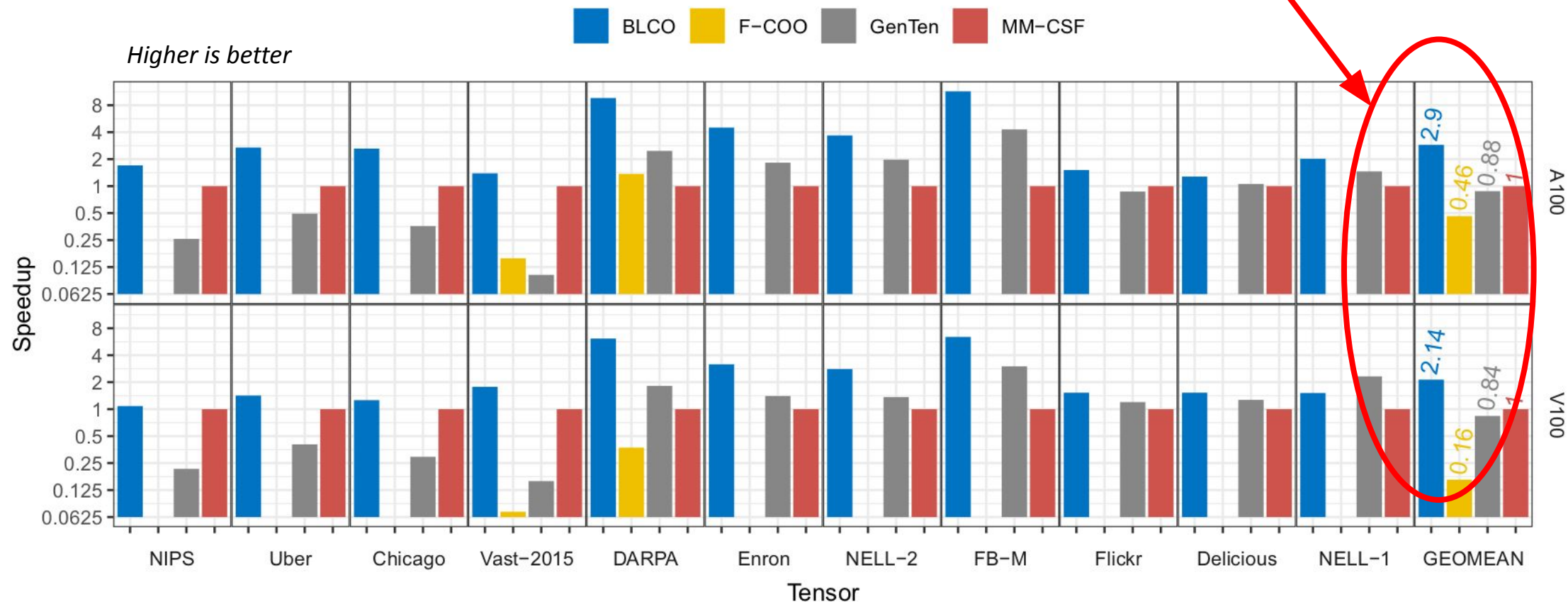
Up to 2.90x speedup over prior state of the art





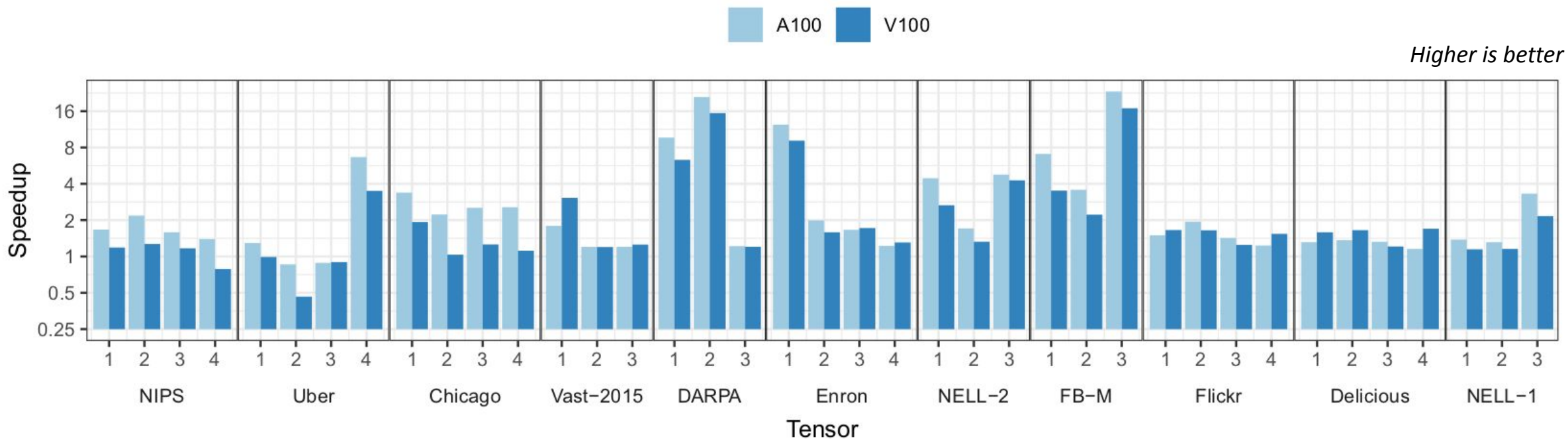
Framework Comparison

Speedup observed across different architectures





State-of-the-Art Comparison



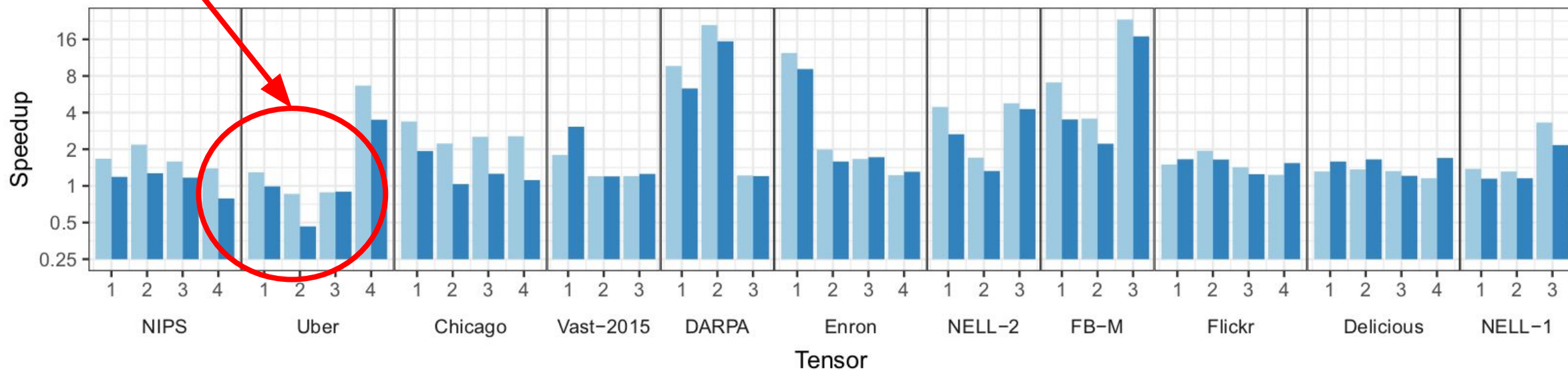


State-of-the-Art Comparison

Worse performance on the smallest of tensors

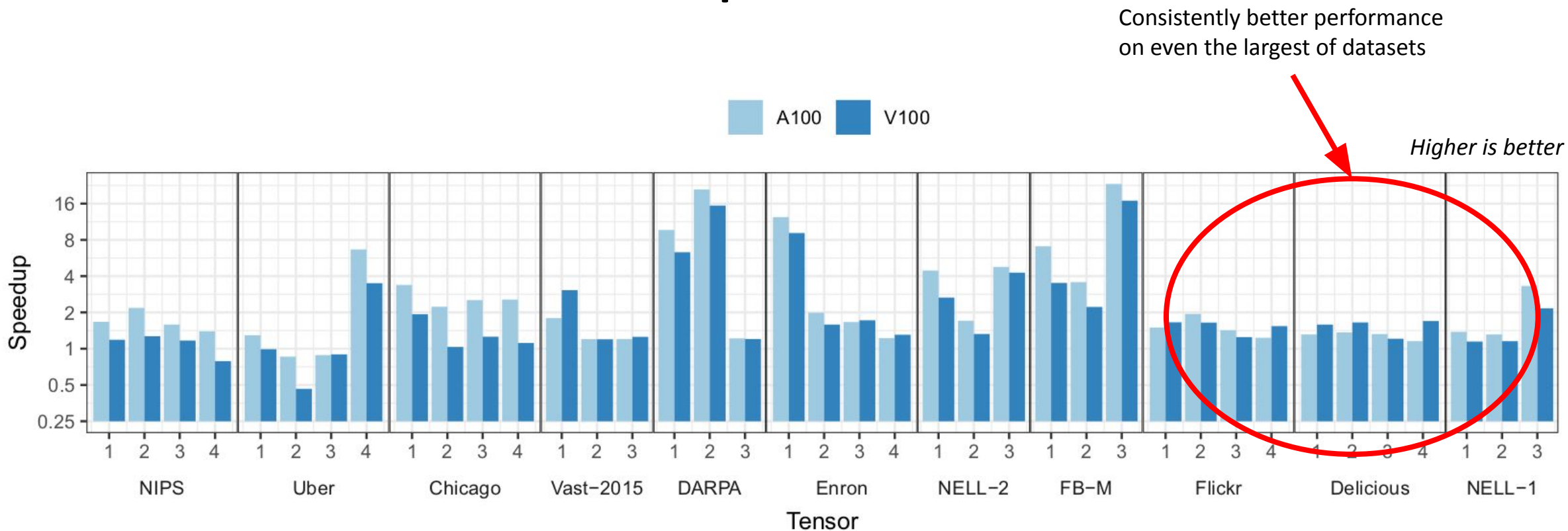
A100 V100

Higher is better





State-of-the-Art Comparison





Memory Traffic Analysis

Data Set	Format	n	Vol ¹	TP ²	Data Set	Format	n	Vol ¹	TP ²
Uber	BLCO	1	2.78	3.60	Enron	BLCO	1	44.82	4.11
		2	2.75	3.61			2	46.23	4.62
		3	2.75	3.53			3	47.88	4.92
		4	2.73	2.77			4	47.22	4.70
	MM-CSF	1	1.68	1.68		MM-CSF	1	41.39	0.31
		2	1.33	2.03			2	62.83	3.16
		3	1.33	1.93			3	37.15	2.29
		4	2.12	0.32			4	37.05	3.01
Vast-2015	BLCO	1	16.91	3.92	NELL-1	BLCO	1	107.5	2.44
		2	16.73	3.77			2	104.5	2.32
		3	13.92	2.90			3	110.7	2.39
	MM-CSF	1	9.19	1.19		MM-CSF	1	123.1	2.21
		2	8.36	1.57			2	118.5	2.19
		3	8.36	1.45			3	122.1	0.86

¹ Memory volume in GB, measured by *l1tex__t_bytes.sum* in Nsight Compute [3]

² Memory throughput in TB/s, calculated by (Vol / total execution time)



Memory Traffic Analysis

MM-CSF
achieves lower
memory traffic
by being more
compressed
overall

Data Set	Format	n	Vol ¹	TP ²	Data Set	Format	n	Vol ¹	TP ²
Uber	BLCO	1	2.78	3.60	Enron	BLCO	1	44.82	4.11
		2	2.75	3.61			2	46.23	4.62
		3	2.75	3.53			3	47.88	4.92
		4	2.73	2.77			4	47.22	4.70
	MM-CSF	1	1.68	1.68		MM-CSF	1	41.39	0.31
		2	1.33	2.03			2	62.83	3.16
		3	1.33	1.93			3	37.15	2.29
		4	2.12	0.32			4	37.05	3.01
Vast-2015	BLCO	1	16.91	3.92	NELL-1	BLCO	1	107.5	2.44
		2	16.73	3.77			2	104.5	2.32
		3	13.92	2.90			3	110.7	2.39
	MM-CSF	1	9.19	1.19		MM-CSF	1	123.1	2.21
		2	8.36	1.57			2	118.5	2.19
		3	8.36	1.45			3	122.1	0.86

¹ Memory volume in GB, measured by *l1tex_t_bytes.sum* in Nsight Compute [3]

² Memory throughput in TB/s, calculated by (Vol / total execution time)



Memory Traffic Analysis

BLCO achieves higher memory throughput with hierarchical conflict resolution (using every level of memory)

Data Set	Format	n	Vol ¹	TP ²	Data Set	Format	n	Vol ¹	TP ²
Uber	BLCO	1	2.78	3.60	Enron	BLCO	1	44.82	4.11
		2	2.75	3.61			2	46.23	4.62
		3	2.75	3.53			3	47.88	4.92
		4	2.73	2.77			4	47.22	4.70
	MM-CSF	1	1.68	1.68		MM-CSF	1	41.39	0.31
		2	1.33	2.03			2	62.83	3.16
		3	1.33	1.93			3	37.15	2.29
		4	2.12	0.32			4	37.05	3.01
Vast-2015	BLCO	1	16.91	3.92	NELL-1	BLCO	1	107.5	2.44
		2	16.73	3.77			2	104.5	2.32
		3	13.92	2.90			3	110.7	2.39
	MM-CSF	1	9.19	1.19		MM-CSF	1	123.1	2.21
		2	8.36	1.57			2	118.5	2.19
		3	8.36	1.45			3	122.1	0.86

¹ Memory volume in GB, measured by *l1tex_t_bytes.sum* in Nsight Compute [3]

² Memory throughput in TB/s, calculated by (Vol / total execution time)



Memory Traffic Analysis

Data Set	Format	n	Vol ¹	TP ²	Data Set	Format	n	Vol ¹	TP ²
Uber	BLCO	1	2.78	3.60	Enron	BLCO	1	44.82	4.11
		2	2.75	3.61			2	46.23	4.62
		3	2.75	3.53			3	47.88	4.92
		4	2.73	2.77			4	47.22	4.70
	MM-CSF	1	1.68	1.68		MM-CSF	1	41.39	0.31
		2	1.33	2.03			2	62.83	3.16
		3	1.33	1.93			3	37.15	2.29
		4	2.12	0.32			4	37.05	3.01
Vast-2015	BLCO	1	16.91	3.92	NELL-1	BLCO	1	107.5	2.44
		2	16.73	3.77			2	104.5	2.32
		3	13.92	2.90			3	110.7	2.39
	MM-CSF	1	9.19	1.19		MM-CSF	1	123.1	2.21
		2	8.36	1.57			2	118.5	2.19
		3	8.36	1.45			3	122.1	0.86

BLCO also has less performance irregularity across the different modes

¹ Memory volume in GB, measured by `l1tex_t_bytes.sum` in Nsight Compute [3]

² Memory throughput in TB/s, calculated by $(\text{Vol} / \text{total execution time})$



Closing Thoughts

- Tensors are becoming prevalent for data analysis
- Efficient tensor formats and algorithms are needed to maximize hardware utilization and throughput
- Our BLCO format outperforms state-of-the-art by up to 33.35x (avg 2.90x) through novel data compression and massively parallel computation
- Future work includes exploring heterogeneous shared-memory computation and extension to distributed-memory platforms



Thank you!