

Decentralized Optimization Algorithms for Large-Scale Deep Neural Network Training

Kun Yuan

DAMO Academy, Alibaba Group

Joint work with Yiming Chen, Pan Pan, Yinghui Xu, Wotao Yin (Alibaba),
Bicheng Ying (UCLA), Hanbin Hu (UCSB), Xinmeng Huang (Upenn),
and Sulaiman A. Alghunaim (Kuwait University)

Aug 5, 2021, Zhejiang University

Contents in the lecture

Introduction to deep neural network (DNN) and various training modes (Part I)

- Single-node training
- Parallel/distributed training
- Decentralized training

Making decentralized algorithms practical for large-scale deep training (Part II)

- Exponential graphs
- Primal-dual decentralized methods
- Multiple gossip loops/Periodic global averaging

Other advanced topics and BlueFog (Part III)

- Large-batch deep training/communication-saving decentralized approaches
- An open source decentralized deep training framework: BlueFog

D-SGD Review

- A network of n nodes (GPUs) collaborate to solve the problem:

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

- Each component $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is local and private to node i
- Random variable ξ_i denotes the local data that follows distribution D_i
- Each local distribution D_i may be different; data heterogeneity
- We consider deep training within high-performance **data-center** clusters

Decentralized SGD

- Decentralized SGD (D-SGD) has the following recursion:

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

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

- Per-iteration communication: $\Omega(d_{\max}) \ll \Omega(n)$ when topology is sparse
- Incurs $\Omega(1)$ comm. overhead on sparse topology (ring or grid)

Decentralized SGD is more communication efficient

Model	Ring-Allreduce	Partial average
ResNet-50	278 ms	150 ms
Bert	1469 ms	567 ms

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

- ResNet-50 has 25.5M parameters; Bert has 300M parameters
- Partial average saves more communication for larger model

Convergence rate: P-SGD v.s. D-SGD

- Convergence comparison (i.i.d data distribution, i.e., $b^2 = 0$):

$$\text{P-SGD : } \frac{1}{T} \sum_{k=1}^T \mathbb{E} \|\nabla f(\bar{x}^{(k)})\|^2 = O\left(\frac{\sigma}{\sqrt{nT}}\right)$$

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

where σ^2 is the gradient noise, and T is the number of iterations.

- D-SGD requires more iteration (i.e., T has to be large enough) to reach that stage due to the extra overhead caused by partial averaging

Transient iterations

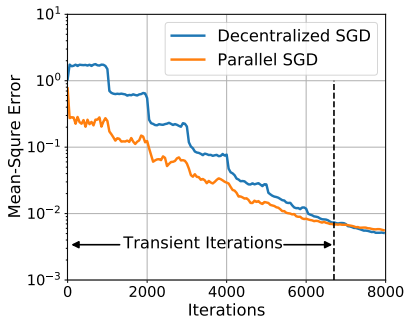
- **Definition** (Pu et al., 2020a): number of iterations before D-SGD achieves linear speedup
- Longer tran. iters. \implies slower convergence than P-SGD
- The transient iteration complexity of D-SGD is

$$\begin{aligned} \text{iid data : } & \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1 - \rho)^{1/3}} \leq \frac{\sigma}{\sqrt{nT}} \implies T = \Omega\left(\frac{\rho^4 n^3}{(1 - \rho)^2}\right) \\ \text{non-iid data : } & \frac{\rho^{2/3} b^{2/3}}{T^{2/3} (1 - \rho)^{2/3}} \leq \frac{\sigma}{\sqrt{nT}} \implies T = \Omega\left(\frac{\rho^4 n^3}{(1 - \rho)^4}\right) \end{aligned}$$

- Sparse topology ($\rho \rightarrow 1$) incurs large tran. iters. complexity

Transient iterations: illustration

Illustration of the tran. iters. on D-SGD over ring (logistic regression)



If the transient stage is too long, we may not be able to achieve linear speedup given the limited time/resource budget

Slower convergence will compensate comm. efficiency

- ImageNet dataset; ResNet-50; 256 V100 GPUs

METHOD	EPOCH	ACC.%	TIME(HRS.)
P-SGD	120	76.26	2.22
D-SGD	120	75.34	1.55

- D-SGD finishes the same epochs faster because it is more comm. efficient
- D-SGD achieves worse accuracy because it converges slower than P-SGD

Slower convergence will compensate comm. efficiency

- ImageNet dataset; ResNet-50; 256 V100 GPUs

METHOD	EPOCH	ACC.%	TIME(HRS.)
P-SGD	120	76.26	2.22
D-SGD	240	76.18	3.03

- When training with more epochs, D-SGD catch up with P-SGD in accuracy; but it takes more wall-clock time than PSGD
- Slower convergence compensates its comm. efficiency

Accelerate D-SGD and make it practical for deep learning

- Recall the transient iteration complexity of D-SGD

$$\text{iid data : } T = \Omega\left(\frac{\rho^4 n^3}{(1 - \rho)^2}\right)$$

$$\text{non-iid data : } T = \Omega\left(\frac{\rho^4 n^3}{(1 - \rho)^4}\right)$$

- Reducing tran. iter. complexity is the key to accelerating D-SGD

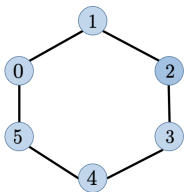
Part II: Making Decentralized SGD Practical for DNN

- Sec. 1. Exponential graphs are provably efficient (Ying et al., 2021)
- Sec. 2. Removing data heterogeneity enhances topology dependence (Huang and Pu, 2021; Yuan and Alghunaim, 2021)
- Sec. 3. Periodic global averaging (Chen et al., 2021)

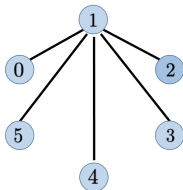
Trade-off between comm. efficiency and convergence rates

- Recall per-iter comm. $\Omega(d_{\max})$ and trans. iters. $\Omega(n^3/(1-\rho)^2)$ (iid data)
- Dense topology: expensive comm. but faster convergence
- Sparse topology: cheap comm. but slower convergence
- What topology shall we use to organize all GPUs?

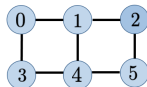
Common topologies



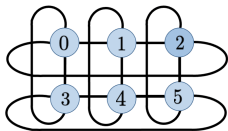
(a) ring



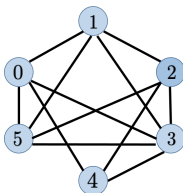
(b) star



(c) 2D-grid



(d) 2D-torus



(e) $\frac{1}{2}$ -random graph (one realization)

Common topologies: comm. cost and tran. iters

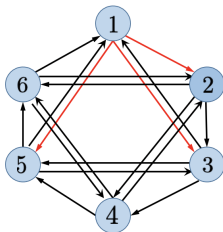
- According to (Nedić et al., 2018), we have

Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$\Omega(2)$	$\Omega(n^7)$
Star	$\Omega(n)$	$\Omega(n^7)$
2D-Grid	$\Omega(4)$	$\Omega(n^5 \log_2^2(n))$
2D-Torus	$\Omega(4)$	$\Omega(n^5)$
$\frac{1}{2}$ -RandGraph	$\Omega(\frac{n}{2})$	$\Omega(n^3)$

- These topologies either have expensive comm. cost or longer tran. iters.
- What topology can enable both cheap comm. and fast convergence?

Static exponential graph

- Static exponential graph (Lian et al., 2017, 2018; Assran et al., 2019) is widely-used in deep training
- Empirically successful but less theoretically understood
- Each node links to neighbors that are $2^0, 2^1, \dots, 2^{\lfloor \log_2(n-1) \rfloor}$ hops away
- In the figure, node 1 connects to 2, 3 and 5.



Weight matrix associated with static exponential graph

- The weight matrix W associated with static exp. graph is defined as

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

- An illustrating example

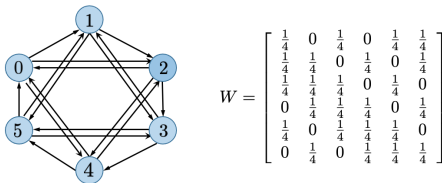


Figure: A 6-node static exponential graph and its associated weight matrix.

Weight matrix over static exponential graph: spectral gap

- Each node has $\lceil \log_2(n) \rceil$ neighbors; per-iter comm. cost is $\Omega(\log_2(n))$
- The following theorem¹ clarifies that $\rho(W^{\text{exp}}) = O(1 - 1/\log_2(n))$; highly non-trivial proofs; requires smart utilization of Fourier transform.

Theorem (Ying et.al., 2021)

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

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

¹B. Ying*, K. Yuan*, Y. Chen*, H. Han, P. Pan, and W. Yin, "Exponential graph is provably efficient for deep training", submitted, 2021

Spectral gap: numerical illustration

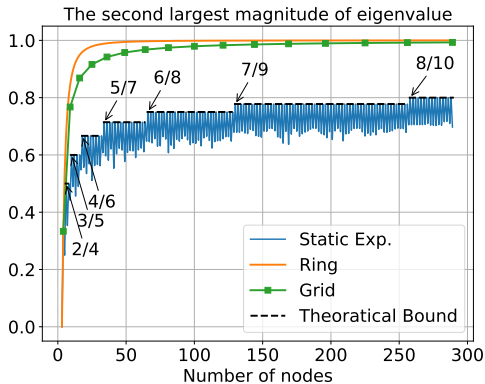


Figure: Illustration of the spectral gaps for ring, grid and static exp. graphs.

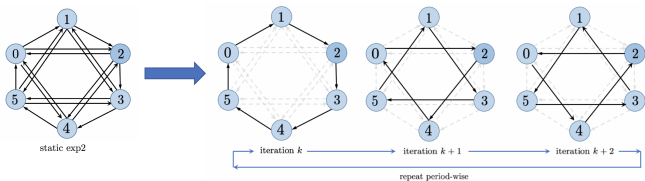
Static exponential graph v.s. other topologies

- Recall D-SGD has tran. iters. $\Omega(n^3/(1 - \rho)^2)$
- With $1 - \rho = O(1/\log_2(n))$, static exp has tran. iters. $\Omega(n^3 \log_2^2(n))$
- Per-iter comm. and tran. iter. of static exp are **nearly best** (up to $\log_2(n)$)

Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$\Omega(2)$	$\Omega(n^7)$
Star	$\Omega(n)$	$\Omega(n^7)$
2D-Grid	$\Omega(4)$	$\Omega(n^5 \log_2^2(n))$
2D-Torus	$\Omega(4)$	$\Omega(n^5)$
$\frac{1}{2}$ -RandGraph	$\Omega(\frac{n}{2})$	$\Omega(n^3)$
Static Exp	$\tilde{\Omega}(1)$	$\tilde{\Omega}(n^3)$

One-peer exponential graph

- Static exponential graph has $\Omega(\log_2(n))$ per-iteration comm.
- Such overhead is still more expensive than ring or grid
- Split exponential graph into a sequence of one-peer realizations (Assran et al., 2019)



- Each realization has $\Omega(1)$ per-iteration communication

One-peer exponential graph: weight matrix

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

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

- An illustrating example

Decentralized SGD over one-peer exponential graph

- The D-SGD recursion over one-peer exponential graph:

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

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

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

- One-loop algorithm; each node has one neighbor; per-iter comm. is $\Omega(1)$
- Since each realization is sparser than static exp., will it enable DSGD with longer transient iterations?

One-peer exp. graphs can achieve periodic exact average

Theorem (PERIODIC GLOBAL-AVERAGING)

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 of one-peer graph is sparser, a [sequence](#) of one-peer graphs will enable effective global averaging.

One-peer exp. graphs can achieve periodic exact average

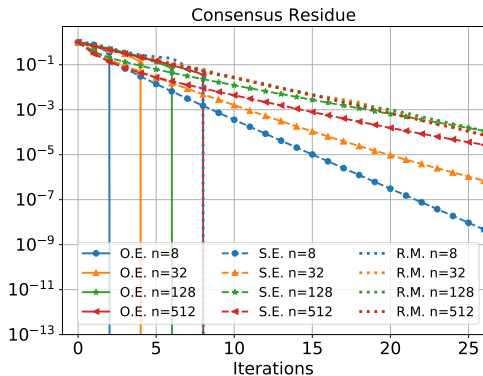


Figure: O.E. graph has periodic global averaging when $\tau = \log_2(n)$ is an integer.

Applying one-peer exp. graphs to DSGD

Assumption

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

Theorem (DSGD CONVERGENCE WITH ONE-PEER EXP.)

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

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

Convergence rate for decentralized **momentum** SGD (DmSGD) with **non-iid data distributions** is also established in (Ying et al., 2021).

Static exp. v.s. one-peer exp.

- Convergence rate for DSGD over static and one-peer exp. graphs

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

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

- DSGD with one-peer exp. converges **as fast as** static exp. in terms of the established bounds; **a surprising result**.
- DSGD with both graphs are with the same tran. iters. $O(n^3 \log_2^2(n))$
- The same results hold for heterogeneous data scenario, and for DmSGD.

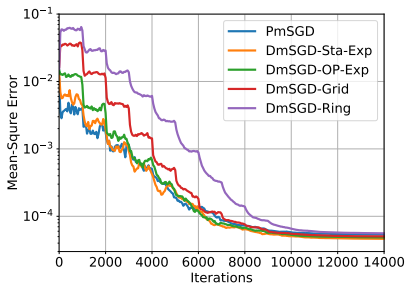
One-peer graph is the state-of-the-art topology

Topology	Per-iter. Comm.	Trans. Iters. (iid scenario)
Ring	$\Omega(2)$	$\Omega(n^7)$
Star	$\Omega(n)$	$\Omega(n^7)$
2D-Grid	$\Omega(4)$	$\Omega(n^5 \log_2^2(n))$
2D-Torus	$\Omega(4)$	$\Omega(n^5)$
$\frac{1}{2}$ -RandGraph	$\Omega(\frac{n}{2})$	$\Omega(n^3)$
Static Exp.	$\tilde{\Omega}(1)$	$\tilde{\Omega}(n^3)$
One-peer Exp.	$\Omega(1)$	$\tilde{\Omega}(n^3)$

- Since one-peer exp. incurs less per-iter comm., it is recommended for DL.

Exponential graphs have shorter transient iterations

Illustration of the tran. iters. on DmSGD for logistic regression.



DmSGD over both exp. graphs converge roughly the same; they are faster than other topologies with 32 nodes.

Experimental results: two metrics

- **Wall-clock time** to finish 90 epochs of training; measures per-iter comm.
- **Validation accuracy** after 90 epochs of training; measures convgt. rate

Image Classification

- ImageNet-1K dataset
- 1.3M training images
- 50K test images
- 1K classes
- DNN Model: ResNet-50
(~25.5M parameters)
- GPU: Tesla V100 clusters
- Framework: Pytorch DDP



D-SGD achieves better linear speedup

Table: Comparison of top-1 validation accuracy(%) and training time (hours).

nodes topology	4(4x8 GPUs)		8(8x8 GPUs)		16(16x8 GPUs)		32(32x8 GPUs)	
	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
Ring	76.16	11.6	76.14	6.5	76.16	3.3	75.62	1.8
one-peer exp.	76.34	11.1	76.52	5.7	76.47	2.8	76.27	1.5

Convergence curves: one-peer exp. v.s. static exp.

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

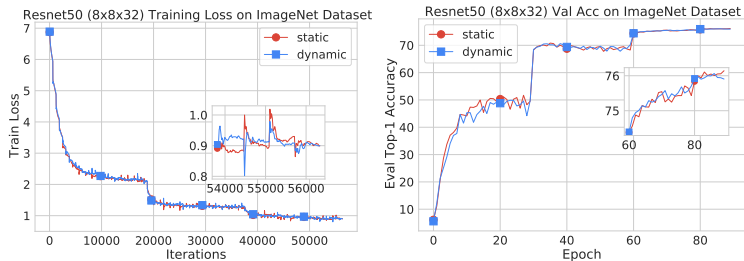


Figure: DmSGD over one-peer exp. converges as fast as over static exp.

Comparing different models/methods: one-peer v.s. static

MODEL TOPOLOGY	RESNET-50			MOBILENET-v2			EFFICIENTNET		
	STATIC	ONE-PEER	DIFF	STATIC	ONE-PEER	DIFF	STATIC	ONE-PEER	DIFF
PARALLEL SGD	76.21	-	-	70.12	-	-	77.63	-	-
VANILLA DMSGD	76.14	76.06	-0.08	69.98	69.81	-0.17	77.62	77.48	-0.14
DMSGD	76.50	76.52	+0.02	69.62	69.98	+0.36	77.44	77.51	+0.07
QG-DMSGD	76.43	76.35	-0.08	69.83	69.81	-0.02	77.60	77.72	+0.12

- setting: ImageNet; $8 \times 8 = 64$ GPUs; diff = o.e - s.e.
- both topo. achieve similar accuracy across different models and algorithms
- accuracy difference is minor (except for MobileNet with DmSGD)
- QG-DmSGD (Lin et al., 2021) and DmSGD can outperform PSGD in ResNet-50 in accuracy

Comparing different tasks: one-peer exp. v.s. static exp.

DATASET MODEL TOPOLOGY	PASCAL VOC				COCO			
	RETINANET		FASTER RCNN		RETINANET		FASTER RCNN	
	STATIC	ONE-PEER	STATIC	ONE-PEER	STATIC	ONE-PEER	STATIC	ONE-PEER
PARALLEL SGD	79.0	-	80.3	-	36.2	-	37.2	-
VANILLA DMSGD	79.0	79.1	80.7	80.5	36.3	36.1	37.3	37.2
DMSGD	79.1	79.0	80.4	80.5	36.4	36.4	37.1	37.0
QG-DMSGD	79.2	79.1	80.8	80.4	36.3	36.2	37.2	37.1

- setting: object detection; $8 \times 8 = 64$ GPUs;
- both topo. achieve similar accuracy across different algorithms in detection

Summary

- Both per-iter comm. and tran. iter. of exp. graphs are nearly best (up to $\log_2(n)$ factors) among known topologies
- While one-peer exp. is sparser, it can converge as fast as staic exp.
- One-peer exponential graph is recommend for decentralized DL

Part II: Making Decentralized SGD Practical for DNN

- Sec. 1. Exponential graphs are provably efficient (Ying et al., 2021)
- Sec. 2. Removing data heterogeneity enhances topology dependence (Huang and Pu, 2021; Yuan and Alghunaim, 2021)
- Sec. 3. Periodic global averaging (Chen et al., 2021)

D-SGD transient iteration complexity review

- Recall the convergence rate of D-SGD for non-convex and non-iid scenario:

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \|\nabla f(x^{(k)})\|^2 = O \left(\frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}} + \frac{\rho^{2/3} b^{2/3}}{T^{2/3} (1-\rho)^{2/3}} \right)$$

where $b^2 > 0$ deteriorates the dependence on network topology $1 - \rho$

- The transient iteration complexity of D-SGD is summarized as

scenario	iid data	non-iid data
strongly-convex	$\Omega\left(\frac{n}{1-\rho}\right)$	$\Omega\left(\frac{n}{(1-\rho)^2}\right)$
generally-convex	$\Omega\left(\frac{n^3}{(1-\rho)^2}\right)$	$\Omega\left(\frac{n^3}{(1-\rho)^4}\right)$
non-convex	$\Omega\left(\frac{n^3}{(1-\rho)^2}\right)$	$\Omega\left(\frac{n^3}{(1-\rho)^4}\right)$

D-SGD transient iteration complexity review

- Can we improve the dependence on topology for non-iid scenario?
- Main idea: remove the influence of b^2 from the convergence rate (Koloskova et al., 2020; Huang and Pu, 2021; Yuan et al., 2020; Yuan and Alghunaim, 2021)²
- Suppose a decentralized method for non-iid scenario can converge as

$$\frac{1}{T} \sum_{k=0}^{T-1} \mathbb{E} \|\nabla f(x^{(k)})\|^2 = O\left(\frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3} \sigma^{2/3}}{T^{2/3} (1-\rho)^{1/3}}\right)$$

it will improve the transient iteration complexity as follows

$$\Omega\left(\frac{\rho^4 n^3}{(1-\rho)^4}\right) \implies \Omega\left(\frac{\rho^4 n^3}{(1-\rho)^2}\right)$$

²K. Yuan and S. A. Alghunaim, "Removing data heterogeneity influence enhances network topology dependence of decentralized SGD", arXiv:2105.08023

How does D-SGD suffer from data heterogeneity?

- For simplicity, we consider the deterministic convex decentralized GD:

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} (x_j^{(k)} - \gamma \nabla f_j(x_j^{(k)})), \quad \forall i \in [n]$$

- Suppose $x_i^{(k)} = x^*$ at iteration k for any $i \in [n]$, it holds that

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

where the last inequality holds because $f_i(x) \neq f(x)$ (data-heterogeneous)

- D-GD cannot stay at x^* ; data heterogeneity incurs oscillation.

How does D-SGD suffer from data heterogeneity?



$$x_i^{(k)} = x^*$$



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

Remove the influence of data-heterogeneity

- EXTRA (Shi et al., 2015) is the first decentralized method to remove the influence of data heterogeneity
- Exact-Diffusion (Yuan et al., 2019) (also known as NIDS (Li et al., 2019) or D^2 (Tang et al., 2018)) improves EXTRA on learning rate stability range
- Gradient-tracking based methods (Xu et al., 2015; Di Lorenzo and Scutari, 2016; Nedic et al., 2017; Qu and Li, 2018; Pu et al., 2020b; Xin and Khan, 2018) remove data heterogeneity, and can be used in more relaxed settings (e.g., asymmetric/directed/time-varying weight matrix)
- All these algorithms can be unified into one decentralized framework (Alghunaim et al., 2020; Xu et al., 2021; Xin et al., 2020a)

Exact-Diffusion

- For Exact-Diffusion, each node run the following recursion in parallel

$\psi_i^{(k+1)} = x_i^{(k)} - \gamma \nabla F(x_i^{(k)}; \xi_i^{(k)})$	(local SGD)
$\phi_i^{(k+1)} = \psi_i^{(k+1)} + x_i^{(k)} - \psi_i^{(k)}$	(bias correction)
$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \phi_j^{(k+1)}$	(partial averaging)

- When correction term $x_i^{(k)} - \psi_i^{(k)}$ is removed from the correction step, Exact-Diffusion reduces to standard D-SGD
- The weight matrix W needs to be symmetric, and satisfies $\lambda_n(W) > -\frac{1}{3}$

How is Exact-Diffusion immune to data heterogeneity?

- Combining all recursions, we achieve the deterministic version

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} \left(2x_i^{(k)} - x_i^{(k-1)} + \gamma(\nabla f(x_i^{(k)}) - \nabla f(x_i^{(k-1)})) \right)$$

- Assume $x_i^{(k-1)} = x_i^{(k)} = x^*$ for any $i \in [n]$, at iteration $k + 1$ we have

$$x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} (2x^* - x^*) = x^*$$

- When initialized from the minimum, Exact-Diffusion can stay there in spite of the data heterogeneity $\nabla f_i(x) \neq \nabla f_j(x)$

Exact-Diffusion convergence

Assumption

(A1) Each local loss function $F(x; \xi_i)$ is L -smooth in terms of x ;

(A2) Each local stochastic gradient is unbiased, and has bounded variance σ^2

(A3) Each local stochastic gradient $g_i^{(k)}$ is independent of each other

(A4) W is positive semi-definite

Theorem (Yuan and Alghunaim (2021))

Under the above assumptions and with appropriate γ , Exact-Diffusion will converge at (S.C. is for strongly-convex and G.C. is for generally-convex)

$$\frac{1}{T+1} \sum_{k=0}^T (\mathbb{E}f(\bar{x}^{(k)}) - f(x^*)) = O\left(\frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3}\sigma^{2/3}}{(1-\rho)^{1/3}T^{2/3}}\right) \quad (\text{G.C.})$$

$$\frac{1}{H_T} \sum_{k=0}^T h_k (\mathbb{E}f(\bar{x}^{(k)}) - f(x^*)) = \tilde{O}\left(\frac{\sigma^2}{nT} + \frac{\rho^2\sigma^2}{(1-\rho)T^2}\right) \quad (\text{S.C.})$$

where h_k is some positive weight and $H_T = \sum_{k=0}^T h_k$.

Convergence comparison: Exact-Diffusion v.s. D-SGD

In the strongly-convex setting,

- The convergence rate comparison:

$$\text{D-SGD} : \tilde{O} \left(\frac{\sigma^2}{nT} + \frac{\rho^2 \sigma^2}{(1-\rho)T^2} + \frac{\rho^2 b^2}{(1-\rho)^2 T^2} \right)$$

$$\text{Exact-Diffusion} : \tilde{O} \left(\frac{\sigma^2}{nT} + \frac{\rho^2 \sigma^2}{(1-\rho)T^2} \right)$$

- The transient iteration complexity comparison (Huang and Pu, 2021; Yuan and Alghunaim, 2021):

$$\text{D-SGD} : \Omega \left(\frac{\rho^2 n}{(1-\rho)^2} \right) \quad \text{Exact-Diffusion} : \Omega \left(\frac{\rho^2 n}{1-\rho} \right)$$

Convergence comparison: Exact-Diffusion v.s. D-SGD

In the generally-convex setting,

- The convergence rate comparison:

$$\text{D-SGD : } O\left(\frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3}\sigma^{2/3}}{(1-\rho)^{1/3}T^{2/3}} + \frac{\rho^{2/3}b^{2/3}}{(1-\rho)^{2/3}T^{2/3}}\right)$$

$$\text{Exact-Diffusion : } O\left(\frac{\sigma}{\sqrt{nT}} + \frac{\rho^{2/3}\sigma^{2/3}}{(1-\rho)^{1/3}T^{2/3}}\right)$$

- The transient iteration comparison (Yuan and Alghunaim, 2021):

$$\text{D-SGD : } \Omega\left(\frac{\rho^4 n^3}{(1-\rho)^4}\right) \quad \text{Exact-Diffusion : } \Omega\left(\frac{\rho^4 n^3}{(1-\rho)^2}\right)$$

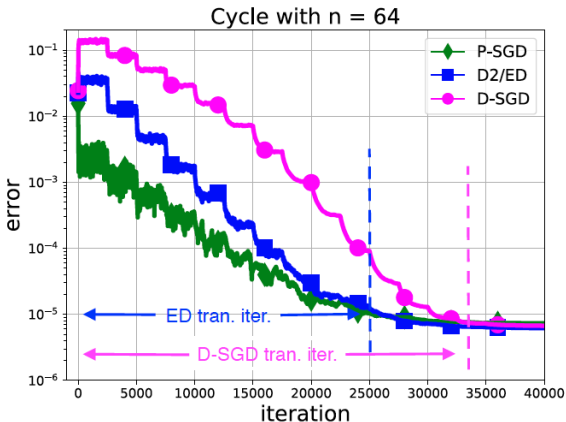
Convergence comparison: Exact-Diffusion v.s. D-SGD

In the non-convex setting,

- Exact-Diffusion can remove data heterogeneity (Tang et al., 2018), but no improved result on network topology dependence was shown
- Gradient-tracking can remove data heterogeneity (Xin et al., 2020b; Zhang and You, 2019; Lu et al., 2019), but no improved result on network topology dependence was shown
- It is still an open question whether data-heterogeneity-corrected methods (such as EXTRA, Exact-Diffusion, and Gradient tracking) can have an improved network topology dependence than P-SGD

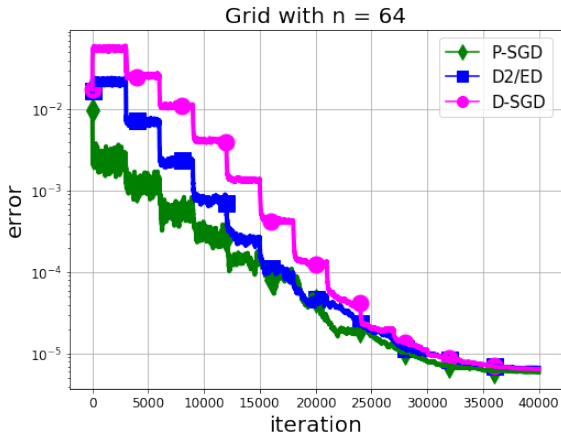
Experiments: Exact-Diffusion v.s. D-SGD

Convex setting: logistic regression problem; non-iid scenario



Convergence comparison: Exact-Diffusion v.s. D-SGD

Strongly-convex setting: least-square problem; non-iid scenario



Convergence comparison: Exact-Diffusion v.s. D-SGD

Deep learning experiments are on-going. No results yet.

Summary

- The data heterogeneity b^2 in D-SGD deteriorates the topology dependence
- EXTRA/Exact-Diffusion/Gradient-tracking can remove the influence of b^2
- Exact-Diffusion improves the topology dependence when b^2 exists.

non-iid scenario	Exact-Diffusion	D-SGD
strongly-convex	$\Omega\left(\frac{\rho^2 n}{1-\rho}\right)$	$\Omega\left(\frac{\rho^2 n}{(1-\rho)^2}\right)$
generally-convex	$\Omega\left(\frac{\rho^4 n^3}{(1-\rho)^2}\right)$	$\Omega\left(\frac{\rho^4 n^3}{(1-\rho)^4}\right)$
non-convex	N.A.	$\Omega\left(\frac{\rho^4 n^3}{(1-\rho)^4}\right)$

Part II: Making Decentralized SGD Practical for DNN

- Sec. 1. Exponential graphs are provably efficient (Ying et al., 2021)
- Sec. 2. Removing data heterogeneity enhances topology dependence (Huang and Pu, 2021; Yuan and Alghunaim, 2021)
- Sec. 3. Periodic global averaging (Chen et al., 2021)

Motivation

- Recall non-convex D-SGD suffers from additional transient iterations

$$\text{homogeneous (iid) data: } \Omega\left(\frac{\rho^4 n^3}{(1-\rho)^2}\right)$$

$$\text{heterogeneous (non-iid) data: } \Omega\left(\frac{\rho^4 n^3}{(1-\rho)^4}\right)$$

- $\rho \rightarrow 1$ will significantly enlarge the transient iteration stage
- Unfortunately, most topologies have $\rho \rightarrow 1$ as n grows
 - Ring: $1 - \rho = O(1/n^2)$;
 - Grid: $1 - \rho = O(1/n)$;
 - Exp.: $1 - \rho = O(1/\log_2(n))$
- We have to alleviate the influence of $1/(1 - \rho)$ in trans. iters. complexity

Per-iteration communication cost

Model	Ring-Allreduce	Partial average
ResNet-50	278 ms	150 ms
Bert	1469 ms	567 ms

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

- While global average takes longer comm. time, it is not too bad
- We can mix partial average with global average (Chen et al., 2021)³.
- In a period of H iterations: run $H - 1$ partial average and 1 global average

³Y. Chen*, K. Yuan*, Y. Zhang, P. Pan, Y. Xu, W. Yin, "Accelerating Gossip SGD with Periodic Global Averaging", ICML 2021

DSGD-PGA: DSGD with Periodic Global Averaging

- DSGD-PGA: accelerate D-SGD with periodic global averaging

$$\begin{aligned}\mathbf{x}_i^{(k+\frac{1}{2})} &= \mathbf{x}_i^{(k)} - \gamma \nabla F(\mathbf{x}_i^{(k)}; \xi_i^{(k+1)}) \\ \mathbf{x}_i^{(k+1)} &= \begin{cases} \frac{1}{n} \sum_{j=1}^n \mathbf{x}_j^{(k+\frac{1}{2})} & \text{If } \text{mod}(k+1, H) = 0 \\ \sum_{j \in \mathcal{N}_i} w_{ij} \mathbf{x}_j^{(k+\frac{1}{2})} & \text{If } \text{mod}(k+1, H) \neq 0 \end{cases}\end{aligned}$$

where H is the global averaging period.

- DSGD-PGA is expected to converge faster than D-SGD.
- DSGD-PGA reduces to D-SGD when $H \rightarrow \infty$
- Similar idea also appeared in topology-changing D-SGD (Koloskova et al., 2020) and SlowMo (Wang et al., 2019)

DSGD-PGA: Transient iteration complexity

- PGA significantly improves the transient stage of D-SGD in the non-convex setting (Chen et al., 2021):

scenario	DSGD-PGA	D-SGD
iid data	$\Omega(\rho^4 n^3 H^2)$	$\Omega(\frac{\rho^4 n^3}{(1-\rho)^2})$
non-iid data	$\Omega(\rho^4 n^3 H^4)$	$\Omega(\frac{\rho^4 n^3}{(1-\rho)^4})$

- PGA bounds $1/(1 - \rho)$ with H ; benefits most for sparse topology

Numerical experiments: D-SGD v.s. DSGD-PGA

Problem: logistic regression problem with non-iid data

Cyclic Topology

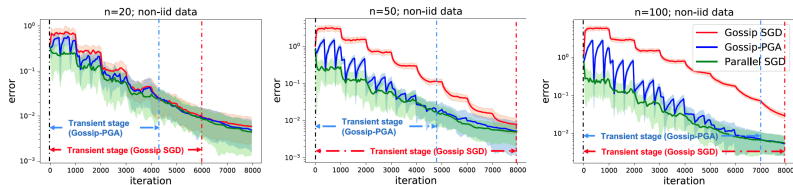


Figure: Transient stage comparison.

DSGD-AGA: D-SGD with Adaptive Global Averaging

- Gossip-AGA avoids the burden of tuning parameters
- An effective period strategy: **more frequent GA in initial stages**
- Intuition: lower consensus variance can speedup convergence

$$\frac{1}{n(T+1)} \sum_{k=0}^T \sum_{i=1}^n \mathbb{E} \|\mathbf{x}_i^{(k)} - \bar{\mathbf{x}}^{(k)}\|^2 \leq \frac{d_1 \gamma^2}{T+1} \sum_{k=0}^T \mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(k)})\|^2 + d_2 \gamma^2$$

Consensus variance gets decreased as $\gamma \rightarrow 0$ and $\mathbb{E} \|\nabla f(\bar{\mathbf{x}}^{(k)})\|^2 \rightarrow 0$

- Adaptive rule: $H^{(\ell)} = \left(\frac{\mathbb{E} f(\bar{\mathbf{x}}^{(0)})}{\mathbb{E} f(\bar{\mathbf{x}}^{(T\ell-1)})} \right)^{\frac{1}{4}} H^{(0)}$;

Experiments on Large-scale Deep Training

Language Modeling:

- Model: BERT-Large (~ 330 M parameters)
- Dataset: Wikipedia (2500M words) and BookCorpus (800M words)
- Hardware: 64 GPUs

Image Classification

Method	Final Loss	Wall-clock Time (hrs)
P-SGD	1.75	59.02
D-SGD	2.17	29.7
D-SGD $\times 2$	1.81	59.7
DSGD-PGA	1.82	35.4
DSGD-AGA	1.77	30.4

Table: Comparison of training loss and training time of BERT training.

- DSGD-AGA achieves similar final loss with $2\times$ speedup

Summary

- Periodic global averaging can improve the transient iteration stage:

$$\Omega\left(\frac{\rho^4 n^3}{(1-\rho)^4}\right) \implies \Omega(\rho^4 n^3 H^4)$$

- PGA benefits most for sparse topology, i.e., $\rho \rightarrow 1$
- Global averaging period H can be adjusted adaptively

Discussion

- We consider deep training within high-performance data-center clusters
- Global averaging conducted by All-reduce has tolerable comm. cost
- For mobile AI or federated learning, global averaging is very expensive
- We can approximate global averaging via multiple partial averaging steps, see [Lu and De Sa, 2021, ICML Outstanding Paper Honorable mention]
- However, multiple partial averaging steps are not recommended for data-center clusters; 3 partial averaging steps may take more wall-clock time than one single global averaging

In Part III, we will

Discuss large-batch decentralized deep training, and

Introduce BlueFog, an open-source library to help deploy decentralized methods into real CPU/GPU clusters

References I

- S. Pu, A. Olshevsky, and I. C. Paschalidis, "Asymptotic network independence in distributed stochastic optimization for machine learning: Examining distributed and centralized stochastic gradient descent," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 114–122, 2020.
- B. Ying, K. Yuan, Y. Chen, H. Hu, P. Pan, and W. Yin, "Exponential graph is provably efficient for deep training," in *Submitted*, 2021.
- K. Huang and S. Pu, "Improving the transient times for distributed stochastic gradient methods," *arXiv preprint arXiv:2105.04851*, 2021.
- K. Yuan and S. A. Alghunaim, "Removing data heterogeneity influence enhances network topology dependence of decentralized sgd," *arXiv preprint arXiv:2105.08023*, 2021.
- Y. Chen, K. Yuan, Y. Zhang, P. Pan, Y. Xu, and W. Yin, "Accelerating gossip sgd with periodic global averaging," *International Conference on Machine Learning*, 2021.

References II

- A. Nedić, A. Olshevsky, and M. G. Rabbat, “Network topology and communication-computation tradeoffs in decentralized optimization,” *Proceedings of the IEEE*, vol. 106, no. 5, pp. 953–976, 2018.
- X. Lian, C. Zhang, H. Zhang, C.-J. Hsieh, W. Zhang, and J. Liu, “Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent,” in *Advances in Neural Information Processing Systems*, 2017, pp. 5330–5340.
- X. Lian, W. Zhang, C. Zhang, and J. Liu, “Asynchronous decentralized parallel stochastic gradient descent,” in *International Conference on Machine Learning*, 2018, pp. 3043–3052.
- M. Assran, N. Loizou, N. Ballas, and M. Rabbat, “Stochastic gradient push for distributed deep learning,” in *International Conference on Machine Learning (ICML)*, 2019, pp. 344–353.
- T. Lin, S. P. Karimireddy, S. U. Stich, and M. Jaggi, “Quasi-global momentum: Accelerating decentralized deep learning on heterogeneous data,” *arXiv preprint arXiv:2102.04761*, 2021.

References III

- A. Koloskova, N. Loizou, S. Boreiri, M. Jaggi, and S. U. Stich, "A unified theory of decentralized sgd with changing topology and local updates," in *International Conference on Machine Learning (ICML)*, 2020, pp. 1–12.
- K. Yuan, S. A. Alghunaim, B. Ying, and A. H. Sayed, "On the influence of bias-correction on distributed stochastic optimization," *IEEE Transactions on Signal Processing*, 2020.
- W. Shi, Q. Ling, G. Wu, and W. Yin, "EXTRA: An exact first-order algorithm for decentralized consensus optimization," *SIAM Journal on Optimization*, vol. 25, no. 2, pp. 944–966, 2015.
- K. Yuan, B. Ying, X. Zhao, and A. H. Sayed, "Exact diffusion for distributed optimization and learning – Part I: Algorithm development," *IEEE Transactions on Signal Processing*, vol. 67, no. 3, pp. 708 – 723, 2019.
- Z. Li, W. Shi, and M. Yan, "A decentralized proximal-gradient method with network independent step-sizes and separated convergence rates," *IEEE Transactions on Signal Processing*, July 2019, early acces. Also available on arXiv:1704.07807.

References IV

- H. Tang, X. Lian, M. Yan, C. Zhang, and J. Liu, " d^2 : Decentralized training over decentralized data," in *International Conference on Machine Learning*, 2018, pp. 4848–4856.
- J. Xu, S. Zhu, Y. C. Soh, and L. Xie, "Augmented distributed gradient methods for multi-agent optimization under uncoordinated constant stepsizes," in *IEEE Conference on Decision and Control (CDC)*, Osaka, Japan, 2015, pp. 2055–2060.
- P. Di Lorenzo and G. Scutari, "Next: In-network nonconvex optimization," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 2, no. 2, pp. 120–136, 2016.
- A. Nedic, A. Olshevsky, and W. Shi, "Achieving geometric convergence for distributed optimization over time-varying graphs," *SIAM Journal on Optimization*, vol. 27, no. 4, pp. 2597–2633, 2017.
- G. Qu and N. Li, "Harnessing smoothness to accelerate distributed optimization," *IEEE Transactions on Control of Network Systems*, vol. 5, no. 3, pp. 1245–1260, 2018.

References V

- S. Pu, W. Shi, J. Xu, and A. Nedić, "Push–pull gradient methods for distributed optimization in networks," *IEEE Transactions on Automatic Control*, vol. 66, no. 1, pp. 1–16, 2020.
- R. Xin and U. A. Khan, "A linear algorithm for optimization over directed graphs with geometric convergence," *IEEE Control Systems Letters*, vol. 2, no. 3, pp. 315–320, 2018.
- S. A. Alghunaim, E. K. Ryu, K. Yuan, and A. H. Sayed, "Decentralized proximal gradient algorithms with linear convergence rates," *IEEE Transactions on Automatic Control*, vol. 66, no. 6, pp. 2787–2794, 2020.
- J. Xu, Y. Tian, Y. Sun, and G. Scutari, "Distributed algorithms for composite optimization: Unified framework and convergence analysis," *IEEE Transactions on Signal Processing*, 2021.
- R. Xin, S. Pu, A. Nedić, and U. A. Khan, "A general framework for decentralized optimization with first-order methods," *Proceedings of the IEEE*, vol. 108, no. 11, pp. 1869–1889, 2020.

References VI

- R. Xin, U. A. Khan, and S. Kar, "An improved convergence analysis for decentralized online stochastic non-convex optimization," *arXiv preprint arXiv:2008.04195*, 2020.
- J. Zhang and K. You, "Decentralized stochastic gradient tracking for non-convex empirical risk minimization," *arXiv preprint arXiv:1909.02712*, 2019.
- S. Lu, X. Zhang, H. Sun, and M. Hong, "Gnsd: A gradient-tracking based nonconvex stochastic algorithm for decentralized optimization," in *2019 IEEE Data Science Workshop (DSW)*. IEEE, 2019, pp. 315–321.
- J. Wang, V. Tantia, N. Ballas, and M. Rabbat, "Slowmo: Improving communication-efficient distributed sgd with slow momentum," *arXiv preprint arXiv:1910.00643*, 2019.