



# Distributed Machine Learning (Part I)

Kun Yuan

Center for Machine Learning Research @ Peking University

## DNN model is notoriously difficult to train

- Highly-nonconvex cost functions; cannot find global minima; trapped into local minimum
- The model size is large, i.e.,  $x \in \mathbb{R}^d$  is of extremely high dimensions
- The size of the dataset is huge

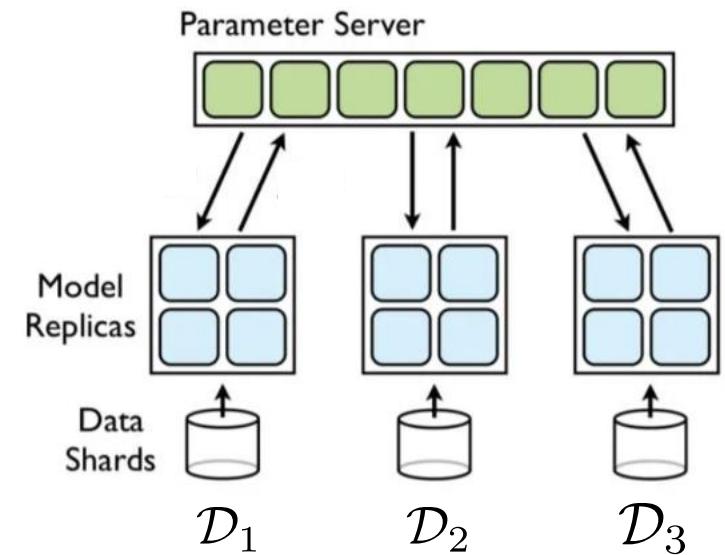
DNN training = Non-convex training + Huge dimensions + Huge dataset

- **Efficient** and **scalable** distributed learning approaches are in urgent need

A network of  $n$  nodes (devices such as GPUs) collaborate to solve the problem:

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad \text{where } f_i(x) = \mathbb{E}_{\xi_i \sim D_i} F(x; \xi_i)$$

- Each component  $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$  is local and private to node  $i$
- Random variable  $\xi_i$  denotes local data that follows distribution  $D_i$
- Each local distribution  $D_i$  is different; data heterogeneity exists
- This talk will sorely focus on the **data parallelism** problem



# Vanilla parallel stochastic gradient descent (PSGD)

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad \text{where} \quad f_i(x) = \mathbb{E}_{\xi_i \sim D_i} F(x; \xi_i).$$

**PSGD**

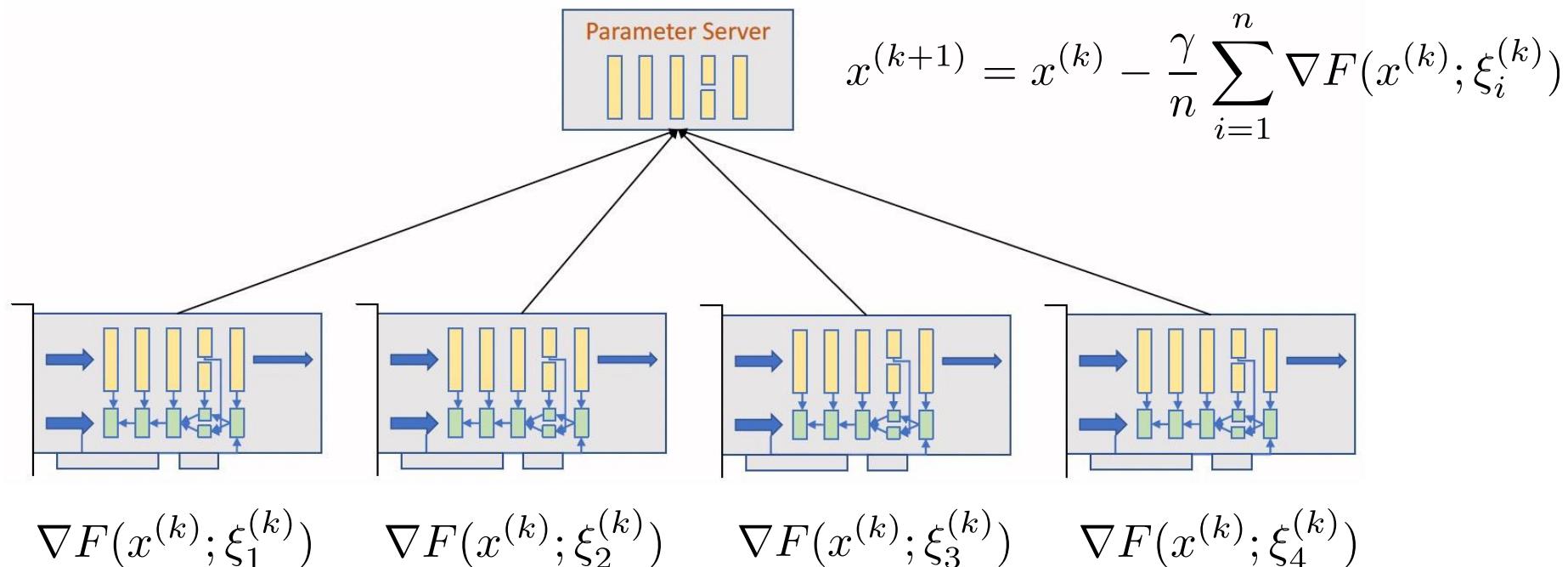
$$g_i^{(k)} = \nabla F(x^{(k)}; \xi_i^{(k)}) \quad (\text{Local compt.})$$

$$x^{(k+1)} = x^{(k)} - \frac{\gamma}{n} \sum_{i=1}^n g_i^{(k)} \quad (\text{Global comm.})$$

- Each node  $i$  samples data  $\xi_i^{(k)}$  and computes gradient  $\nabla F(x^{(k)}; \xi_i^{(k)})$
- All nodes synchronize (i.e. globally average) to update model  $x$  per iteration

# Vanilla parallel stochastic gradient descent (PSGD)

- Federated learning typically implements PSGD using parameter server



- LLM training within data-centers implements PSGD using Ring-Allreduce

## Assumption [PSGD assumption]

- (1) Each  $f_i(x)$  is  $L$ -smooth, i.e.,  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$  for any  $x, y$ ;
- (2) Each stochastic gradient is unbiased and has bounded variance, i.e.,

$$\mathbb{E}[\nabla F(x; \xi_i)] = \nabla f_i(x), \quad \mathbb{E}\|\nabla F(x; \xi_i) - \nabla f_i(x)\|^2 \leq \sigma^2$$

## Theorem [PSGD convergence]

Under the above assumptions and with proper  $\gamma$ , PSGD converges as follows

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E}\|\nabla f(x^{(k)})\|^2 = \mathcal{O}\left(\frac{\sigma}{\sqrt{nT}}\right)$$

[K. Yuan, Lecture 6: stochastic gradient descent, PKU Class “Optimization for deep learning”, check Kun Yuan’s website]

## Theorem [PSGD convergence]

Under the above assumptions and with proper  $\gamma$ , PSGD converges as follows

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(x^{(k)})\|^2 = \mathcal{O}\left(\frac{\sigma}{\sqrt{nT}}\right)$$

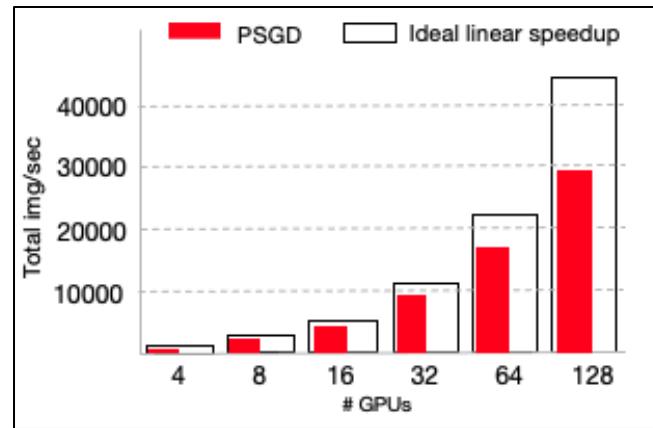
- This implies that to achieve an  $\epsilon$ -accurate solution, PSGD needs

$$\frac{\sigma}{\sqrt{nT}} \leq \epsilon \implies T \geq \frac{\sigma^2}{n\epsilon^2}$$

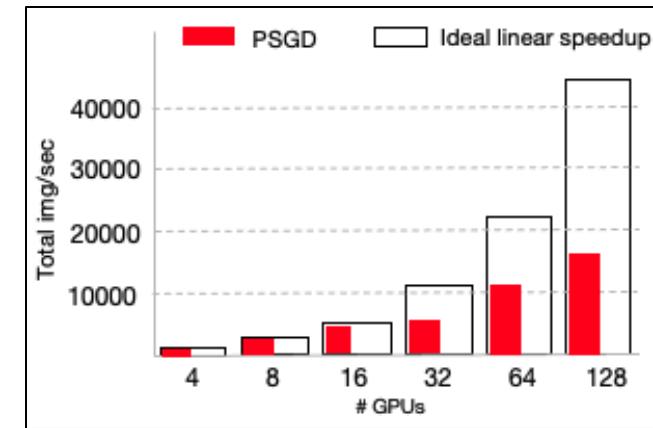
iterations, which decreases linearly with number of nodes  $n$ ; this is called **linear speedup**

# PSGD cannot achieve ideal scalability due to comm. overhead

- PSGD **cannot** achieve ideal linear speedup in throughput due to comm. overhead
- Larger comm-to-compt ratio leads to worse performance in PSGD



Small comm.-to-compt. ratio



Large comm.-to-compt. ratio

Ring-Allreduce is used in each experiment

- How can we reduce the communication overhead in PSGD?

- Global average incurs  $O(n)$  comm. overhead; proportional to network size n

## [Decentralized communication]

- Each node interacts with the server at every iteration; proportional to iteration numbers

## [Lazy communication]

- Each node sends a full model (or gradient) to the server; proportional to dimension d

## [Compressed communication]

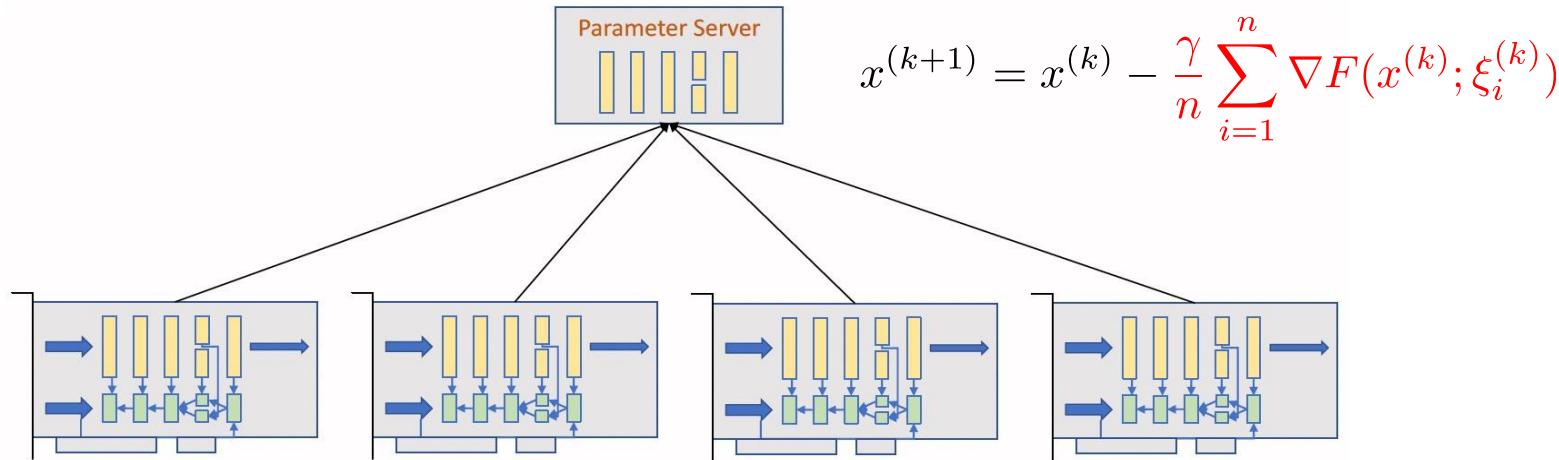
- and more (asynchronous communication; robust communication against Byzantine nodes, etc.)

## PART 03-1

---

### Decentralized SGD: communication efficiency

## Parameter-server architecture



- Global average using PS incurs  $O(n)$  comm. overhead; **proportional to network size n**
- Global average using Ring-Allreduce incurs  $O(n)$  latency; **proportional to network size n**
- When network size n is large (say  $n \geq 10000$ ), PSGD suffers severe communication overhead

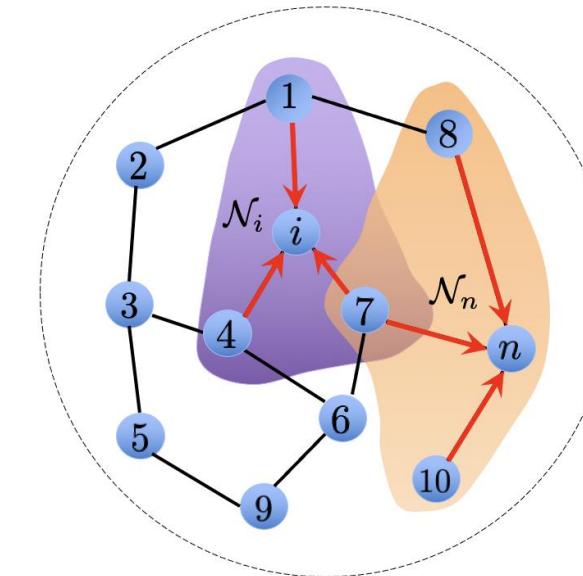
# Decentralized SGD (DSGD)

- To break  $O(n)$  comm. overhead, we replace global average with partial average

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{(k+\frac{1}{2})} \quad (\text{Partial averaging})$$

- DSGD = local SGD update + partial averaging [1, 2]
- $\mathcal{N}_i$  is the set of neighbors at node  $i$ ;  $w_{ij}$  scales information from  $j$  to  $i$
- Incurs  $O(d_{\max})$  comm. overhead per iteration where  $d_{\max} = \max_i \{|\mathcal{N}_i|\}$  is the graph maximum degree

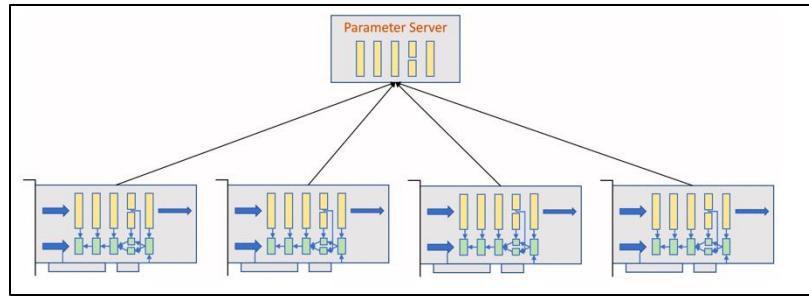


[1] C. G. Lopes and A. H. Sayed, "Diffusion least-mean-squares over adaptive networks", ICASSP, 2007

[2] A. Nedich and A. Ozdaglar, " Distributed subgradient methods for multi-agent optimization", IEEE TAC, 2009

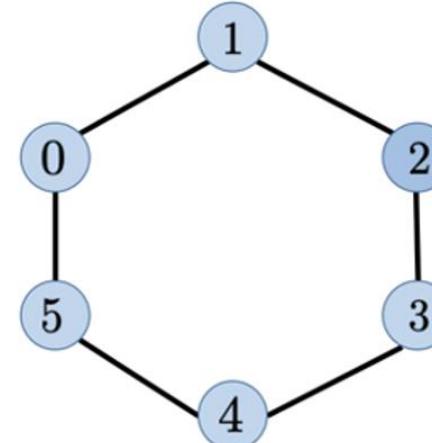
# DSGD is more communication-efficient than PSGD

- Incurs  $O(1)$  comm. overhead on **sparse** topologies; much less than global average  $O(n)$

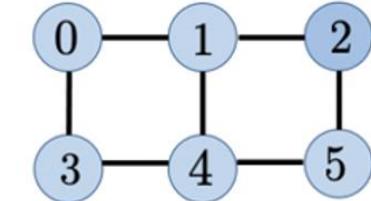


Global averaging over centralized network

comm. overhead  $O(n)$



Partial averaging over ring or grid



comm. overhead  $O(1)$

# DSGD is more communication-efficient than PSGD

---

- A real experiment on a 256-GPUs cluster [1]

Model	Ring-Allreduce	Partial average
ResNet-50 (25.5M)	278 ms	150 ms

Table. Comparison of per-iter comm. time in terms of runtime with 256 GPUs

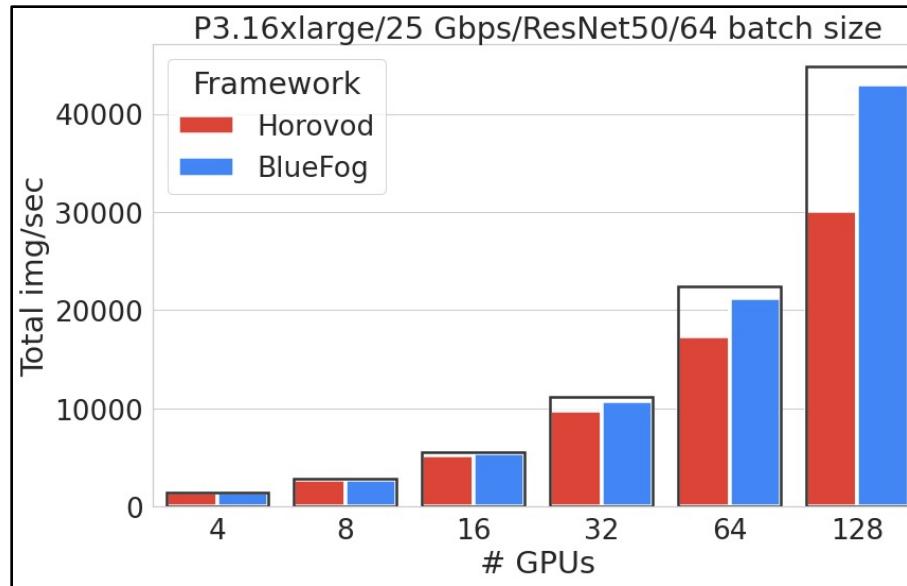
- DSGD saves more communications per iteration for larger models

---

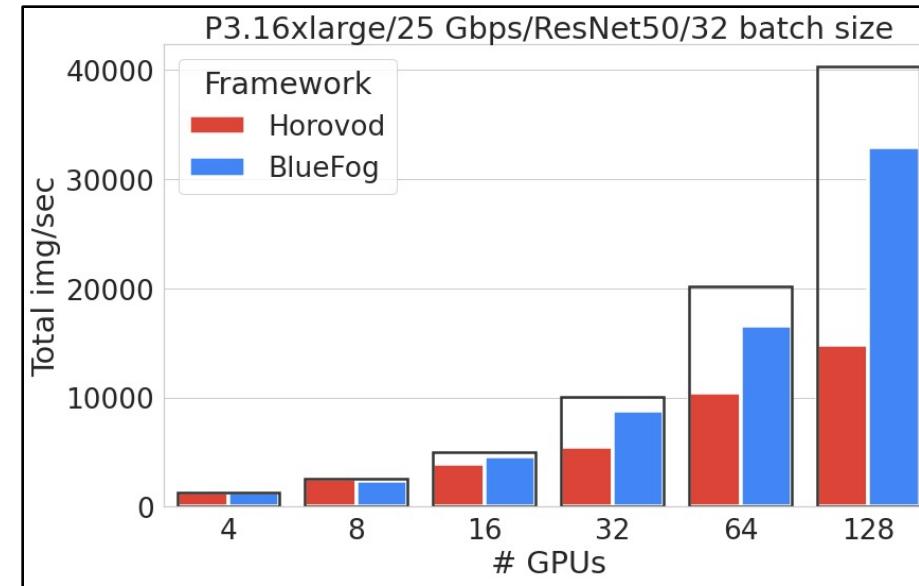
[1] Y. Chen, K. Yuan, Y. Zhang, P. Pan, Y. Xu, and W. Yin, ``Accelerating Gossip SGD with Periodic Global Averaging'', ICML 2021

# DSGD is more communication-efficient than PSGD

- DSGD (BlueFog) has **better scalability** than PSGD (Horovod) due to its small comm. overhead



Small comm.-to-compt. ratio



Large comm.-to-compt. ratio

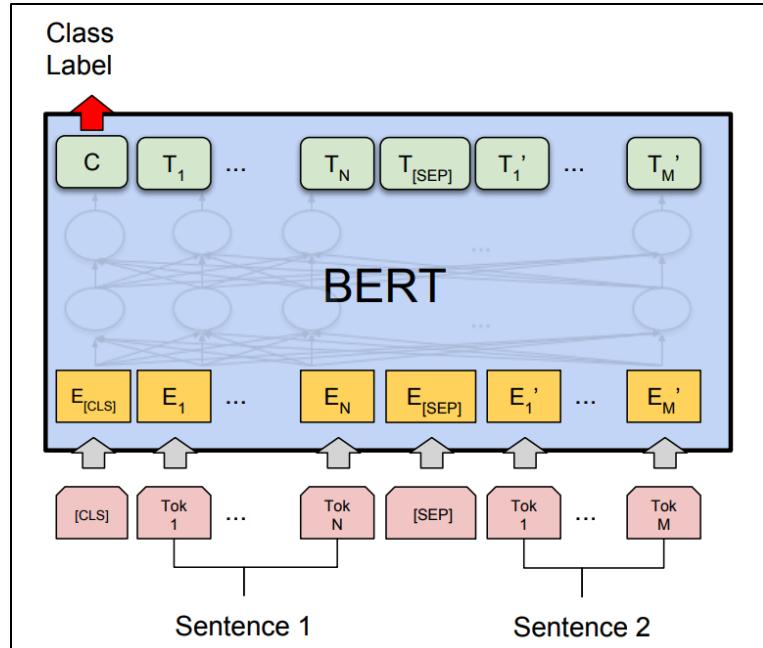
# DSGD saves more wall-clock time without severely hurting performance



nodes topology	4(4x8 GPUs) acc.	4(4x8 GPUs) time	8(8x8 GPUs) acc.	8(8x8 GPUs) time	16(16x8 GPUs) acc.	16(16x8 GPUs) time	32(32x8 GPUs) acc.	32(32x8 GPUs) time
P-SGD	76.32	11.6	76.47	6.3	76.46	3.7	76.25	2.2
Decentralized	<b>76.34</b>	<b>11.1</b>	<b>76.52</b>	<b>5.7</b>	<b>76.47</b>	<b>2.8</b>	<b>76.27</b>	<b>1.5</b>

DSGD shows very impressive linear speedup performance and saves more time than PSGD!

# DSGD saves more wall-clock time without severely hurting performance



Model: BERT-Large (330M parameters)

Dataset: Wikipedia (2500M words) and  
BookCorpus (800M words)

Hardware: 64 GPUs

Table. Comparison in loss and training time [1]

Method	Final Loss	Wall-clock Time (hrs)
P-SGD	1.75	59.02
D-SGD	1.77	30.4

[1] Y. Chen, K. Yuan, Y. Zhang, P. Pan, Y. Xu, and W. Yin, ``Accelerating Gossip SGD with Periodic Global Averaging'', ICML 2021

# The impressive performance of DSGD requires many efforts

---



- Vanilla DSGD cannot achieve strong performance against PSGD
- This lecture will highlight the efforts we devoted to making DSGD practical for real applications

## PART 03-2

---

### Decentralized SGD: convergence property

# Does decentralized SGD converge? Yes!

- Recall the PSGD and DSGD recursions

## Parallel SGD

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

$$x_i^{(k+1)} = \frac{1}{n} \sum_{j=1}^n x_j^{(k+\frac{1}{2})} \quad (\text{Global averaging})$$

## Decentralized SGD

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{(k+\frac{1}{2})} \quad (\text{Partial averaging})$$

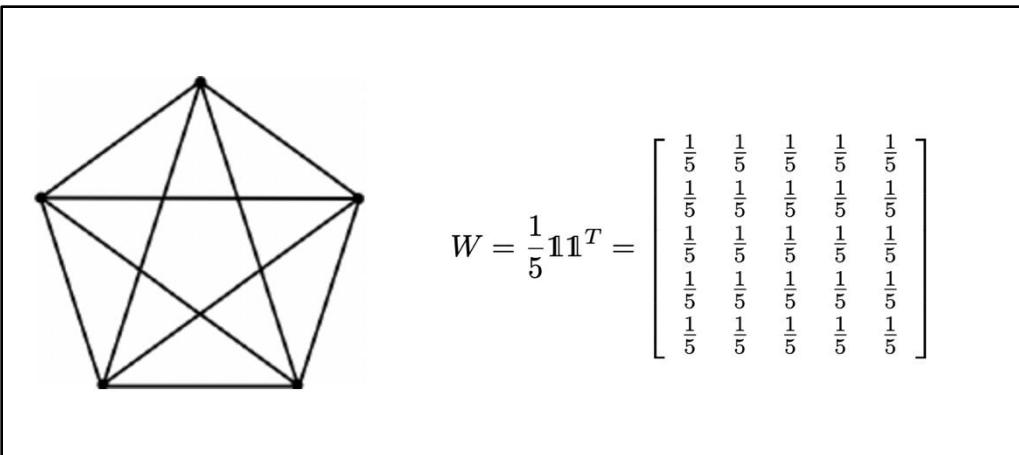
- DSGD will converge if the **partial average** asymptotically converge to the **global average**
- This argument is true if the weight matrix  $W = [w_{ij}]_{i=1,j=1}^n \in \mathbb{R}^{n \times n}$  is **doubly-stochastic**

$$W\mathbf{1}_n = \mathbf{1}_n \quad \text{and} \quad \mathbf{1}_n^T W = \mathbf{1}_n^T$$

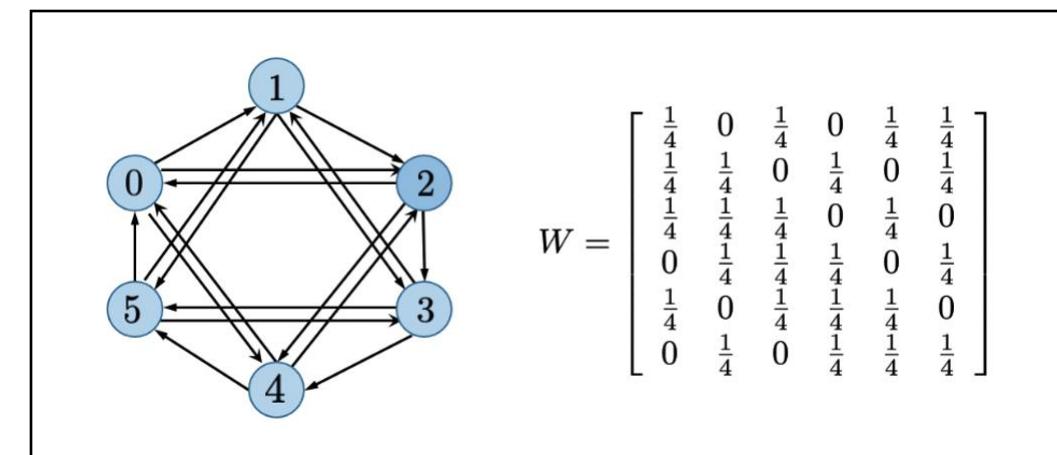
# Doubly-stochastic weight matrix is easy to construct

- Doubly-stochastic weight matrix is easy to construct over all undirected graphs and some directed graphs

Fully-connected weight matrix



Exponential weight matrix



# A numerical illustration

- We use examples to test whether partial average will converge to the global average

- Partial average:  $x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{(k)}$  for  $k = 0, 1, \dots, T$

Partial average over **ring** graph

Iter T	x1	x2	x3	x4	x5
0	0	1	2	3	4
1	1.67	1	2	3	2.33
5	1.95	1.93	2.00	2.07	2.05
10	1.998	1.996	2.000	2.000	2.002
20	1.999	1.999	2.000	2.000	2.000
50	2.000	2.000	2.000	2.000	2.000

Global average

Partial average over **exponential** graph

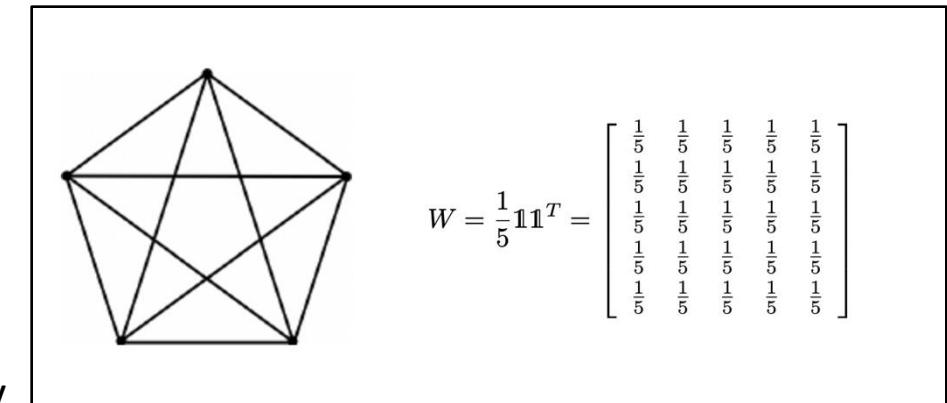
Iter T	x1	x2	x3	x4	x5
0	0	1	2	3	4
1	1.75	1.5	2.5	2.25	2.0
5	2.002	2.001	2.00	1.999	1.998
10	1.999	1.999	2.000	2.000	2.002
20	2.000	2.000	2.000	2.000	2.000

Global average

# Network topology determines the consensus rate

- The above numerical example implies that
  - Partial average **asymptotically converges** to global average
  - Network topology determines how fast** that partial average will converge to global average
- We introduce quantity  $\rho$  to gauge the graph connectivity

$$\rho = \|W - \frac{1}{n}\mathbf{1}\mathbf{1}^T\|_2 \in (0, 1) \text{ where } W = [w_{ij}] \in \mathbb{R}^{n \times n}$$



- Well-connected topology has  $\rho \rightarrow 0$ , e.g. fully-connected topology
- Sparingly-connected topology has  $\rho \rightarrow 1$ , e.g. ring has  $\rho = O(1 - \frac{1}{n^2})$

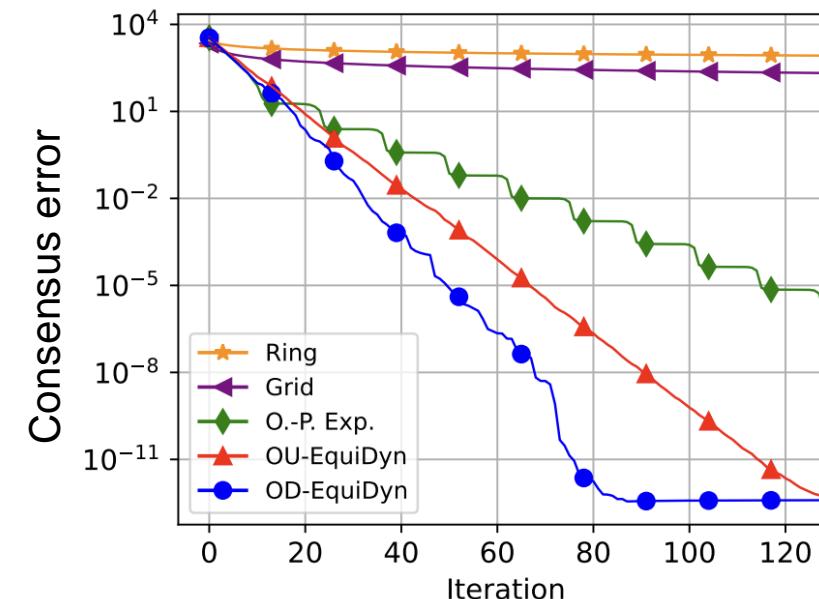
# Network topology determines the consensus rate

- With graph connectivity indicator  $\rho$ , it can be proved that

$$\|k\text{-th partial average} - \text{global average}\| \leq C\rho^k$$

Based on the above fact,  $\rho$  is also referred to as **consensus rate**

Network topology	Consensus rate $\rho$
Ring	$O(1 - \frac{1}{n^2})$
Grid	$O(1 - \frac{1}{n \ln(n)})$
Torus	$O(1 - \frac{1}{n})$
ExpoGraph	$O(1 - \frac{1}{\ln(n)})$
GeoMedian	$O(1 - \frac{\ln(n)}{n})$
Erdos-Renyi	$O(1)$
EquiGraph	$O(1)$



## Assumption [DSGD assumption]

- (1) Each  $f_i(x)$  is  $L$ -smooth, i.e.,  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$  for any  $x, y$ ;
- (2) Each stochastic gradient is unbiased and has bounded variance, i.e.,

$$\mathbb{E}[\nabla F(x; \xi_i)] = \nabla f_i(x), \quad \mathbb{E}\|\nabla F(x; \xi_i) - \nabla f_i(x)\|^2 \leq \sigma^2$$

- (3) Gradient dissimilarity can be upper bounded, i.e.,

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2 \leq b^2$$

Same assumptions  
as PSGD

Data heterogeneity  
assumption

- Quantity  $b^2$  gauges the magnitude of data heterogeneity; Larger  $b^2$  implies worse heterogeneity
- If  $b^2 = 0$ , we have  $\nabla f_i(x) = \nabla f(x)$  for any  $i$ , implying that all distributions  $\mathcal{D}_i$  are **homogeneous**

## Theorem [DSGD convergence]

Under the above assumptions and with proper  $\gamma$ , DSGD converges as follows [1]

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}} + \frac{\rho^{2/3} b^{2/3}}{T^{2/3} (1-\rho)^{2/3}} \right)$$

where  $\bar{x}^k = \frac{1}{n} \sum_{i=1}^n x_i^{(k)}$ .

- As iteration  $T \rightarrow \infty$ , DSGD converges to a stationary solution
- **Network topology** influences the convergence; Sparse topology ( $\rho \rightarrow 1$ ) results in slower convergence
- **Data heterogeneity** influences the convergence; large hetero.  $b^2$  results in slower convergence

[1] A. Koloskova, et. al., "A unified theory of decentralized sgd with changing topology and local updates", ICML 2020

# PSGD v.s. DSGD

---

$$\textbf{PSGD} \quad \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(x^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} \right)$$

$$\textbf{DSGD} \quad \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} \right)$$


---

- DSGD convergence is tight; utilizing fully-connected topology ( $\rho = 0$ ) reduces to PSGD

## asymptotic rate

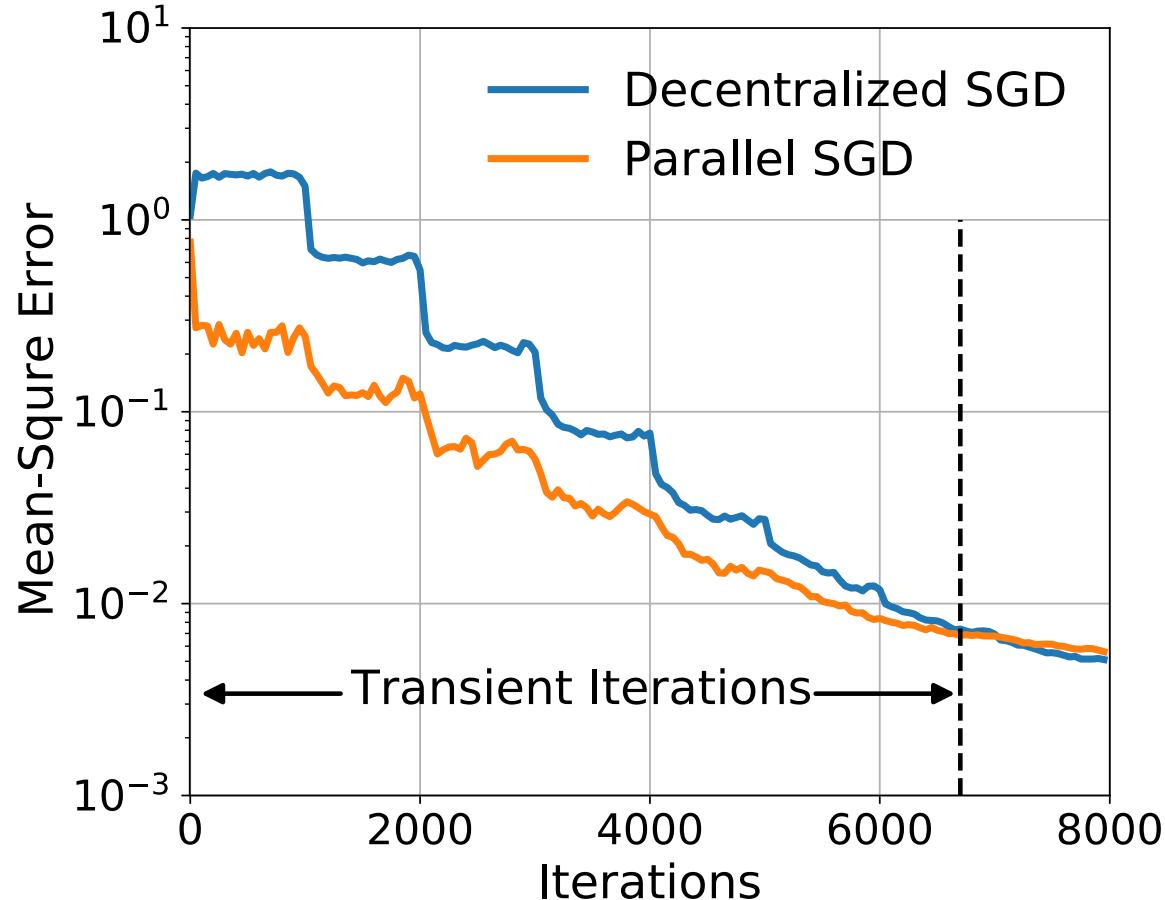
$$\text{PSGD} \quad \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(x^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} \right)$$

Extra overhead

$$\text{DSGD} \quad \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3}(1-\rho)^{1/3}} + \frac{\rho^{2/3} b^{2/3}}{T^{2/3}(1-\rho)^{2/3}} \right)$$

- DSGD convergence is tight; utilizing fully-connected topology ( $\rho = 0$ ) reduces to PSGD
- D-SGD can asymptotically converge as fast as P-SGD when  $T \rightarrow \infty$ ; the first term dominates; reach **linear speedup** asymptotically
- But DSGD **requires more iterations** to reach that stage due to the overhead caused by partial average

# PSGD v.s. DSGD



- DSGD converges asymptotically as fast as PSGD when  $T$  is sufficiently large
- DSGD has to experience transient iterations to catch up with PSGD
- In real practice, we may not be able to run very large iterations due to resource or time constraints
- We need to accelerate DSGD and reduce its transient iterations

# Transient iteration complexity

---

- Transient iteration complexity gauges the non-asymptotic stage
- **Definition:** number of iterations before a distributed algorithm achieves its linear speedup stage

DSGD convergence:  $\frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}} + \frac{\rho^{2/3} b^{2/3}}{T^{2/3} (1-\rho)^{2/3}} \right)$

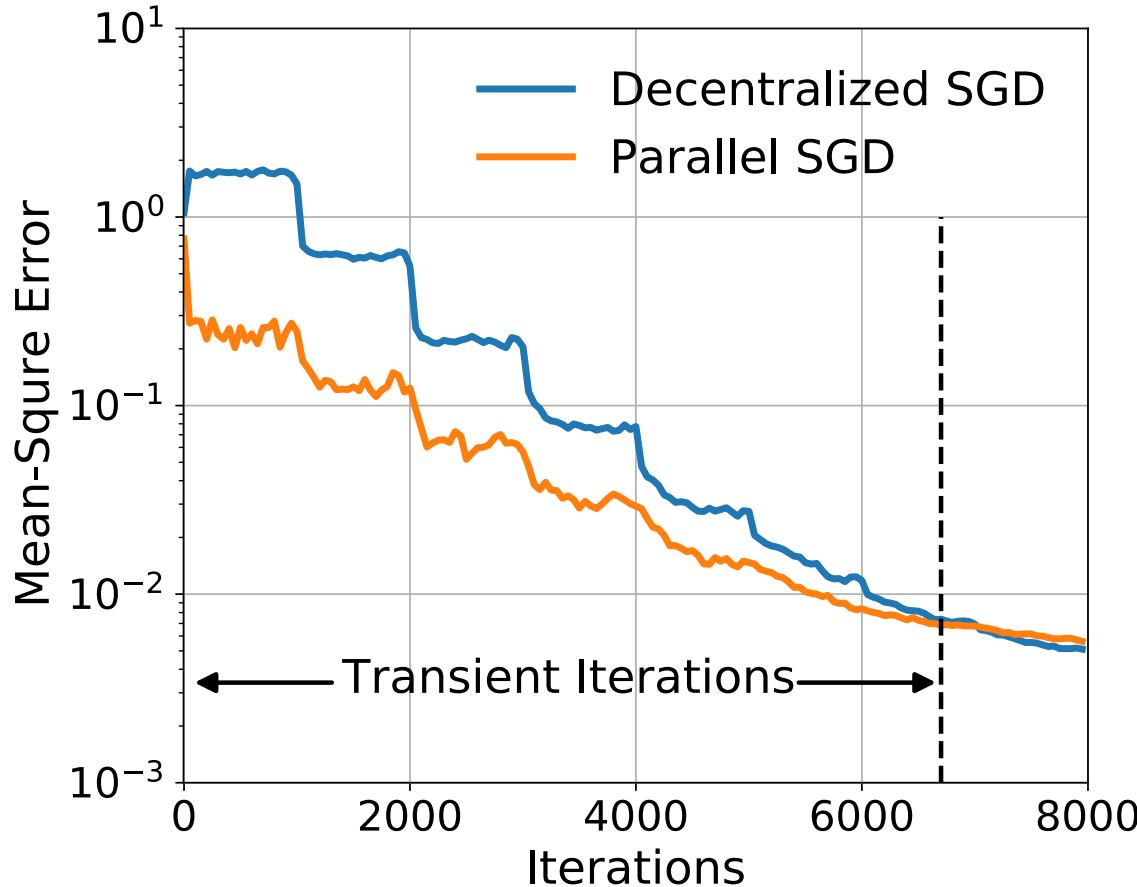
Transient iterations:

$$\frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}} \leq \frac{\sigma}{\sqrt{nT}}$$

$$\frac{\rho^{2/3} b^{2/3}}{T^{2/3} (1-\rho)^{2/3}} \leq \frac{\sigma}{\sqrt{nT}}$$
]

$$\mathcal{O} \left( \frac{\rho^4 n^3}{\sigma^2 (1-\rho)^2} + \frac{\rho^4 n^3 b^4}{\sigma^6 (1-\rho)^4} \right)$$

# Major factors that influence the transient iteration complexity



DSGD tran. iters.  $\mathcal{O} \left( \frac{\rho^4 n^3}{\sigma^2 (1 - \rho)^2} + \frac{\rho^4 n^3 b^4}{\sigma^6 (1 - \rho)^4} \right)$

- **Large data heterogeneity**  $b^2$  will significantly enlarge the transient iterations
- **Sparsely-connected network** ( $\rho \rightarrow 1$ ) will significantly enlarge the transient iterations
- Data heterogeneity and network topology will amplify the influence of each other

# A brief summary

- DSGD utilizes **partial average** in model updates

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

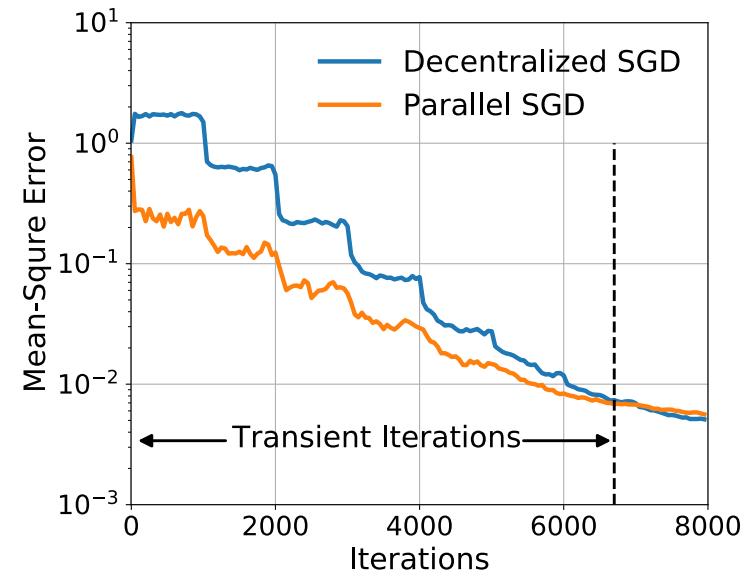
$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{(k+\frac{1}{2})} \quad (\text{Partial averaging})$$

- Partial average asymptotically converges to the global average; makes DSGD converge

- Achieve asymptotical linear speedup rate, but needs transient iterations

$$\mathcal{O} \left( \frac{\rho^4 n^3}{\sigma^2 (1-\rho)^2} + \frac{\rho^4 n^3 b^4}{\sigma^6 (1-\rho)^4} \right)$$

- In next few slides, we will reduce transient iterations by improving data heterogeneity and network topology



## PART 03-3

---

**Decentralized SGD: removing data heterogeneity**

# D-SGD suffers from data heterogeneity

- Recall the distributed stochastic optimization problem

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad \text{where} \quad f_i(x) = \mathbb{E}_{\xi_i \sim D_i} F(x; \xi_i).$$

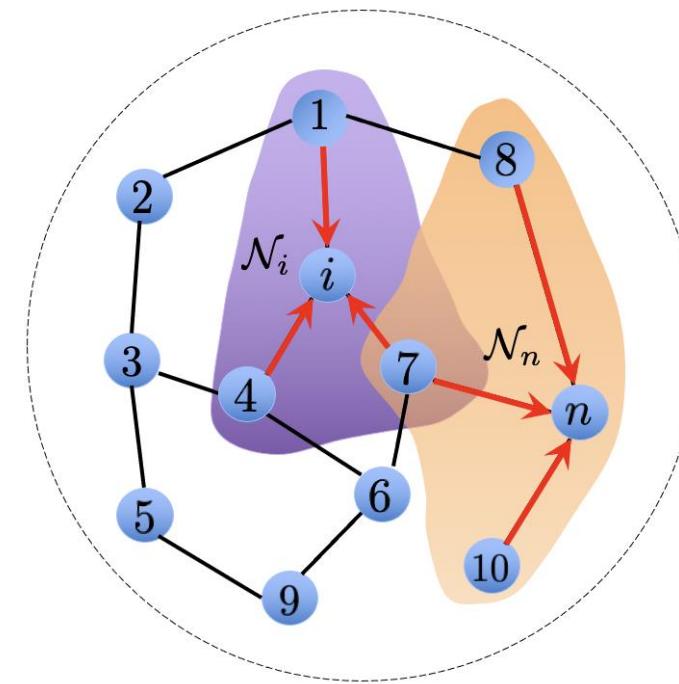
- The D-SGD algorithm iterates as follows

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{(k+\frac{1}{2})} \quad (\text{Partial averaging})$$

- We can write the above algorithm into one line

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left( x_j^{(k)} - \gamma \nabla F(x_j^{(k)}; \xi_j^{(k)}) \right)$$



# D-SGD suffers from data heterogeneity

---

- Suppose no gradient noise exists and all nodes stay at the stationary solution, i.e.,

$$\nabla F(x_i^{(k)}; \xi_i^{(k)}) = \nabla f(x_i^{(k)}) \quad \text{and} \quad x_1^{(k)} = x_2^{(k)} = \dots = x_n^{(k)} = x^*$$

- The D-SGD algorithm becomes

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} (x^* - \gamma \nabla f_j(x^*)) = x^* - \gamma \sum_{j \in \mathcal{N}_i} w_{ij} \nabla f_j(x^*)$$

- In homogeneous-data scenario where  $f_i(x) = f(x)$  for any node  $i$ , we have  $\nabla f_i(x^*) = \nabla f(x^*) = 0$  and

$$x_i^{(k+1)} = x^* - \gamma \sum_{j \in \mathcal{N}_i} w_{ij} \nabla f_j(x^*) = x^*$$

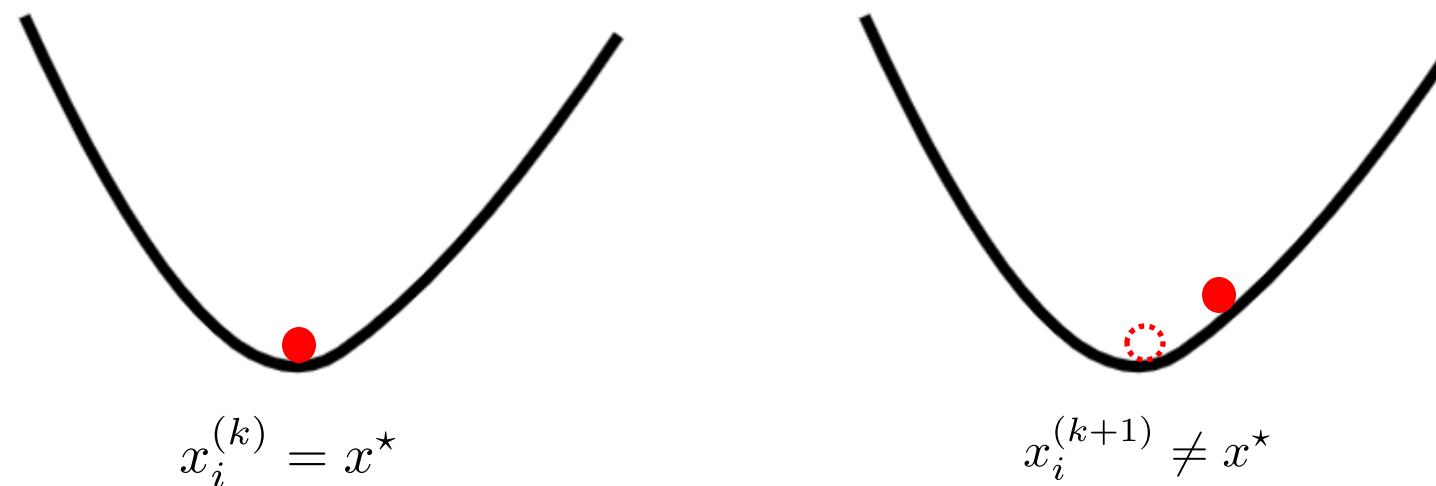
Therefore, the stationary solution is stable in homogeneous-data scenario

# D-SGD suffers from data heterogeneity

- In heterogeneous-data scenario where  $\nabla f_i(x) \neq \nabla f(x)$  for any node i, we have  $\nabla f_i(x^*) \neq \nabla f(x^*)$  and

$$x_i^{(k+1)} = x^* - \gamma \sum_{j \in \mathcal{N}_i} w_{ij} \nabla f_j(x^*) \neq x^* - \gamma \nabla f(x^*) = x^*$$

Therefore, the stationary solution is NOT stable in heterogeneous-data scenario



- This bias is caused by the algorithmic structure, not the gradient noise

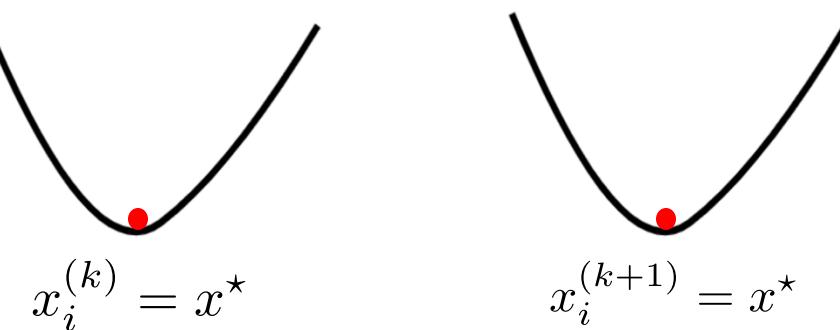
# A crude idea to remove data heterogeneity

- The fundamental reason for the bias is  $\nabla f_i(x) \neq (1/n) \sum_{i=1}^n \nabla f_i(x)$  due to data heterogeneity
- To remove the bias caused by data heterogeneity, an impractical but intuitive algorithm is

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left( x_j^{(k)} - \gamma \nabla f_j(x_j^{(k)}) \right) \longrightarrow x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left( x_j^{(k)} - \frac{1}{n} \sum_{j=1}^n \nabla f_j(x_j^{(k)}) \right)$$

- Assume all nodes initialize at  $x_i^{(k)} = x^*$ , we find  $x_i^{(k+1)}$  still remains at  $x^*$

$$\begin{aligned} x_i^{(k+1)} &= \sum_{j \in \mathcal{N}_i} w_{ij} \left( x^* - \frac{1}{n} \sum_{j=1}^n \nabla f_j(x^*) \right) \\ &= x^* - \frac{1}{n} \sum_{j=1}^n \nabla f_j(x^*) = x^* \end{aligned}$$



# Track the global gradient

---

- How to get the dynamic global average  $\frac{1}{n} \sum_{j=1}^n \nabla f_j(x_j^{(k)})$  in a decentralized manner ?
- **Decentralized tracking** (also known as dynamic average consensus [1])

$$y_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left( y_j^{(k)} + \nabla f_j(x_j^{(k+1)}) - \nabla f_j(x_j^{(k)}) \right) \quad \text{where} \quad y_i^{(0)} = \nabla f_i(x_i^{(0)})$$

- When  $\nabla f_i(x_i^{(k)})$  converges to a fixed point (such as  $\nabla f_i(x^\star)$ ) or oscillates slowly, we have [1]

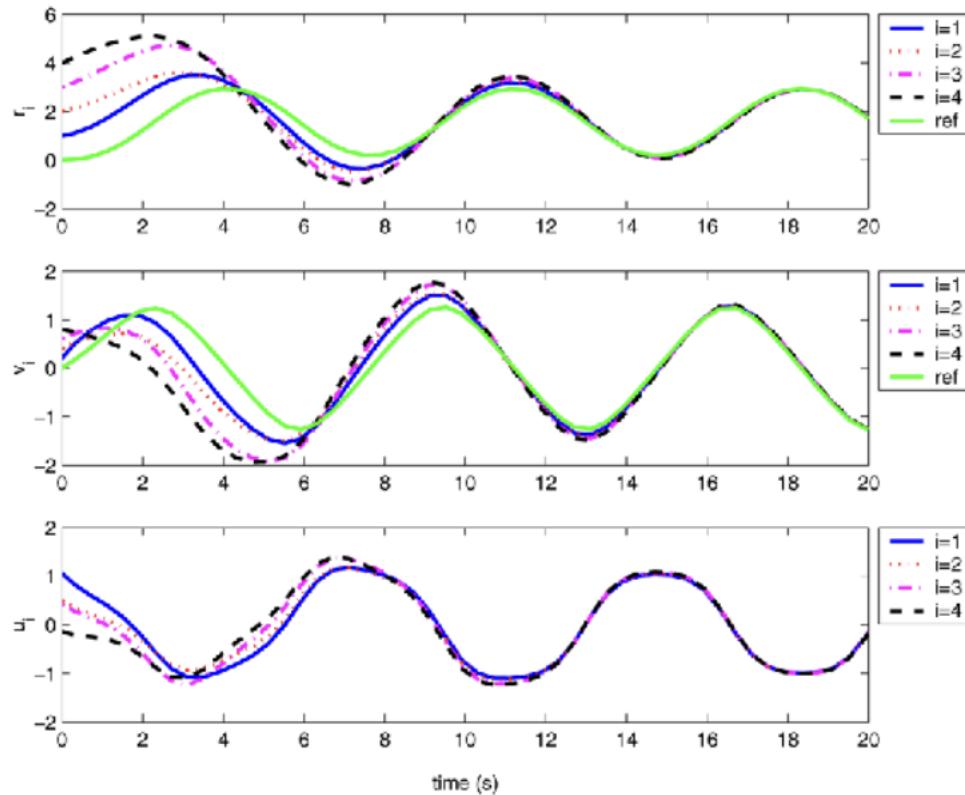
$$y_i^{(k)} \rightarrow \frac{1}{n} \sum_{j=1}^n \nabla f_j(x_j^{(k)}), \quad \forall i \in [n]$$

---

[1] M. Zhu and S. Martinez, "Discrete-time dynamic average consensus", Automatica, 2010

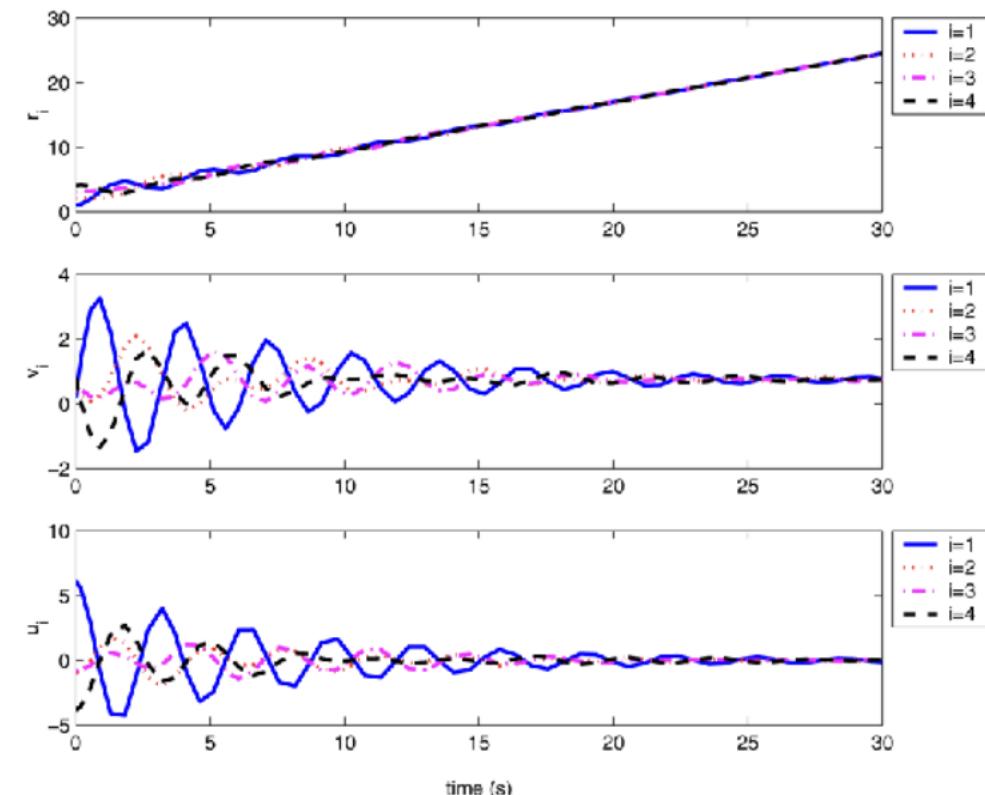
# Track the global gradient: illustration

Dynamic gradient has small oscillation



**Green: globally averaged gradient**

Dynamic gradient converges to a stationary point



Simulation results are from [Ren IEEE TAC 2007]

# Gradient tracking algorithm

- By dynamically track the globally-averaged gradient, we can achieve a new decentralized algorithm [1-3]

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left( x_j^{(k)} - \gamma y_j^{(k)} \right) \quad (\text{DSGD})$$

$$y_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left( y_j^{(k)} + \nabla F(x_j^{(k+1)}; \xi_j^{(k+1)}) - \nabla F(x_j^{(k+1)}; \xi_j^{(k)}) \right) \quad (\text{Track gradient})$$

where  $y_i^{(0)} = \nabla F(x_i^{(0)}; \xi_i^{(0)})$  in the initialization

- Since  $y_i^{(k)}$  is supposed to converge to  $\frac{1}{n} \sum_{i=1}^n \nabla F(x_i^{(k)}; \xi_i^{(k)})$ , this algorithm is named as **Gradient Tracking**

---

[1] J. Xu et. al., "Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes", IEEE CDC 2015

[2] P Di Lorenzo and G Scutari, "Next: In-network nonconvex optimization", IEEE TSIPN, 2016

[3] A Nedic, A Olshevsky, W Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs", SIOPT, 2017

# Gradient tracking: convergence

## Assumption [GT assumption]

- (1) Each  $f_i(x)$  is  $L$ -smooth, i.e.,  $\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$  for any  $x, y$ ;
- (2) Each stochastic gradient is unbiased and has bounded variance, i.e.,

$$\mathbb{E}[\nabla F(x; \xi_i)] = \nabla f_i(x), \quad \mathbb{E}\|\nabla F(x; \xi_i) - \nabla f_i(x)\|^2 \leq \sigma^2$$

## Theorem [GT convergence]

Under the above assumptions and with proper  $\gamma$ , Gradient tracking converges as

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E}\|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O}\left(\frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3}\sigma^{2/3}}{T^{2/3}(1-\rho)^{1/3}}\right)$$

[1] S. Alghunaim and K. Yuan, “A Unified and Refined Convergence Analysis for Non-Convex Decentralized Learning”, IEEE TSP 2022

## Theorem [GT convergence]

Under the above assumptions and with proper  $\gamma$ , Gradient tracking converges as

$$(GT) \quad \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}} \right)$$

- Recall DSGD with heterogeneous data converges as follows

$$(DSGD) \quad \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = \mathcal{O} \left( \frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}} + \frac{\rho^{2/3} b^{2/3}}{T^{2/3} (1-\rho)^{2/3}} \right)$$

- Gradient tracking eliminates the influence of data heterogeneity and shorten the transient stage

$$\text{DSGD: } \mathcal{O} \left( \frac{\rho^4 n^3}{\sigma^2 (1-\rho)^2} + \frac{\rho^4 n^3 b^4}{\sigma^6 (1-\rho)^4} \right) \quad \longrightarrow \quad \text{Gradient tracking: } \mathcal{O} \left( \frac{\rho^4 n^3}{\sigma^2 (1-\rho)^2} \right)$$

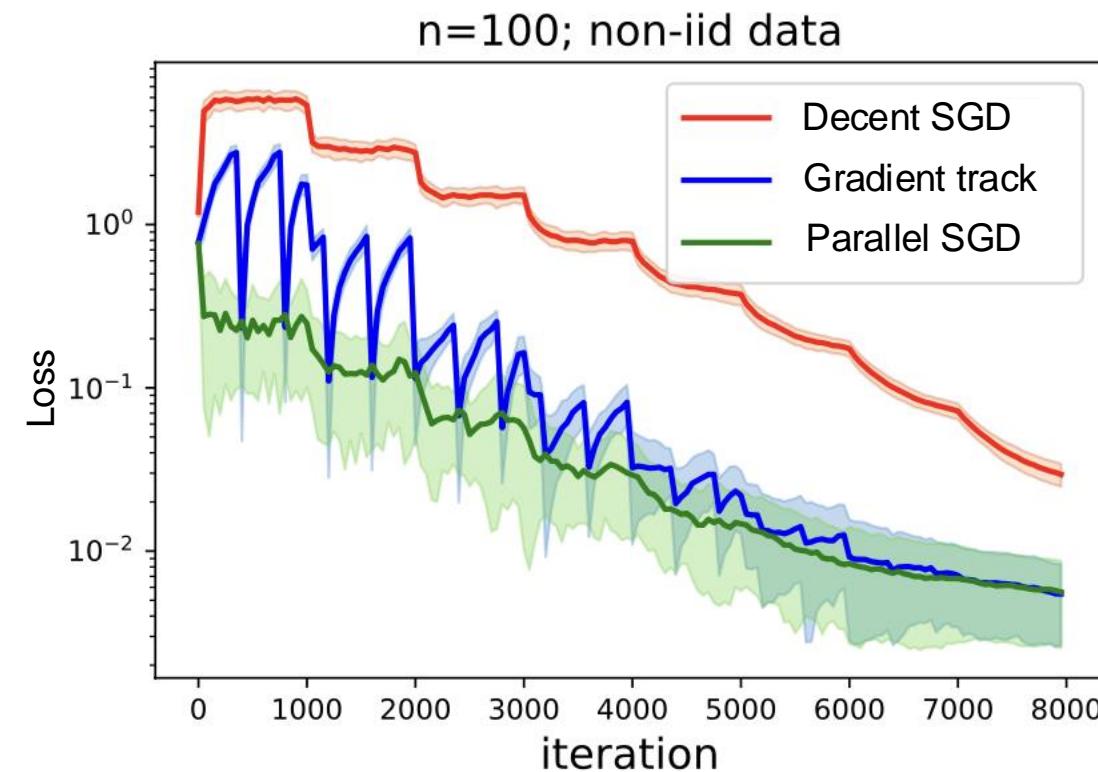
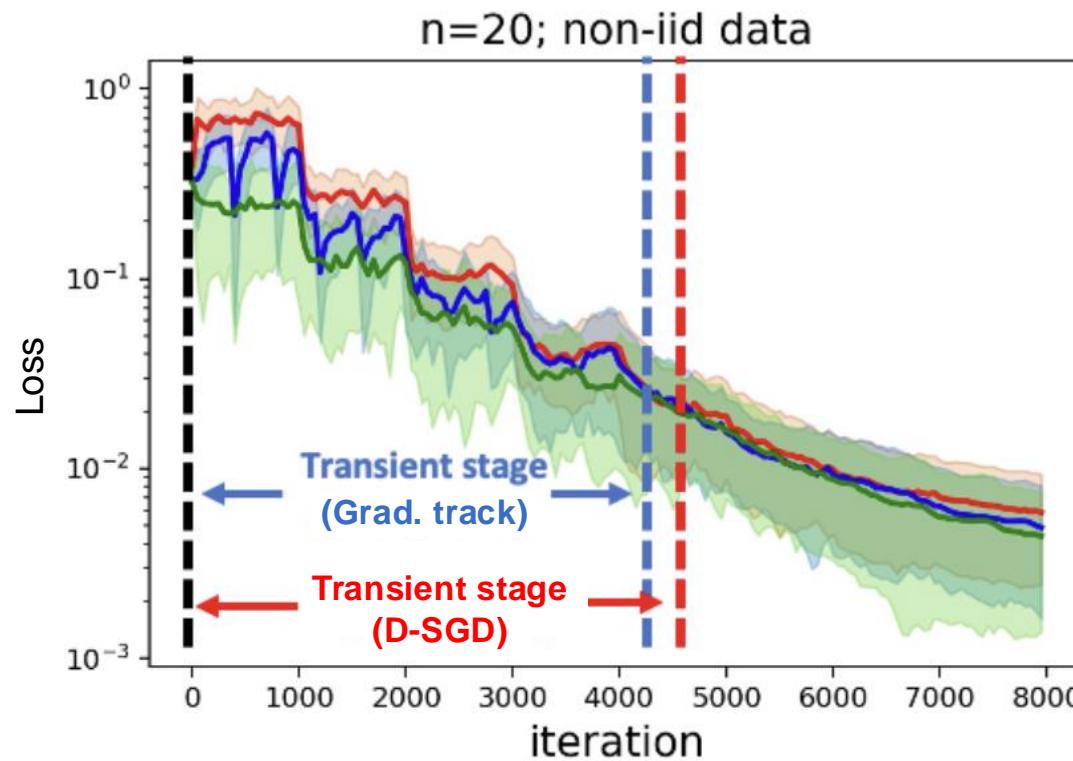
## Empirical studies

$$\text{DSGD: } \mathcal{O}\left(\frac{\rho^4 n^3}{\sigma^2(1-\rho)^2} + \frac{\rho^4 n^3 b^4}{\sigma^6(1-\rho)^4}\right)$$

$$\text{Gradient tracking: } \mathcal{O}\left(\frac{\rho^4 n^3}{\sigma^2(1-\rho)^2}\right)$$

- DNN training over ring topology with **non-iid data**

$$1 - \rho = O(n^{-2})$$



- Gradient tracking has much shorter transient stage than D-SGD for sparse topology

## PART 03-4

---

### Decentralized SGD: topology design

# Trade-off in topology design

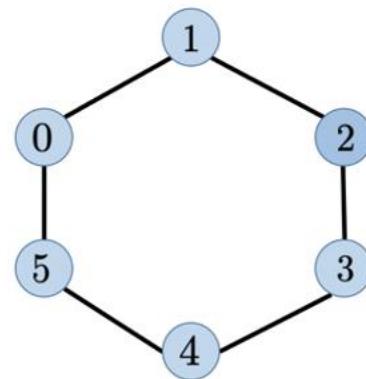


- Given a topology, its maximum degree  $d_{\max}$  decides communication overhead per iteration
- Its connectivity  $\rho = \|W - \frac{1}{n}\mathbf{1}\mathbf{1}^T\|_2 \in (0, 1)$  decides the length of transient stage  $\mathcal{O}(n^3/(1 - \rho)^2)$
- Trade-off between per-iteration communication and transient iteration complexity

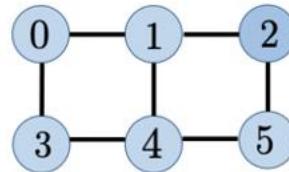
	Sparse topology	Dense topology
per-iter comm.	✓	✗
trans. iter. complexity	✗	✓

# What topology to use?

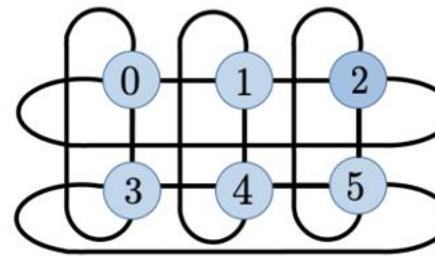
- Shall we use these common topologies to organize all nodes?



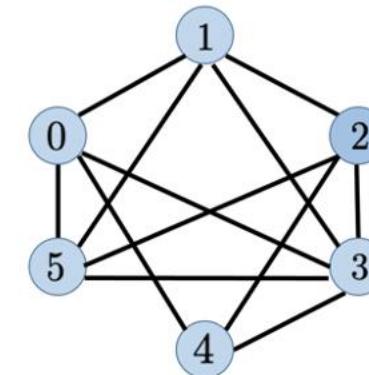
Ring



Grid



Torus



Erdos-Renyi Random

$$\rho = \mathcal{O}(1 - \frac{1}{n^2})$$

$$\mathcal{O}(1 - \frac{1}{n \ln(n)})$$

$$\mathcal{O}(1 - \frac{1}{n})$$

$$\mathcal{O}(\rho)$$

$\rho \in (0, 1)$  independent of n

# What topology to use?

---

- Communication cost v.s. transient iteration complexity in DSGD

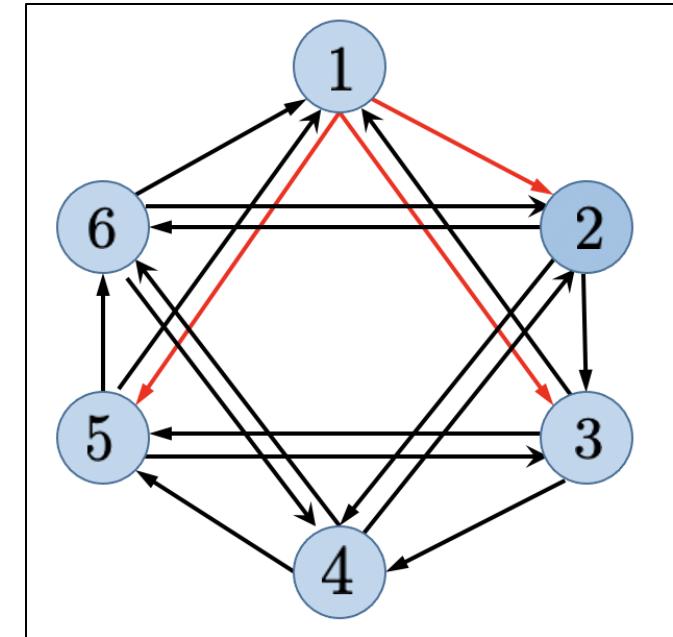
Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$O(1)$	$O(n^7)$
2D-Grid	$O(1)$	$\tilde{O}(n^5)$
2D-Torus	$O(1)$	$O(n^5)$
$\frac{1}{2}$ -RandGraph	$O(n)$	$O(n^3)$

The smaller both comm. cost and tran. Iters. are, the better

- These topologies either have expensive communication cost or longer transient stage
- Is there any topology that enables both cheap communication and fast convergence?**

# Static exponential graph: topology and per-iteration comm.

- Each node links to neighbors that are  $2^0, 2^1, \dots, 2^{\lfloor \log_2(n-1) \rfloor}$  away [1,2]
- In the figure, node 1 connects to node 2, 3 and 5.
- Each node has  $\lceil \log_2(n) \rceil$  neighbors; per-iter comm. cost is  $O(\log_2(n))$
- Empirically successful in deep training but less theoretically understood



[1] M. Assran et. al., “Stochastic Gradient Push for Distributed Deep Learning”, ICML 2018

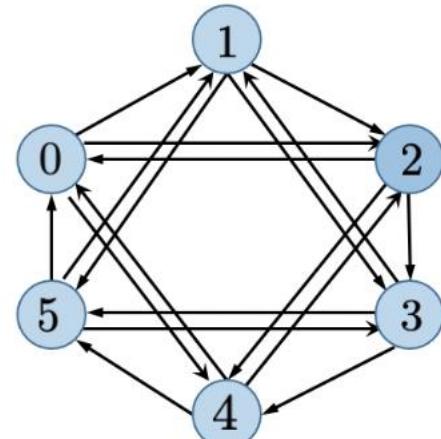
[2] B. Ying, K. Yuan, Y. Chen, H. Hu, P. Pan, W. Yin, Exponential Graph is Provably Efficient for Decentralized Deep Training, NeurIPS 2021

# Static exponential graph: weight matrix

- The weight matrix associated with exponential graph is defined as (doubly stochastic)

$$w_{ij}^{\text{exp}} = \begin{cases} \frac{1}{\lceil \log_2(n) \rceil + 1} & \text{if } \log_2(\text{mod}(j - i, n)) \text{ is an integer or } i = j \\ 0 & \text{otherwise.} \end{cases}$$

- An illustrating example:



$$W = \begin{bmatrix} \frac{1}{4} & 0 & \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} & 0 & \frac{1}{4} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} & 0 \\ 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 & \frac{1}{4} \\ \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & 0 \\ 0 & \frac{1}{4} & 0 & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix}$$

# Static exponential graph: connectivity

- Is the static exponential graph well connected?

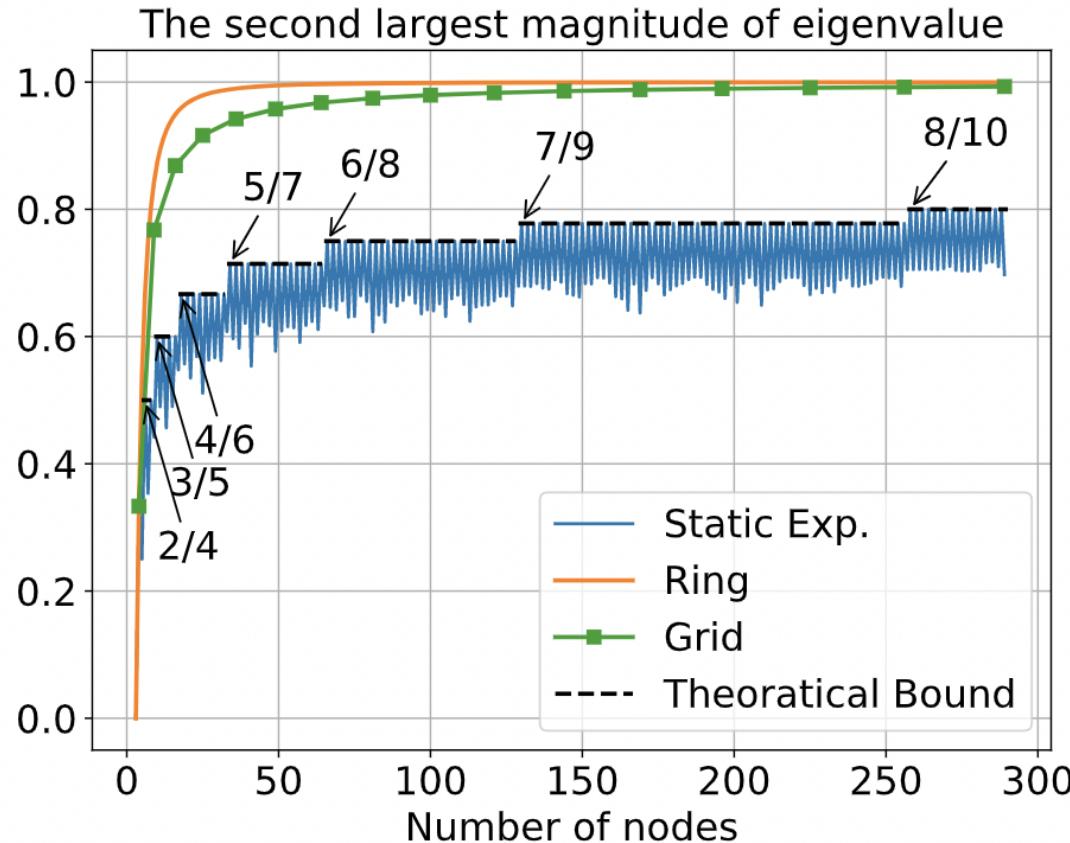
**Theorem.** Let  $\tau = \lceil \log_2(n) \rceil$ , and  $\rho = \|W - \frac{1}{n}\mathbf{1}\mathbf{1}^T\|_2$  be the spectral gap. It holds that

$$\begin{cases} \rho = 1 - \frac{2}{\tau + 1}, & \text{when } n \text{ is even} \\ \rho < 1 - \frac{2}{\tau + 1}, & \text{when } n \text{ is odd} \end{cases}$$

- This theorem implies that exponential graph has  $\rho(W) = O(1 - 1/\log_2(n))$
- Highly non-trivial proofs; requires smart utilization of Fourier transform

B. Ying, K. Yuan, Y. Chen, H. Hu, P. Pan, W. Yin, Exponential Graph is Provably Efficient for Decentralized Deep Training, NeurIPS 2021

# Static exponential graph: illustration of the spectral gap



- Our theoretical bound is very tight
- Spectral gap increases slowly when  $n$  grows

# Static exponential graph: transient iterations in DSGD

---

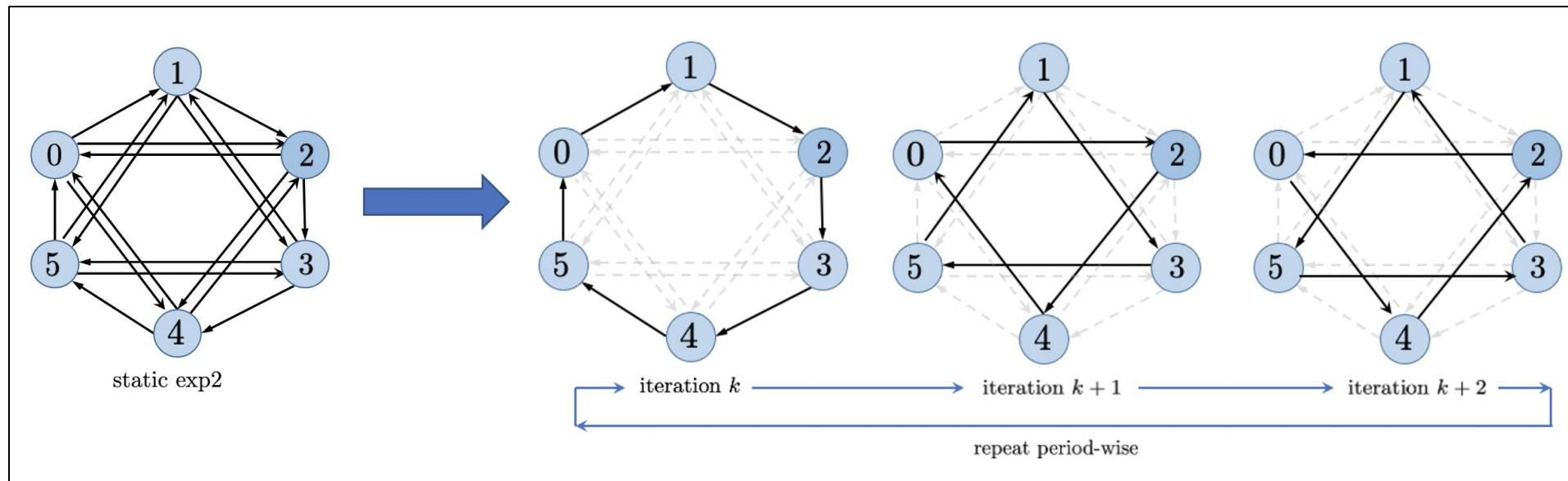
- Recall DSGD has transient iteration complexity  $O(n^3/(1 - \rho)^2)$  in iid scenarios
- With  $\rho(W) = O(1 - 1/\log_2(n))$ , exponential graphs have tran. iters. as  $O(n^3 \log_2^2(n))$
- Per-iteration communication and transient iteration complexity are **nearly the best** (up to  $\log_2(n)$ )

Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$O(1)$	$O(n^7)$
2D-Grid	$O(1)$	$\tilde{O}(n^5)$
2D-Torus	$O(1)$	$O(n^5)$
$\frac{1}{2}$ -RandGraph	$O(n)$	$O(n^3)$
Static Exp	$\tilde{O}(1)$	$\tilde{O}(n^3)$

- **Can we achieve even better topology?**

# One-peer exponential graph: topology

- **Split** exponential graph into a sequence of one-peer realizations;
- Each node has **exactly one** neighbor per iteration
- **$O(1)$  per-iteration communication**; cheaper than ring (2 neighbors) or grid (3 or 4 neighbors)

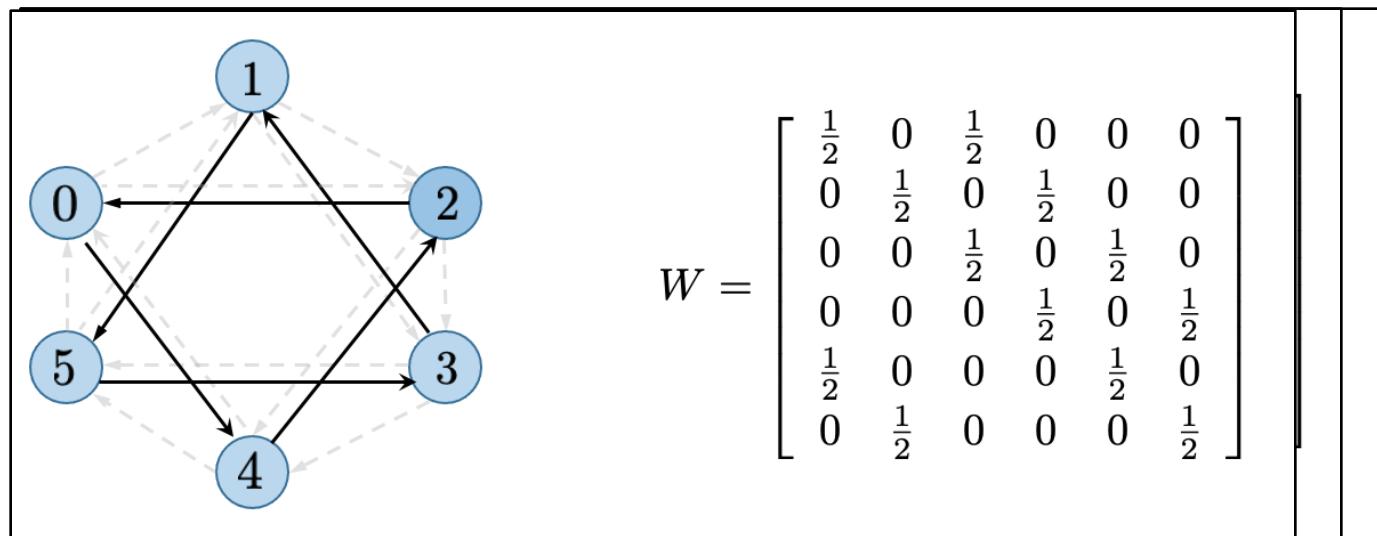


# One-peer exponential graph: weight matrix

- We let  $\tau = \lceil \log_2(n) \rceil$ . The weight matrix  $W^{(k)}$  is defined as

$$w_{ij}^{(k)} = \begin{cases} \frac{1}{2} & \text{if } \log_2(\text{mod}(j - i, n)) = \text{mod}(k, \tau) \\ \frac{1}{2} & \text{if } i = j \\ 0 & \text{otherwise.} \end{cases}$$

- An illustrating example



# DSGD over one-peer exponential graph

Sample  $W^{(k)}$  over one-peer exponential graph

$$x_i^{(k+\frac{1}{2})} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)}) \quad (\text{Local update})$$

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij}^{(k)} x_j^{(k+\frac{1}{2})} \quad (\text{Partial averaging})$$

- DSGD with **time-varying** weight matrix;
- Per-iteration communication cost is  **$O(1)$** ; very efficient

## One-peer exponential graph: Periodic exact average

- While one-peer exponential graph is **sparse**, it is **effective** to aggregate information

### Finite-time exact convergence

**Theorem.** Suppose  $\tau = \log_2(n)$  is a positive integer. It holds that

$$W^{(k+\ell)} W^{(k+\ell-1)} \dots W^{(k+1)} W^{(k)} = \frac{1}{n} \mathbf{1} \mathbf{1}^T$$

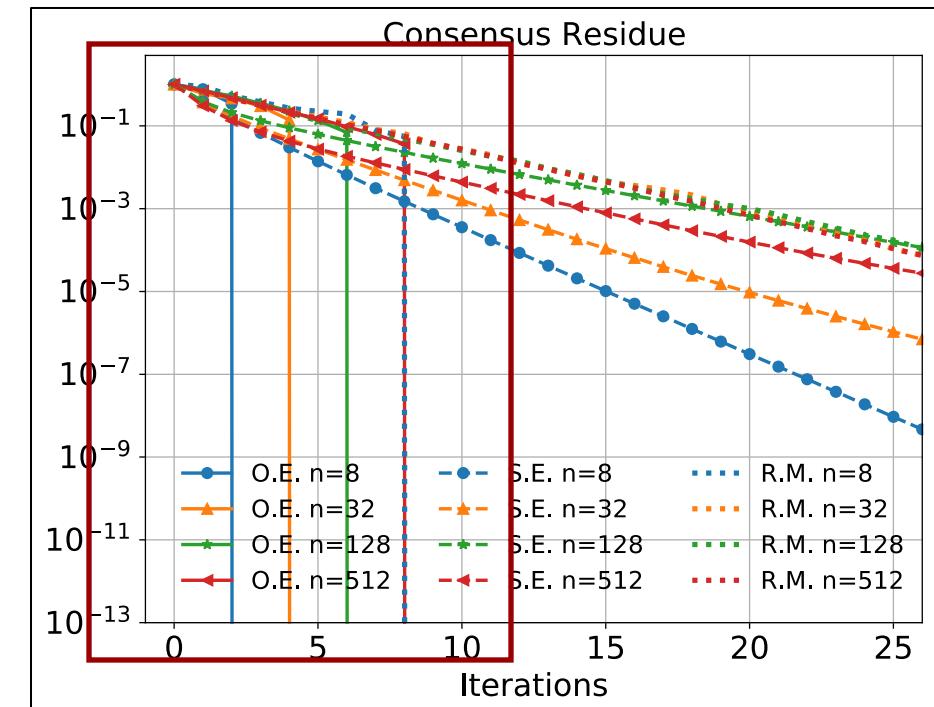
for any integer  $k \geq 0$  and  $\ell \geq \tau - 1$ .

- While each realization is sparser, a sequence (with length  $\tau$ ) of one-peer graphs will enable effective global averaging

# One-peer exponential graph: Periodic exact average



- We examine  $\left\| \frac{1}{n} \mathbf{1} \mathbf{1}^T x - \prod_{k=0}^T W^{(k)} x \right\|$  for a vector  $x$
- $\frac{1}{n} \mathbf{1} \mathbf{1}^T x$  is the global average
- $\prod_{k=0}^T W^{(k)} x$  is the partial average after  $T$  iterations
- One-peer exp. achieves global average after  $\log_2(n)$  iters.



# Apply one-peer exponential graph to DSGD

**Assumption** (1) Each  $f_i(x)$  is  $L$ -smooth; (2) Each gradient noise is unbiased and has bounded variance  $\sigma^2$ ; (3) Each local distribution  $D_i$  is identical (iid)

**Theorem** Under the above assumptions and with  $\gamma = O(1/\sqrt{T})$ , let  $\tau = \log_2(n)$  be an integer, DSGD with one-peer exponential graph will converge at

$$\frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = O\left( \frac{\sigma}{\sqrt{nT}} + \underbrace{\frac{\sigma^{2/3} \log_2^{1/3}(n)}{T^{2/3}}}_{\text{extra overhead}} \right)$$

Novel analysis; require new tricks to utilize periodic exact average to establish tight convergence

# Static v.s. one-peer exponential graph (iid scenario)

---

- Convergence rate of DSGD with iid dataset over static and one-peer exponential graphs are

Static exp.  $O\left(\frac{\sigma}{\sqrt{nT}} + \frac{\sigma^{2/3}}{T^{2/3}(1-\rho)^{1/3}}\right)$  (where  $1-\rho = O(1/\log_2(n))$ )

One-peer exp.  $O\left(\frac{\sigma}{\sqrt{nT}} + \frac{\sigma^{2/3} \log_2^{1/3}(n)}{T^{2/3}}\right)$

- DSGD with one-peer exp. converges **as fast as** static exp.; **a surprising result.**
- DSGD with both graphs are with the **same** transient iteration complexity  $O(n^3 \log_2^2(n))$
- The communication cost saving in one-peer exponential graph is a **free lunch**

# Compare one-peer exp. with other topologies

---

Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$O(1)$	$O(n^7)$
2D-Grid	$O(1)$	$\tilde{O}(n^5)$
2D-Torus	$O(1)$	$O(n^5)$
$\frac{1}{2}$ -RandGraph	$O(n)$	$O(n^3)$
Static Exp	$\tilde{O}(1)$	$\tilde{O}(n^3)$
One-Peer Expo	$O(1)$	$\tilde{O}(n^3)$

We recommend using one-peer exponential graph in deep training.

# Experiments in deep training (image classification)



ImageNet-1K dataset

1.3M training images

50K test images

1K classes

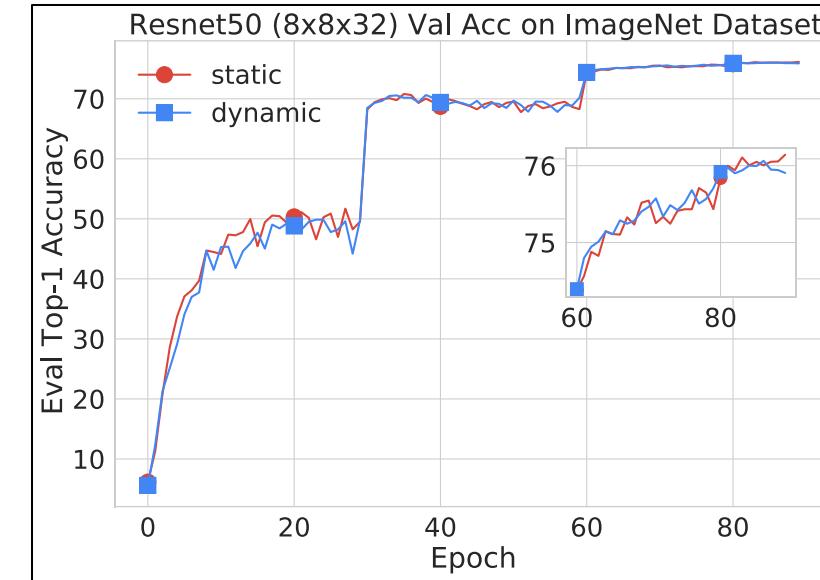
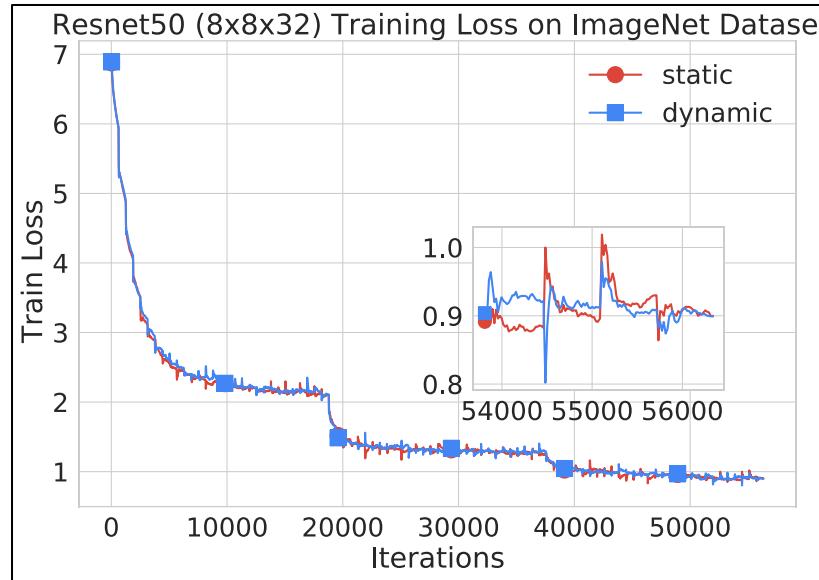
DNN model: ResNet-50 (25.5M parameters)

GPU: Up to 256 Tesla V100 GPUs

- **Wall-clock time** to finish 90 epochs of training; measures per-iter communication
- **Validation accuracy** after 90 epochs of training; measures convergence rate

# One peer is not slower than static exponential graph

Image classification: ResNet-50 for ImageNet;  $8 \times 8 = 64$  GPUs.



One-peer and exponential graphs converge **roughly the same**; but one-peer is more comm. efficient

# DSGD over one-peer Exp. achieves better linear speedup



nodes	4(4x8 GPUs)		8(8x8 GPUs)		16(16x8 GPUs)		32(32x8 GPUs)	
topology	acc.	time	acc.	time	acc.	time	acc.	time
P-SGD	76.32	11.6	76.47	6.3	76.46	3.7	76.25	2.2

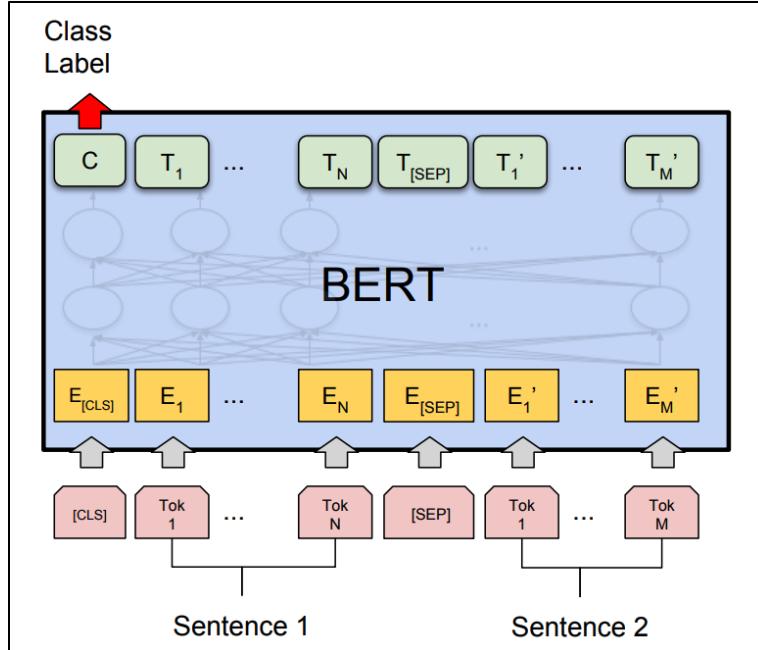
DSGD over ring has more efficient comm. than PSGD; **suffers from performance degradation**

DSGD over one-peer exp. graph is more comm.-efficient **without performance degradation**

---

[YYC+21] B. Ying, K. Yuan, Y. Chen, H. Hu, P. Pan, and W. Yin, ``Exponential Graph is Provably Efficient for Deep Training'', NeurIPS 2021

# Experiments in deep training (language modeling)



Model: BERT-Large (330M parameters)

Dataset: Wikipedia (2500M words) and  
BookCorpus (800M words)

Hardware: 64 GPUs

Table. Comparison in loss and training time [CYZ+21]

Method	Final Loss	Wall-clock Time (hrs)
P-SGD	1.75	59.02
D-SGD	1.77	30.4

[CYZ+21] Y. Chen, K. Yuan, Y. Zhang, P. Pan, Y. Xu, and W. Yin, "Accelerating Gossip SGD with Periodic Global Averaging", ICML 2021

# A brief summary

---

- Exponential graphs are both sparse and effective. They are nearly best up to logarithm terms
- One-peer exponential graph is even sparser without hurting effectiveness

Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$O(1)$	$O(n^7)$
2D-Grid	$O(1)$	$\tilde{O}(n^5)$
2D-Torus	$O(1)$	$O(n^5)$
$\frac{1}{2}$ -RandGraph	$O(n)$	$O(n^3)$
Static Exp	$\tilde{O}(1)$	$\tilde{O}(n^3)$
One-Peer Expo	$O(1)$	$\tilde{O}(n^3)$

## However ...

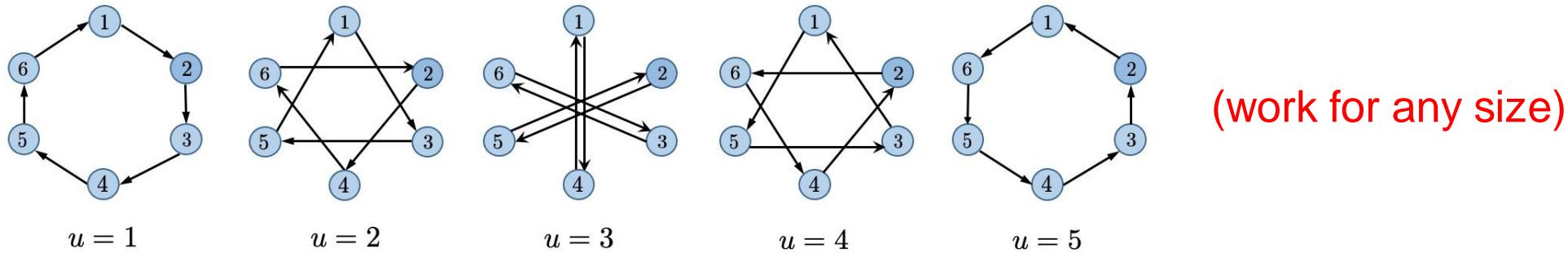
---

- Periodic exact average for one-peer exp. **only holds** when network size **n is a power of 2**
- Not known when one-peer exp. performs well when n is **not** a power of 2
- Not known whether the transient iteration  $O(n^3 \log_2^2(n))$  can be further improved to  $O(n^3)$

Can we develop topologies that

- have  $O(1)$  per-iteration communication cost;
- enable DSGD to converge with  $O(n^3)$  transient iteration complexity;
- and are valid for any network size n?

# More advanced topologies: EquiTopo graph



Topology	Per-iter. Comm.	Trans. Iter. (iid scenario)
Ring	$O(1)$	$O(n^7)$
2D-Grid	$O(1)$	$\tilde{O}(n^5)$
2D-Torus	$O(1)$	$O(n^5)$
$\frac{1}{2}$ -RandGraph	$O(n)$	$O(n^3)$
Static Exp	$\tilde{O}(1)$	$\tilde{O}(n^3)$
One-Peer Expo	$O(1)$	$\tilde{O}(n^3)$
O.-P. EquiDyn	$O(1)$	$O(n^3)$ (in expectation)

Z. Song, et. al., “Communication-Efficient Topologies for Decentralized Learning with  $O(1)$  Consensus Rate”, NeurIPS 2022

# More advanced topologies: Base-k graph

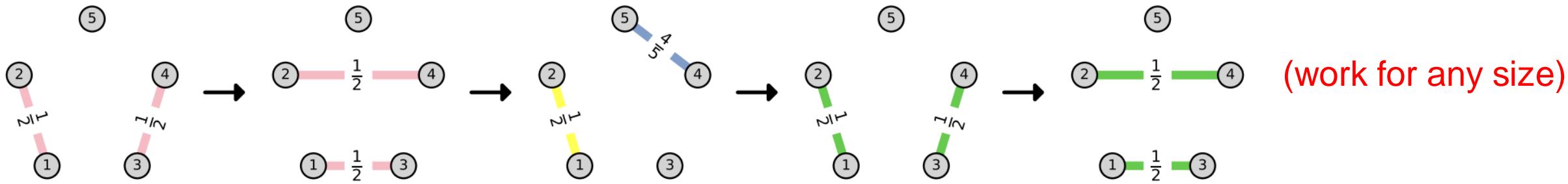


Table 1: Comparison among different topologies with  $n$  nodes. The definition of the consensus rate and finite-time convergence is shown in Sec. 3.

Topology	Consensus Rate	Connection	Maximum Degree	#Nodes $n$
Ring [28]	$1 - O(n^{-2})$	Undirected	2	$\forall n \in \mathbb{N}$
Torus [28]	$1 - O(n^{-1})$	Undirected	4	$\forall n \in \mathbb{N}$
Exp. [43]	$1 - O((\log_2(n))^{-1})$	Directed	$\lceil \log_2(n) \rceil$	$\forall n \in \mathbb{N}$
1-peer Exp. [43]	$O(\log_2(n))$ -finite time conv.	Directed	1	A power of 2
1-peer Hypercube [31]	$O(\log_2(n))$ -finite time conv.	Undirected	1	A power of 2
<b>Base-<math>(k + 1)</math> Graph (ours)</b>	$O(\log_{k+1}(n))$ -finite time conv.	Undirected	$k$	$\forall n \in \mathbb{N}$

Y. Takezawa, et. al., "Beyond Exponential Graph: Communication-Efficient Topologies for Decentralized Learning via Finite-time Convergence", NeurIPS 2023

## PART 03-5

---

### BlueFog: An open-source and high-performance python library

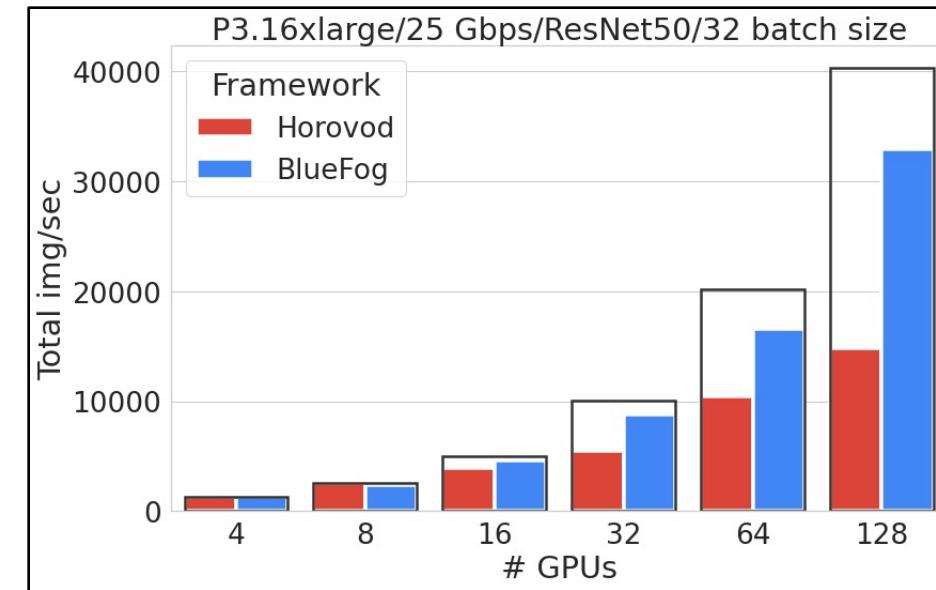
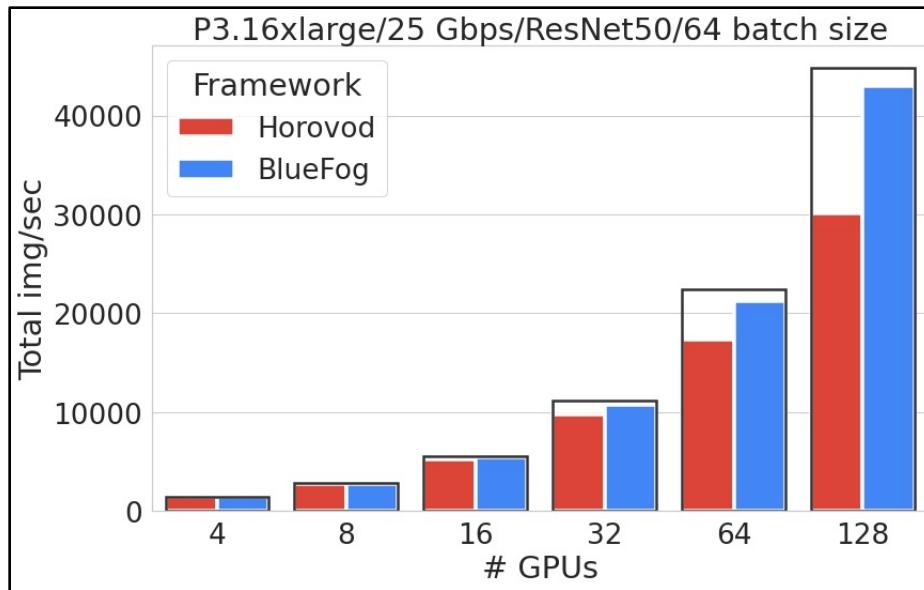


<https://github.com/Bluefog-Lib/bluefog>

- An open-source library to support decentralized communication in optimization and deep learning
- High-performance
- Easy-to-use

# High-performance

- BlueFog has larger throughput than Horovod (the SOTA DL system implementing PSGD) [YYH+21]



- All our research progresses are involved in BlueFog

[YYH+21] B. Ying, K. Yuan, H. Hu, Y. Chen, and W. Yin, ``BlueFog: Make Decentralized Algorithms Practical for Optimization and machine learning'', arXiv:2111.04287 [GitHub site: [github.com/Bluefog-Lib/bluefog](https://github.com/Bluefog-Lib/bluefog)]

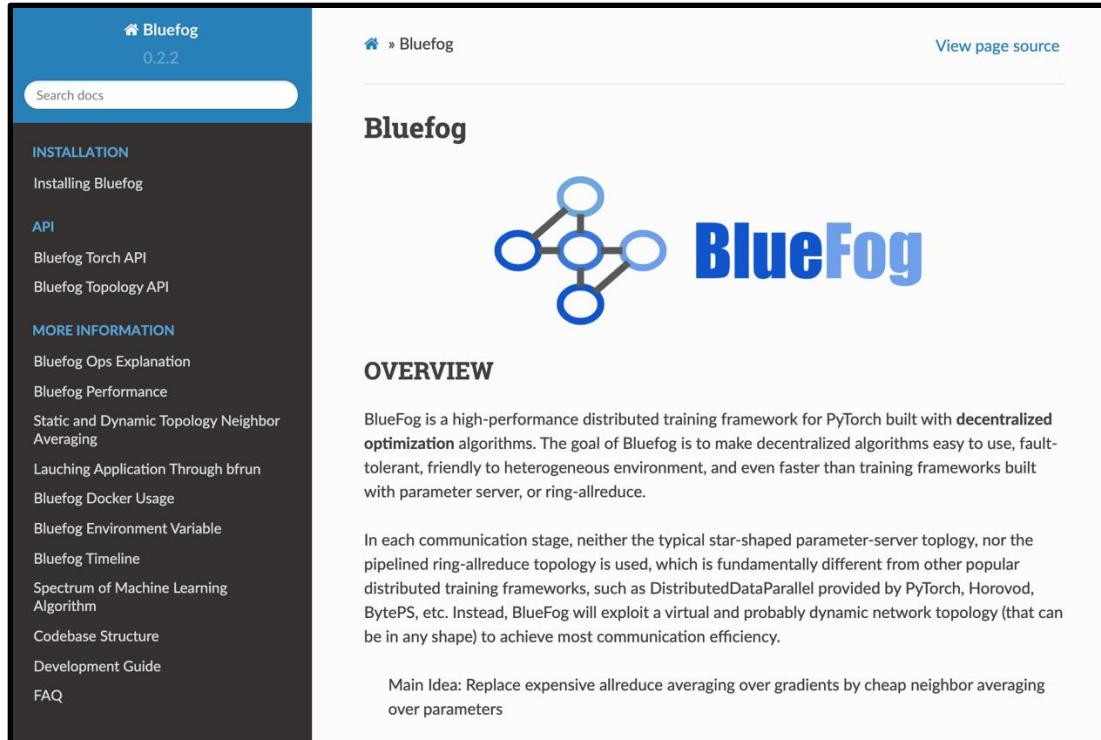
- Writing codes for decentralized methods is as easy as writing equations

## Decentralized least-square algorithms

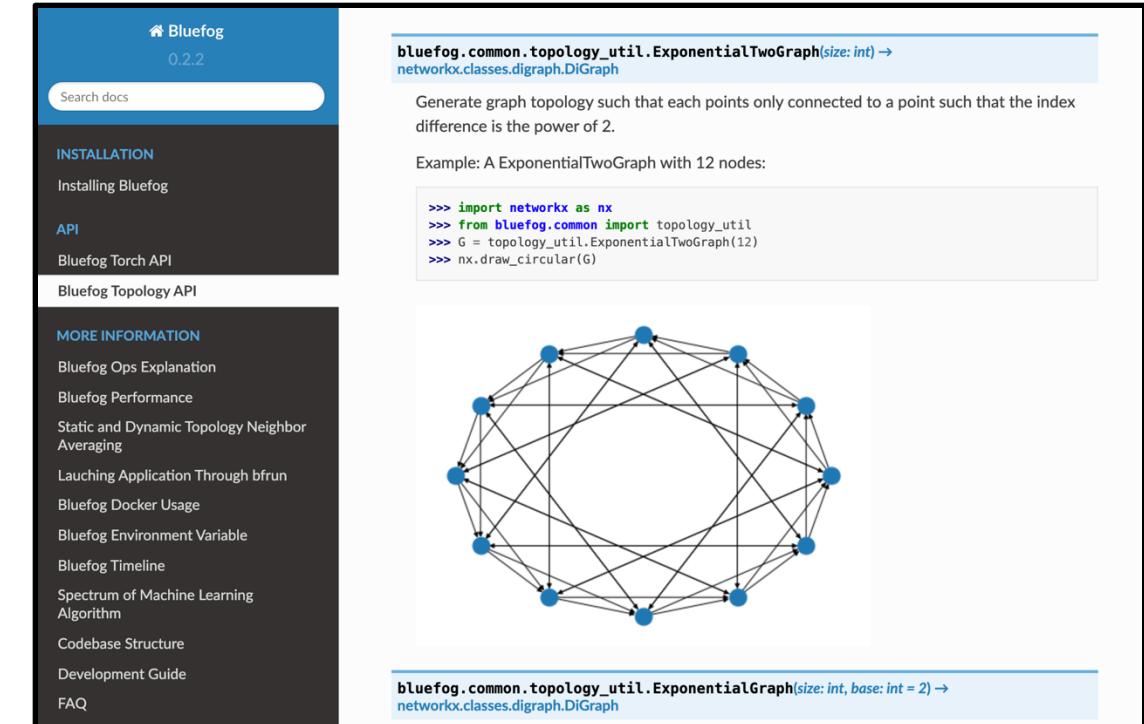
$$y_i^{(k)} = x_i^{(k)} - \gamma A_i^T (A_i x_i^{(k)} - b_i)$$
$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} y_j^{(k)}$$

```
1 import bluefog.torch as bf
2 bf.init()  # Initialize the BlueFog
3
4 # Set topology as static exponential graph.
5 G = bf.ExponentialTwoGraph(bf.size())
6 bf.set_topology(G)
7
8 # DGD implementation
9 for ite in range(maxite):
10     grad_local = A.t().mm(A.mm(x) - b)  # compute local grad
11     y = x - gamma * grad_local           # local update
12     x = bf.neighbor_allreduce(y)         # partial averaging
```

## Abundant documents



The screenshot shows the Bluefog documentation homepage. The top navigation bar includes links for "Search docs", "INSTALLATION", "API", and "MORE INFORMATION". The main content features a network graph icon and the text "BlueFog" in large blue letters. Below this is the "OVERVIEW" section, which contains a detailed description of the framework's purpose and how it differs from other distributed training frameworks. It also includes a "Main Idea" section explaining the neighbor averaging technique.



The screenshot shows a specific documentation page for the `ExponentialTwoGraph` class. The top navigation bar is identical to the homepage. The main content includes a code snippet demonstrating how to generate such a graph, followed by a detailed explanation of the topology. Below this is another code snippet showing an example of creating a graph with 12 nodes. The bottom section features a complex network graph diagram with many nodes and connections.

## Detailed tutorials

### Contents

#### 1 Preliminary

Learn how to write your first "hello world" program over the real multi-CPU system with BlueFog.

#### 2 Average Consensus Algorithm

Learn how to achieve the globally averaged consensus among nodes in a decentralized manner.

#### 3 Decentralized Gradient Descent

Learn how to solve a general distributed (possibly stochastic) optimization problem in a decentralized manner.

#### 4 Decentralized Gradient Descent with Bias-Correction

Learn how to accelerate your decentralized (possibly stochastic) optimization algorithms with various bias-correction techniques.

#### 5 Decentralized Optimization over directed and time-varying networks

Learn how to solve distributed optimization in a decentralized manner if the connected topology is directed or time-varying.

#### 6 Asynchronous Decentralized Optimization

Learn how to solve a general distributed optimization problem with asynchronous decentralized algorithms.

#### 7 Decentralized Deep Learning

Learn how to train a deep neural network with decentralized optimization algorithms.

#### 2.1.3 Initialize BlueFog and test it

All contents in this section are displayed in Jupyter notebook, and all experimental examples are written with BlueFog and iParallel. Readers not familiar with how to run BlueFog in ipython notebook environment is encouraged to read Sec. [HelloWorld section] first. In the following codes, we will initialize BlueFog and test whether it works normally.

The output of `rc.ids` should be a list from 0 to the number of processes minus one. The number of processes is the one you set in the `ibrun start -np {X}`.

```
In [1]: import ipyparallel as ipp  
rc = ipp.Client(profile="bluefog")  
rc.ids
```

Let each agent import necessary modules and then initialize BlueFog. You should be able to see the printed information like:

```
[stdout:0] Hello, I am 1 among 4 processes  
...  
[stdout:1] Hello, I am 2 among 4 processes  
...  
[stdout:3] Hello, I am 3 among 4 processes  
...  
[stdout:4] Hello, I am 4 among 4 processes
```

```
In [2]: %%px  
import numpy as np  
import bluefog.torch as bf  
import torch  
from bluefog.common import topology_util  
import networkx as nx  
  
bf.init()  
print(f"Hello, I am {bf.rank()} among {bf.size()} processes")
```

Push seed to each agent so that the simulation can be reproduced.

```
In [3]: dview = rc[:] # A DirectView of all engines  
dview.block = True  
  
# Push the data into all workers  
# `dview.push({'seed': 2021}, block=True)`  
# Or equivalently  
dview['seed'] = 2021
```

After running the following code, you should be able to see the printed information like

```
[stdout:0] I received seed as value: 2021
```

# Decentralized SGD final summary

---



- Decentralized algorithms save remarkable communication compared to centralized ones
- By removing data heterogeneity, we can significantly improve the transient iterations
- Sparse and effective topologies make decentralized optimization practical for deep training
- In real GPU clusters, decentralized SGD show strong scaling performance without hurting the accuracy performance



# Thank you!

**Kun Yuan homepage:** <https://kunyuan827.github.io/>

**BlueFog homepage:** <https://github.com/Bluefog-Lib/bluefog>

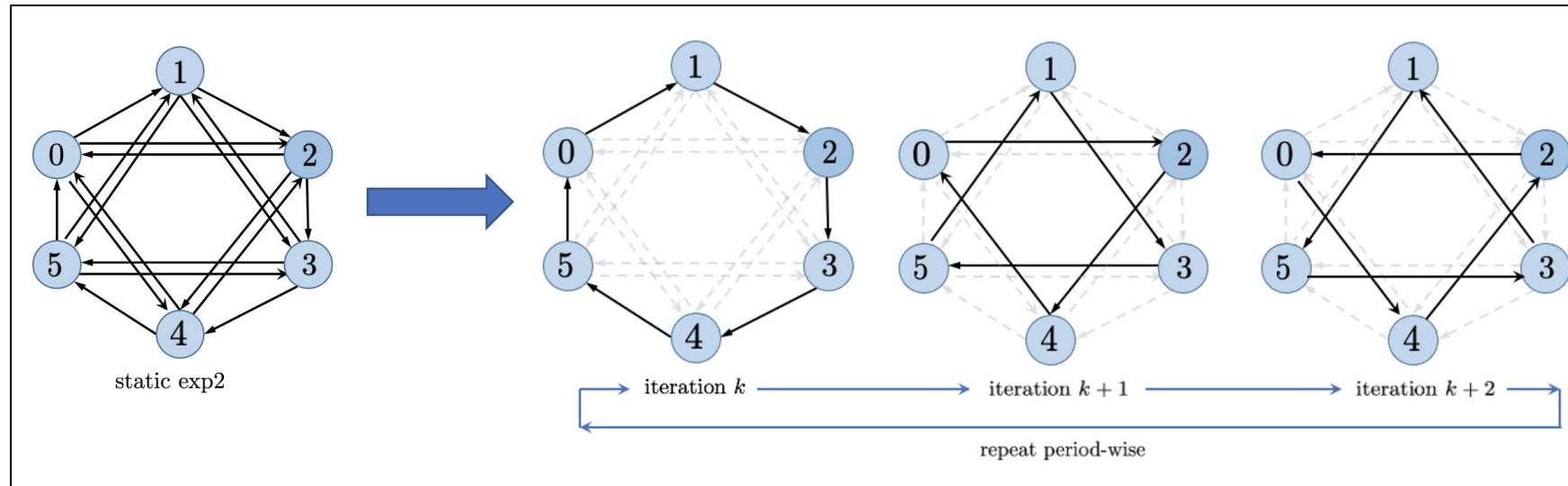
# Backup Slides

---

**EquiTopo graphs are new state-of-the-art**

# Why does exponential graph suffer $\log(n)$ deterioration?

- Exponential graphs are still not well-connected



- For example, node 0 never sends messages to nodes 3 and 5
- We need to develop topologies that every pair of nodes is connected in positive probability

# Basis weight matrix and graph

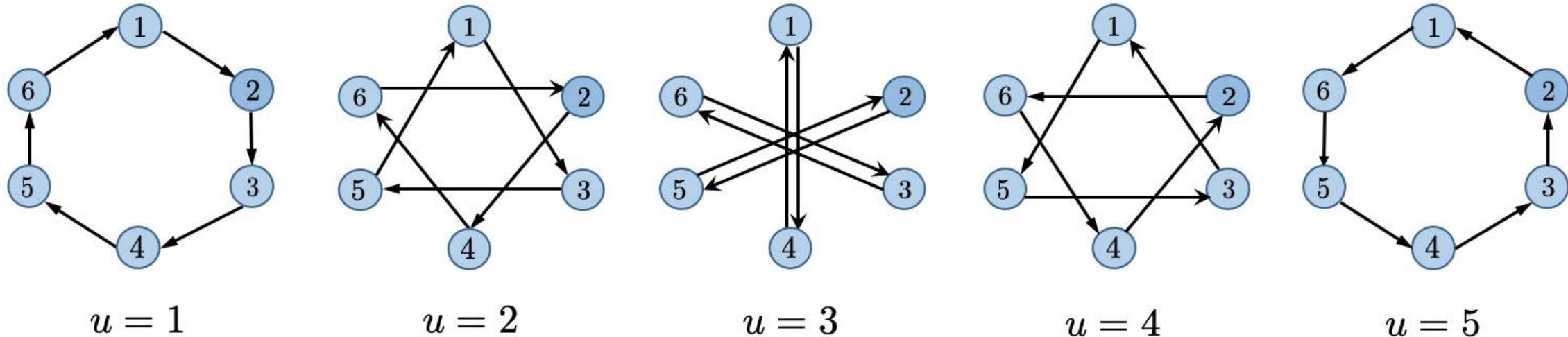
**Definition** Given a graph of size  $n$ , we introduce a set of doubly stochastic *basis matrices*  $\{A^{(u,n)}\}_{u=1}^{n-1}$ , where  $A^{(u,n)} = [a_{ij}^{(u,n)}] \in \mathbb{R}^{n \times n}$  with

$$a_{ij}^{(u,n)} = \begin{cases} \frac{n-1}{n}, & \text{if } i = (j + u) \bmod n, \\ \frac{1}{n}, & \text{if } i = j, \\ 0, & \text{otherwise.} \end{cases}$$

Their associated graphs  $\{\mathcal{G}(A^{(u,n)})\}_{u=1}^{n-1}$  are called *basis graphs*.

# Basis weight matrix and graph: illustration

The set of basis graphs  $\{\mathcal{G}(A^{(u)})\}_{u=1}^5$  for  $n = 6$



- Each basis graph  $\mathcal{G}(A^{(u)})$  is an **one-peer** graph;  $O(1)$  per-iteration communication overhead
- Each edge in a basis graph has the same **label difference**

---

## Generate EquiDyn realization $W^{(k)}$

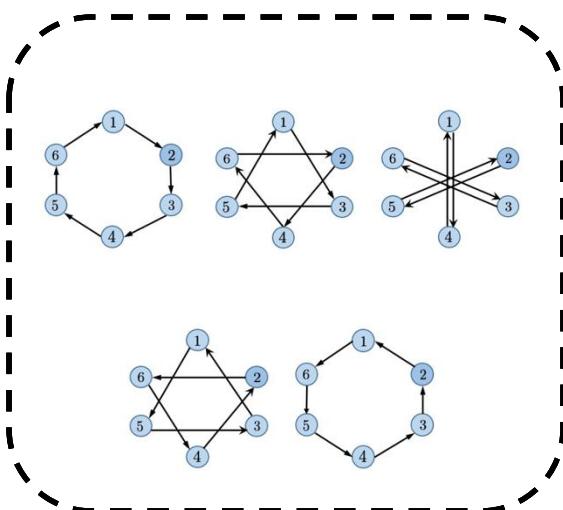
---

Pick  $v_k$  from uniform distribution over the basis index set  $[n - 1]$

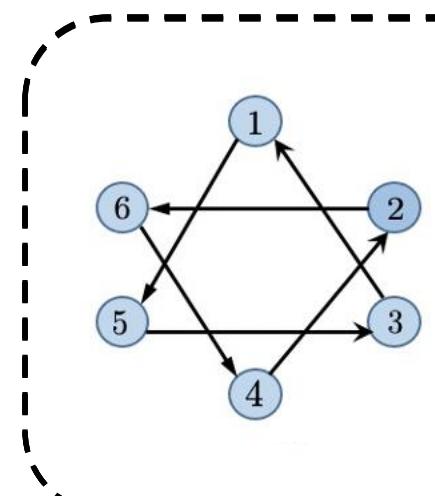
Produce basis matrix  $A^{(v_k)}$  according to the definition

Generate weight matrix  $W^{(k)} = (1 - \eta)I + \eta A^{(v_k)}$

---



Sample data  $k=2$



$$W = \begin{bmatrix} \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 0 & 0 & 0 & \frac{1}{2} \end{bmatrix}$$

# OP-EquiDyn v.s. OP-Exponential

---



**OP-Exp**

OP-EquiDyn

Sampled in a **cyclic** manner

Nodes with **exponential** label differences can be connected

Sampled in a **random** manner

Nodes with **any** label differences can be connected

# OP-EquiDyn has a network-size-independent spectral gap

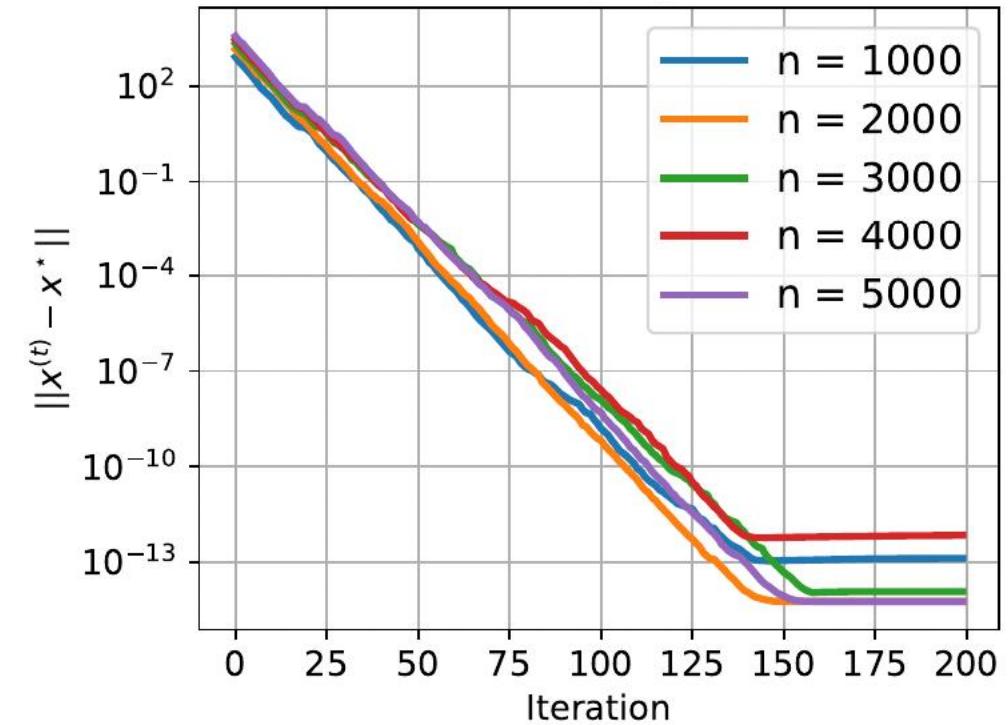
**Theorem.** Let the one-peer directed weight matrix  $W^{(k)}$  be generated by the above EquiDyn algorithm. If we let  $\eta = 1/2$ , it then holds that

$$\rho = \mathbb{E} \|W^{(k)} - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T\|_2 \leq \frac{\sqrt{2}}{2}$$

- Such spectral gap is **independent of network size**, and holds for **any size n**
- Recall that DSGD has transient iteration complexity  $O(n^3(1 - \rho)^{-2})$
- Substituting  $\rho = \sqrt{2}/2$ , DSGD over OP-EquiDyn has tran. iters.  $O(n^3)$

# OP-EquiDyn has a network-size-independent spectral gap

- We examine  $\|\frac{1}{n}\mathbf{1}\mathbf{1}^T x - \prod_{k=0}^T W^{(k)}x\|$  for a vector  $x$
- OP-EquiDyn has a network-size-independent rate
- While network size increases, consensus rate remains almost unchanged



# OP-EquiDyn achieves new SOTA results



DSGD with different network topology

Topology	Per-iter.	Comm.	Trans.	Iters. (iid scenario)
Ring		$O(1)$		$O(n^7)$
2D-Grid		$O(1)$		$\tilde{O}(n^5)$
2D-Torus		$O(1)$		$O(n^5)$
$\frac{1}{2}$ -RandGraph		$O(n)$		$O(n^3)$
Static Exp		$\tilde{O}(1)$		$\tilde{O}(n^3)$
One-Peer Expo		$O(1)$		$\tilde{O}(n^3)$
O.-P. EquiDyn		$O(1)$		$O(n^3)$

- OP-EquiDyn achieves  $O(1)$  comm.,  $O(n^3)$  transient iteration complexity, and holds for any size n
- Since DSGD has a transient complexity as  $O(n^3(1 - \rho)^{-2})$ , the order  $O(n^3)$  cannot be improved

# OP-EquiDyn can also accelerate other decentralized methods

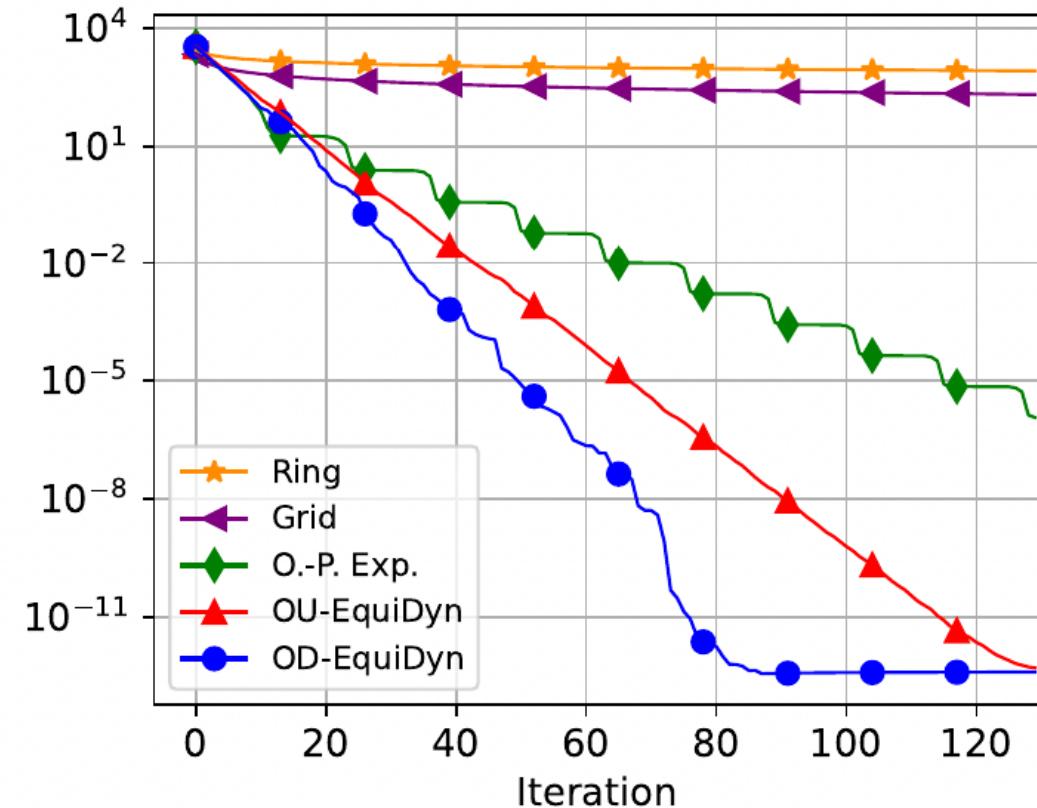
## Gradient tracking

$$\begin{aligned}\mathbf{x}_i^{(t+1)} &= \sum_{j=1}^n w_{ij}^{(t)} (\mathbf{x}_j^{(t)} - \gamma \mathbf{y}_j^{(t)}); \\ \mathbf{y}_i^{(t+1)} &= \sum_{j=1}^n w_{ij}^{(t)} \mathbf{y}_j^{(t)} + \mathbf{g}_i^{(t+1)} - \mathbf{g}_i^{(t)}, \quad \mathbf{y}_i^{(0)} = \mathbf{g}_i^{(0)}.\end{aligned}$$

Topology	Per-iter Comm.	Convergence Rate	Trans. Iters.
Ring	$\Theta(1)$	$\mathcal{O}\left(\frac{\sigma}{\sqrt{nT}} + \frac{n^2\sigma^{\frac{2}{3}}}{T^{\frac{2}{3}}} + \frac{n^4}{T}\right)$	$\mathcal{O}(n^{15})$
Torus	$\Theta(1)$	$\mathcal{O}\left(\frac{\sigma}{\sqrt{nT}} + \frac{n\sigma^{\frac{2}{3}}}{T^{\frac{2}{3}}} + \frac{n^2}{T}\right)$	$\mathcal{O}(n^9)$
Static Exp.	$\Theta(\ln(n))$	$\mathcal{O}\left(\frac{\sigma}{\sqrt{nT}} + \frac{\ln(n)\sigma^{\frac{2}{3}}}{T^{\frac{2}{3}}} + \frac{\ln^2(n)}{T}\right)$	$\mathcal{O}(n^3 \ln^6(n))$
O.-P. Exp.	1	$\mathcal{O}\left(\frac{\sigma}{\sqrt{nT}} + \frac{\ln(n)\sigma^{\frac{2}{3}}}{T^{\frac{2}{3}}} + \frac{\ln^2(n)}{T}\right)$	$\mathcal{O}(n^3 \ln^6(n))$
OD (OU)-EquiDyn	1	$\mathcal{O}\left(\frac{\sigma}{\sqrt{nT}} + \left(\frac{\sigma}{T}\right)^{\frac{2}{3}} + \frac{1}{T}\right)$	$\mathcal{O}(n^3)$

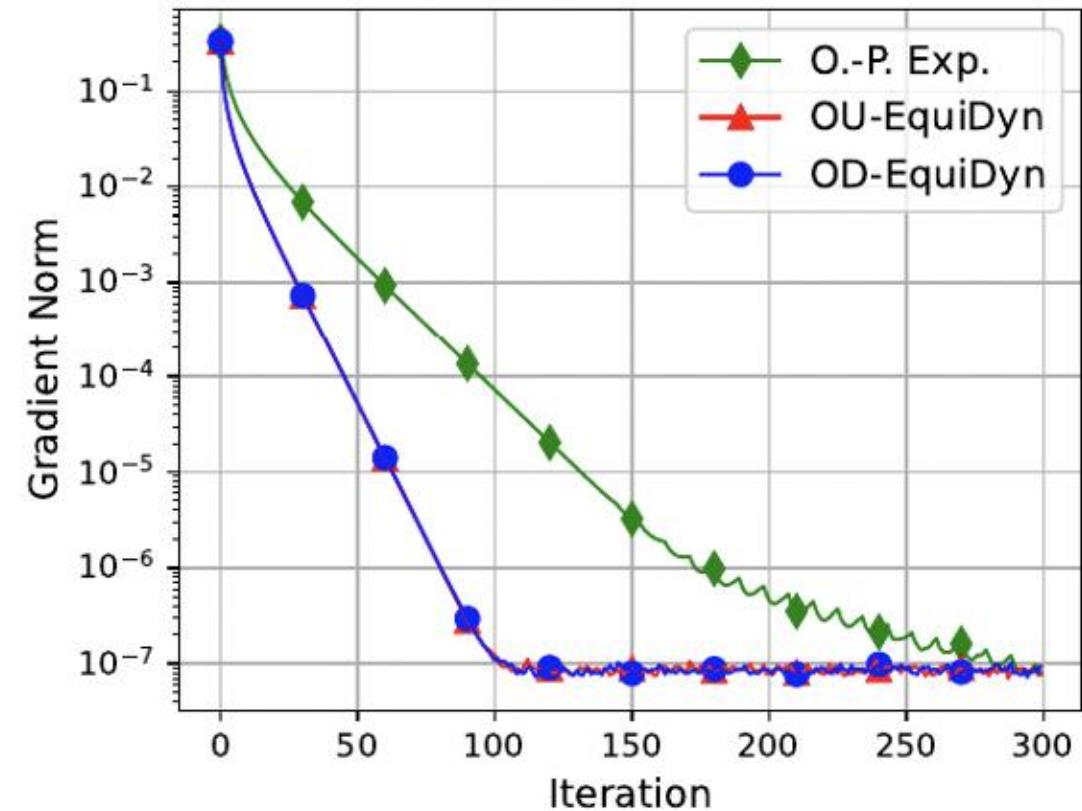
# Experiments: compare with other topologies

- We examine  $\left\| \frac{1}{n} \mathbf{1} \mathbf{1}^T x - \prod_{k=0}^T W^{(k)} x \right\|$  for a vector  $x$
- Network size is 4900
- EquiDyn converges the fastest



# Experiments: gradient tracking with different topologies

- We use GT to solve logistic regression with non-convex regularizes
- Network size is 300
- GT with EquiDyn converges faster than OP-Exp



# Experiments: deep learning experiments

- EquiTopo graph has many variants, i.e., OU-EquiDyn supports undirected graphs
- EquiTopo graph outperforms other common topologies with 17 GPUs

Topology	MNIST Acc.	CIFAR-10 Acc.
Centralized SGD	98.34	91.76
Ring	98.32	91.25
Static Exp.	98.31	91.48
O.-P. Exp.	98.17	90.86
D-EquiStatic	98.29	<b>92.01</b>
U-EquiStatic	98.26	91.74
OD-EquiDyn	<b>98.39</b>	91.44
OU-EquiDyn	98.12	91.56

# Summary

---

Can we develop topologies that

- have  $O(1)$  per-iteration communication cost;
- enable DSGD to converge with  $O(n^3)$  transient iteration complexity;
- and are valid for any network size n?

**One-peer EquiDyn is the answer!**

## More results

---

- One-peer EquiDyn performs well for any network size  $n$  **in expectation**
- It occasionally performs badly with non-zero probability
- We further develop a new algorithm DSGD-CECA that is guaranteed to perform well in all trials

L. Ding, K. Jin, B. Ying, K. Yuan, and W. Yin, “DSGD-CECA: Decentralized SGD with Communication-Optimal Exact Consensus Algorithm”, ICML 2023