# Accelerating Distributed Inference of Sparse Deep Neural Networks via Mitigating the Straggler Effect

**Mohammad Hasanzadeh Mofrad**, Rami Melhem,
Yousuf Ahmad and Mohammad Hammoud
Github Repository: https://github.com/hmofrad/DistSparseDNN

Carnegie
Mellon
University
Qatar

IEEE
HPEC

# Dense and Sparse Neural Networks

- **Deep Neural Networks (DNNs)** are pervasive
  - Speech processing
  - Friend suggestion
  - Autonomous driving
  - Item recommendation

- 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
  - Comparable accuracy

- Sparse DNNs core kernel is **Sparse Matrix-matrix Multiplication (SpMM)**

# 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
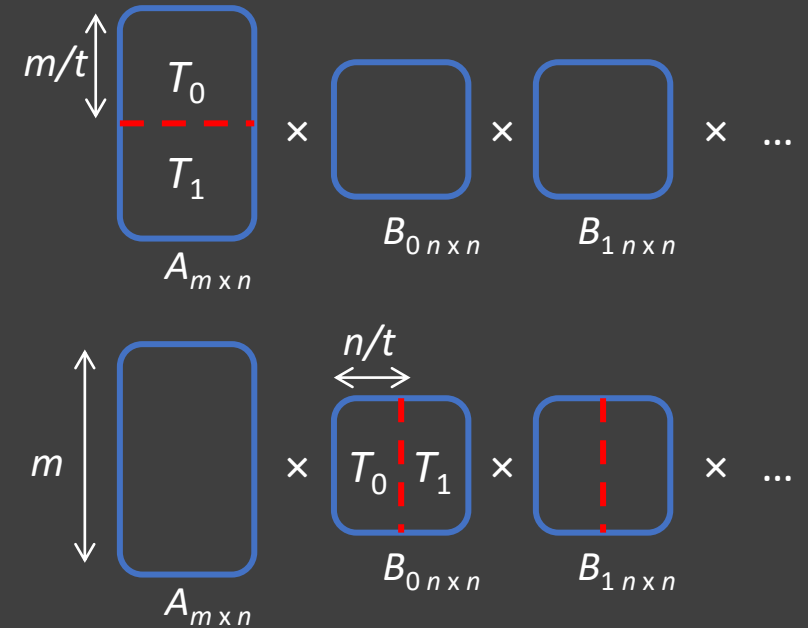
- **Model Parallelism**
  - Vertical 1D-Column partitioning of network (B)
    - $t$ network partitions
    - **−** Strict synchronization, **+** better L1 and L3 utilization

- **Manager-worker** (single queue of $m*t$ i.e. $m>>t$)

- **Work-Stealing** ($t$ queues of $m$ tiles)

- SpMM algorithm: Two-step right SpMM with CSC
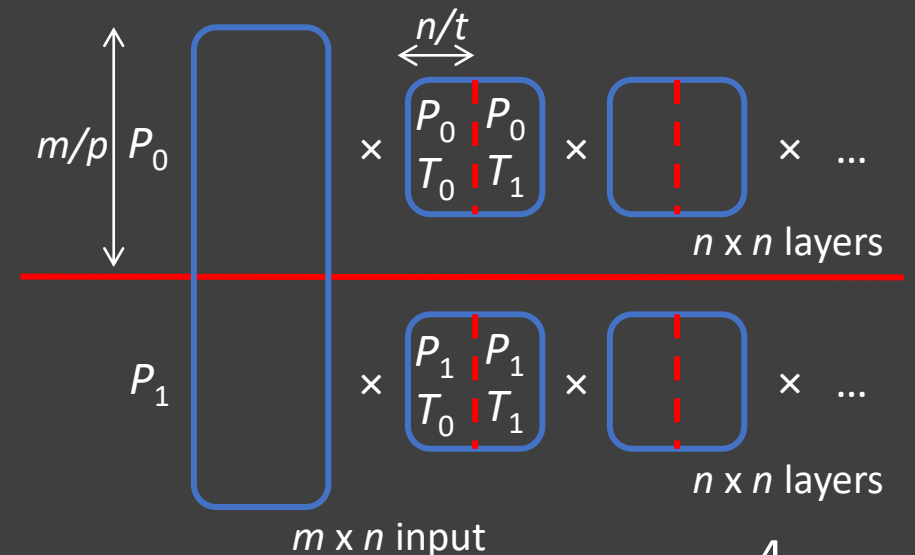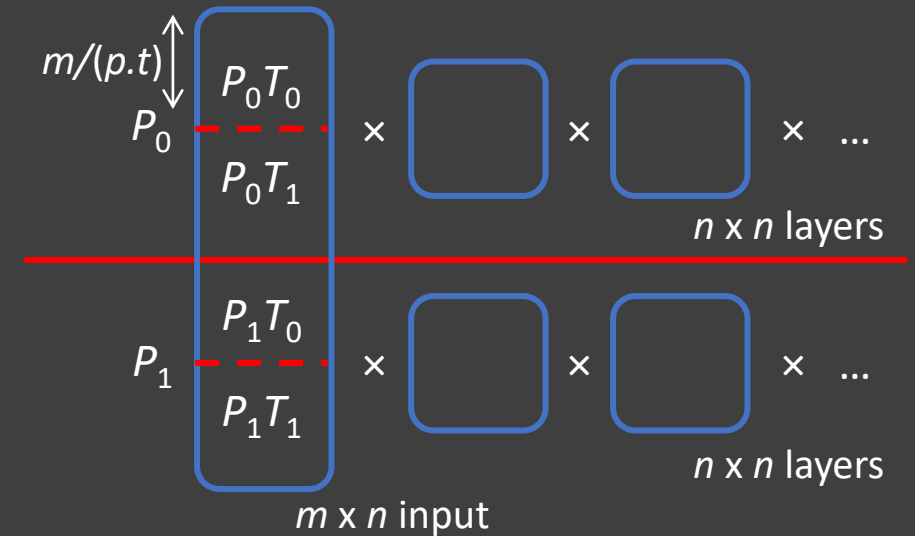
# Distributed (Sparse) DNN Inference

- **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

- Data * Manager-worker, and Data * Work-stealing

- Network is replicated for each process

- No communication is happening among processes



$m/(p.t)$

$P_0$   $P_0 T_0$   $P_0 T_1$   $\times$   $\times$   $\times$   ...

$n \times n$ layers

$P_1$   $P_1 T_0$   $P_1 T_1$   $\times$   $\times$   $\times$   ...

$n \times n$ layers

$m \times n$ input

$n/t$

$m/p$   $P_0$   $\times$   $P_0 T_0$ $P_0 T_1$   $\times$   $\times$   ...

$n \times n$ layers

$P_1$   $\times$   $P_1 T_0$ $P_1 T_1$   $\times$   $\times$   ...

$n \times n$ layers

$m \times n$ input

4

# Data-then-model Parallelism

- Due to imbalance Data parallelism suffers from straggler effect
- Imperfect solution is **hashing**

- 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
- Zero data movement
- No contention
- Less synchronization cost

# Experiments

| Dataset (RadixNet Sparse DNN, MNIST) | | | | | |
|---|---|---|---|---|---|
| Input | | Network | | | |
| | | Each Layer | | All Layers | |
| Size $_{m \times n}$ | NNZ | Size $_{n \times n}$ | NNZ | L | NNZ |
| 60K x 1K | 6.3M | 1K x 1K | 32K | 120 | 3.9M |
| | | | | 480 | 15.7M |
| | | | | 1920 | 62.9M |
| 60K x 4K | 25M | 4K x 4K | 131K | 120 | 15.7M |
| | | | | 480 | 62.7M |
| | | | | 1920 | 251M |
| 60K x 16K | 99M | 16K x 16K | 524K | 120 | 62.9M |
| | | | | 480 | 251M |
| | | | | 1920 | 1B |
| 60K x 65K | 392M | 65K x 65K | 21 M | 120 | 251M |
| | | | | 480 | 1B |
| | | | | 1920 | 4B |

| Node Specification (16) | |
|---|---|
| CPU | 28-core @ 2.6 GHZ |
| Memory | 192 GB |
| OS | Linux |
| MPI | Intel |
| Network | Intel Omni-path Fabric |

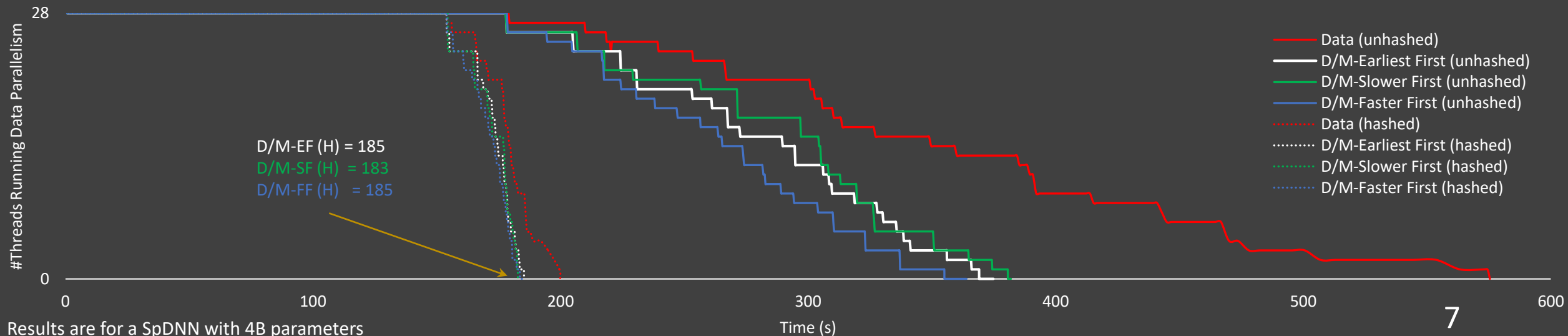| Parallelisms |
|---|
| Data Parallelism |
| Model Parallelism |
| Data-then-model Parallelism |
| Manager-worker Parallelism |
| Work-stealing Parallelism |

6

# 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



D/M-EF (H) = 185
D/M-SF (H)  = 183
D/M-FF (H)  = 185

Legend:
- Data (unhashed)
- D/M-Earliest First (unhashed)
- D/M-Slower First (unhashed)
- D/M-Faster First (unhashed)
- Data (hashed)
- D/M-Earliest First (hashed)
- D/M-Slower First (hashed)
- D/M-Faster First (hashed)

Y-axis: #Threads Running Data Parallelism (0 to 28)
X-axis: Time (s) (0 to 600)

Results are for a SpDNN with 4B parameters

# Experiments: Scalability

## Questions?

**Strong Cluster Scaling (392M, 4B)**



Legend:
- Data*Work-Sharing
- Data*Work-Stealing
- Data*Model
- Data*Data
- Data*Data/Model

Y-axis: Time (s) — 0, 300, 600, 900, 1200, 1500
2109

X-axis: (#Cores), (#Nodes)
(28), (1)   (56), (2)   (112), (4)   (224), (8)   (448), (16)

**Weak Scaling**



Legend:
- Data*Work-Sharing
- Data*Work-Stealing
- Data*Model
- Data*Data
- Data*Data/Model

Y-axis: Time (s) — 0, 50, 100, 150, 200, 250

X-axis: (#Input Nonzeros, #Layers Nonzeros), (#Nodes)
(6.3M,62.9M), (2)   (25M,251M), (4)   (98.8M,1B), (8)   (392M,4B), (16)

**Strong Data Scaling**



Legend:
- Data*Work-Sharing
- Data*Work-Stealing
- Data*Model
- Data*Data
- Data*Data/Model

Y-axis: Time (s) — 0, 50, 100, 150, 200, 250

X-axis: (#Input Nonzeros, #Layers Nonzeros), (#Nodes)
(6.3M,62.9M), (16)   (25M,251M), (16)   (98.8M,1B), (16)   (392M,4B), (16)