



A Mathematics-Inspired Learning-to-Optimize Framework for Decentralized Optimization

Kun Yuan (袁 坤)

Center for Machine Learning Research @ Peking University

Oct. 21, 2024

Joint work with



Yutong He
(Peking U.)



Qiulin Shang
(Peking U.)



Xinmeng Huang
(UPenn)



Jialin Liu
(CFU)



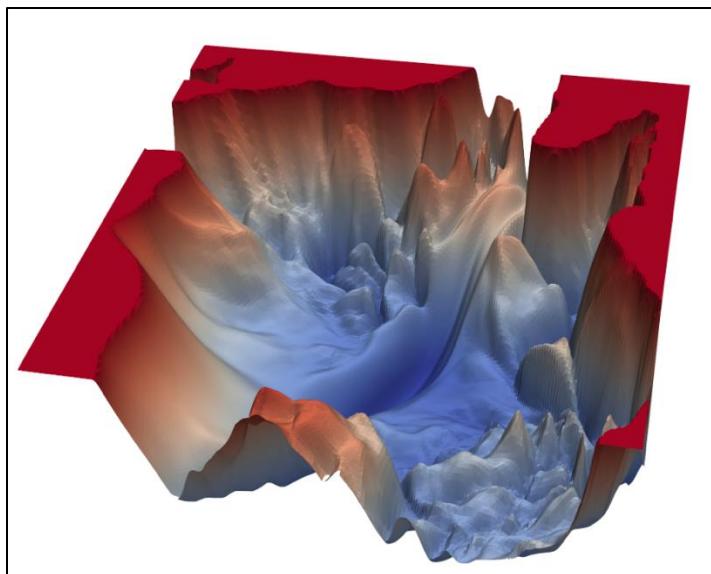
PART 01

Basics and Motivation

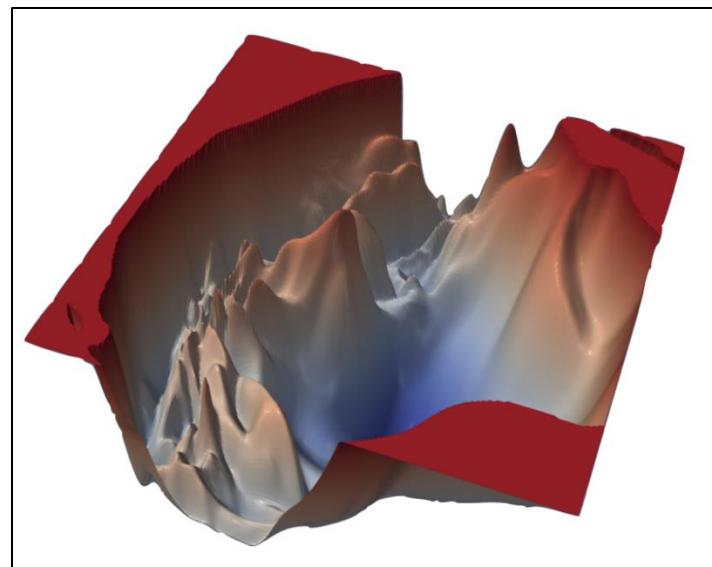
DNN model is notoriously difficult to train

- Highly **non-convex** cost functions; cannot find global minima; trapped into local minimum

VGG-56



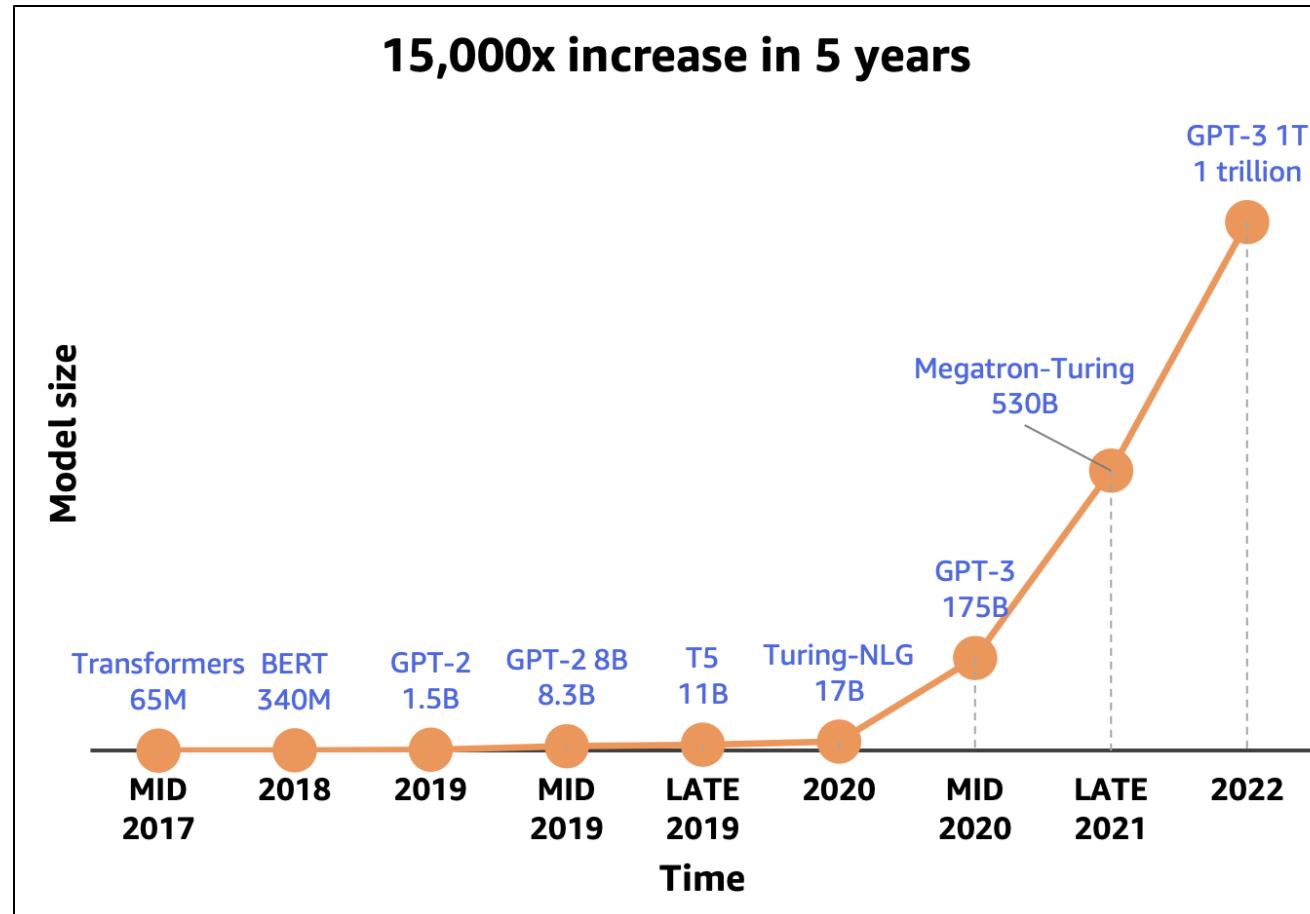
VGG-110



[Visualizing the loss landscape of neural nets]

DNN model is notoriously difficult to train

- The model size is **large**, i.e., $x \in \mathbb{R}^d$ is of extremely high dimensions



[Train and deploy large language models on Amazon SageMaker]

DNN model is notoriously difficult to train

- The size of the dataset is **huge**

Dataset	Sampling prop.	Epochs	Disk size
CommonCrawl	67.0%	1.10	3.3 TB
C4	15.0%	1.06	783 GB
Github	4.5%	0.64	328 GB
Wikipedia	4.5%	2.45	83 GB
Books	4.5%	2.23	85 GB
ArXiv	2.5%	1.06	92 GB
StackExchange	2.0%	1.03	78 GB

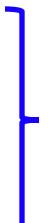
 Crawled data from websites; in both high quality and low quality
 High-quality data

Table 1: **Pre-training data.** Data mixtures used for pre-training, for each subset we list the sampling proportion, number of epochs performed on the subset when training on 1.4T tokens, and disk size. The pre-training runs on 1T tokens have the same sampling proportion.

[LLaMA: Open and Efficient Foundation Language Models]

DNN model is notoriously difficult to train



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

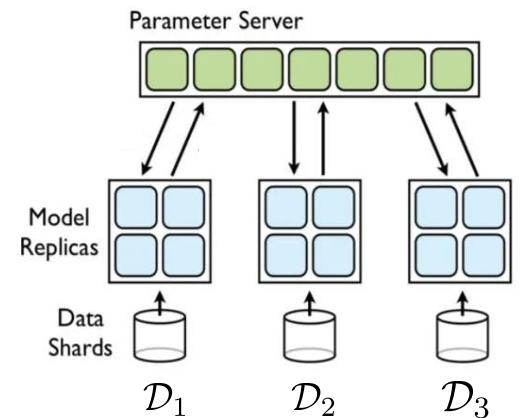
Efficient and **scalable** distributed optimization approaches are in urgent need

Distributed learning

- This talk focuses on data-parallel approaches to process **massive** datasets;
- 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 ξ_i denotes the local data that follows distribution D_i
- Each local distribution D_i is different; data heterogeneity exists



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

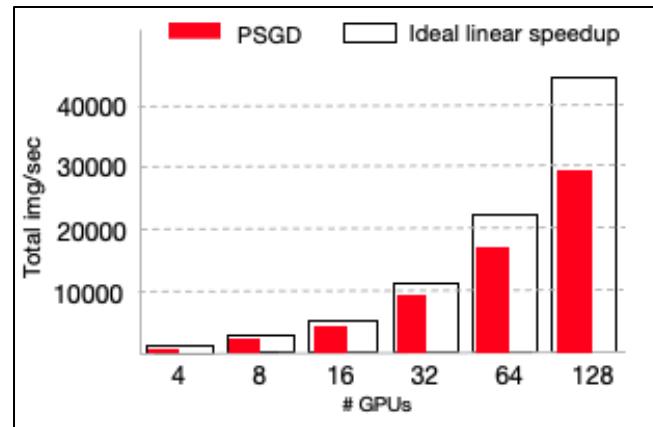
$$g_i^{(k)} = \nabla F(x^{(k)}; \xi_i^{(k)}) \quad (\text{Local compt.})$$

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

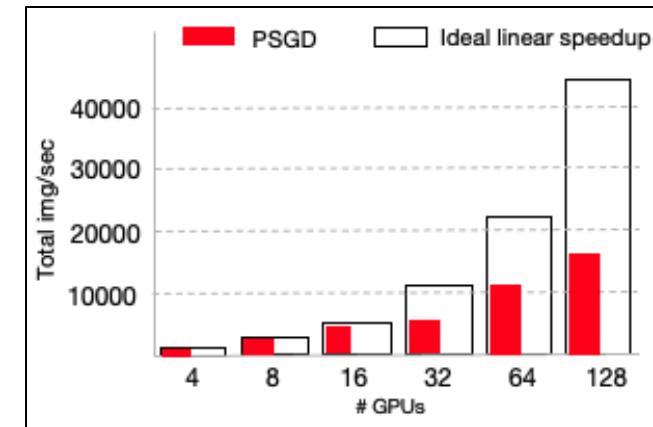
- Each node i samples data $\xi_i^{(k)}$ and computes gradient $\nabla F(x^{(k)}; \xi_i^{(k)})$
- All nodes synchronize (i.e. globally average) to update model x per iteration
- The global average takes significant communication overhead (esp. when n is large)

PSGD cannot achieve ideal scalability due to comm. overhead

- PSGD **cannot** achieve ideal linear speedup in throughput due to comm. overhead in global average
- 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?

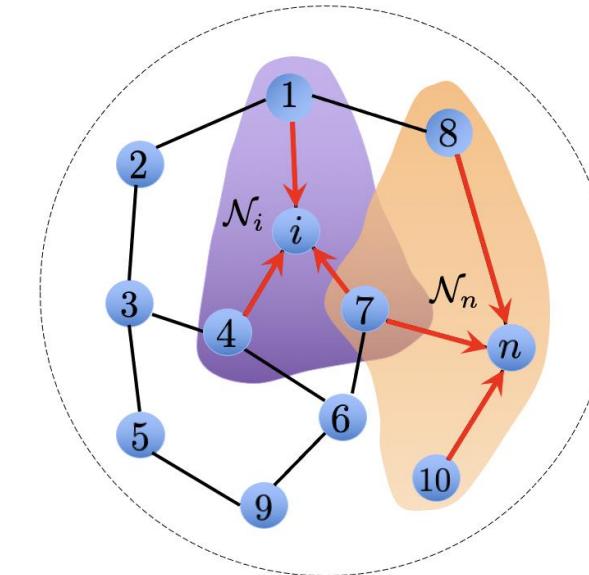
Decentralized SGD (DSGD)

- To reduce 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

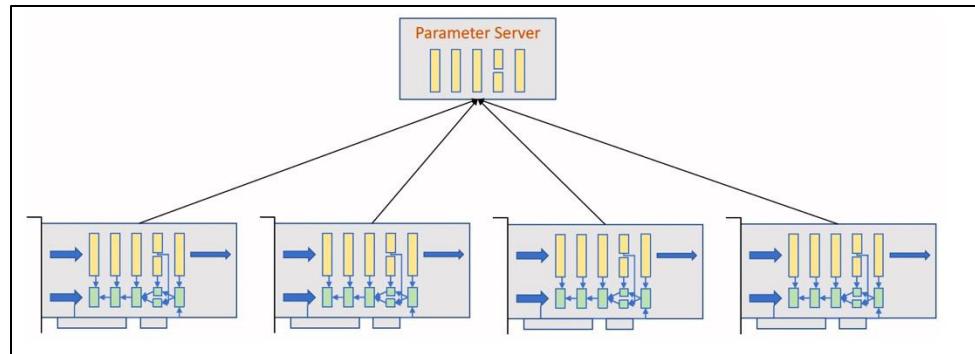


[1] C. G. Lopes and A. H. Sayed, "Diffusion least-mean-squares over adaptive networks", ICASSP, 2007

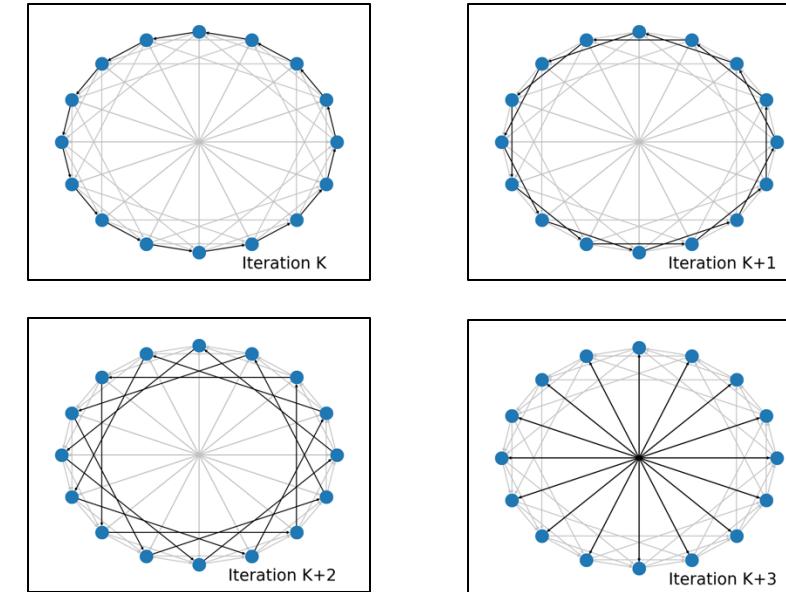
[2] A. Nedich and A. Ozdaglar, " Distributed subgradient methods for multi-agent optimization", IEEE TAC, 2009

DSGD is more communication-efficient than PSGD

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



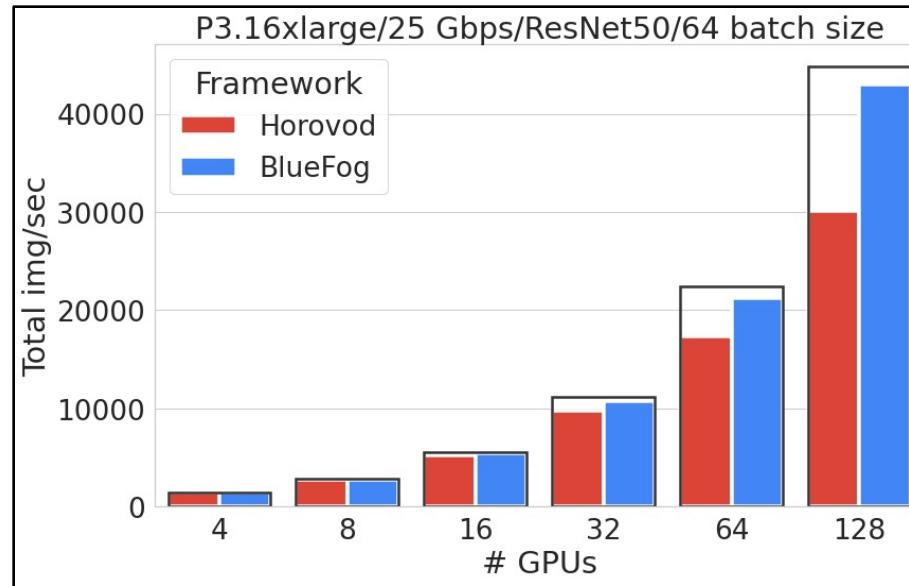
Global averaging has $O(n)$ comm. overhead



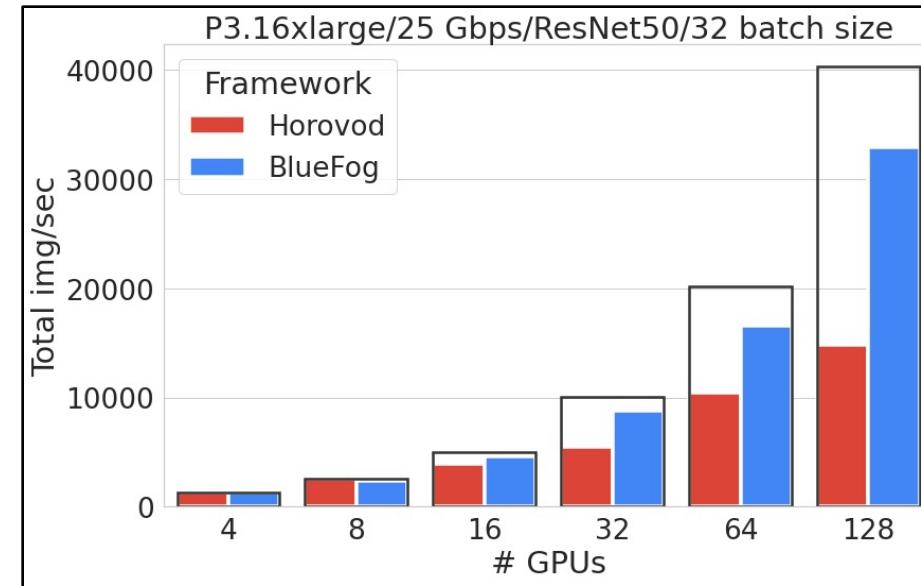
Partial averaging has $O(1)$ comm. overhead

DSGD is more communication-efficient than PSGD

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



Small comm.-to-compt. ratio



Large comm.-to-compt. ratio

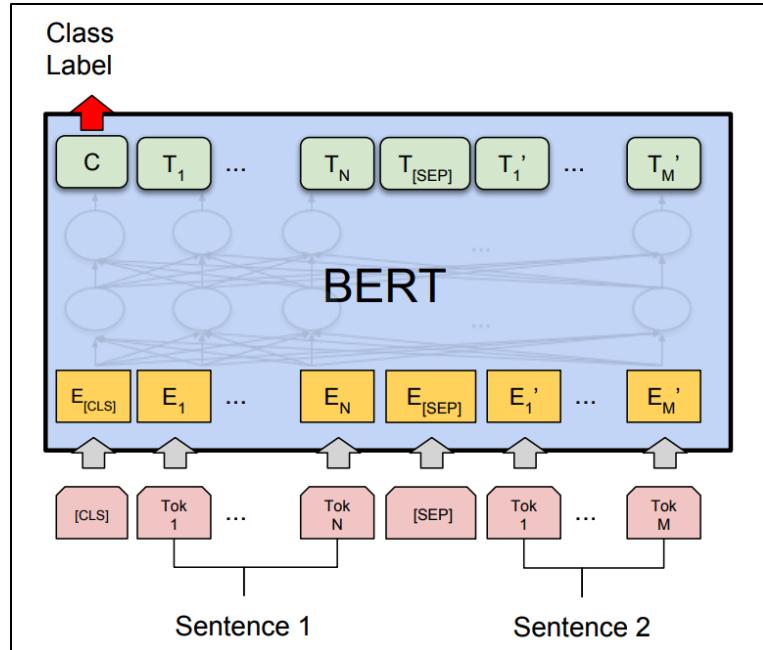
DSGD saves more wall-clock time without severely hurting performance



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	76.34	11.1	76.52	5.7	76.47	2.8	76.27	1.5

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

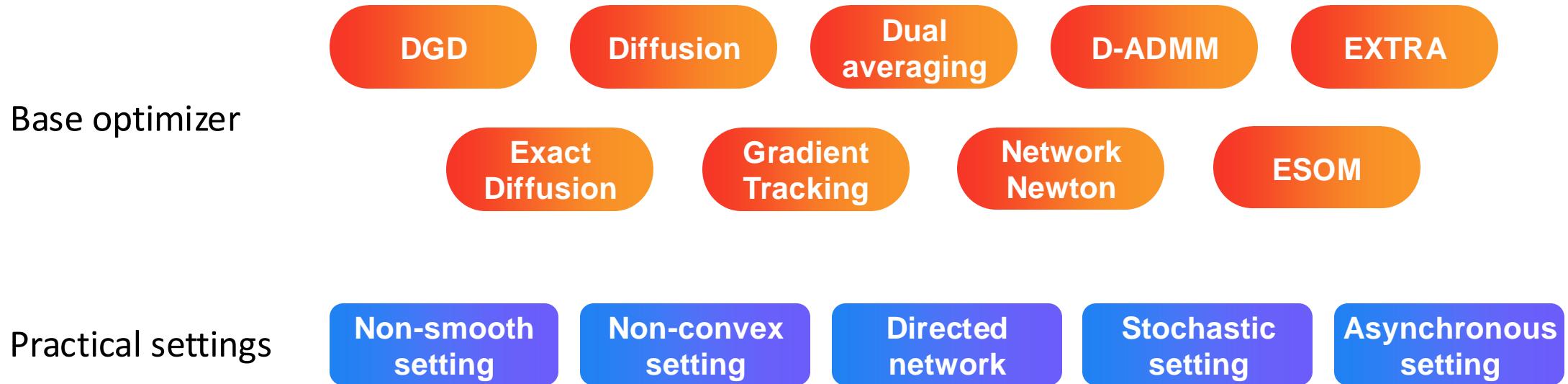
Table. Comparison in loss and training time [1]

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

Many well-known decentralized algorithms have been proposed

- Since 2009, many well-known decentralized algorithms have been proposed in literature



- The combination of base optimizer and practical settings gives rise to massive algorithms
- All these algorithms are **hand-crafted**

Hand-crafted algorithms

$$f(x) \in \mathcal{F}$$

Function space



Expert design



$$x^{(k+1)} = x^{(k)} - \gamma \nabla f(x^{(k)})$$

Hand-crafted optimizers

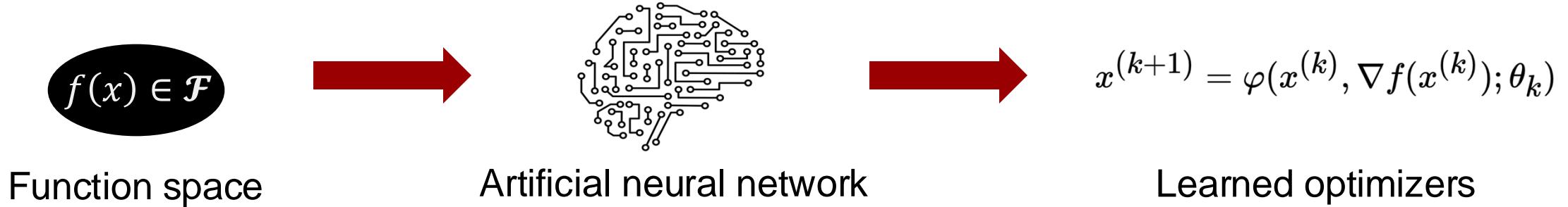
- 😊 Works well for a **wide range** of target problems
- 😢 Excessively focused on worst-case problems, overlooking the problem-specific structures
- 😢 Hyperparameters to **tune**
- 😢 Hand-crafted algorithms can be **sub-optimal**

Can we develop algorithms automatically?

PART 02

Learn-to-Optimize

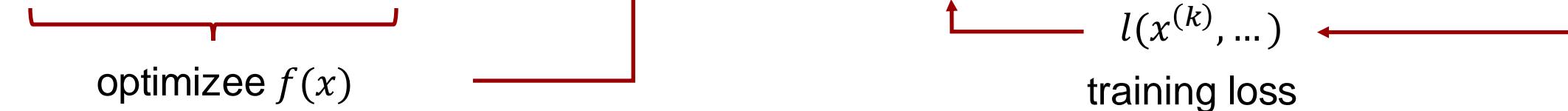
Learning to optimize



- :() Works well on a **specialized** domain; not on all problems
- :) Trained on real problems, learned to utilize problem-specific structures
- :) **Free** from hyperparameter tuning
- :) **Faster** than hand-crafted optimizers

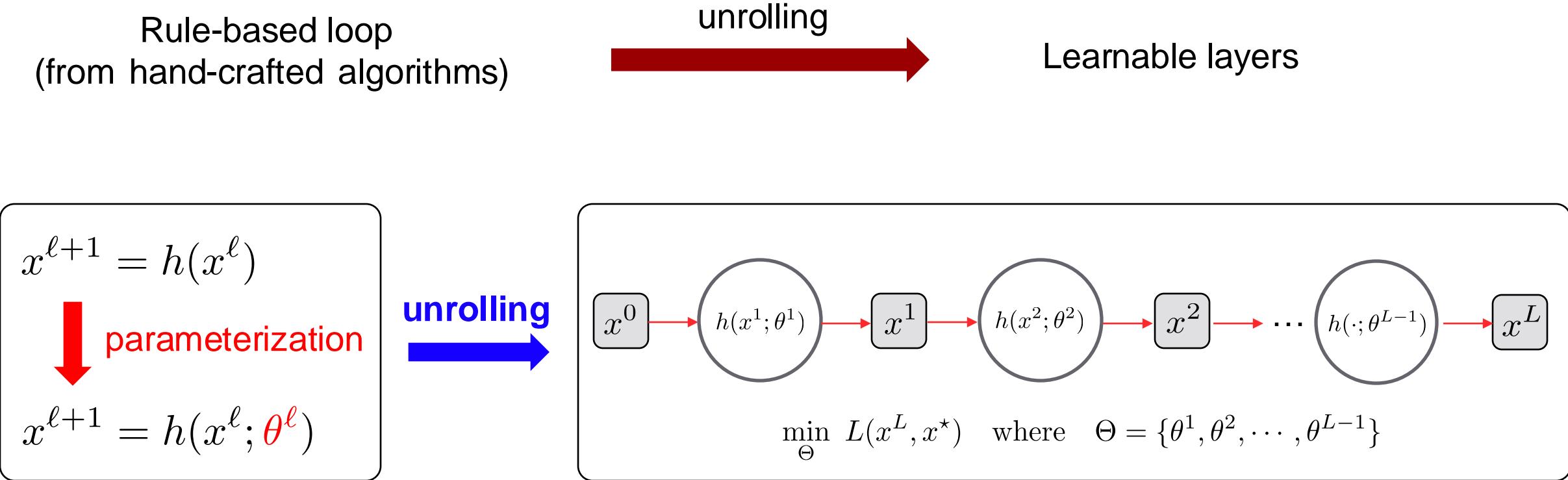
How to learn an optimizer?

$$\begin{aligned}
 & \min \|A_1x - b_1\|^2 \\
 & \min \|A_2x - b_2\|^2 \\
 & \min \|A_3x - b_3\|^2 \\
 & \dots
 \end{aligned}
 \quad
 \begin{aligned}
 & \min (1/N) \sum_{n=1}^N \ln (1 + \exp(-b_{1,n}a_{1,n}^T x)) \\
 & \min (1/N) \sum_{n=1}^N \ln (1 + \exp(-b_{2,n}a_{2,n}^T x)) \\
 & \dots
 \end{aligned}$$



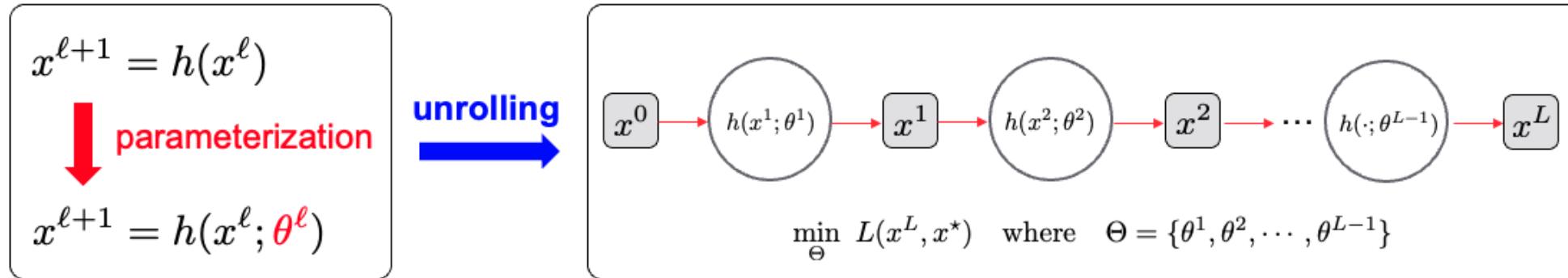
- Training data: target problem functions to optimize, called **optimizees**.
- Learned optimizer: neural network models that map x^{k-1} to x^k
- Loss function: performance of the learned optimizer on training optimizees
- Two major approaches: **algorithm unrolling** and **generic L2O**

Algorithm unrolling: principle



K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. In Proceedings of the 27th international conference on international conference on machine learning, pages 399–406, 2010.

Algorithm unrolling: Pros and Cons



- 😊 Perform **well** on in-distribution tasks
- 😊 **Compact** search space; easy to train
- 😢 Based on **hand-crafted** algorithms; not learning new algorithms but only improve them
- 😢 **Hardly** generalize to longer iterations due to **storage issues**

- Generic L2O typically learns a recurrent neural network (RNN) as the iterative update rule without a hand-crafted base algorithm.

Specific rule:

$$\theta_{t+1} = \theta_t - \alpha_t \nabla f(\theta_t)$$



parameterization

Abstract rule:

$$\theta_{t+1} = \theta_t + g_t(\nabla f(\theta_t), \phi)$$



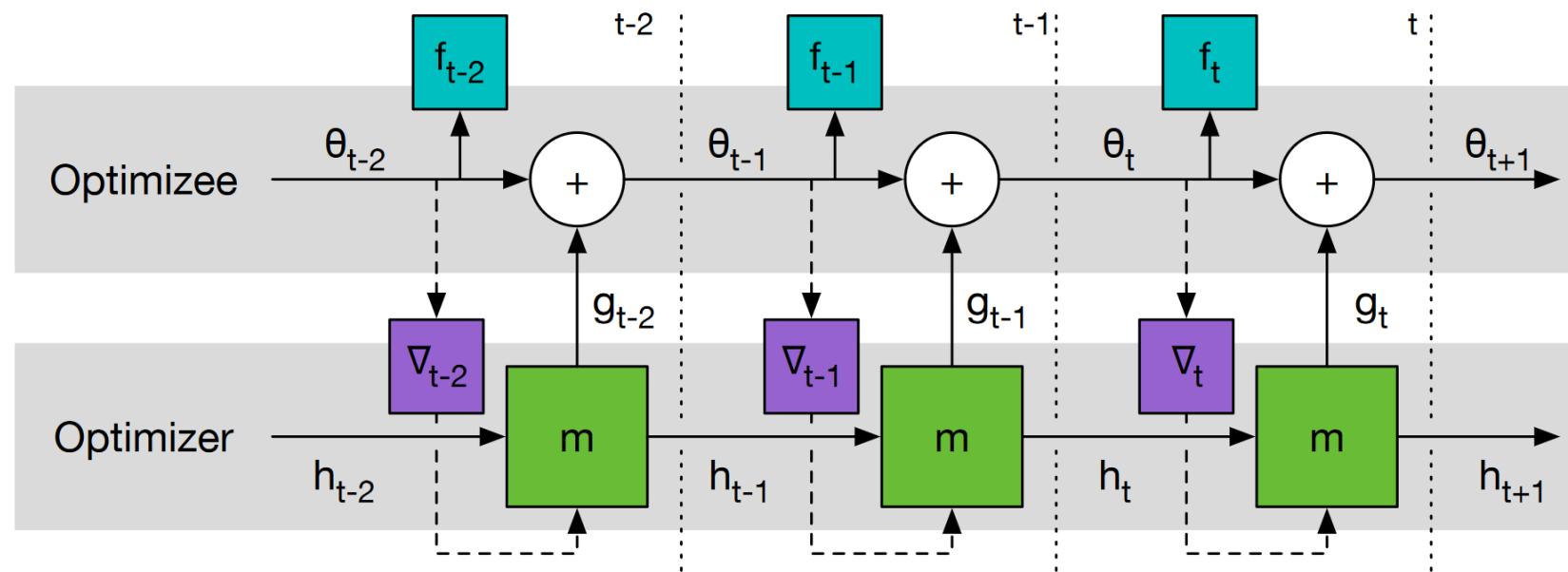
RNN parameterization

Learning to learn with recurrent neural networks

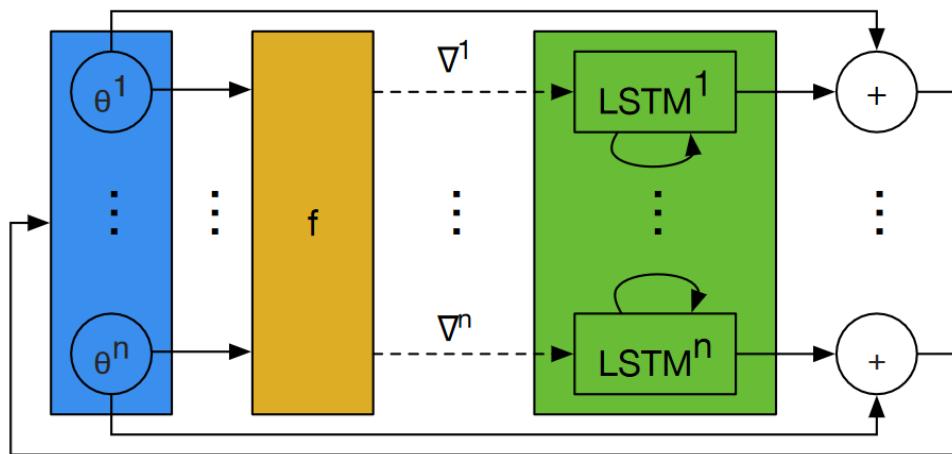
$$\theta_{t+1} = \theta_t + g_t ,$$

$$\begin{bmatrix} g_t \\ h_{t+1} \end{bmatrix} = m(\nabla_t, h_t, \phi) .$$

Very general; do not specialize to any concrete algorithms such as momentum GD, adaptive GD, etc.



Coordinate-wise LSTM optimizer



- LSTMs for different coordinates have **shared parameters** but **different hidden state**.
- Coordinate-wise LSTM optimizer has **fewer trainable parameters**, and is **more flexible** to problems with different dimensions.

- 😊 RNN shares the same parameters; **Memory efficient**; Generalize to **longer** iterations
- 😊 Do **not depend on** a hand-crafted algorithm
- 😢 **Vast search space**; difficult to learn a good update rule
- 😢 **Cannot handle constraints**

Unrolling

- 😊 Perform **well** on in-distribution tasks
- 😊 **Compact** search space; easy to train
- 😢 Based on **hand-crafted** algorithms
- 😢 Memory **expensive**

Generic L2O

- 😊 **Memory efficient**; fit for long iterations
- 😊 **No need** for hand-crafted algorithms
- 😢 **Cannot handle constraints**
- 😢 **Vast** search space; **difficult** to train

Both approaches show **weak generalization** on **out-of-distribution** problems!

Can we develop an effective L2O algorithm for decentralized optimization that is

- **Memory efficient:** the learned algorithm can run arbitrarily large iterations
- **Easy to train:** the search space is compact
- **Enforcing consensus constraint:** ensure convergence to the consensus solution among nodes
- **Strong generalization:** can generalize well to out-of-distribution problems

PART 03

A math-inspired L2O framework for decentralized optimization

A trivial decentralized L2O framework

Consider a distributed optimization problem:

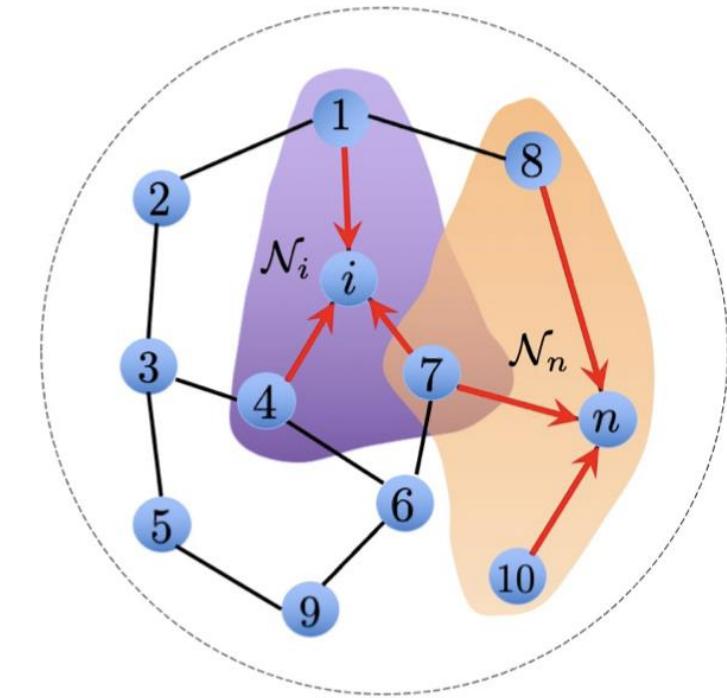
$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x) \quad \text{where} \quad f_i(x) = \mathbb{E}_{\xi_i \sim D_i} F(x; \xi_i)$$

Recall the generic L2O approach:

$$x^{k+1} = x^k + g^k(\nabla f(x^k); \theta), \quad k = 0, 1, \dots, K-1$$

If we directly follow the generic L2O approach, we arrive at

$$x_i^{k+1} = x_i^k + \boxed{g_i^k} \left(\{x_j^k, \nabla f_j(x_j^k)\}_{j \in \mathcal{N}(i) \cup \{i\}}; \boxed{\theta_i} \right), \quad \forall i \in \mathcal{V},$$



takes gradients and variables
from neighbors

each node i learns an RNN
parameterized by θ_i

$$x_i^{k+1} = x_i^k + g_i^k \left(\{x_j^k, \nabla f_j(x_j^k)\}_{j \in \mathcal{N}(i) \cup \{i\}}; \theta_i \right), \quad \forall i \in \mathcal{V},$$

- The search space consisting of all combinations between the gradients and variables from neighbors is extremely **large**; very **difficult** to train
- Each node updates its own variables independently; **no mechanism to ensure** convergence to the optimal consensus solution $x_1 = \dots = x_n = x^*$
- **Weak** generalization is observed in experiments

Can we find more structures in the update rule $g_i(\cdot)$

Base update rules

We reformulate the unconstrained optimization problem into the following constrained form:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times d}} \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}_i) + r(\mathbf{x}_i), \quad \text{s.t.} \quad \mathbf{x}_i = \mathbf{x}_j, \forall \{i, j\} \in \mathcal{E}.$$

The corresponding Lagrangian function is:

$$\mathcal{L}\left(X, \{\nu_{i,j}\}_{\{i,j\} \in \mathcal{E}}\right) = \frac{1}{n} \sum_{i=1}^n (f_i(x_i) + r(x_i)) + \sum_{\{i,j\} \in \mathcal{E}} \langle \nu_{i,j}, x_i - x_j \rangle,$$

Solving the above function using the primal-dual algorithm yields:

$$\begin{aligned} x_i^{k+1} &= x_i^k - \frac{\gamma}{n} (\nabla f_i(x_i^k) + g_i^{k+1} + y_i^k), \quad g_i^{k+1} \in \partial r(x_i^{k+1}) \\ y_i^{k+1} &= y_i^k + 2\gamma \sum_{j \in \mathcal{N}(i)} (x_i^{k+1} - x_j^{k+1}) \end{aligned}$$

Base update rules

$$x_i^{k+1} = x_i^k - \frac{\gamma}{n} (\nabla f_i(x_i^k) + g_i^{k+1} + y_i^k), \quad g_i^{k+1} \in \partial r(x_i^{k+1}) \quad (1)$$

$$y_i^{k+1} = y_i^k + 2\gamma \sum_{j \in \mathcal{N}(i)} (x_i^{k+1} - x_j^{k+1}) \quad (2)$$

- Update (2) implies **consensus**. With $x_i^k \rightarrow x_i^*$ and $y_i^k \rightarrow y_i^*$ we have

$$x_i^* = \left(\frac{1}{|\mathcal{N}(i)|} \right) \sum_{j \in \mathcal{N}(i)} x_j^*, \quad \forall i \in \{1, \dots, n\} \quad \longrightarrow \quad \mathbf{x}_1^* = \dots = \mathbf{x}_n^*.$$

- Update (1) implies **optimality**. Since $x_i^k \rightarrow x_i^*$, $y_i^k \rightarrow y_i^*$ and $\sum_{i=1}^n y_i^* = 0$ we have

$$\nabla f_i(\mathbf{x}^*) + \partial r(\mathbf{x}^*) + \mathbf{y}_i^* \ni 0 \quad \longrightarrow \quad \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}^*) + \partial r(\mathbf{x}^*) \ni 0,$$

Base update rules

$$x_i^{k+1} = x_i^k - \frac{\gamma}{n} (\nabla f_i(x_i^k) + g_i^{k+1} + y_i^k), \quad g_i^{k+1} \in \partial r(x_i^{k+1})$$

$$y_i^{k+1} = y_i^k + 2\gamma \sum_{j \in \mathcal{N}(i)} (x_i^{k+1} - x_j^{k+1})$$

- Inspired by the above primal-dual algorithm, we propose a new L2O update rule

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^k - \mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k; \boldsymbol{\theta}_{i,1}), \quad \mathbf{g}_i^{k+1} \in \partial r(\mathbf{x}_i^{k+1}),$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + \mathbf{s}_i^k(\{\mathbf{x}_i^{k+1} - \mathbf{x}_j^{k+1}\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,2}),$$

- Compared with the following trivial update rule

$$x_i^{k+1} = x_i^k + g_i^k \left(\{x_j^k, \nabla f_j(x_j^k)\}_{j \in \mathcal{N}(i) \cup \{i\}}; \theta_i \right), \quad \forall i \in \mathcal{V},$$

We introduce the **dual variable \mathbf{y}** to enforce consensus constraints

Base update rules

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^k - \mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k; \boldsymbol{\theta}_{i,1}), \quad \mathbf{g}_i^{k+1} \in \partial r(\mathbf{x}_i^{k+1})$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + s_i^k \left(\left\{ \mathbf{x}_i^{k+1} - \mathbf{x}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,2} \right)$$

- However, the above update rules does not fully utilize neighboring info to update x
- We thus make a small adjustment to ensure that neighboring info benefits both x and y

$$\mathbf{z}_i^{k+1} = \mathbf{x}_i^k - \mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k; \boldsymbol{\theta}_{i,1}), \quad \mathbf{g}_i^{k+1} \in \partial r(\mathbf{z}_i^{k+1})$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + s_i^k \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,2} \right)$$

$$\mathbf{x}_i^{k+1} = \mathbf{z}_i^{k+1} - \mathbf{u}_i^k \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,3} \right)$$

Base update rule

Mathematics-inspired update rules

- Does the base update rule converge to the consensus and optimal solution?

No theoretical guarantees. There exists bad m_i, s_i and u_i to deviate the convergence

- What conditions does a good base update rule have to satisfy?
- Can these conditions provide guidance for the development of m_i, s_i and u_i ?

Mathematics-inspired update rules

$$\begin{aligned}
 \mathbf{z}_i^{k+1} &= \mathbf{x}_i^k - \mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k; \boldsymbol{\theta}_{i,1}), \quad \mathbf{g}_i^{k+1} \in \partial r(\mathbf{z}_i^{k+1}) \\
 \mathbf{y}_i^{k+1} &= \mathbf{y}_i^k + s_i^k \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,2} \right) \\
 \mathbf{x}_i^{k+1} &= \mathbf{z}_i^{k+1} - \mathbf{u}_i^k \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,3} \right)
 \end{aligned}$$

If (x_i^k, y_i^k, z_i^k) stays at the optimal solution (x^*, y_i^*, z_i^*) , the next iterate $(x_i^{k+1}, y_i^{k+1}, z_i^{k+1})$ should be fixed

Condition 1 (FIXED POINT): For any $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + r(\mathbf{x})$, $\mathbf{g}_i^* \in \partial r(\mathbf{x}^*)$ and $\mathbf{y}_i^* = -\nabla f_i(\mathbf{x}^*) - \mathbf{g}_i^*$, it holds for any $i \in \mathcal{V}$ that

$$\lim_{k \rightarrow \infty} \mathbf{m}_i^k(\nabla f_i(\mathbf{x}^*), \mathbf{g}_i^*, \mathbf{y}_i^*) = \lim_{k \rightarrow \infty} \mathbf{s}_i^k \left(\{\mathbf{0}_d\}_{j \in \mathcal{N}(i)} \right) = \lim_{k \rightarrow \infty} \mathbf{u}_i^k \left(\{\mathbf{0}_d\}_{j \in \mathcal{N}(i)} \right) = \mathbf{0}_d$$

Mathematics-inspired update rules

$$\mathbf{z}_i^{k+1} = \mathbf{x}_i^k - \mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k; \boldsymbol{\theta}_{i,1}), \quad \mathbf{g}_i^{k+1} \in \partial r(\mathbf{z}_i^{k+1})$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + s_i^k \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,2} \right)$$

$$\mathbf{x}_i^{k+1} = \mathbf{z}_i^{k+1} - \mathbf{u}_i^k \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,3} \right)$$

Any fixed points should converge to the optimal primal and dual solution (x^*, y_i^*, z_i^*)

Condition 2 (GLOBAL CONVERGENCE): For any sequences generated by the base update rules, *there exists* $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + r(\mathbf{x})$, $\mathbf{y}_i^* \in -\nabla f_i(\mathbf{x}^*) - \partial r(\mathbf{x}^*)$ such that

$$\lim_{k \rightarrow \infty} \mathbf{x}_i^k = \lim_{k \rightarrow \infty} \mathbf{z}_i^k = \mathbf{x}^*, \quad \lim_{k \rightarrow \infty} \mathbf{y}_i^k = \mathbf{y}_i^*, \quad \forall i \in \mathcal{V}$$

Mathematics-inspired update rules

Conditions 1 and 2 provide guidance to develop explicit structures for base update rule

Theorem 1 (MATHEMATICS-INSPIRED STRUCTURE): Given $f_i \in \mathcal{F}_L(\mathbb{R}^d)$, $r \in \mathcal{F}(\mathbb{R}^d)$ and base update rules $\{\mathbf{m}_i^k, \mathbf{s}_i^k, \mathbf{u}_i^k\}_{k=0}^\infty$ with $\mathbf{m}_i^k \in \mathcal{D}_C(\mathbb{R}^{3d})$, $\mathbf{s}_i^k, \mathbf{u}_i^k \in \mathcal{D}_C(\mathbb{R}^{|\mathcal{N}(i)|d})$, if Conditions 1 and 2 hold, there exist $\mathbf{P}_i^k, \mathbf{P}_{i,j,1}^k, \mathbf{P}_{i,j,2}^k \in \mathbb{R}^{d \times d}$ and $\mathbf{b}_{i,1}^k, \mathbf{b}_{i,2}^k, \mathbf{b}_{i,3}^k \in \mathbb{R}^d$ satisfying

$$\begin{aligned}\mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k) &= \mathbf{P}_i^k(\nabla f_i(\mathbf{x}_i^k) + \mathbf{g}_i^{k+1} + \mathbf{y}_i^k) + \mathbf{b}_{i,1}^k \\ \mathbf{s}_i^k\left(\left\{\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}\right\}_{j \in \mathcal{N}(i)}\right) &= \sum_{j \in \mathcal{N}(i)} \mathbf{P}_{i,j,1}^k\left(\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}\right) + \mathbf{b}_{i,2}^k \\ \mathbf{u}_i^k\left(\left\{\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}\right\}_{j \in \mathcal{N}(i)}\right) &= \sum_{j \in \mathcal{N}(i)} \mathbf{P}_{i,j,2}^k\left(\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}\right) + \mathbf{b}_{i,3}^k\end{aligned}$$

We name it as **structured update rule**. $P_i^k, P_{i,j,1}^k, P_{i,j,2}^k$ and $b_{i,1}^k, b_{i,2}^k$ and $b_{i,3}^k$ are parameters to learn.

MiLoDo update rules

- To improve computational efficiency, we make the P matrices diagonal.
- Inspired by Theorem 1, which shows that the bias terms $b_{i,1}^k, b_{i,2}^k, b_{i,3}^k$ vanish asymptotically, we eliminate these terms for simplicity.
- Specifically, we set:

$$\mathbf{P}_i^k = \text{Diag}(\mathbf{p}_i^k), \mathbf{P}_{i,j,1}^k = \text{Diag}(\mathbf{p}_{i,j,1}^k), \mathbf{P}_{i,j,2}^k = \text{Diag}(\mathbf{p}_{i,j,2}^k), \mathbf{b}_{i,1}^k = \mathbf{b}_{i,2}^k = \mathbf{b}_{i,3}^k = \mathbf{0}_d,$$

where $\mathbf{p}_i^k, \mathbf{p}_{i,j,1}^k, \mathbf{p}_{i,j,2}^k \in \mathbb{R}_+^d$, and $\mathbf{p}_{i,j,1}^k = \mathbf{p}_{j,i,1}^k, \forall \{i, j\} \in \mathcal{E}$

MiLoDo update rules

Finally, we reach the MiLoDo update rules

$$\boldsymbol{z}_i^{k+1} = \text{prox}_{r, \text{Diag}(\boldsymbol{p}_i^k)} \left(\boldsymbol{x}_i^k - \boxed{\boldsymbol{p}_i^k} \odot (\nabla f_i(\boldsymbol{x}_i^k) + \boldsymbol{y}_i^k) \right)$$

$$\boldsymbol{y}_i^{k+1} = \boldsymbol{y}_i^k + \sum_{j \in \mathcal{N}(i)} \boxed{\boldsymbol{p}_{i,j,1}^k} \odot (\boldsymbol{z}_i^{k+1} - \boldsymbol{z}_j^{k+1})$$

$$\boldsymbol{x}_i^{k+1} = \boldsymbol{z}_i^{k+1} - \sum_{j \in \mathcal{N}(i)} \boxed{\boldsymbol{p}_{i,j,2}^k} \odot (\boldsymbol{z}_i^{k+1} - \boldsymbol{z}_j^{k+1})$$

MiLoDo: Mathematics-inspired Learning-to-optimize framework for Decentralized optimization

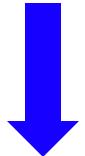
Theoretical guarantees for MiLoDo

Any fixed points of MiLoDo update rules are the primal and dual optimal solutions.

Theorem 2 (EXACT CONVERGENCE): Assume $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ is strongly connected, $\{f_i\}_{i \in \mathcal{V}} \subset \mathcal{F}_L(\mathbb{R}^d)$, $r \in \mathcal{F}(\mathbb{R}^d)$ and there exists $0 < m < M < \infty$ such that $[p_i^k]_l \geq m$, $m \leq [p_{i,j,1}^k]_l \leq M$, $|[p_{i,j,2}^k]_l| \leq M$ for all $k \geq 0$ and $1 \leq l \leq d$. Here $[\mathbf{x}]_l$ denotes the l -th coordinate of the vector \mathbf{x} . If a sequence generated by (22)-(24) with initialization $\mathbf{y}_i^0 = \mathbf{0}_d$ converges to $\{\mathbf{x}_i^*, \mathbf{y}_i^*, \mathbf{z}_i^*\}_{i=1}^n$, then this limit must be the primal and dual optimal solutions to problem (2). In other words, there exists $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) + r(\mathbf{x})$ such that $\mathbf{x}_i^* = \mathbf{z}_i^* = \mathbf{x}^*$ and $\mathbf{y}_i^* \in -\nabla f_i(\mathbf{x}^*) - \partial r(\mathbf{x}^*)$ holds.

① Generic L2O

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^k + g_i^k \left(\{ \mathbf{x}_j^k, \nabla f_j(\mathbf{x}_j^k) \}_{j \in \mathcal{N}(i) \cup \{i\}}; \boldsymbol{\theta}_i \right)$$

 Inspired by primal-dual algorithm

② Base rule

$$\begin{aligned}\mathbf{z}_i^{k+1} &= \mathbf{x}_i^k - \mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k; \boldsymbol{\theta}_{i,1}), \\ \mathbf{y}_i^{k+1} &= \mathbf{y}_i^k + s_i^k \left(\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,2} \right) \\ \mathbf{x}_i^{k+1} &= \mathbf{z}_i^{k+1} - \mathbf{u}_i^k \left(\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \}_{j \in \mathcal{N}(i)}; \boldsymbol{\theta}_{i,3} \right)\end{aligned}$$

 Fixed-point cond.

 Global converge

④ MiLoDo

$$\begin{aligned}\mathbf{z}_i^{k+1} &= \text{prox}_{r, \text{Diag}(\mathbf{p}_i^k)} (\mathbf{x}_i^k - \mathbf{p}_i^k \odot (\nabla f_i(\mathbf{x}_i^k) + \mathbf{y}_i^k)) \\ \mathbf{y}_i^{k+1} &= \mathbf{y}_i^k + \sum_{j \in \mathcal{N}(i)} \mathbf{p}_{i,j,1}^k \odot (\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}) \\ \mathbf{x}_i^{k+1} &= \mathbf{z}_i^{k+1} - \sum_{j \in \mathcal{N}(i)} \mathbf{p}_{i,j,2}^k \odot (\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1})\end{aligned}$$

 Simplification

③ Math-inspired rule

$$\begin{aligned}\mathbf{m}_i^k(\nabla f_i(\mathbf{x}_i^k), \mathbf{g}_i^{k+1}, \mathbf{y}_i^k) &= \mathbf{P}_i^k(\nabla f_i(\mathbf{x}_i^k) + \mathbf{g}_i^{k+1} + \mathbf{y}_i^k) + \mathbf{b}_{i,1}^k \\ \mathbf{s}_i^k \left(\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \}_{j \in \mathcal{N}(i)} \right) &= \sum_{j \in \mathcal{N}(i)} \mathbf{P}_{i,j,1}^k (\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}) + \mathbf{b}_{i,2}^k \\ \mathbf{u}_i^k \left(\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \}_{j \in \mathcal{N}(i)} \right) &= \sum_{j \in \mathcal{N}(i)} \mathbf{P}_{i,j,2}^k (\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}) + \mathbf{b}_{i,3}^k\end{aligned}$$

Introduce math structures to L2O framework to narrow down search space and promote generalization

Generic L2O vs MiLoDo

Generic L2O

$$x_i^{k+1} = x_i^k + g_i^k \left(\{x_j^k, \nabla f_j(x_j^k)\}_{j \in \mathcal{N}(i) \cup \{i\}}; \theta_i \right)$$

MiLoDo

$$z_i^{k+1} = \text{prox}_{r, \text{Diag}(\mathbf{p}_i^k)} \left(\mathbf{x}_i^k - \mathbf{p}_i^k \odot (\nabla f_i(\mathbf{x}_i^k) + \mathbf{y}_i^k) \right),$$

$$\mathbf{y}_i^{k+1} = \mathbf{y}_i^k + \sum_{j \in \mathcal{N}(i)} \mathbf{p}_{i,j,1}^k \odot (z_i^{k+1} - z_j^{k+1}),$$

$$x_i^{k+1} = z_i^{k+1} - \sum_{j \in \mathcal{N}(i)} \mathbf{p}_{i,j,2}^k \odot (z_i^{k+1} - z_j^{k+1}).$$

- **Large** search space; **difficult** to train
- Consensus and optimality are **NOT** ensured
- **Weak** out-of-distribution generalization

- **Compact** search space; **easy** to train
- Consensus and optimality are **ensured**
- **Strong** out-of-distribution generalization

PART 04

Parameterization and Training

$$\boldsymbol{z}_i^{k+1} = \text{prox}_{r, \text{Diag}(\boldsymbol{p}_i^k)} \left(\boldsymbol{x}_i^k - \boldsymbol{p}_i^k \odot (\nabla f_i(\boldsymbol{x}_i^k) + \boldsymbol{y}_i^k) \right),$$

$$\boldsymbol{y}_i^{k+1} = \boldsymbol{y}_i^k + \sum_{j \in \mathcal{N}(i)} \boldsymbol{p}_{i,j,1}^k \odot (\boldsymbol{z}_i^{k+1} - \boldsymbol{z}_j^{k+1}),$$

$$\boldsymbol{x}_i^{k+1} = \boldsymbol{z}_i^{k+1} - \sum_{j \in \mathcal{N}(i)} \boldsymbol{p}_{i,j,2}^k \odot (\boldsymbol{z}_i^{k+1} - \boldsymbol{z}_j^{k+1}).$$

How to learn $p_i^k, p_{i,j,1}^k, p_{i,j,2}^k$ with deep neural network?

$$\begin{aligned}\mathbf{p}_i^k, \mathbf{h}_{M,i}^{k+1} &= \phi_{M,i}(\nabla f(\mathbf{x}_i^k), \mathbf{y}_i^k, \mathbf{h}_{M,i}^k; \boldsymbol{\theta}_{M,i}), \\ \left\{ \tilde{\mathbf{p}}_{i,j,1}^k \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{S,i}^{k+1} &= \phi_{S,i}\left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{S,i}^k; \boldsymbol{\theta}_{S,i} \right), \\ \left\{ \mathbf{p}_{i,j,2}^k \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{U,i}^{k+1} &= \phi_{U,i}\left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{U,i}^k; \boldsymbol{\theta}_{U,i} \right),\end{aligned}$$

- We parameterize $\mathbf{p}_i^k, \mathbf{p}_{i,j,1}^k, \mathbf{p}_{i,j,2}^k$ using three local LSTM networks $\phi_{M,i}, \phi_{S,i}, \phi_{U,i}$.
- LSTM hidden states $\mathbf{h}_{M,i}^k, \mathbf{h}_{S,i}^k, \mathbf{h}_{U,i}^k$ are randomly initialized, and parameters $\boldsymbol{\theta}_{M,i}, \boldsymbol{\theta}_{S,i}, \boldsymbol{\theta}_{U,i}$ are learnable.

$$\nabla f_i(\mathbf{x}_i^k), \mathbf{y}_i^k, \{\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}\}_{j \in \mathcal{N}(i)}$$

Optimization

$$\begin{aligned}\mathbf{z}_i^{k+1} &= \text{prox}_{r, \text{Diag}(\mathbf{p}_i^k)} \left(\mathbf{x}_i^k - \mathbf{p}_i^k \odot (\nabla f_i(\mathbf{x}_i^k) + \mathbf{y}_i^k) \right), \\ \mathbf{y}_i^{k+1} &= \mathbf{y}_i^k + \sum_{j \in \mathcal{N}(i)} \mathbf{p}_{i,j,1}^k \odot (\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}), \\ \mathbf{x}_i^{k+1} &= \mathbf{z}_i^{k+1} - \sum_{j \in \mathcal{N}(i)} \mathbf{p}_{i,j,2}^k \odot (\mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1}).\end{aligned}$$

Learning

$$\mathbf{p}_i^k, \mathbf{h}_{M,i}^{k+1} = \phi_{M,i}(\nabla f(\mathbf{x}_i^k), \mathbf{y}_i^k, \mathbf{h}_{M,i}^k; \boldsymbol{\theta}_{M,i}),$$

$$\left\{ \tilde{\mathbf{p}}_{i,j,1}^k \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{S,i}^{k+1} = \phi_{S,i} \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{S,i}^k; \boldsymbol{\theta}_{S,i} \right),$$

$$\left\{ \mathbf{p}_{i,j,2}^k \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{U,i}^{k+1} = \phi_{U,i} \left(\left\{ \mathbf{z}_i^{k+1} - \mathbf{z}_j^{k+1} \right\}_{j \in \mathcal{N}(i)}, \mathbf{h}_{U,i}^k; \boldsymbol{\theta}_{U,i} \right),$$

(solve the target optimization problem)

(learn the parameters $p_i, p_{i,j,1}, p_{i,j,2}$)

$$\mathbf{p}_i^k, \{\tilde{\mathbf{p}}_{i,j,1}^k\}_{j \in \mathcal{N}(i)}, \{\mathbf{p}_{i,j,2}^k\}_{j \in \mathcal{N}(i)}$$

- The parameters to train is $\Theta = \{\theta_{M,i}, \theta_{S,i}, \theta_{U,i}\}_{i=1}^n$
- We evaluate and refine the performance of the optimizer over the initial K steps on a batch of training problems

$$\min_{\Theta} \mathcal{L}_K(\Theta, \mathcal{F}_B) := \frac{1}{|\mathcal{F}_B|} \sum_{f \in \mathcal{F}_B} \left[\frac{1}{K} \sum_{k=1}^K f(\bar{x}^k) \right].$$

where $\bar{x}^k = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i^k$ is the global average of all local variables

- We typically set $K = 100$ during the training phrase

Pre-training + fine-tuning

Inspired by LLM fine-tuning techniques, we propose the following innovative training strategy:

1. Special Initialization:

- Initialize MiLoDo parameters with traditional decentralized optimization algorithms instead of random values. (e.g., EXTRA, Exact Diffusion)

2. Mixed Dataset Pre-Training:

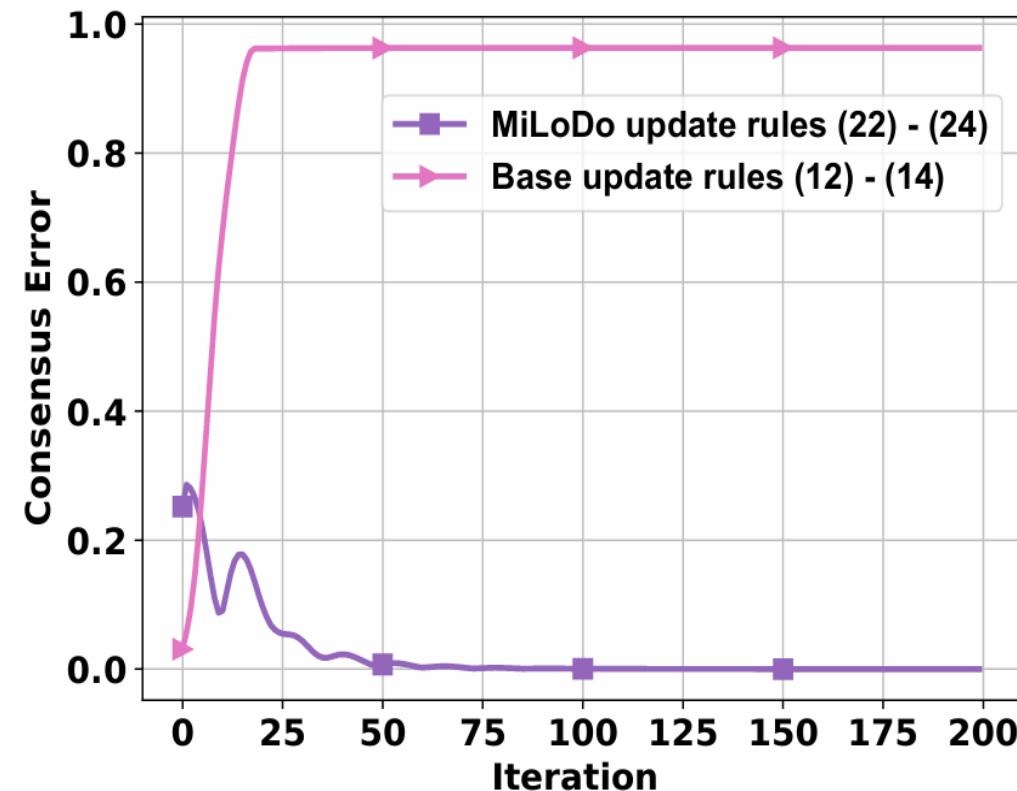
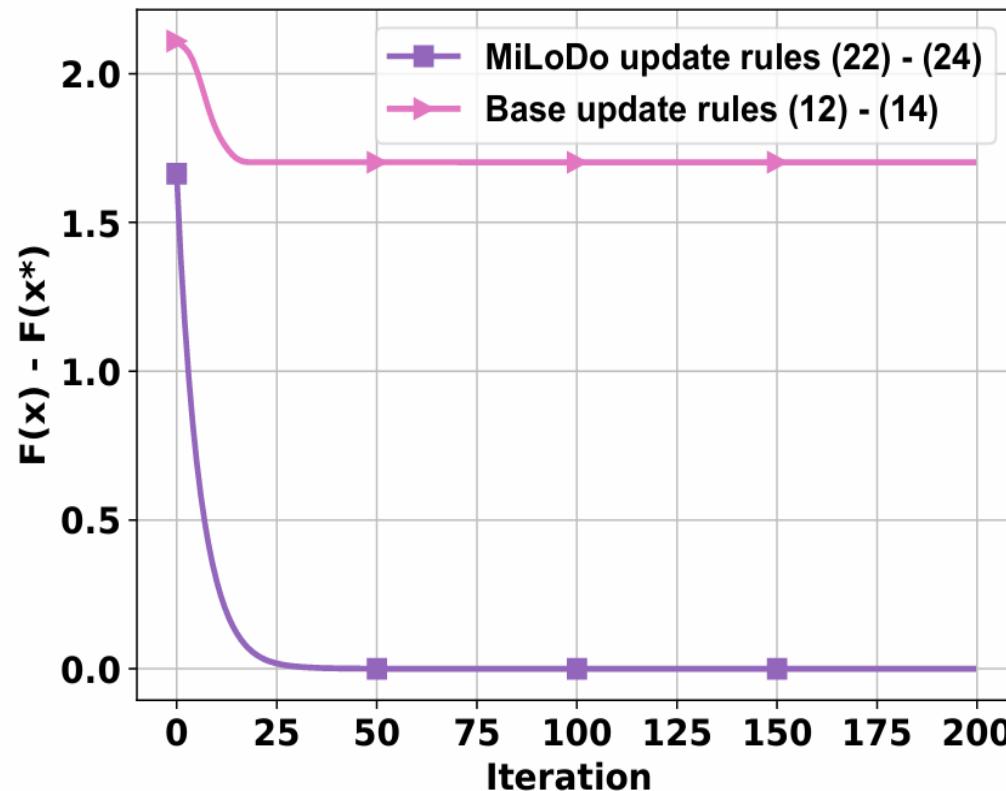
- **Pre-train** MiLoDo on a mixed dataset of problems with various dimensions to learn general optimization strategies, enhancing the generalization ability.
- Then, **fine-tune** on specific tasks to improve adaptability and performance.

PART 05

Experimental results

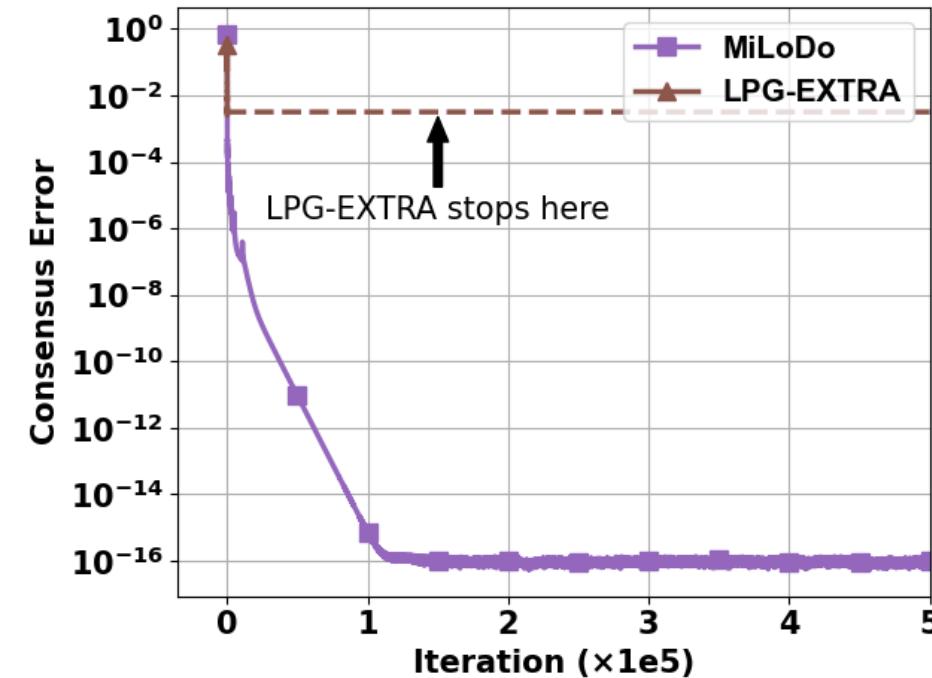
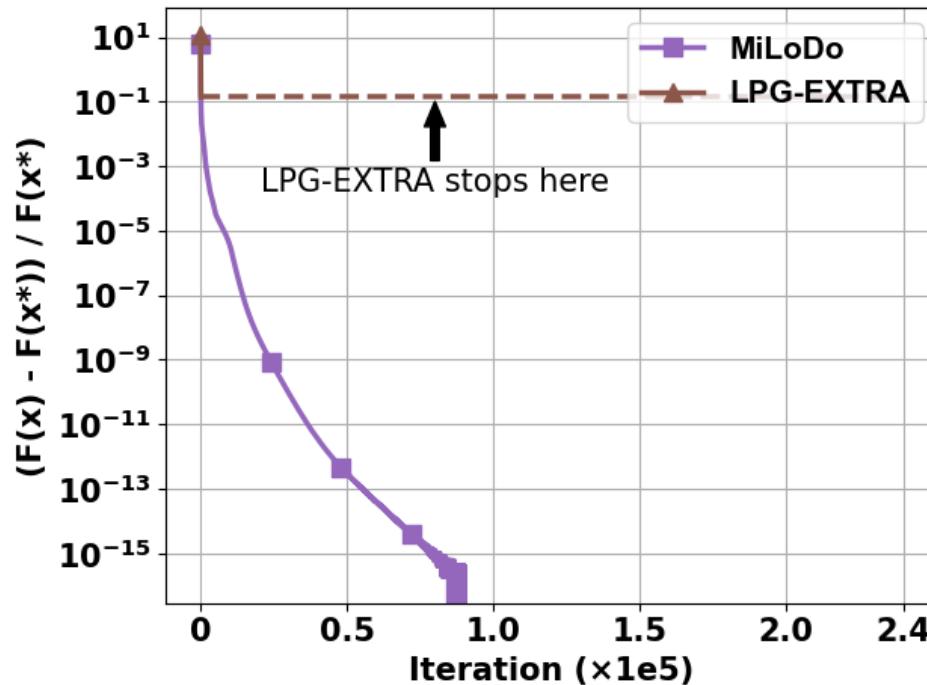
MiLoDo vs Generic L2O

LASSO problem (10,10,5,0)



MiLoDo vs algorithm unrolling

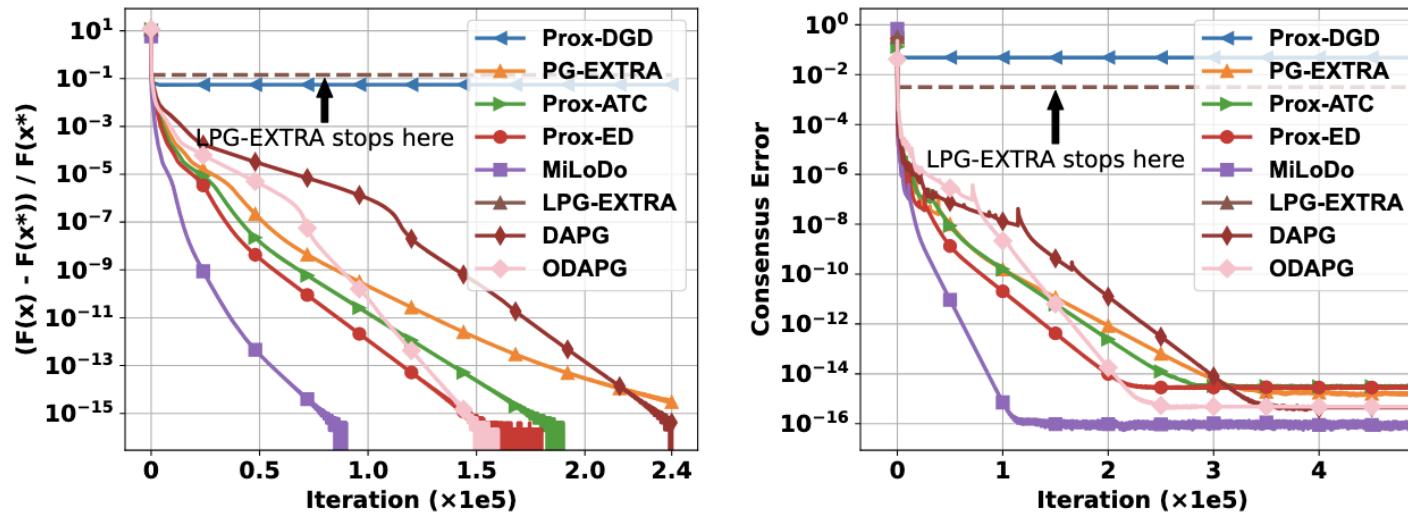
LPG-EXTRA: PG-EXTRA with neural network unrolling; can only run 100 iterations



H. Wang, Y. Shen, Z. Wang, D. Li, J. Zhang, K. B. Letaief, and J. Lu. Decentralized statistical inference with unrolled graph neural networks, IEEE CDC, 2021

Significant acceleration

In tests on LASSO and logistic optimization problems, MiLoDo achieves approximately **1.8 to 2.5 times speedup** compared to other hand-crafted optimization algorithms.

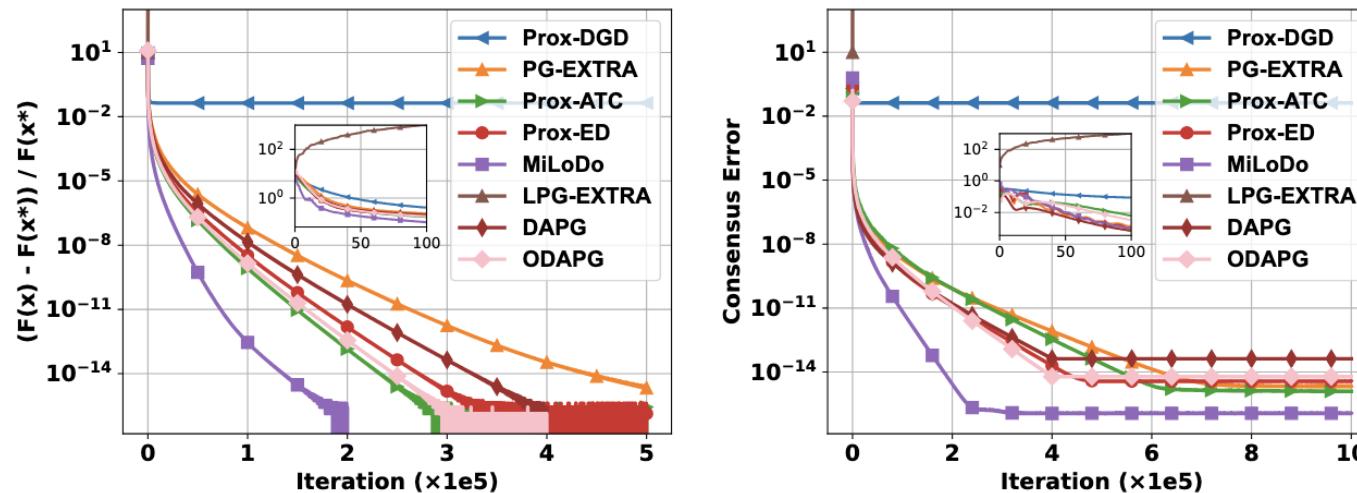


The first learned optimizer that can beat man-crafted optimizer in decentralized optimization!

Figure 2: MiLoDo-optimizer trained on synthetic LASSO($10, 300, 10, 0.1$) and tested on unseen LASSO($10, 300, 10, 0.1$) instances.

Generalization across different dimensions

The MiLoDo optimizer, trained on low-dimensional problems, successfully **generalizes to high-dimensional problems** during inference, demonstrating exceptional scalability.

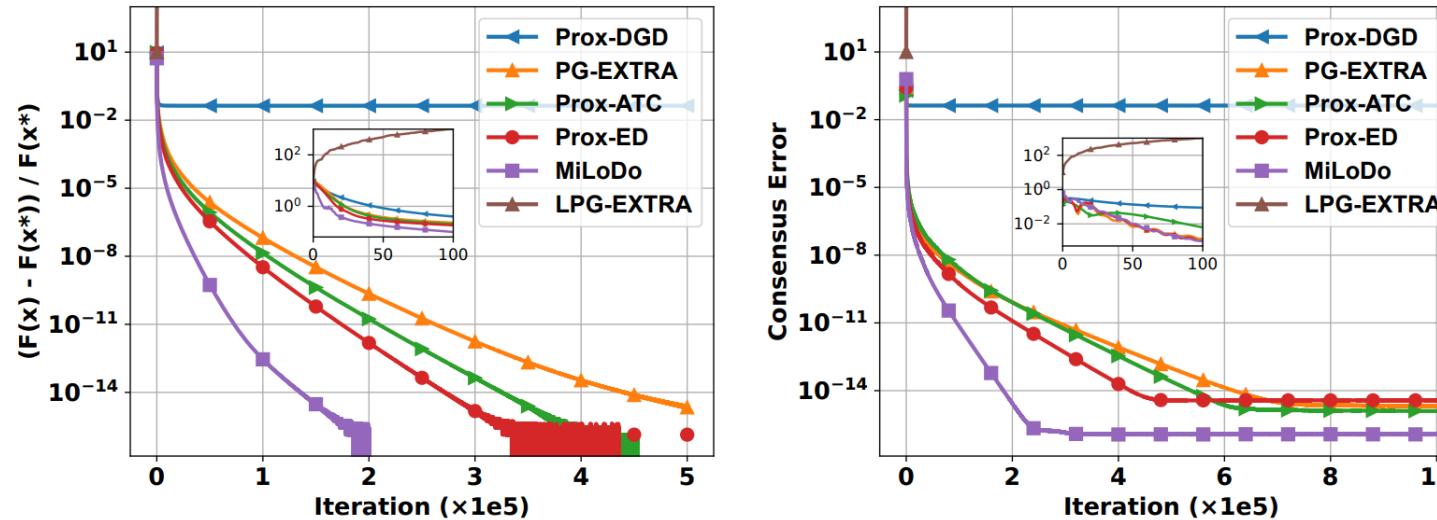


Generalization from $d = 300$ to $d = 30000$

Figure 3: MiLoDo-optimizer trained on synthetic LASSO($10, 300, 10, 0.1$) and tested on synthetic LASSO($10, 30000, 1000, 0.1$).

Generalization from short training to long inference

The MiLoDo optimizer, limited to **100 iterations** during training, can extend to over **100,000 iterations** during inference. This exceptional scalability significantly enhances its practical value.

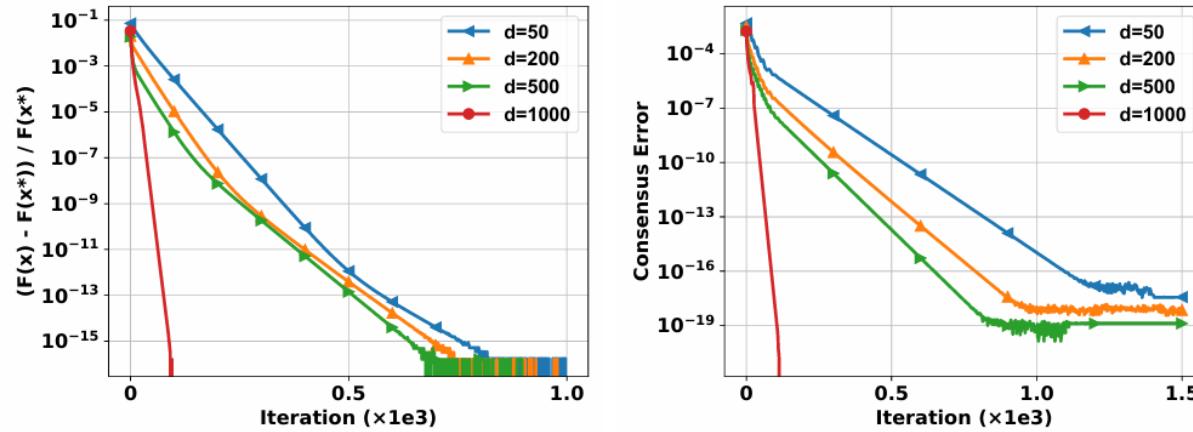


Trained to operate for 100 iterations, but performs well for 10000 iterations when solving unseen problems !

Figure 3: MiLoDo-optimizer trained on synthetic LASSO($10, 300, 10, 0.1$) and tested on synthetic LASSO($10, 30000, 1000, 0.1$).

Generalization across different tasks

MiLoDo generalizes across task types. Trained on **LASSO problems**, it accurately converges on **logistic regression tasks** of various dimensions.

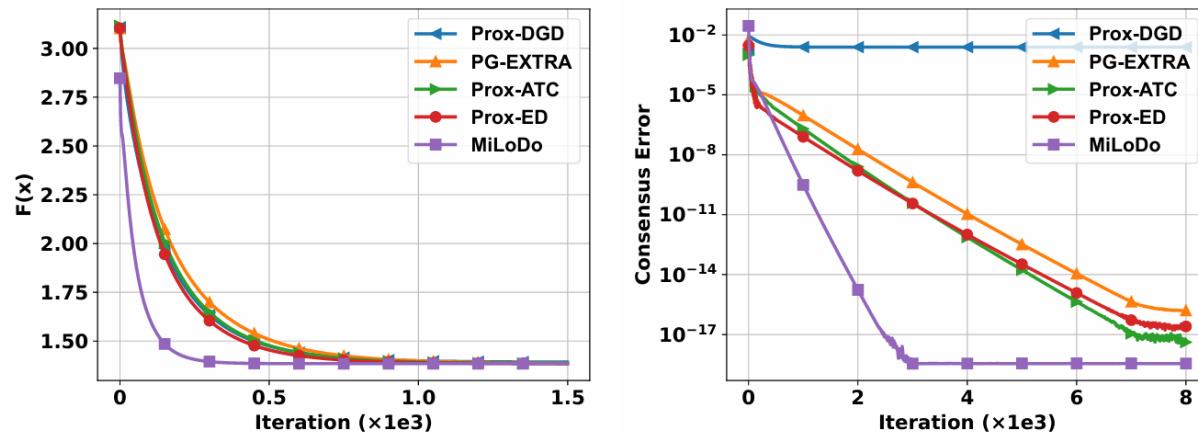


**Trained to LASSO but
generalizes to logistic
regression!**

Figure 5: MiLoDo-optimizer trained on meta training set and tested on Logistic($10, d, 100, 0.1$) with $d \in \{50, 200, 500, 1000\}$.

Out-of-Distribution Testing

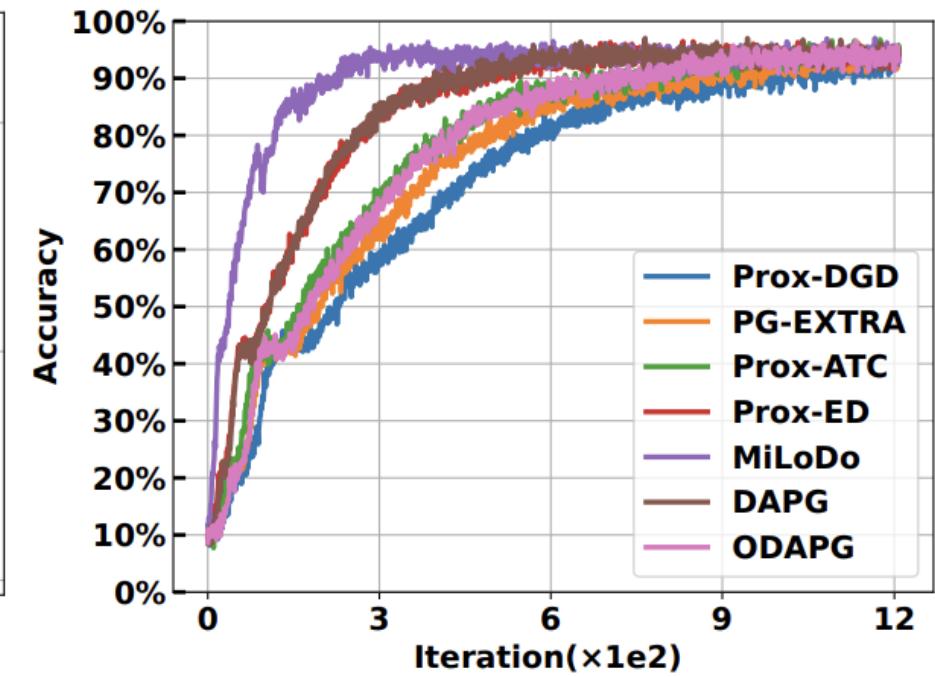
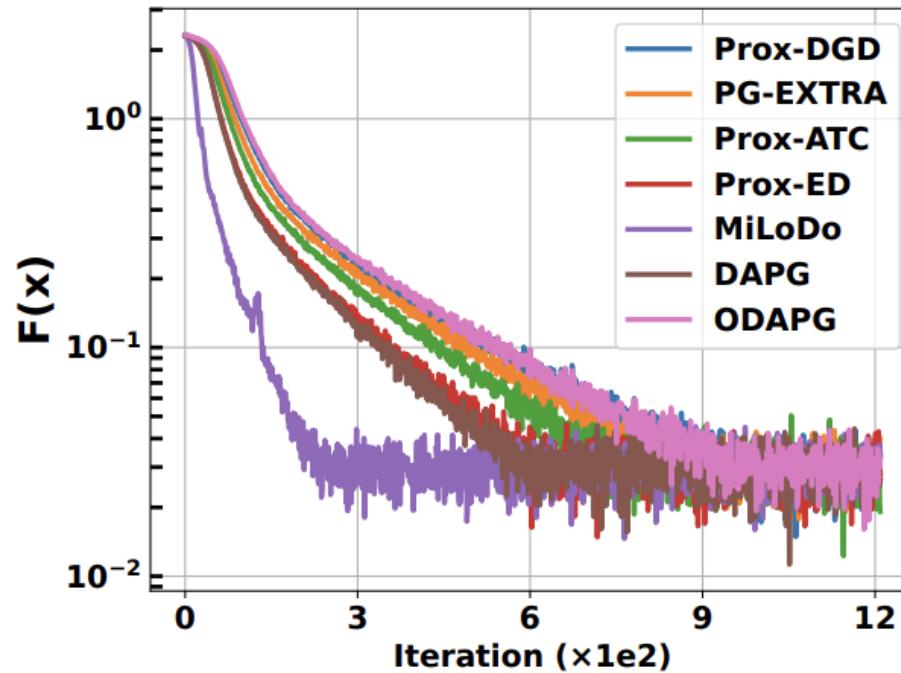
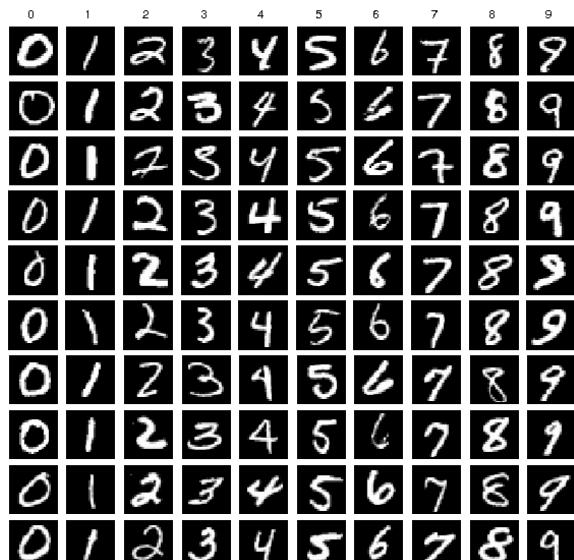
Trained on **synthetic datasets**, MiLoDo shows clear advantages when tested on the challenging **real-world BSDS500 dataset**.



**Trained on synthetic datasets
but generalizes to real dataset !**

Figure 4: MiLoDo-optimizer trained on meta training set and tested on LASSO(10, 200, 10, 0.05) with real dataset BSDS500[29].

Multi-layer neural network (MLP)



ResNet on Cifar-10

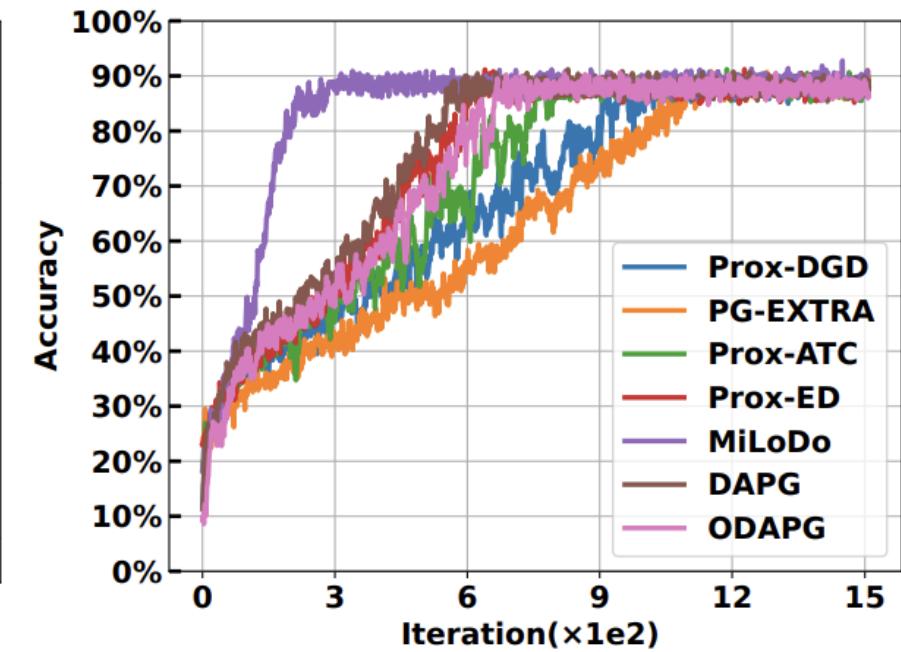
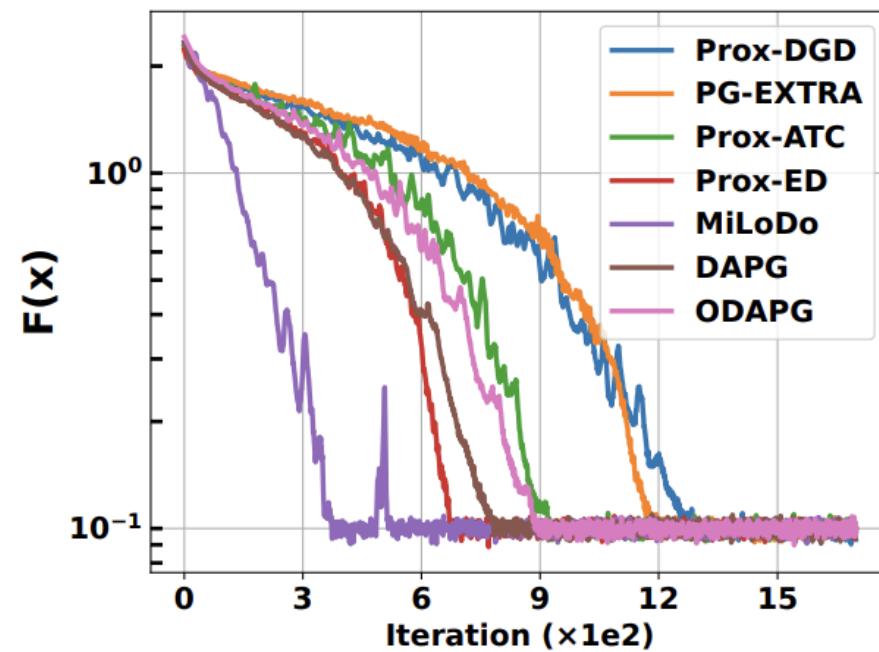
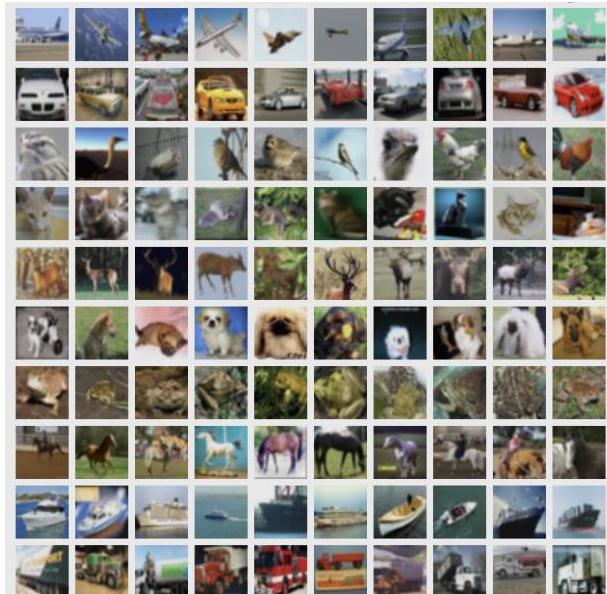


Table 1: Runtime Comparison on LASSO(10, 30000, 1000, 0.1).

	Stopping condition: $\text{Gap} < 10^{-7}$		Stopping condition: $\text{Gap} < 10^{-15}$		
	Time/Iters	Iters	Total Time	Iters	Total Time
MiLoDo	5.91 ms	2.45e+04	144.80 s	1.62e+05	957.42 s
Prox-ED	4.99 ms	6.22e+04	310.38 s	3.07e+05	1531.93 s

- Increases the per-iteration computational cost by 18.4%
- Achieves a convergence speedup by 1.6x ~ 2.1x

Saves the total time significantly

Summary

Can we develop an effective L2O algorithm for decentralized optimization that is

- **Memory efficient? Yes! MiLoDo avoids unrolling the algorithms**
- **Easy to train? Yes! MiLoDo introduces math structures and narrow down search space**
- **Enforcing consensus constraint? Yes! MiLoDo has theoretically guarantees**
- **Strong generalization? Yes! The math principles are critical to ensure strong generalization**

Take-home message



Mathematical structures are fundamental in Learning-to-Optimize!

Future work

- How to generalize across different network sizes and topologies?
- How to develop L2O framework for decentralized **stochastic** optimization which is crucial in deep learning tasks?
- We find the potential of pretraining + finetuning in this work. Can we develop foundational large models for decentralized optimization?

Recent work on optimization for Large language models (LLMs)



Subspace Optimization for Large Language Models with Convergence Guarantees

[https://arxiv.org/pdf/2410.11289](https://arxiv.org/pdf/2410.11289.pdf)

Enhancing Zeroth-Order Fine-Tuning for Language Models with Low-rank Structures

[https://arxiv.org/pdf/2410.07698](https://arxiv.org/pdf/2410.07698.pdf)

Significantly save memory cost to pre-train or fine-tune LLMs



Thank you!

Paper address: <https://arxiv.org/pdf/2410.01700>

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