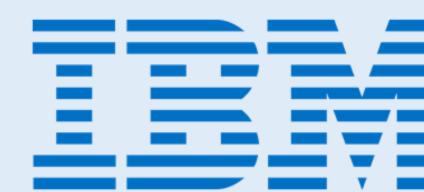
# Can Decentralized Algorithms Outperform Centralized Algorithms? Paper ID 2767

A Case Study for Decentralized Parallel Stochastic Gradient Descent

## Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu











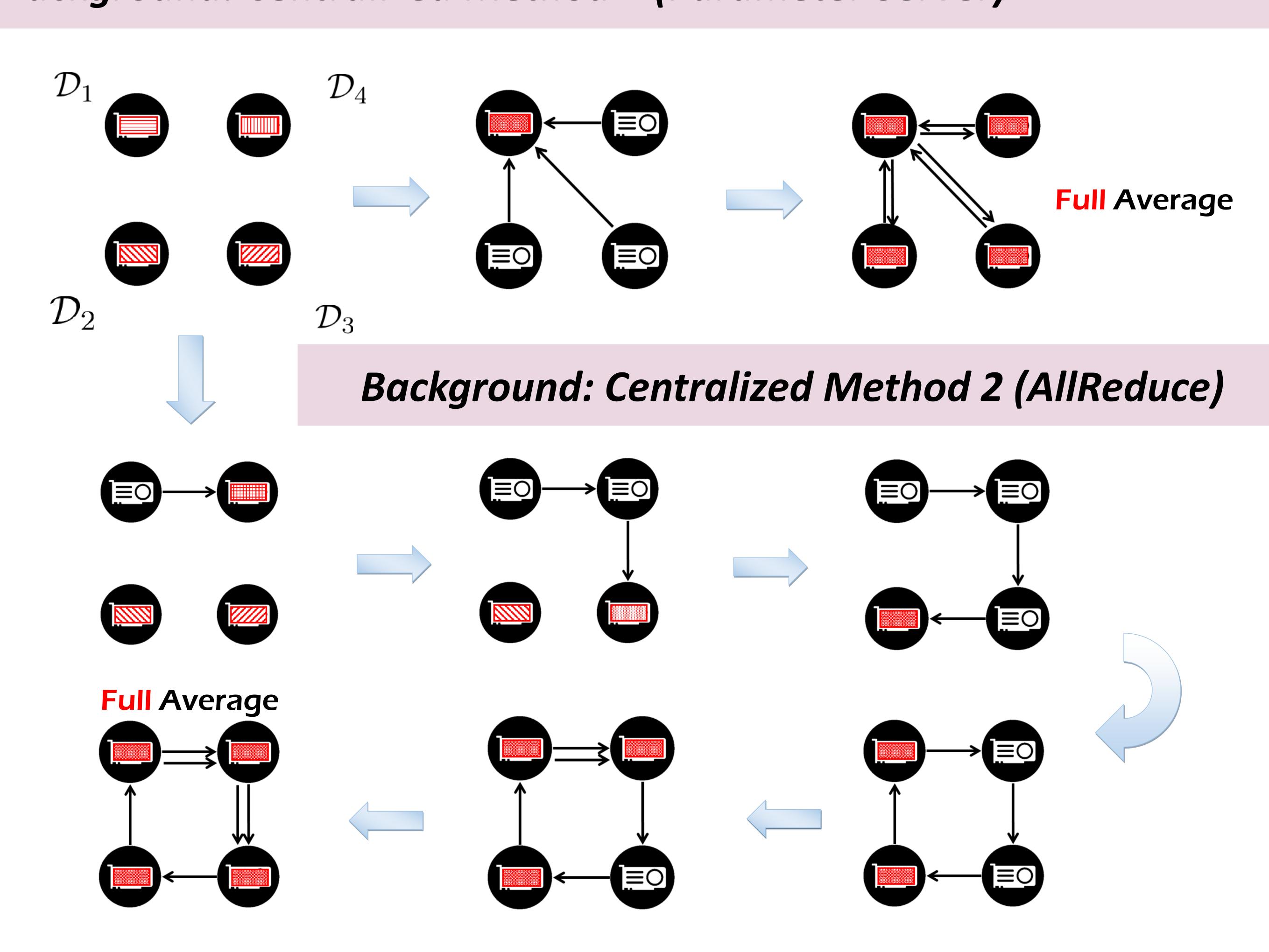
#### Executive summary

- In large scale machine learning, instead of using only one machine, we distribute data into multiple machine and let them collaborate on solving the following optimization problem:

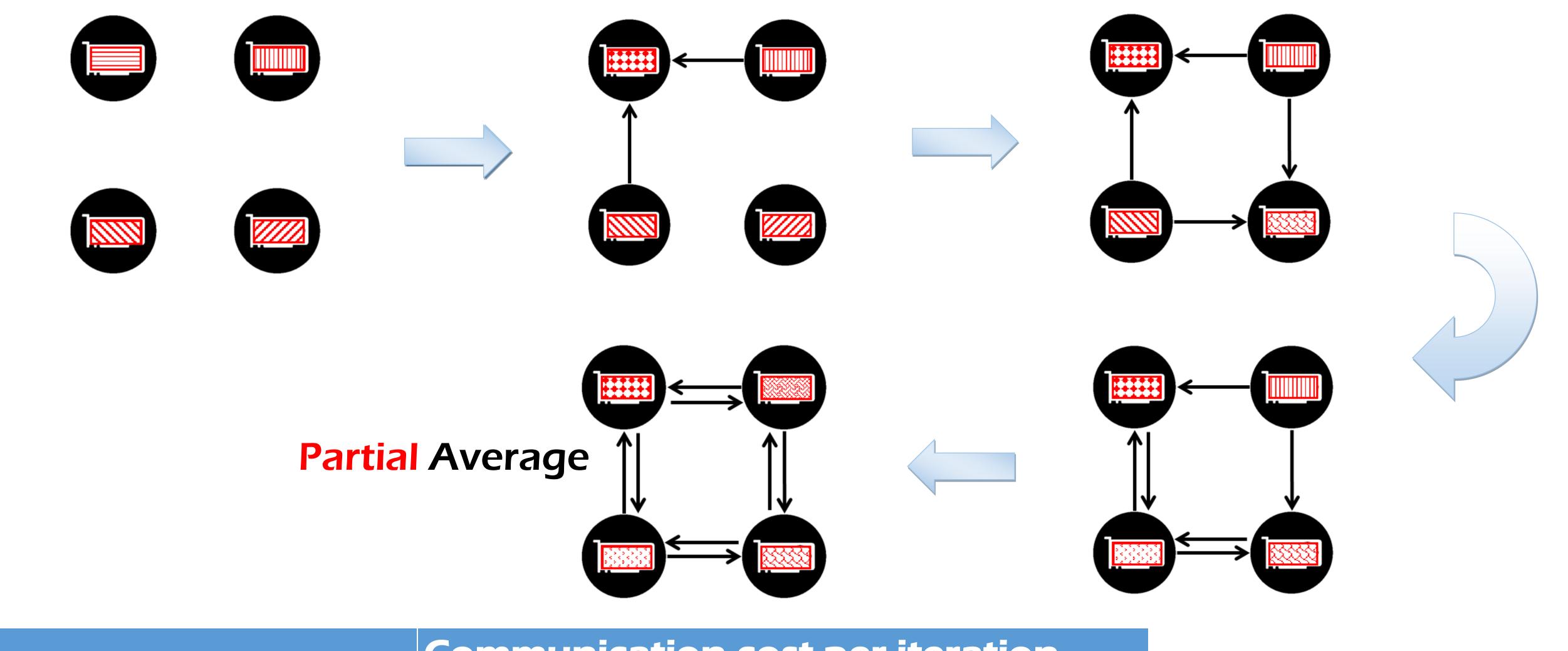
$$\min_{\mathbf{x}} \quad f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^{N} \underbrace{\mathbb{E}_{a \sim \mathcal{D}_i} F_i(\mathbf{x}; a)}_{=:f_i(\mathbf{x})}$$

- Based on our case study, the decentralized algorithm may not need more iterations to converge than its centralized counterpart
- Decentralized algorithms outperform centralized algorithms for networks with low bandwidth and high latency

#### Background: Centralized Method 1 (Parameter Server)



#### Our Proposal: Decentralized Method



Communication cost per iteration Single Machine PS O(N \* alpha + NB \* beta)

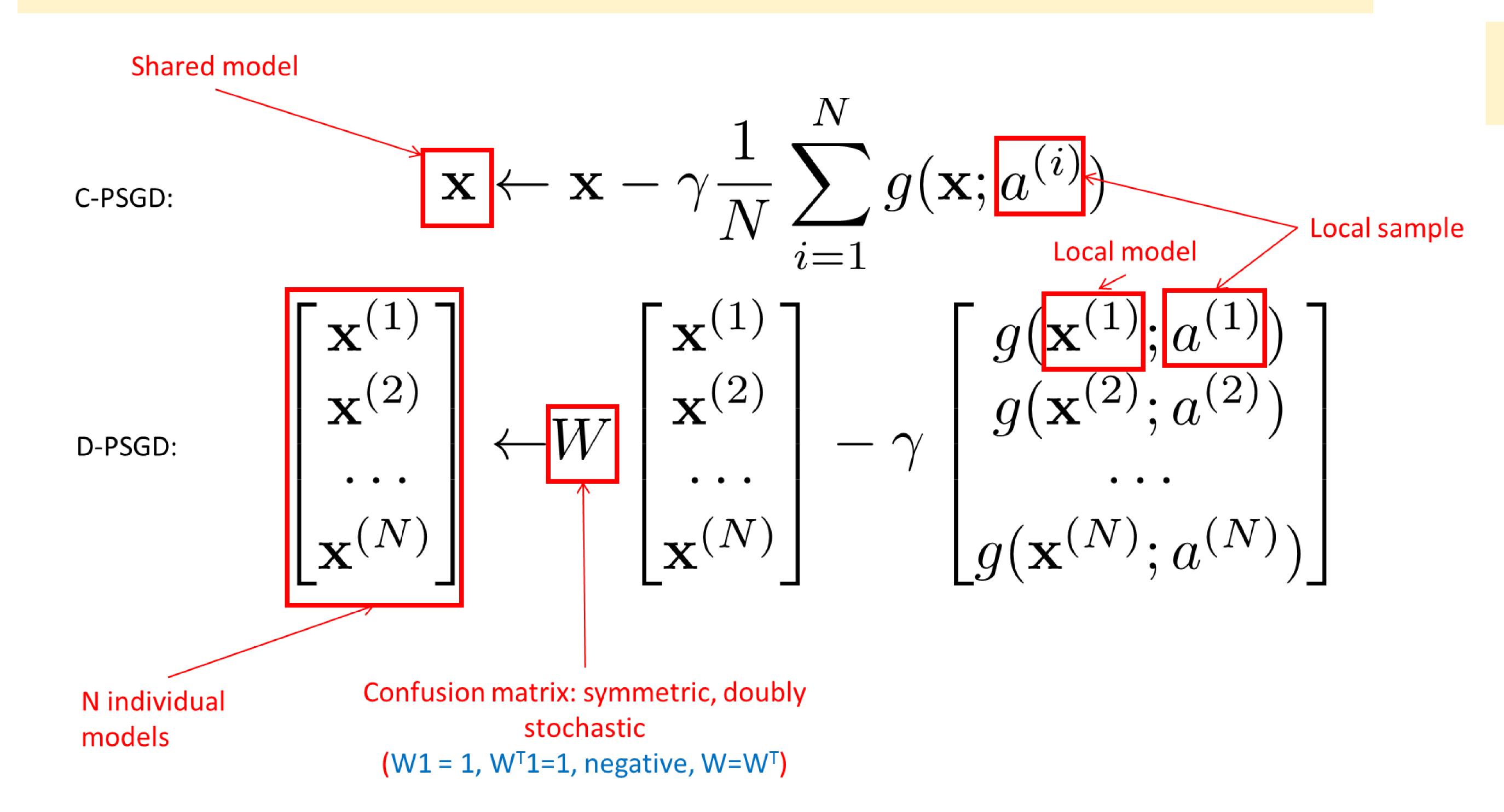
Decentralization (ring) O(alpha + B \* beta)

AllReduce

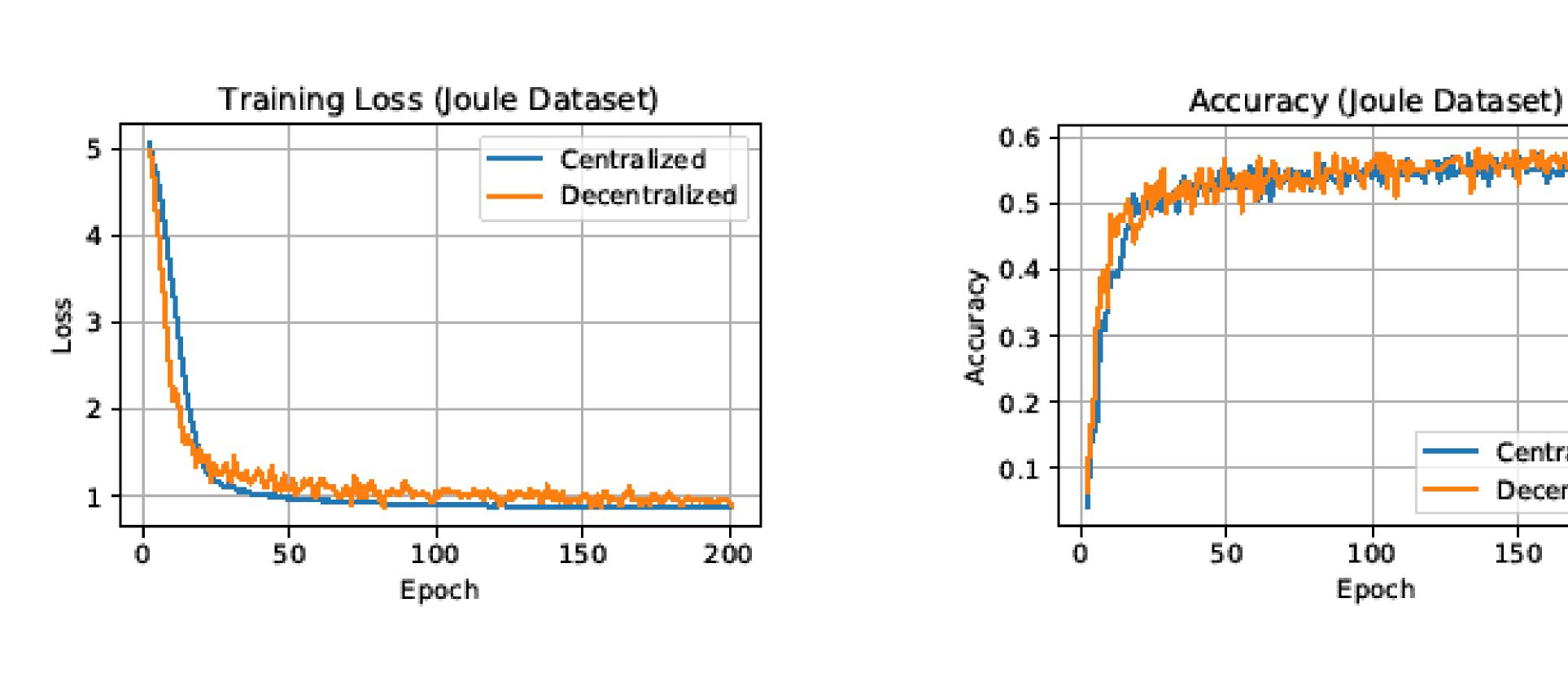
O(log N \* alpha + B \* beta)

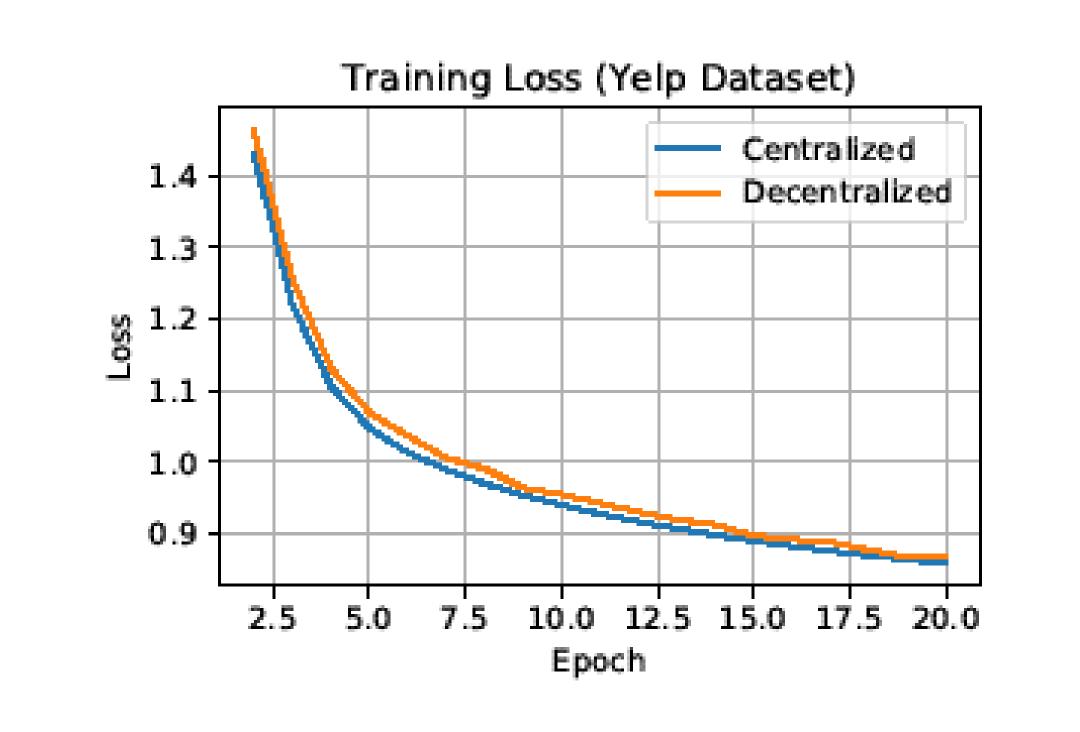
- alpha: latency per message - beta: transfer time per byte
- N: # workers - B: # bytes of the message

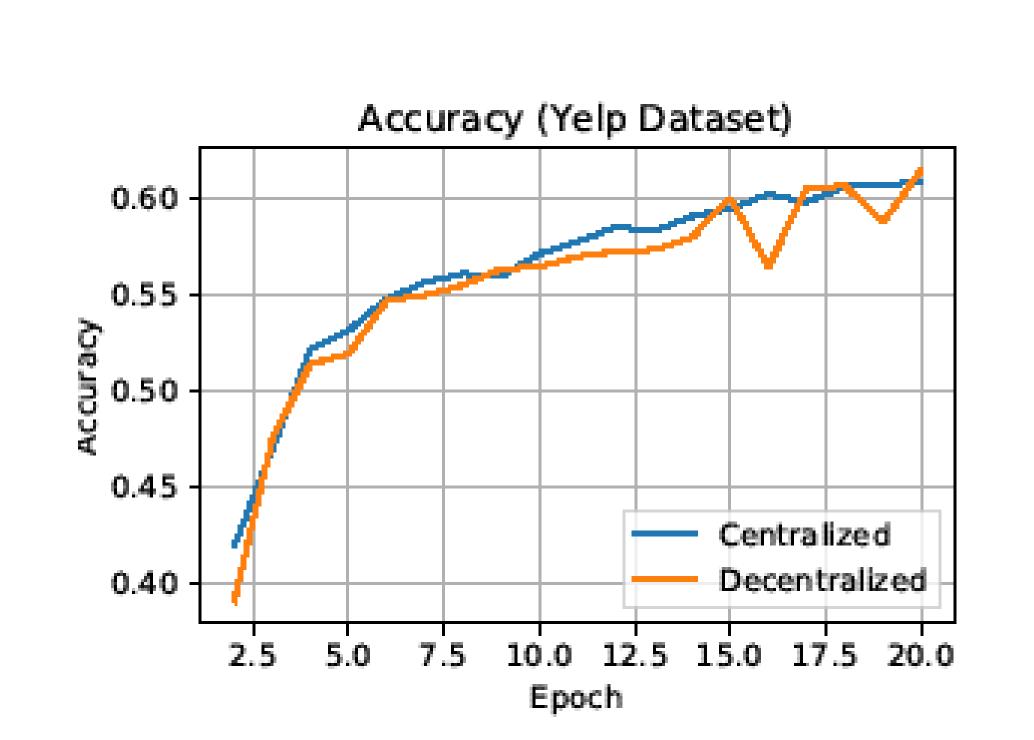
#### Algorithm (Decentralized Parallel SGD)



### Evaluation: Proprietary dataset and model (IBM Watson Natural Language Classifier)







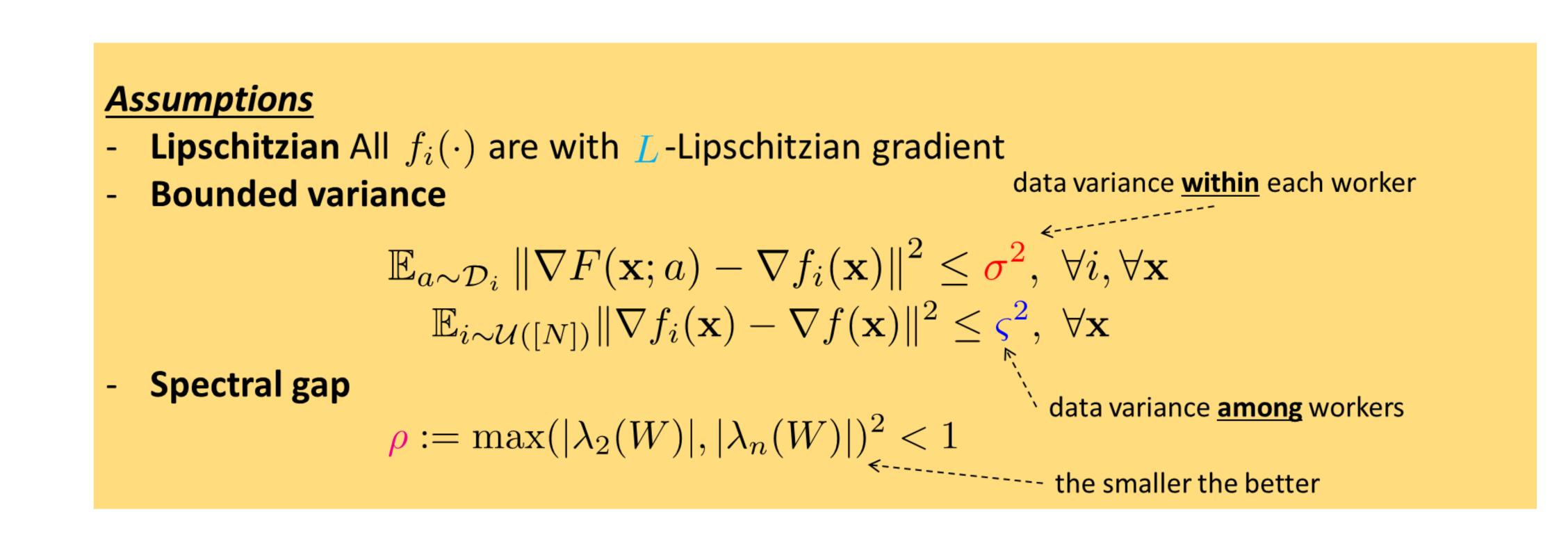
#### Future work

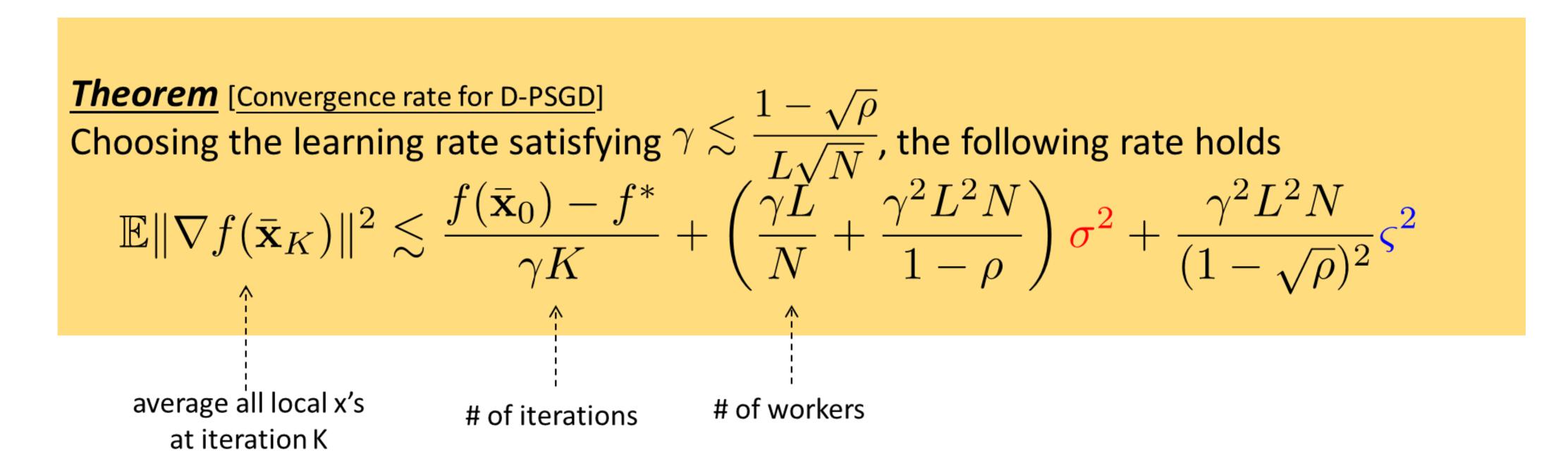
- Asynchronous parallelism for decentralized algorithms - Investigate new topologies to improve communication efficiency

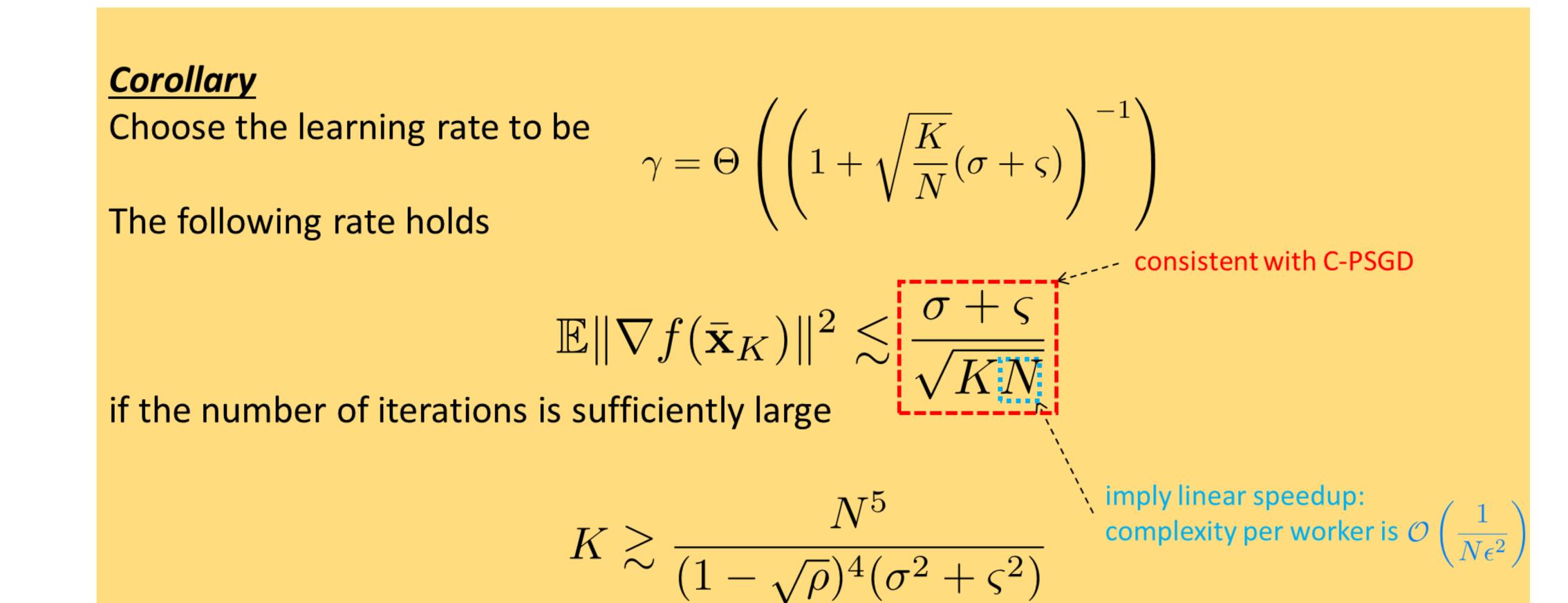


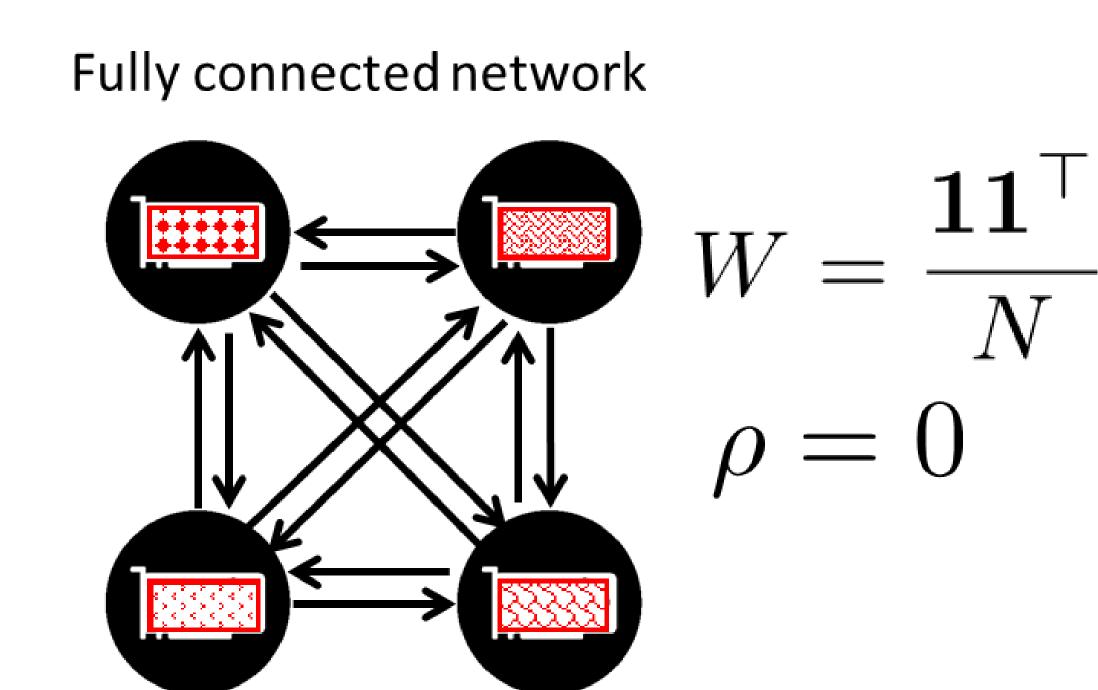
This Paper

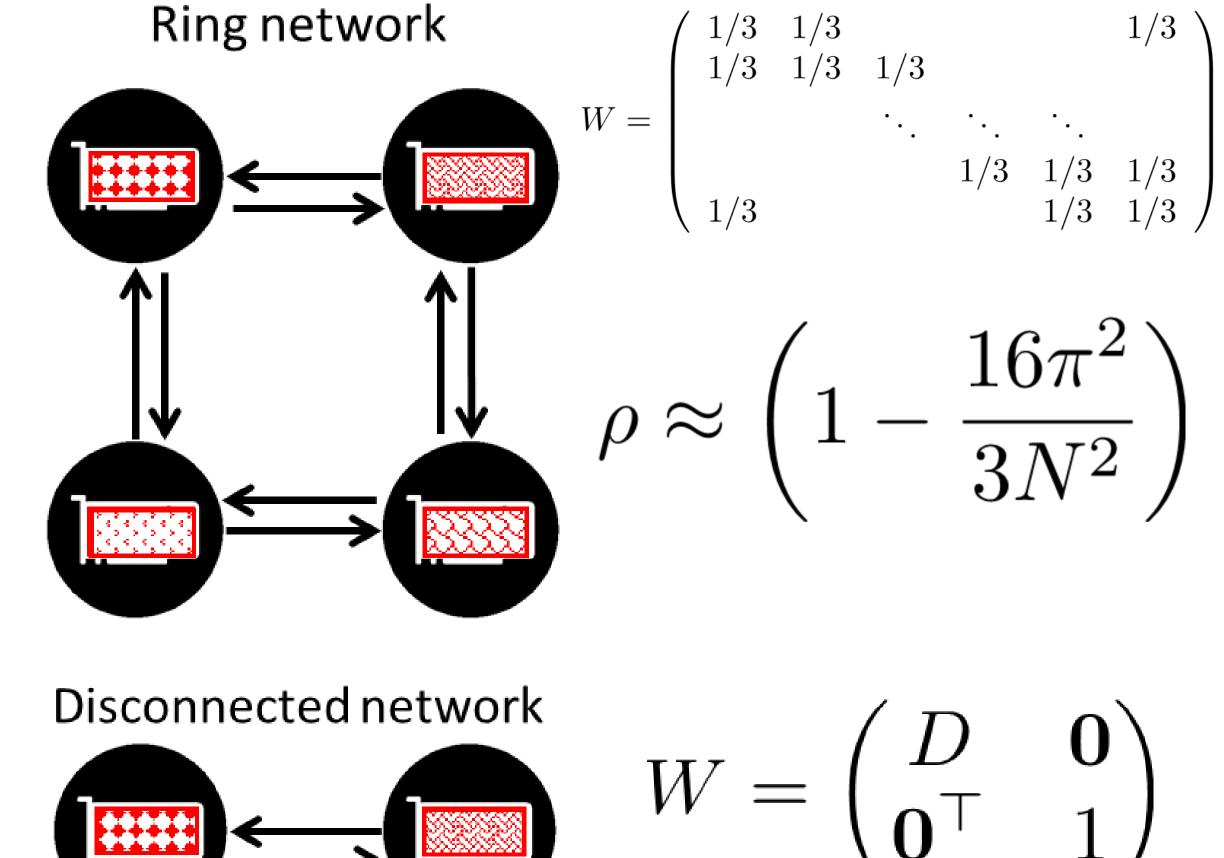
#### Theoretical Results

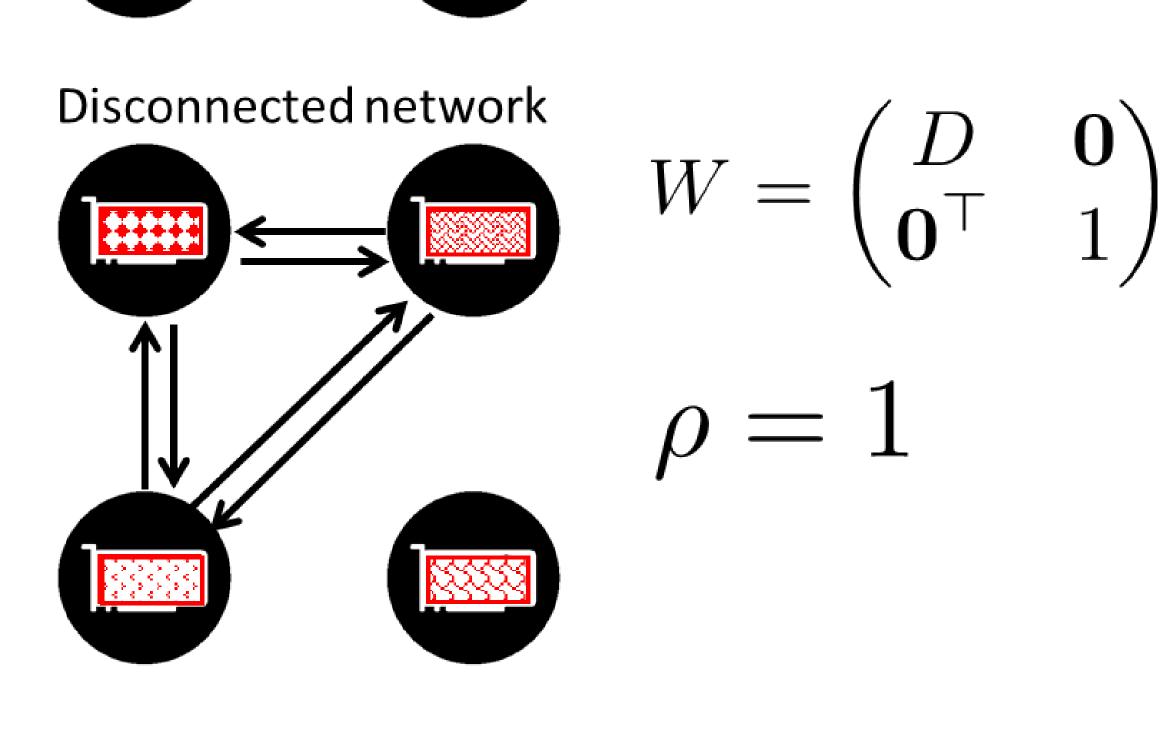












### Evaluation: Public dataset and model (CIFAR10/ResNet)

