## **Dissertation Defense**

**Distributed Sparse** Computing and Communication for **Big Graph Analytics** and Deep Learning

by

#### Mohammad Hasanzadeh Mofrad

Date: Friday, October 30, 2020 Time: 9:30 AM EST

Committee:

Dr. Rami Melhem (advisor), Univ. of PittsburghDr. Alexandros Labrinidis, Univ. of PittsburghDr. John Lange, Univ. of PittsburghDr. Balaji Palanisamy, Univ. of PittsburghDr. Mohammad Hammoud, CMU-Qatar



#### Overview: Distribute Sparse Computing and Communication



### Motivation: Graph Analytics

- **Big data** is a multi-billion-dollar industry
  - \$230 Billion by 2025 www.prnewswire.com
- Relational data represented by graphs is dominating the big data industry
- Graphs are ubiquitous
  - Google PageRank
    - Internet graph
  - Uber routing's engine
    - Commute graph
  - Amazon item recommendation
    - Bipartite items/users' graph
  - Facebook friend suggestion
    - Relationship graph
- Majority of these graphs are sparse



#### Background: Duality Between Graphs & Sparse Matrices

- The relational data represented by graphs is often hyper-sparse
- Linear algebra primitives backed by sparse matrix formats are new alternatives to their graph-based versions
- 1. Compressed Sparse Column (CSC)
- 2. Compressed Sparse Row (CSR)
- 3. Doubly Compressed Sparse Column (DCSC)
- 4. Doubly Compressed Sparse Column (DCSR)
- 5. Gather, Apply, and Scatter (GAS) abstraction
- 6. Sparse Matrix Vector (SpMV)

0.1 0.2 0.2 0.4 0.9 0.3 0.9 0.5	Src       Dst       Wgt         1       1       0.1         1       3       0.2         1       4       0.4         2       3       0.3         2       4       0.5         4       1       0.9         5       1       0.3         5       4       0.8	0       1       2       3       4       5         0       -       -       -       -         1       .1       .2       .4       -         2       -       -       -       -         3       -       -       -       -         4       .9       -       -       -         5       .3       -       .8       -
Graph <b>G</b> ( <i>n</i> vertice	es) Adjacency List	Adjacency Matrix <b>A</b> <sub>nxn</sub>
	Graph Theory	Linear Algebra
Representation	Graph	Adjacency Matrix
Data structure	Adjacency List	CSC <sup>1</sup> , CSR <sup>2</sup> , DCSC <sup>3</sup> , DCSR <sup>4</sup>
Computing Model	Fan-in/ Fan-out Operations	GAS <sup>5</sup>
Primitive		SpMV <sup>6</sup>
Scalability	Graph Partitioning	Sparse Matrix Partitioning $4$

### State-of-the-art Comparison: Graph Analytics

System	Computing Model	Platform	Parallelism Model	Computation & Communication	Matrix Compression	Authors
GraphPad	Linear Algebra	HPC/Cloud	MPI+OpenMP		DCSC	2016, Intel
Gemini	Frontier-based	HPC	MPI+OpenMP	NUMA-aware CPU	CSR/CSC ( <mark>Dual</mark> )	2016, Tsinghua University
CombBLAS	Linear Algebra	HPC	MPI+OpenMP		DCSC	2018, LBNL
LA3	Linear Algebra	Cloud	MPI+OpenMP	Optimized for cloud communication	DCSC	2018, CMUQ
LAGraph	Graph Theory	HPC	OpenMP		CSR	2020, Texas A&M

• State-of-the-art big graph analytics systems are not 100% mature yet!

### Motivation: Deep Learning

- Deep Neural Networks (DNNs) are pervasive
  - Speech processing
  - Post Translation
  - Autonomous driving
  - Knowledge Discovery



- The core kernel behind inference/training of DNNs is **Dense matrix-matrix** multiplication ( $C_{m \times n} = A_{m \times n} \times B_{n \times n}$ )
- Sparse DNNs are new alternative to dense DNNs with
  - Less time and space complexities
  - Better or comparable accuracy
  - Avoids overfitting
- Sparse DNNs core kernel is Sparse Matrix-matrix Multiplication (SpMM)

# Background: Duality Between Linear Algebra & Deep Learning

- Neural networks connections
  - Triplet format graph theory: (*i*, *j*, *w*<sub>*ij*</sub>)
- Primitive: Matrix-matrix multiplication
  - $C = A \times B$
- Iterative matrix-matrix multiplication for DNN
  - $C_{l+1} = h(A_l \times B_l + b)$
- For Sparse neural networks, the primitive is
  - **SpMM** (Sparse Matrix-matrix multiplication)



#### State-of-the-art Comparison: Deep Learning

System	Computing Model	Platform	Parallelism Model	Computation & Communication	Matrix Compression	Authors
TensorFlow	Linear Algebra	Commodity/ HPC/ <mark>Cloud</mark>	Data parallelism (Keras + CPU/GPU/TPU)	Horovod	Dense	2015 – present, Google
PyTorch	Linear Algebra	Commodity/ HPC	Data parallelism (MPI/ Gloo + OpenMP/ CUDA)	Horovod	Dense	2016 – present, Facebook
MXNet	Linear Algebra	Commodity/ HPC	Data parallelism (MPI4Py/ Gloo + OpenMP/ CUDA)	Horovod/ TCP/IP socket	Dense	2020 – present, Amazon
DeepSpeed	Linear Algebra	Commodity/ HPC/ Cloud	Data parallelism (MPI/ NCCL + X/ CUDA)	Communication Overlapping	Sparse	2019 – present, Microsoft

• Deep Learning frameworks are still evolving ...



#### Scalability

			<ul><li>(1) A distributed sparse data structure &amp; algorithm</li></ul>
Distributed	GraphTap, Cluster 2019 (1) Graphite, VLDB 2020 (2)	DistSpDNN, HPEC 2020 (3) DistSpDNN, HPEC 2020 (4)	(2) A new parallelism model suitable for comp./comm intensive analytics
			(3) Hashing for accelerating DNN inference
Multicore	GraphTap, Cluster 2019 (1)	SpDNN, HPEC 2019 (3)	(4) A new parallelism for mitigating the straggler effect
	<i>SpMV</i> (Graph Analytics)	<i>SpMM</i> (Deep Learning)	Operator (application)



	SpMV (Graph Analytics)	SpMM (Deep Learning)	
Multicore	GraphTap, Cluster 2019	SpDNN, HPEC 2019	• Operator (application)
Distributed	GraphTap, Cluster 2019 Graphite, VLDB 2020	DistSpDNN, HPEC 2020 DistSpDNN, HPEC 2020	

## Compressed Sparse Column (CSC)





#### Doubly Compressed Sparse Column (DCSC)





DCSC SpMV Diagram ( $y = \overline{A} \oplus . \otimes x$ )

X

0

1

2

3

Δ

5

*n* x 1

**DCSC Data Structures** 

#### Triply Compressed Sparse Column (TCSC)

CSC space requirement is n + 2nnz + 1 where
nnz is the number of nonzero entries
+ Sequential column-major access (No indirection for SpMV)
- Length of JA and, x and y vectors equals n



# GraphTap: Distributed SpMSpV<sup>2</sup> Using TCSC for Graph Analytics

- Graph computing model: GAS
  - Gather (input vector), Apply (SpMspV<sup>2</sup>), and Scatter (output vector)
- Scalability: 2D matrix partitioning
  - Tile size translates into computation volume
  - Vector size translates into communication volume
- Compression: TCSC
  - TCSC's less indirections means less computation time
  - TCSC's smaller vectors means less accumulation time
- Communication:
  - TCSC's smaller vectors means less communication time



#### Experiments

Graph Processing Systems						Dataset	t	
	System	Сс	Compression Primitive		Graph	V	<i>E</i>	Туре
Graphite		TCS	5C	SpMSpV <sup>2</sup> GAS	UK'05 (UK5)	39.4 M	0.93 B	Web
LA3 (VLD	B, 2018)	DCS	SC	SpMV GAS	IT'04 (IT4)	41.2 M	1.15 B	Web
GraphPad (IPDPS, 2016)		CSC		SpMV GAS	Twitter (TWT)	41.6 M	1.46 B	Social
					GSH'15 (G15)	68.6 M	1.8 B	Web
Nod	le Specification (32)		Grap	h Applications	UK'06 (UK6)	80.6 M	2.48 B	Web
CPU	28-core @ 2.6 GHZ		PageRank (Pl	R)	UK Union (UKU)	133 M	5.5 B	Web
Memory	192 GB		Single Source	e Shortest Path (SSSP)	Rmat26 (R26)	67.1 M	1.07 B	Synthetic
OS	Linux		Breadth First	: Search (BFS)	Rmat27 (R27)	134 M	2.14 B	Synthetic
MPI	Intel		Connected C	omponent (CC)	Rmat28 (R28)	268 M	4.29 B	Synthetic
Network	Intel Omni-path Fabric	2			Rmat29 (R29)	536 M	8.58 B	Synthetic
	•				Rmat30 (R30)	1.07 B	17.1 B	Synthetic

#### Experiments: CSC, DCSC, and TCSC Comparison



#### Experiments: GraphTap, GraphPad, and LA3 Comparison (weak scaling)





18

R28



#### MPI+X Parallelism Model vs. the New MPI\*X Parallelism Model

- MPI+X uses process-based partitioning and hence
  - When computing each tile, multiple threads are forked & joined
  - Only MPI processes are the communication endpoints
  - Synchronization points are designed for MPI processes

- MPI\*X uses thread-based partitioning and so
  - Threads stay alive
  - Communication is done by threads (better overlapping)
  - Synchronization is done by threads



20

The Process-based 2D Matrix Partitioning and Placement for scalability (e.g., *p*=4)

- **Processed-based 2D partitioning**: Given *p* processes, a matrix is partitioned into a *p* x *p* grid of tiles
  - Each process **computes** *p* tiles
- **2D cyclic placement (GraphPad, IPDPS 16)**: Given *p* processes,  $\sqrt{p}$  processes are placed in each row/column group
  - Process **communication** is not part of the placement
- **2D Staggered partitioning (LA3, VLDB 18)**: Like 2D cyclic, however, unique processes are placed in diagonal tiles
  - Process communication is democratized by diagonal processes (O(2p)+O(p^2) time)



2D-process-based partitioning



**2D-Cyclic Placement** 



#### MPI+X Parallelism Model for Multicore HPC Systems (e.g., p=4)

- MPI+X parallelism model utilizes
  - MPI to achieve horizontal scaling
    - One MPI process per machine
  - A threading library (X), e.g, OpenMP, or Pthread to achieve **vertical scaling**

 MPI+X uses process-based partitioning & placement



*p* x *p* (4x4)

**MPI+X with Processed-based 2D-Staggered** 

*m* is the number of sub-partitions

## The New Thread-based 2D Matrix Partitioning and Placement for scalability (e.g., *p*=4, *t*=2)

- Thread-based 2D partitioning: Given *p* processes & *t* threads, a matrix is partitioned into a (*p.t*)x(*p.t*) grid of tiles
  - *p.t* global threads
  - Each global thread **computes** *t* tiles



*n*/(p.t)

 $(p.t) \ge (p.t) (8 \ge 8)$ 

2D Thread-based Partitioning

## The New Thread-based 2D Matrix Partitioning and Placement for scalability (e.g., *p*=4, *t*=2)

- Thread-based 2D partitioning: Given *p* processes & *t* threads, a matrix is partitioned into a (*p.t*)x(*p.t*) grid of tiles
  - *p.t* global threads
  - Each global thread **computes** *t* tiles
- New Thread-based 2D Staggered placement (2DT-Staggereed):
  - Thread communication is democratized by diagonal threads (O(2p.t) time)

A							
$T_0$	$T_1$	$T_0$	$T_1$	$T_0$	$T_1$	$T_0$	$T_1$
<i>t</i> <sub>2</sub>	$T_3$	$T_2$	<b>T</b> <sub>3</sub>	<i>†</i> <sub>2</sub>	<b>T</b> <sub>3</sub>	<i>T</i> <sub>2</sub>	<b>T</b> <sub>3</sub>
$t_4$	$t_{5}$	$T_4$	$t_{5}$	$t_4$	$T_5$	$T_4$	$T_5$
$t_6$	<b>t</b> 77	$T_6$	<b>†</b> 7	$T_6$	<b>†</b> 7	$T_6$	$t_7$
<i>t</i> <sub>2</sub>	$T_{3}$	$T_2$	$T_3$	$t_2$	$T_3$	$T_2$	$T_3$
$t_4$	$t_{5}$	$t_4$	$t_{5}$	$t_4$	$t_{5}$	$t_4$	$t_{5}$
$T_6$	<b>†</b> <sub>7</sub>	$T_6$	<b>†</b> 7	$T_6$	<b>Ť</b> <sub>7</sub>	$T_6$	Ť <sub>7</sub>
$T_0$	$t_1$	$T_0$	$t_1$	$t_0$	$T_1$	$T_0$	$T_1$

 $(p.t) \ge (p.t) (8 \ge 8)$ 

2D Thread-based Placement (Global Thread Ids)

#### The New MPI\*X Parallelism Model (e.g., *p*=4 and *t*=2)

#### • MPI\*X parallelism model utilizes

- Threads to achieve diagonal scaling
  - t global threads are grouped together to form a process
- Global thread ids are combined to form process and local thread assignments
- One MPI process per CPU socket and one thread per core
  - NUMA-aware

				4			
$P_0T_0$	$P_1T_0$	$P_0T_0$	<i>P</i> <sub>1</sub> <i>T</i> <sub>0</sub>	$P_0T_0$	<i>P</i> <sub>1</sub> <i>T</i> <sub>0</sub>	$P_0T_0$	$P_1T_0$
$P_2T_0$	$P_{3}T_{0}$	$P_{2}T_{0}$	$P_{3}T_{0}$	$P_{2}T_{0}$	<i>P</i> <sub>3</sub> <i>T</i> <sub>0</sub>	$P_{2}T_{0}$	$P_3T_0$
$P_0 T_1$	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>0</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>0</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>0</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>
$P_{2}T_{1}$	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>2</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>2</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>2</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>
$P_{2}T_{0}$	$P_{3}T_{0}$	$P_{2}T_{0}$	$P_{3}T_{0}$	$P_{2}T_{0}$	$P_{3}T_{0}$	$P_{2}T_{0}$	$P_{3}T_{0}$
$P_{0}T_{1}$	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>0</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>0</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>0</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>1</sub> <i>T</i> <sub>1</sub>
$P_{2}T_{1}$	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>2</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>2</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>2</sub> <i>T</i> <sub>1</sub>	<i>P</i> <sub>3</sub> <i>T</i> <sub>1</sub>
$P_0T_0$	$P_1T_0$	$P_0T_0$	<i>P</i> <sub>1</sub> <i>T</i> <sub>0</sub>	$P_0T_0$	$P_1T_0$	$P_0T_0$	$P_1T_0$

 $(p.t) \times (p.t) \ (8 \times 8)$ 

2D Thread-based Placement (Local Thread Ids)

#### The **Graphite** HPC Graph Analytics System

#### • MPI\*X

- Thread-based partitioning & placement (O(2*pt*) time versus previous O(*p*^2))
- Diagonal scaling
- TCSC (Triply Compressed Sparse Column) (GraphTap, Cluster 19)

#### NUMA-aware

- Computation (CPU and memory affinity)
- Communication (MPI shared-memory transport)
- Asynchronous computation & communication
  - Broadcasting (log(*p*) time) when possible

#### Experiments

	Graph Processing Systems							
System	Parallelism	Unit	Partitioning	NUMA-aware	Compression	Primitive		
Graphite	MPI*X	Threads	2DT-Staggered	CPU/Memory	TCSC	SpMSpV <sup>2</sup> GAS		
LA3 (VLDB, 2018)	MPI+X	Processes	2D-Staggered	N/A	DCSC	SpMV GAS		
GraphPad (IPDPS, 2016)	MPI+X	Processes	2D-Cyclic	N/A	CSC	SpMV GAS		
Gemini (USENIX, 2016)	MPI+X	Processes	1D-Row	Memory	CSC/CSR	Push/Pull		

Nod	e Specification (20)	Graph Applications	Datasets				
CPU	28-core @ 2.6 GHZ	PageRank (PR)	Graph	<i>V</i>	<i>E</i>	Туре	
			UK'05 (UK5)	39.4 M	0.93 B	Web	
Memory	192 GB	Single Source Shortest Path (SSSP)	IT'04 (IT4)	41.2 M	1.15 B	Web	
OS	Linux	Breadth First Search (BFS)	Twitter (TWT)	41.6 M	1.46 B	Social	
			GSH'15 (G15)	68.6 M	1.8 B	Web	
MPI	Intel	Connected Component (CC)	UK'06 (UK6)	80.6 M	2.48 B	Web	
Network	Intel Omni-path Fabric		UK Union (UKU)	133 M	5.5 B	Web	
			Rmat26 (R26)	67.1 M	1.07 B	Synthetic	
			Rmat27 (R27)	134 M	2.14 B	Synthetic	
			Rmat28 (R28)	268 M	4.29 B	Synthetic	
			Rmat29 (R29)	536 M	8.58 B	Synthetic	

## Experiments: MPI\*X and MPI+X comparison (scalability)



#### Results (Weak Scaling)



29

#### GraphTap & Graphite vs. State-of-the-art

System	Computing Model	Platform	Parallelism Model	Computation & Communication	Matrix Compression	Authors
GraphPad	Linear Algebra	HPC/Cloud	MPI+OpenMP		DCSC <sup>2</sup>	2016, Intel
Gemini	Frontier-based	HPC	MPI+OpenMP	NUMA-aware CPU	CSR/CSC ( <mark>Dual</mark> )	2016, Tsinghua University
CombBLAS	Linear Algebra	HPC	MPI+OpenMP		DCSC	2018, LBNL
LA3	Linear Algebra	Cloud	MPI+OpenMP	Optimized for cloud communication	DCSC	2018, CMUQ
LAGraph	Graph Theory	НРС	OpenMP		CSR	2020, Texas A&M
GraphTap	Linear Algebra	HPC	MPI		CSC/DCSC/ TCSC	2019, UPitt
Graphite	Linear Algebra	HPC	MPI*X	Asynchronous Comm. Overlapping of Comm. NUMA-aware	TCSC	2020, UPitt



## Sources of Sparsity in DNNs<sub>1</sub>

A. Input can be sparse

- B. Network can be sparse
  - Being pruned during/after training
  - Following a predefined sparse architecture
    - TensorFlow Model Optimization
  - Propagating weights can create sparsity  $\frac{1}{2}$
  - Due to either input/network sparsity
  - Due to the activation function



Datasets

CIFAR-10 CIFAR-100

NNZ Distribution

Zero

IMDB

Nonzero

Ratio (%) 5<sup>.0</sup>

0

MNIST

Fashion

MNIST



## Multithreaded Single Machine (Sparse) DNN Inference

#### Data Parallelism

- Horizontal 1D-Row partitioning of input (A)
  - *t* input partitions where *t* is the number of threads
- + No synchronization, stragglers, bad L3 utilization



#### Model Parallelism

- Vertical 1D-Column partitioning of network (B)
  - *t* network partitions
- Strict synchronization, + better L1 & L3 utilization
- Each tile is stored using a compressed sparse format



### **Distributed** (Sparse) DNN Inference

- Data \* Data parallelism
  - Horizontal 1D-Row partitioning of input (A)
  - *p.t* input partitions
  - *p* is the number of processes
- Data \* Model Parallelism
  - vertical 1D-Column partitioning of network (B)
  - *p.t* network partitions
- Network is replicated for each process
- No communication is happening among processes



## Neural Network Hashing to Achieve Load Balance & Locality

nput (A

- Hashing mitigates data parallelism stragglers
- Input hashing:
  - Hashes rows of A
  - + Load balance in data parallelism
- Network hashing:
  - Hashes columns of A & rows of B
  - + Better access pattern in model parallelism in favor of the cache
- Input + network hashing:
  - Hashes rows of A
  - Hashes columns of A & rows of B
  - Hashes columns of B
  - + Both above advantages
  - + Pseudo sequential access pattern (log(n))



#### Experiments

	Datase	t (Radix	Net Sparse	DNN, I	MNIST	)
	Input			Netwo	ork	
#	Sizo		Each La	yer	All I	ayers
#	512e <sub>m x n</sub>	ININZ	Size <sub>n x n</sub>	NNZ	L	NNZ
$A_0$					120	3.9M
$A_1$	60K x 1K	6.3M	1K x 1K	32K	480	15.7M
$A_2$					1920	62.9M
$B_0$					120	15.7M
$B_1$	60K x 4K	25M	4K x 4K	131K	480	62.7M
B <sub>2</sub>					1920	251M
C <sub>0</sub>					120	62.9M
$C_1$	60K x 16K	99M	16K x 16K	524K	480	251M
C <sub>2</sub>					1920	1B
$D_0$					120	251M
$D_1$	60K x 65K	392M	65K x 65K	21 M	480	1B
$D_2$					1920	4B

Parallelism/		Right Multiplication	Left Multiplication
Network	Intel Omn	i-path Fabric	
MPI	Intel		
OS	Linux		
Memory	192 GB		
CPU	28-core @	2.6 GHZ	
Nod	e Specifica	tion (32)	

Parallelism/ Multiplication	Right Multiplication (CSC)	Left Multiplication (CSR)
Data Parallelism	Yes	Yes
Model Parallelism	Yes	N/A

### Experiments: Comparison of SpMMs (Single)



Runtime Variation (No Hashing)



L1 & L3 Utilization (Input + DNN Hashing)





Results are for D<sub>2</sub>

#### Experiments: Comparison of SpMMs (Distributed)

Strong Cluster Scaling



Results are for D<sub>2</sub>



#### Motivation: Data Versus Model Parallelism

- Due to imbalance Data parallelism suffers from straggler effect
  - Hashing alleviates the imbalance problem of data parallelism



#### Data-then-model Parallelism

- Switch from data to model parallelism and turning idle threads into additional processing power to mitigate the effect of stragglers
- Lazy load balancing by reusing idle threads
- Less synchronization cost





## Experiments

	Datase	t (Radix	Net Sparse	DNN, I	MNIST	)
Input			DNN			
#	Sizo	n NNZ	Each Layer		All Layers	
#	# SIZe <sub>mxn</sub>		Size <sub>n x n</sub>	NNZ	L	NNZ
$A_0$					120	3.9M
$A_1$	60K x 1K	6.3M	1K x 1K	32K	480	15.7M
$A_2$					1920	62.9M
B <sub>0</sub>					120	15.7M
$B_1$	60K x 4K	25M	4K x 4K	131K	480	62.7M
B <sub>2</sub>					1920	251M
C <sub>0</sub>					120	62.9M
$C_1$	60K x 16K	99M	16K x 16K	524K	480	251M
C <sub>2</sub>					1920	1B
$D_0$					120	251M
$D_1$	60K x 65K	392M	65K x 65K	21 M	480	1B
$D_2$					1920	4B

Node Specification (16)				
CPU	28-core @ 2.6 GHZ			
Memory	192 GB			
OS	Linux			
MPI	Intel			
Network	Intel Omni-path Fabric			
Parallelisms				
	Parallelisms			
Data Para	Parallelisms allelism			
Data Para Model Pa	Parallelisms allelism arallelism			
Data Para Model Pa Data-the	Parallelisms allelism arallelism n-model Parallelism			
Data Para Model Pa Data-the Manager	Parallelisms allelism arallelism n-model Parallelism r-worker Parallelism			
Data Para Model Pa Data-the Manager Work-ste	Parallelisms allelism arallelism n-model Parallelism r-worker Parallelism ealing Parallelism			

#### Experiments: Data-then-model Parallelism Thread Scheduling Algorithms

#### • Locking Mechanism Runtime

- Threads are either
  - Worker (active) or helper (newly idle)
- An idle thread enlists into an idle queue
- A working thread probes the idle queue
  - Helper threads get recruited by a working one

#### • Advantages

- + Overloadable with scheduling strategies
  - Earliest first, slower first, and faster first
- + Fully decentralized and asynchronous
- + Minimal lock contention
  - condition variables + locks
- + Elastic: adding/removing threads on the fly
  - zero data movement



#### Scalability Experiment



Weak Scaling







## SpDNN & DistSpDNN vs. State-of-the-art

System	Computing Model	Platform	Parallelism Model	Computation & Communication	Matrix Compression	Authors
TensorFlow	Linear Algebra	Commodity/ HPC/ <mark>Cloud</mark>	Data parallelism (Keras + CPU/GPU/TPU)	Horovod	Dense	2015 – present, Google
PyTorch	Linear Algebra	Commodity/ HPC	Data parallelism (MPI/ Gloo + OpenMP/ CUDA)	Horovod	Dense	2016 – present, Facebook
MXNet	Linear Algebra	Commodity/ HPC	Data parallelism (MPI4Py/ Gloo + OpenMP/ CUDA)	Horovod/ TCP/IP socket	Dense	2020 – present, Amazon
DeepSpeed	Linear Algebra	Commodity/ HPC/ Cloud	Data parallelism (MPI/ NCCL + X/ CUDA)	Communication overlapping	Sparse	2019 – present, Microsoft
SpDNN	Linear Algebra	Multicore	Model (OpenMP)		Sparse	2019, UPitt
DistSpDNN	Linear Algebra	HPC	Data-then-model (MPI*OpenMP)	NUMA-aware Memory re(allocation)	Sparse, 2-step SpMM	2020, UPitt

## Summary & Conclusion (1)

- Sparse Matrix Vector (SpMV)
  - **TCSC** that reduces time, space, and communication complexities (GraphTap, Cluster 2019)
    - Co-compression reduces the communication volume
  - MPI\*X that deems threads as basic unit of computation, computation, and synchronization (Graphite, VLDB 2020)
    - Better scalability requires laying out the computation/communication path beforehand
- Sparse Matrix Matrix Multiplication (SpMM)
  - Neural network hashing for alleviating the load imbalance and offering cache locality (DistSpDNN, HPEC 2020)
    - Proper hashing can offer super-linear speedup (log(n) time)
  - Data-then-model parallelism for mitigating the straggler effect (DistSpDNN, HPEC 2020)
    - Sparse computations require new parallelism that takes account of the runtime imbalance

#### Summary & Conclusion (2)



Implementations are available at <a href="https://github.com/hmofrad">https://github.com/hmofrad</a>

Acknowledgment!

Questions?