# Communication Efficient Distributed Training

Ji Liu, Ph.D.
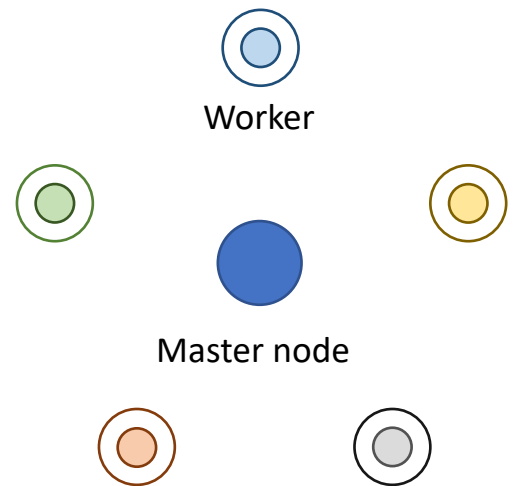
# Objective

model

$$\min_{\boldsymbol{x}} \quad f(\boldsymbol{x}) := \frac{1}{n} \sum_{i=1}^{n} \mathbb{E}_{\boldsymbol{a}^{(i)} \sim \mathcal{D}_i} F(\boldsymbol{x}; \boldsymbol{a}^{(i)})$$

# of workers

samples on worker i

Worker

Master node

All functions are assumed to be L-Lipschitzian

Centralized distributed learning

How to reduce communication cost?

# Summary

Compression



Asynchronization



Decentralization



Compression
+
Decentralization

# Compression

# Algorithm

$\mathbf{C}(\cdot)$ Compression operator (maybe randomized)

$$g^{(i)} := \nabla F(\boldsymbol{x}; \boldsymbol{a}^{(i)})$$

(SGD) $\quad \mathbf{x} \leftarrow \mathbf{x} - \gamma \bar{\mathbf{g}}$



**Worker 1**

$g^{(1)}$  $\bar{g}$

**Server**

$\bar{g}$  $\bar{g}$

$g^{(2)}$  $g^{(3)}$

**Worker 2**  **Worker 3**

$$\bar{\mathbf{g}} = \frac{1}{3}(\mathbf{g}^{(1)} + \mathbf{g}^{(2)} + \mathbf{g}^{(3)}) \qquad \text{(Standard)}$$

Exchange 2N full vectors

$$\bar{\mathbf{g}} = \frac{1}{3}(\mathbf{C}(\mathbf{g}^{(1)}) + \mathbf{C}(\mathbf{g}^{(2)}) + \mathbf{C}(\mathbf{g}^{(3)})) \qquad \text{(Single compression)}$$

Exchange N(1+c) full vectors

$$\bar{\mathbf{g}} = \mathbf{C}\left(\frac{1}{3}(\mathbf{C}(\mathbf{g}^{(1)}) + \mathbf{C}(\mathbf{g}^{(2)}) + \mathbf{C}(\mathbf{g}^{(3)}))\right) \text{(Double compression)}$$

Exchange 2cN full vectors

# Unfortunately

To ensure convergence, it should satisfy $\mathbb{E}(\mathbf{C}(\mathbf{x})) = \mathbf{x}$

Early methods only work for $\mathbf{C}(\cdot)$ compression operator
- *Randomized quantization   (unbiased)*
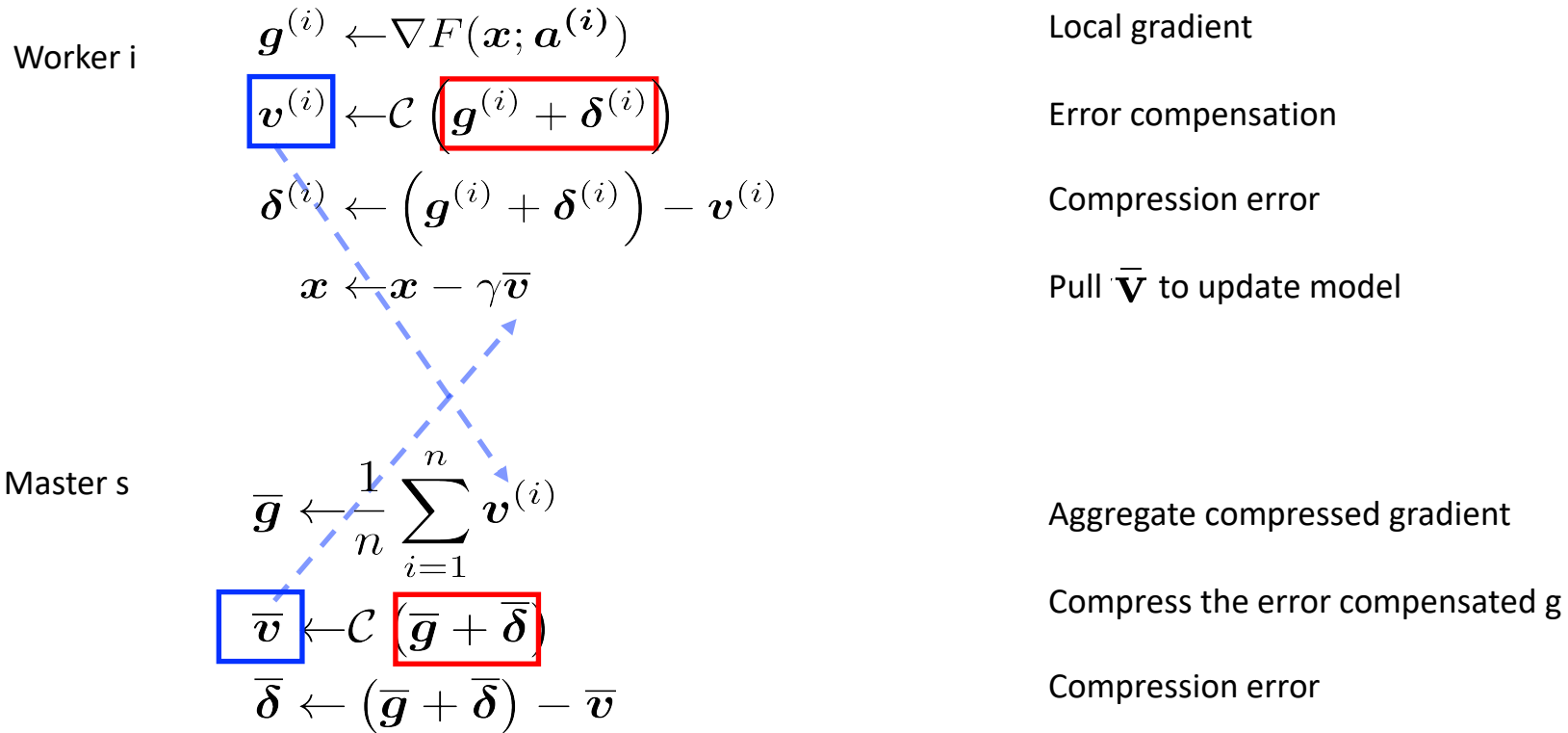- ~~*Randomized quantization   (biased)*~~
- ~~*1bit quantization*~~
- ~~*Clipping*~~
- ~~*Top-k sparsification*~~

Can we relax it to allow more aggressive or even arbitrary compression?

# Double Squeeze: Error Compensated SGD

Worker i

$$\boldsymbol{g}^{(i)} \leftarrow \nabla F(\boldsymbol{x}; \boldsymbol{a}^{(i)})$$

Local gradient

$$\boxed{\boldsymbol{v}^{(i)}} \leftarrow \mathcal{C}\left(\boxed{\boldsymbol{g}^{(i)} + \boldsymbol{\delta}^{(i)}}\right)$$

Error compensation

$$\boldsymbol{\delta}^{(i)} \leftarrow \left(\boldsymbol{g}^{(i)} + \boldsymbol{\delta}^{(i)}\right) - \boldsymbol{v}^{(i)}$$

Compression error

$$\boldsymbol{x} \leftarrow \boldsymbol{x} - \gamma \overline{\boldsymbol{v}}$$

Pull $\overline{\mathbf{v}}$ to update model

Master s

$$\overline{\boldsymbol{g}} \leftarrow \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{v}^{(i)}$$

Aggregate compressed gradient

$$\boxed{\overline{\boldsymbol{v}}} \leftarrow \mathcal{C}\left(\boxed{\overline{\boldsymbol{g}} + \overline{\boldsymbol{\delta}}}\right)$$

Compress the error compensated g

$$\overline{\boldsymbol{\delta}} \leftarrow \left(\overline{\boldsymbol{g}} + \overline{\boldsymbol{\delta}}\right) - \overline{\boldsymbol{v}}$$

Compression error

# Intuition

Essential updating rule of DoubleSqueeze (SGD alike)

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_t - \gamma \overline{\boldsymbol{g}}_t + \gamma(\hat{\boldsymbol{\delta}}_t - \hat{\boldsymbol{\delta}}_{t-1})$$

$$\overline{\boldsymbol{g}}_t = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{g}_t^i$$

$$\hat{\boldsymbol{\delta}}_t = \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{\delta}_t^{(i)} + \overline{\boldsymbol{\delta}}_t$$

C-SGD (Uncompressed)

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_0 - \gamma \sum_{s=0}^{t} \overline{\boldsymbol{g}}_s$$

Naive Compressed C-SGD

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_0 - \gamma \sum_{s=0}^{t} \overline{\boldsymbol{g}}_s + \gamma \sum_{s=0}^{t} \hat{\boldsymbol{\delta}}_s$$

DoubleSqueeze

$$\boldsymbol{x}_{t+1} = \boldsymbol{x}_0 - \gamma \sum_{s=0}^{t} \overline{\boldsymbol{g}}_s + \gamma \hat{\boldsymbol{\delta}}_t$$  **Much smaller**

# Convergence

Assumption
$$\mathbb{E}[\|\mathbf{C}(\mathbf{x}) - \mathbf{x}\|^2] \leq \sigma'^2$$

Convergence rates

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left(\|\nabla f(\overline{\boldsymbol{x}}_t)\|^2\right) \lesssim \frac{1}{T} + \frac{\sigma}{\sqrt{nT}}$$

SGD

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left(\|\nabla f(\overline{\boldsymbol{x}}_t)\|^2\right) \lesssim \frac{1}{T} + \frac{\sigma}{\sqrt{nT}} + \boxed{\frac{\sigma'}{\sqrt{T}}}$$

C-SGD (C(.) needs to be unbiased)

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left(\|\nabla f(\overline{\boldsymbol{x}}_t)\|^2\right) \lesssim \frac{1}{T} + \frac{\sigma}{\sqrt{nT}} + \boxed{\left(\frac{\sigma'}{T}\right)^{\frac{2}{3}}}$$
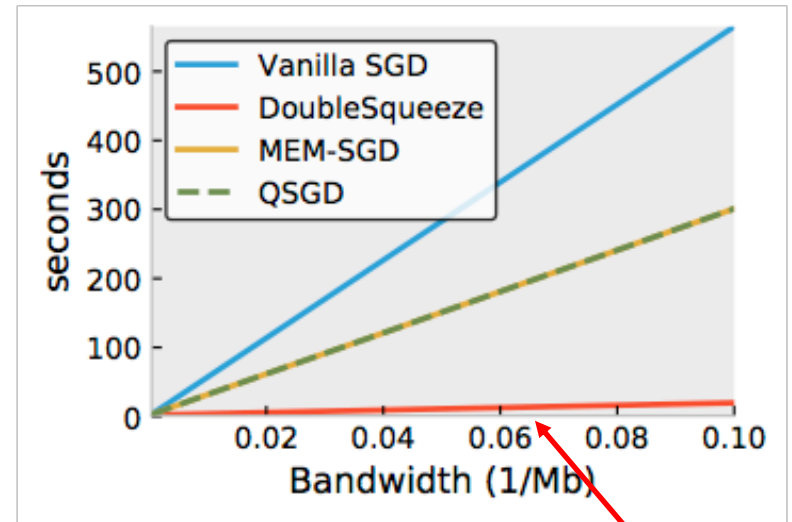
Double squeeze
EC-SGD
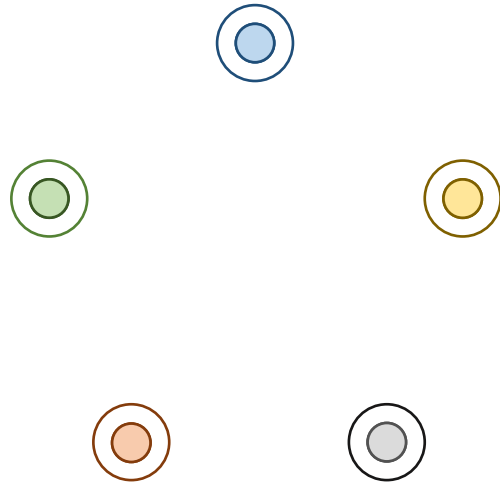
better

# Experiments

ResNet-18. CIFAR-10. 8 workers



Iteration (epoch) is consistent with SGD
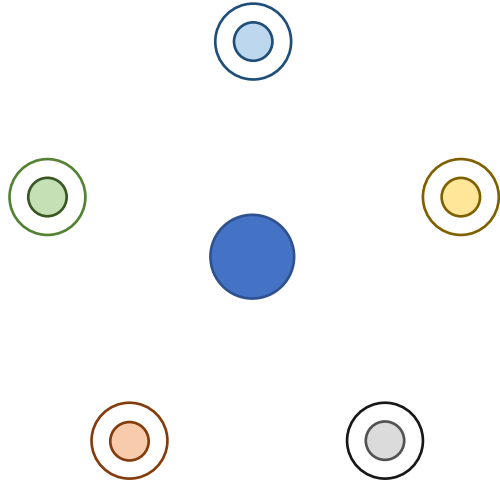


Running time in each iteration is faster
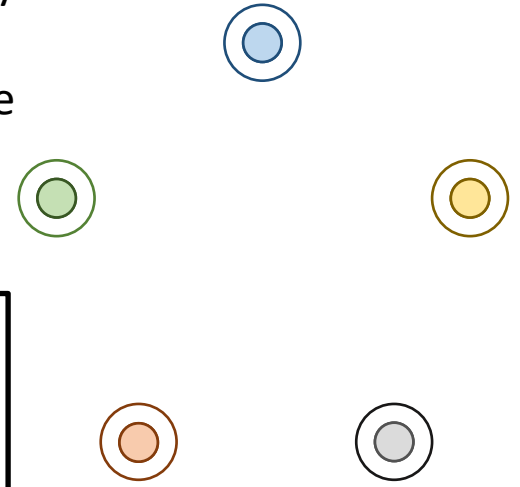
# Decentralization

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

How does the **decentralized** approach compare to the **centralized** approach?

Centralized communication
（fully exchanged）

O(N * alpha + NB * beta)

Decentralized communication
(partially exchanged)

O(alpha + B * beta)

Centralized-SGD:

$$\boxed{\boldsymbol{x}} \leftarrow \boldsymbol{x} - \gamma \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{g}(\boldsymbol{x}; \boxed{\boldsymbol{a}^{(i)}})$$

shared model

local sample

Centralized-SGD:

$$\boldsymbol{x} \leftarrow \boldsymbol{x} - \gamma \frac{1}{n} \sum_{i=1}^{n} \boldsymbol{g}(\boldsymbol{x}; \boldsymbol{a}^{(i)})$$

Decentralized-SGD:

Local model   Local sample

$$\begin{bmatrix} \boldsymbol{x}^{(1)} \\ \boldsymbol{x}^{(2)} \\ \dots \\ \boldsymbol{x}^{(n)} \end{bmatrix} \leftarrow W \left( \begin{bmatrix} \boldsymbol{x}^{(1)} \\ \boldsymbol{x}^{(2)} \\ \dots \\ \boldsymbol{x}^{(n)} \end{bmatrix} - \gamma \begin{bmatrix} \boldsymbol{g}(\boldsymbol{x}^{(1)}; \boldsymbol{a}^{(1)}) \\ \boldsymbol{g}(\boldsymbol{x}^{(2)}; \boldsymbol{a}^{(2)}) \\ \dots \\ \boldsymbol{g}(\boldsymbol{x}^{(n)}; \boldsymbol{a}^{(n)}) \end{bmatrix} \right)$$
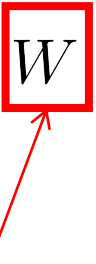
n individual models

Average the local model with neighbor's, e.g.,

$$\mathbf{x}^{(2)} \leftarrow \frac{1}{3} \sum_{i=1,2,3} \left( \mathbf{x}^{(i)} - \gamma \mathbf{g}^{(i)} \right)$$

Decentralized SGD

$$
\begin{bmatrix} \boldsymbol{x}^{(1)} \\ \boldsymbol{x}^{(2)} \\ \cdots \\ \boldsymbol{x}^{(n)} \end{bmatrix} \leftarrow \boxed{W} \left( \begin{bmatrix} \boldsymbol{x}^{(1)} \\ \boldsymbol{x}^{(2)} \\ \cdots \\ \boldsymbol{x}^{(n)} \end{bmatrix} - \gamma \begin{bmatrix} \boldsymbol{g}(\boldsymbol{x}^{(1)}; \boldsymbol{a}^{(1)}) \\ \boldsymbol{g}(\boldsymbol{x}^{(2)}; \boldsymbol{a}^{(2)}) \\ \cdots \\ \boldsymbol{g}(\boldsymbol{x}^{(n)}; \boldsymbol{a}^{(n)}) \end{bmatrix} \right)
$$

weight matrix: symmetric, doubly stochastic
(W1 = 1, W$^T$1=1, nonnegative, W=W$^T$)

ring network

$$
W = \begin{pmatrix} 1/3 & 1/3 & & & & 1/3 \\ 1/3 & 1/3 & 1/3 & & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1/3 & 1/3 & 1/3 \\ 1/3 & & & & 1/3 & 1/3 \end{pmatrix}
$$

## Assumptions

- **Lipschitzian** All $f_i(\cdot)$ are with $L$-Lipschitzian gradient     data variance **within** each worker
- **Bounded variance**

$$\mathbb{E}_{\boldsymbol{a}\sim\mathcal{D}_i}\|\nabla F(\boldsymbol{x};\boldsymbol{a}) - \nabla f_i(\boldsymbol{x})\|^2 \leq \sigma^2,\ \forall i,\ \forall \boldsymbol{x}$$

$$\|\nabla f_i(\boldsymbol{x}) - \nabla f(\boldsymbol{x})\|^2 \leq \zeta^2,\ \forall i,\ \forall \boldsymbol{x}$$

data variance **among** workers

### *Assumptions*

- **Lipschitzian** All $f_i(\cdot)$ are with $L$-Lipschitzian gradient
- **Bounded variance**

$$\mathbb{E}_{\boldsymbol{a} \sim \mathcal{D}_i} \|\nabla F(\boldsymbol{x}; \boldsymbol{a}) - \nabla f_i(\boldsymbol{x})\|^2 \leq \sigma^2, \; \forall i, \; \forall \boldsymbol{x}$$

$$\|\nabla f_i(\boldsymbol{x}) - \nabla f(\boldsymbol{x})\|^2 \leq \zeta^2, \; \forall i, \; \forall \boldsymbol{x}$$

- **Spectral gap**

$$\rho := \max_{j \geq 2} |\lambda_j(W)|$$

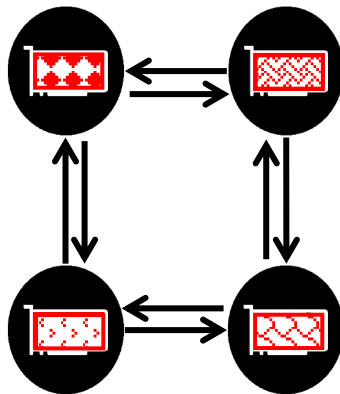Measure how fast the information can spread across the network

Fully connected network

Ring network

Disconnected network

$$W = \frac{\mathbf{1}\mathbf{1}^\top}{N}$$
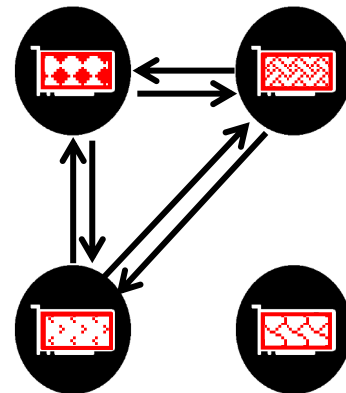
$$W = \begin{pmatrix} 1/3 & 1/3 & & & & 1/3 \\ 1/3 & 1/3 & 1/3 & & & \\ & & \ddots & \ddots & \ddots & \\ & & & 1/3 & 1/3 & 1/3 \\ 1/3 & & & & 1/3 & 1/3 \end{pmatrix}$$

$$W = \begin{pmatrix} D & \mathbf{0} \\ \mathbf{0}^\top & 1 \end{pmatrix}$$

$$\rho = 0$$

$$\rho \approx \left(1 - \frac{16\pi^2}{3N^2}\right)$$
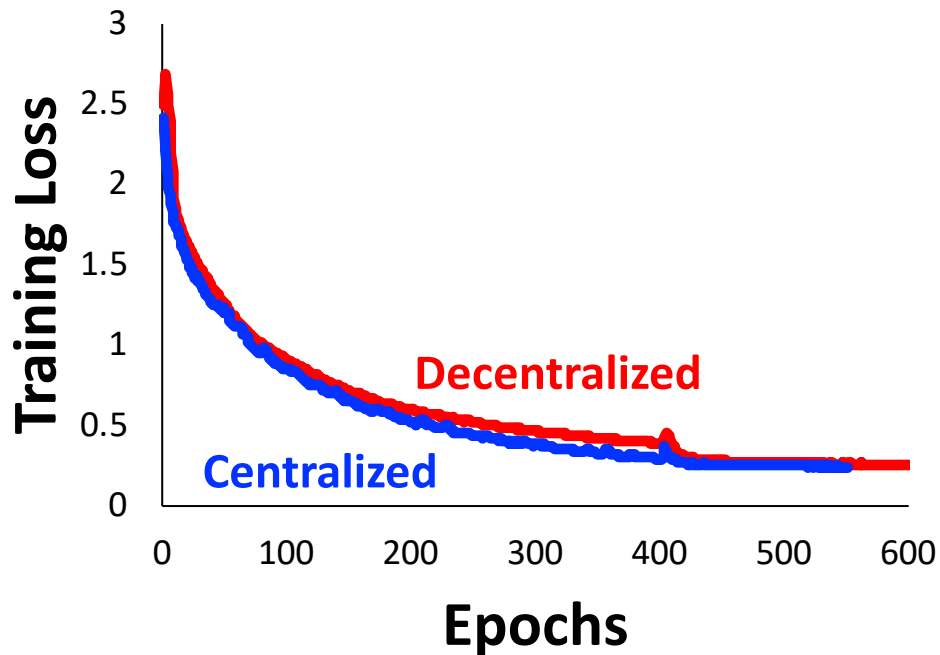
$$\rho = 1$$

**_Theorem [DSGD]_** Choose the learning rate approximately. When T is sufficiently large, we have

$$\frac{1}{T}\sum_{t=1}^{T}\mathbb{E}\left(\|\nabla f(\overline{\boldsymbol{x}}_t)\|^2\right) \lesssim \frac{1}{T} + \frac{\sigma}{\sqrt{nT}} + \left(\frac{\zeta\rho}{T(1-\rho)}\right)^{\frac{2}{3}}$$
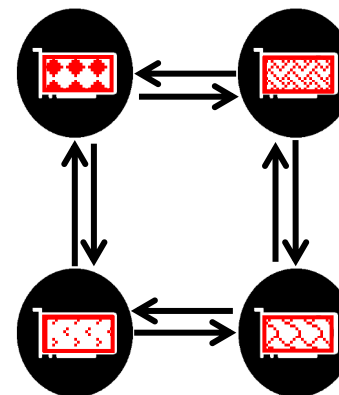
Average of local models

Convergence rate of CSGD

Cost of using decentralized communication (minor)
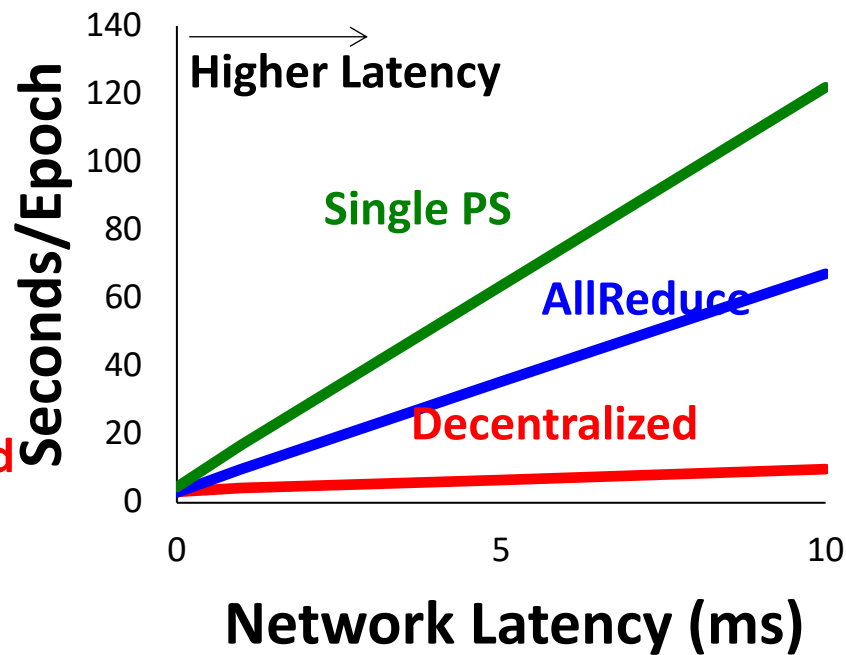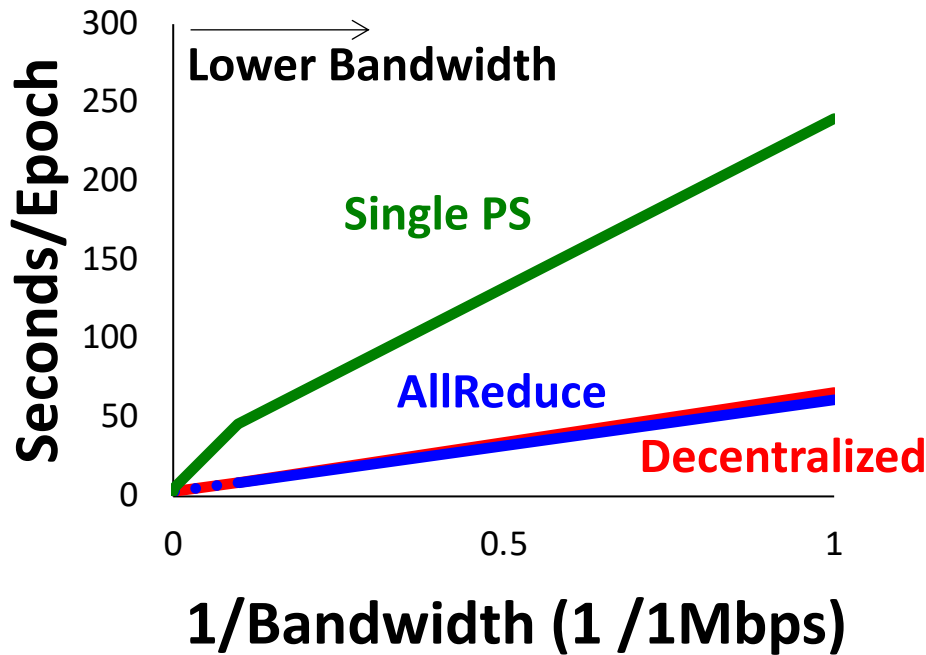
DECENTRALIZED METHOD
Ring Topology

Training Loss vs Epochs chart showing Decentralized (red) and Centralized (blue) curves.

Centralized includes PS and AllReduce!

100 GPUs

ResNet

CIFAR10

Decentralized algorithms **outperform** centralized algorithms for networks with low bandwidth and high latency
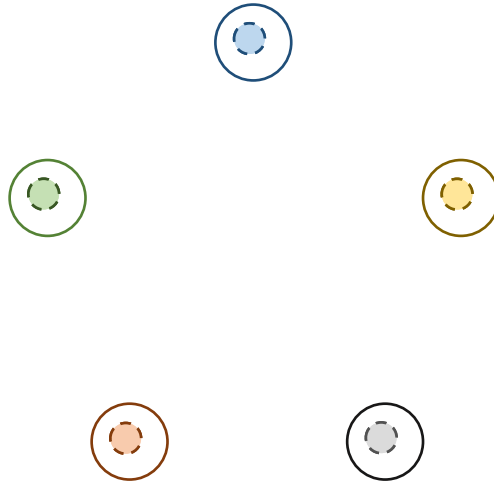
**Theoretical** view | Decentralized-SGD achieves the same convergence rate as Centralized-PSGD

**Practical** view | When the network is with high latency, decentralized communication can outperform its centralized counterpart.

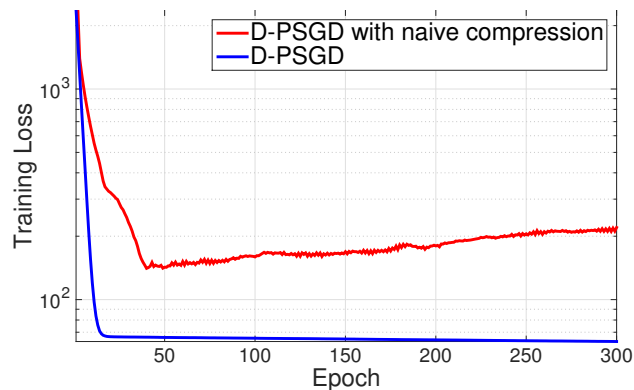# Compression + Decentralization

# Naïve compression does not work

Can we further reduce the communication cost?

Naïve compression for D-SGD

$$\boldsymbol{x}_{t+1}^{(i)} = \sum_{j} W_{ij} \mathcal{C}\left(\boldsymbol{x}_t^{(j)}\right) - \gamma \nabla F(\boldsymbol{x}_t^{(i)}; \boldsymbol{a}^{(i)})$$

# DCD-SGD

Store a copy of its neighbors' models

$$\hat{\boldsymbol{x}}_{t+1}^{(i)} = \sum_j W_{ij} \boxed{\boldsymbol{x}_t^{(j)}} - \gamma \nabla F(\boldsymbol{x}_t^{(i)}; \boldsymbol{a}^{(i)})$$

$$\boldsymbol{x}_{t+1}^{(i)} = \hat{\boldsymbol{x}}_t^{(i)} + \boxed{\mathcal{C}\left(\hat{\boldsymbol{x}}_t^{(i)} - \boldsymbol{x}_t^i\right)}$$

Compress the difference and send to its neighbors

$$\sup_{\boldsymbol{x}} \frac{\mathbb{E}\|\mathcal{C}(\boldsymbol{x}) - \boldsymbol{x}\|^2}{\|\boldsymbol{x}\|^2} \leq \alpha^2$$

$$\mathbb{E}(\mathcal{C}(\mathbf{x})) = \mathbf{x}$$

$$\frac{1}{T} \sum_{t=1}^{T} \mathbb{E}\left(\|\nabla f(\overline{\boldsymbol{x}}_t)\|^2\right) \lesssim \frac{1}{T} + \frac{\sigma(1 + \alpha^2)}{\sqrt{nT}} + \frac{\zeta^{\frac{2}{3}}(1 + \alpha^2)}{T^{\frac{2}{3}}}$$
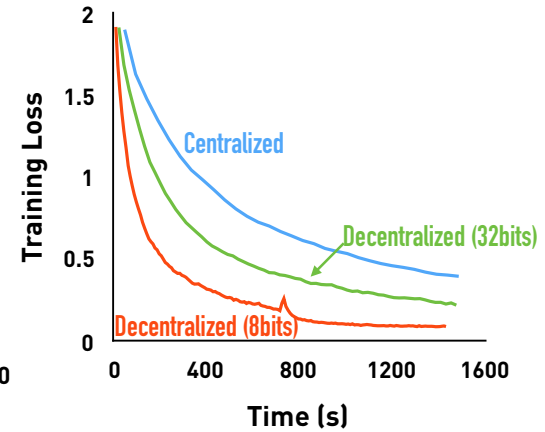
Consistent with D-SGD

*[NIPS 2018]*

(b) Time vs Training Loss
Bandwidth = 1.4Gbps,
Latency = 0.13ms

(c) Time vs Training Loss
Bandwidth = 1.4Gbps,
Latency = 20ms

(d) Time vs Training Loss
Bandwidth = 5Mbps,
Latency = 20ms

**Two Issues** of DCD-SGD:

Require  $\mathbb{E}(\mathcal{C}(\mathbf{x})) = \mathbf{x}$

Diverges when using 4-bit compression in most cases

Can we fix it by using error compression strategy?

$$X_{t+1} = (X_t - \gamma G_t)W$$
$$= X_t - \gamma G_t + \boxed{(X_t - \gamma G_t)}(W - I)$$

**Share this with Error Compensation**

DCD-SGD + Error Compensation

**Local:**

Error Compensation

$$\boxed{\boldsymbol{v}_{t+1}^{(i)}} = \mathcal{C}\left(\boldsymbol{x}_t^{(i)} - \gamma \boldsymbol{g}_t^{(i)} + \boxed{\boldsymbol{\delta}_t^{(i)}}\right)$$

$$\boldsymbol{\delta}_{t+1}^{(i)} = \boldsymbol{v}_{t+1}^{(i)} - \left(\boldsymbol{x}_t^{(i)} - \gamma \boldsymbol{g}_t^{(i)} + \boldsymbol{\delta}_t^{(i)}\right)$$

**Communicate:**

$$\boldsymbol{x}_{t+1}^{(i)} = \boldsymbol{x}_t^{(i)} - \gamma \boldsymbol{g}_t^{(i)} + \boxed{\eta} \sum_{j \in \mathcal{N}_i} (W_{ij} - I_{ij}) \boxed{\boldsymbol{v}_{t+1}^{(j)}}$$

Control the compression error explicitly

**DCD-SGD**

$$\sup_{\boldsymbol{x}} \frac{\mathbb{E}\|\mathcal{C}(\boldsymbol{x}) - \boldsymbol{x}\|^2}{\|\boldsymbol{x}\|^2} \le \alpha^2$$

$$\mathbb{E}(\mathcal{C}(\mathbf{x})) = \mathbf{x}$$

**Fails for 4-bit compression**

$$\mathcal{O}\left(\frac{1}{T} + \frac{\sigma(1 + \alpha^2)}{\sqrt{nT}} + \frac{\zeta^{\frac{2}{3}}(1 + \alpha^2)}{T^{\frac{2}{3}}}\right)$$
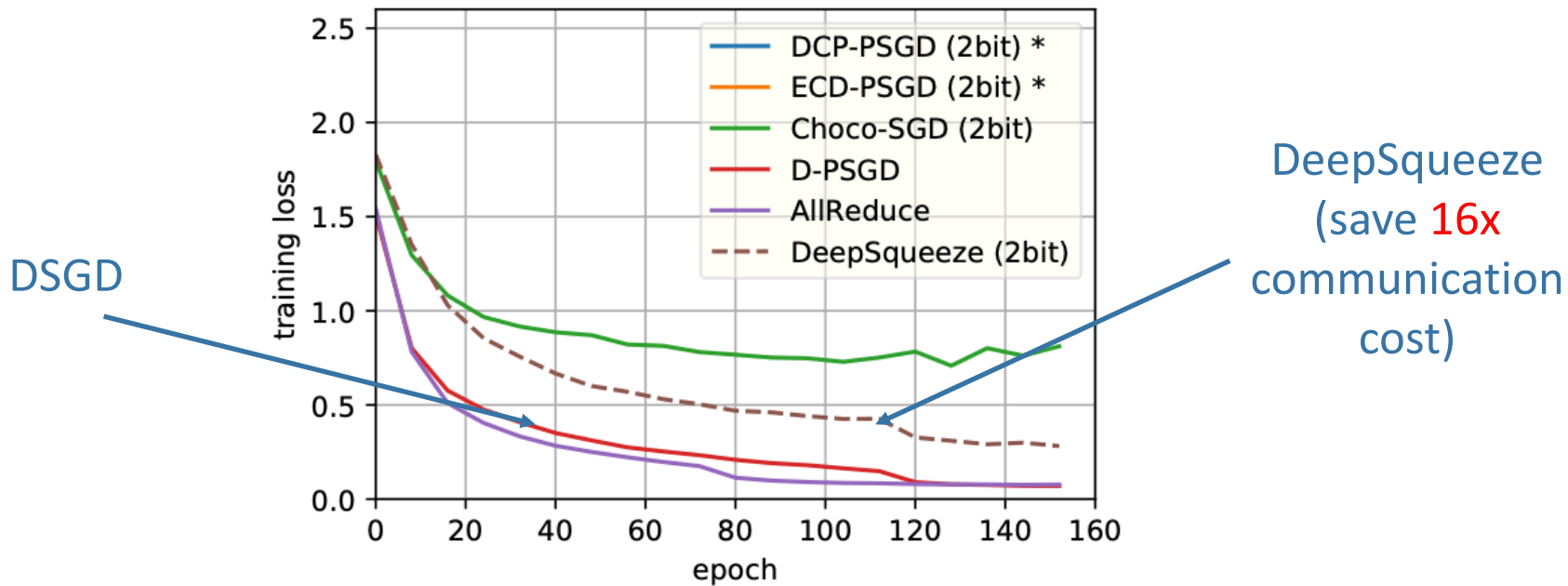
**DeepSqueeze**

$$\sup_{\boldsymbol{x}} \frac{\mathbb{E}\|\mathcal{C}(\boldsymbol{x}) - \boldsymbol{x}\|^2}{\|\boldsymbol{x}\|^2} \le \alpha^2$$

**Compression can be biased**

**Robust to 2-bit compression**

$$\mathcal{O}\left(\frac{1}{T} + \frac{\sigma\left(1 + \frac{\alpha^2\sqrt{n}}{\sqrt{T}}\right)}{\sqrt{nT}} + \frac{\zeta^{\frac{2}{3}}(1 + \alpha^2)}{T^{\frac{2}{3}}}\right)$$

# Experiments



DSGD

DeepSqueeze
(save 16x
communication
cost)

# Summary



Compression

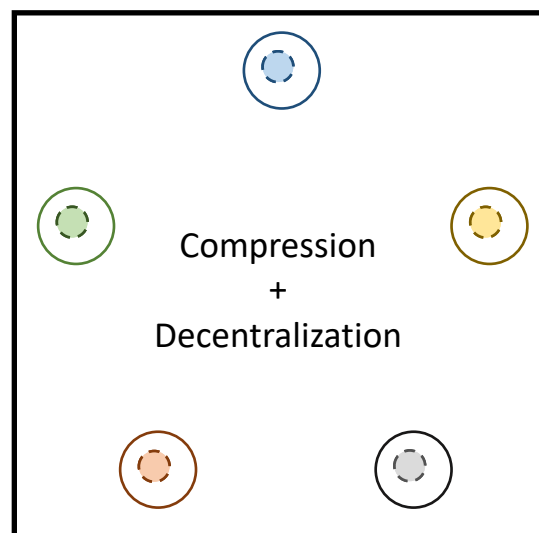Asynchronization

Decentralization

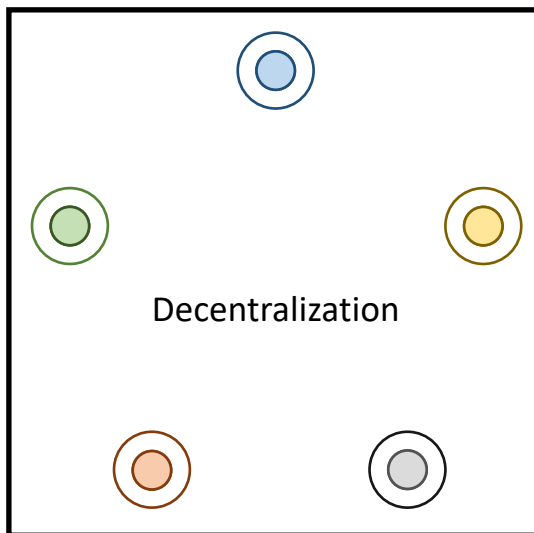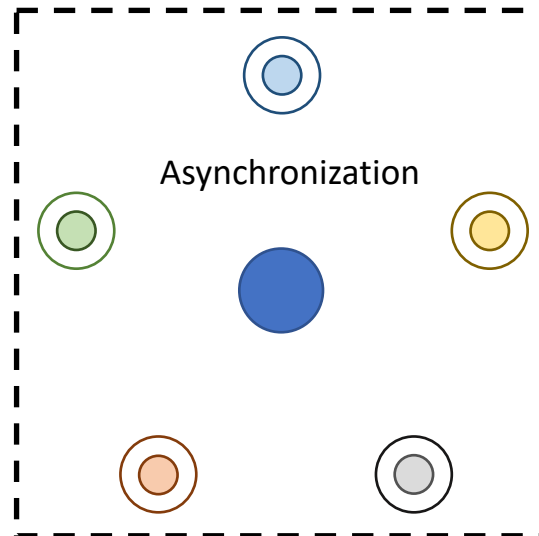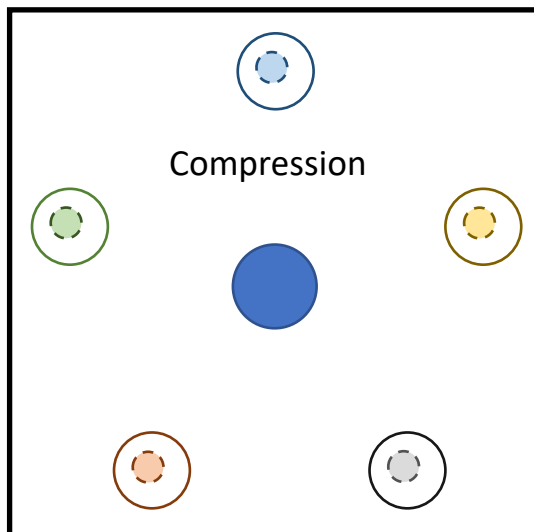Compression
+
Decentralization

**Ji Liu**
University of Rochester and
Kuaishou Inc., USA
ji.liu.uwisc@gmail.com

**Ce Zhang**
ETH Zurich, Switzerland
ce.zhang@inf.ethz.ch

*Coming Soon.*

**now**
the essence of knowledge
Boston — Delft

# Q&A