

# **FedCBO: Reaching Group Consensus in Clustered Federated Learning and Robustness to Backdoor Adversarial Attacks**

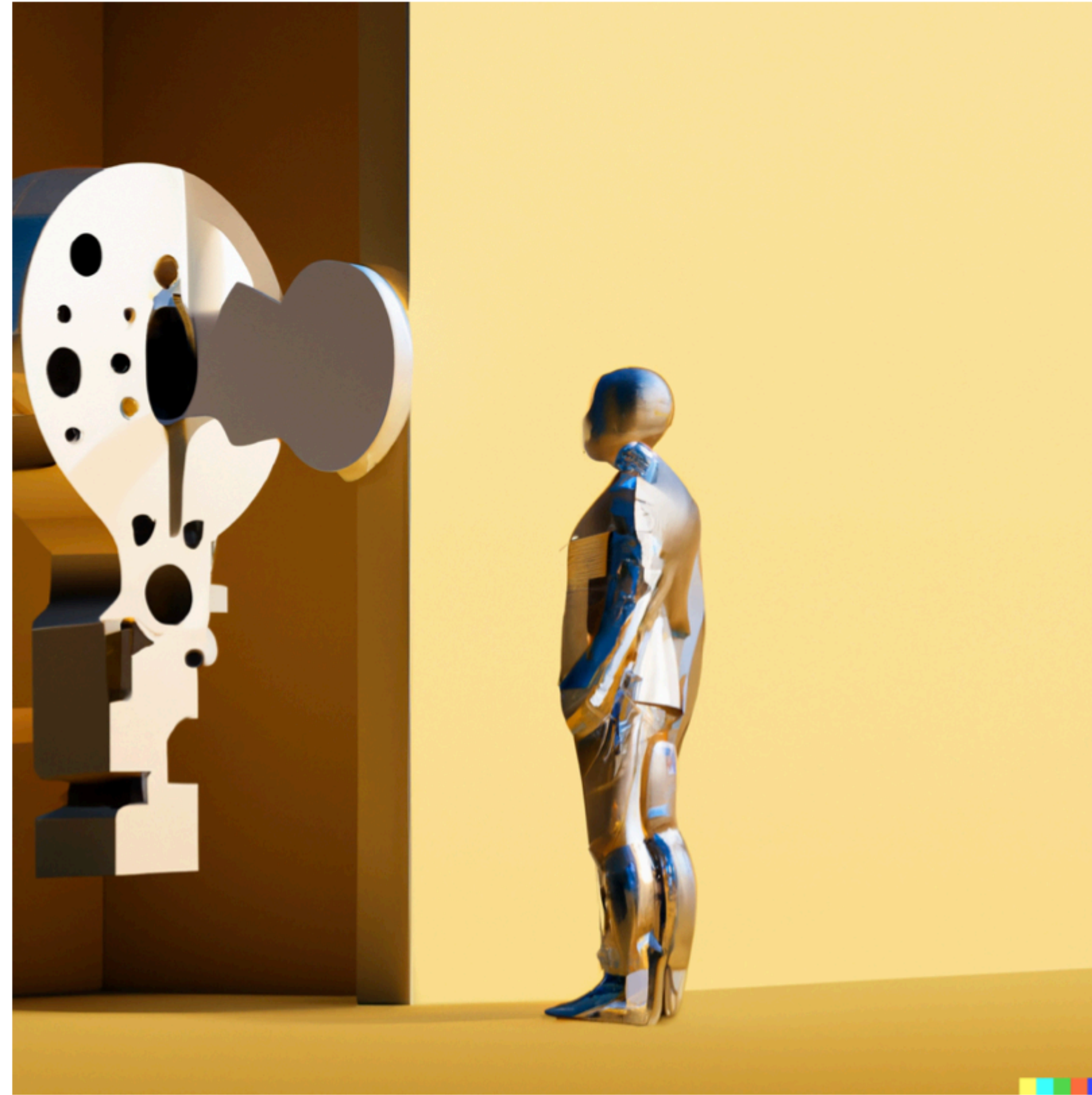
Nicolás García Trillos  
University of Wisconsin-Madison

ICERM  
May 2024

# Based on

- J.A. Carrillo, **NGT**, S. Li, Y. Zhu "*FedCBO: Reaching Group Consensus in Clustered Federated Learning through Consensus-based Optimization.*" <https://arxiv.org/abs/2305.02894>
- **NGT**, S. Li, K. Riedl, Y. Zhu “*Bilevel Consensus-based Optimization and applications to backdoor attack robustness in Clustered Federated Learning.*” (In preparation).

# The success stories of AI



**From Dall-e:** generate an image representing the success stories of AI.

# However...

UNIVERSITÄT BONN

UNIVERSITY STUDYING RES

02. July 2019

## Building trust in artificial intelligence

Interdisciplinary team from IT, philosophy and law defines priorities for the certification of AI



**Safety** (e.g., systems should be robust to perturbations of data)

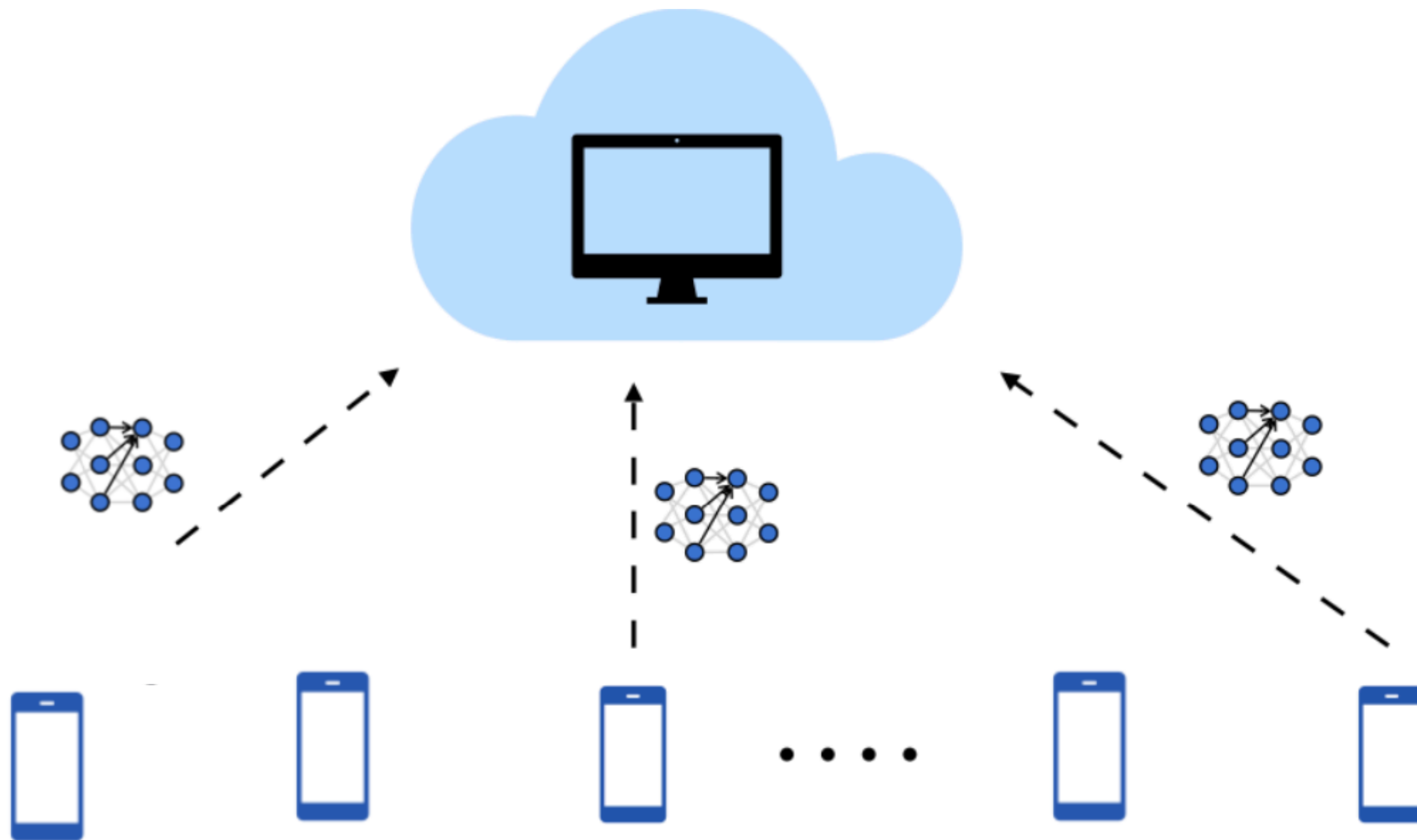
**Security** (e.g., systems should be robust to adversarial agents)

**Trust** (e.g., privacy, fairness, watermarks for generated images)

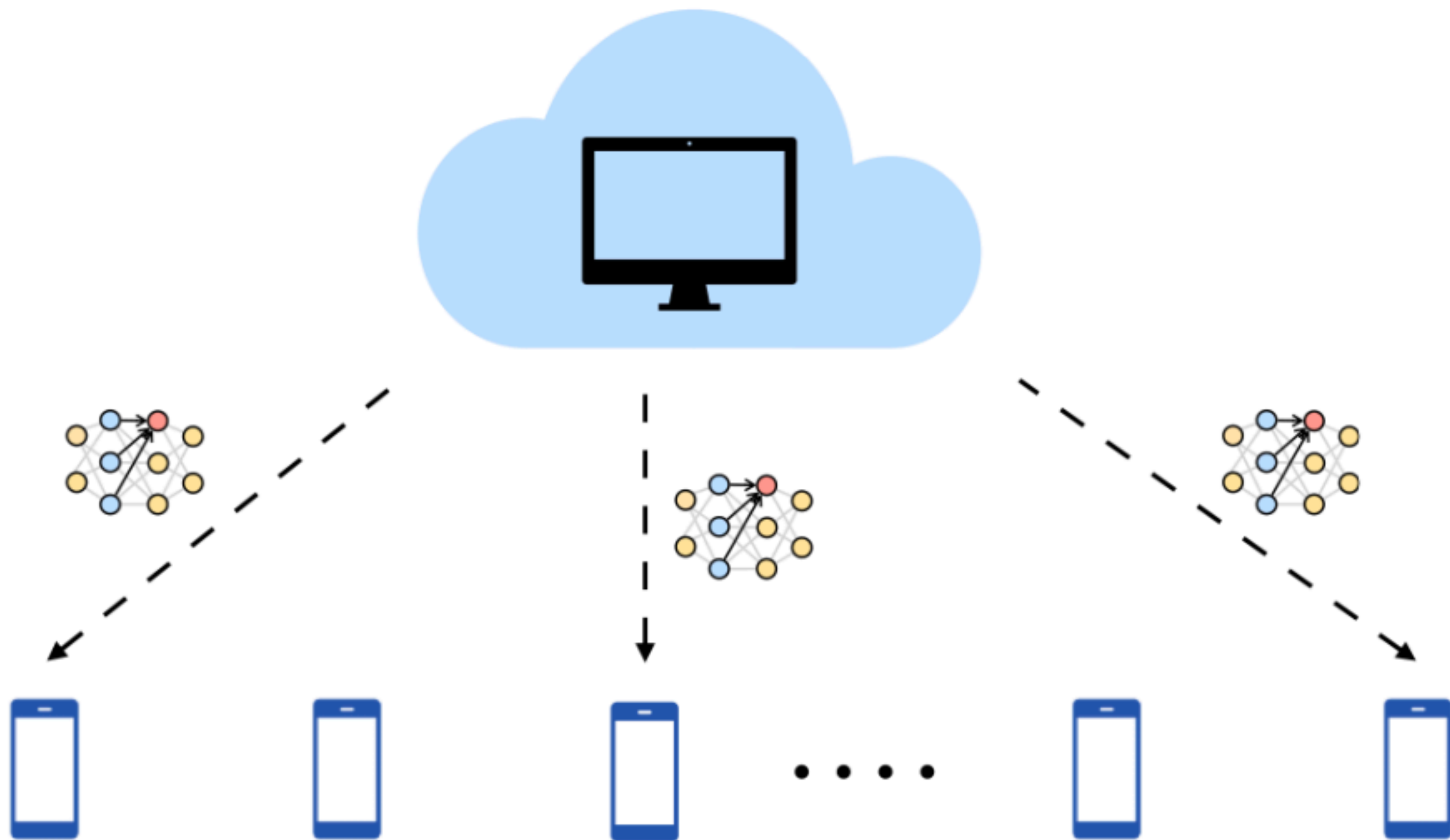
# Outline

1. Introduction: federated learning and clustered federated learning
2. Federated learning through consensus based optimization (CBO)
3. Backdoor attacks.

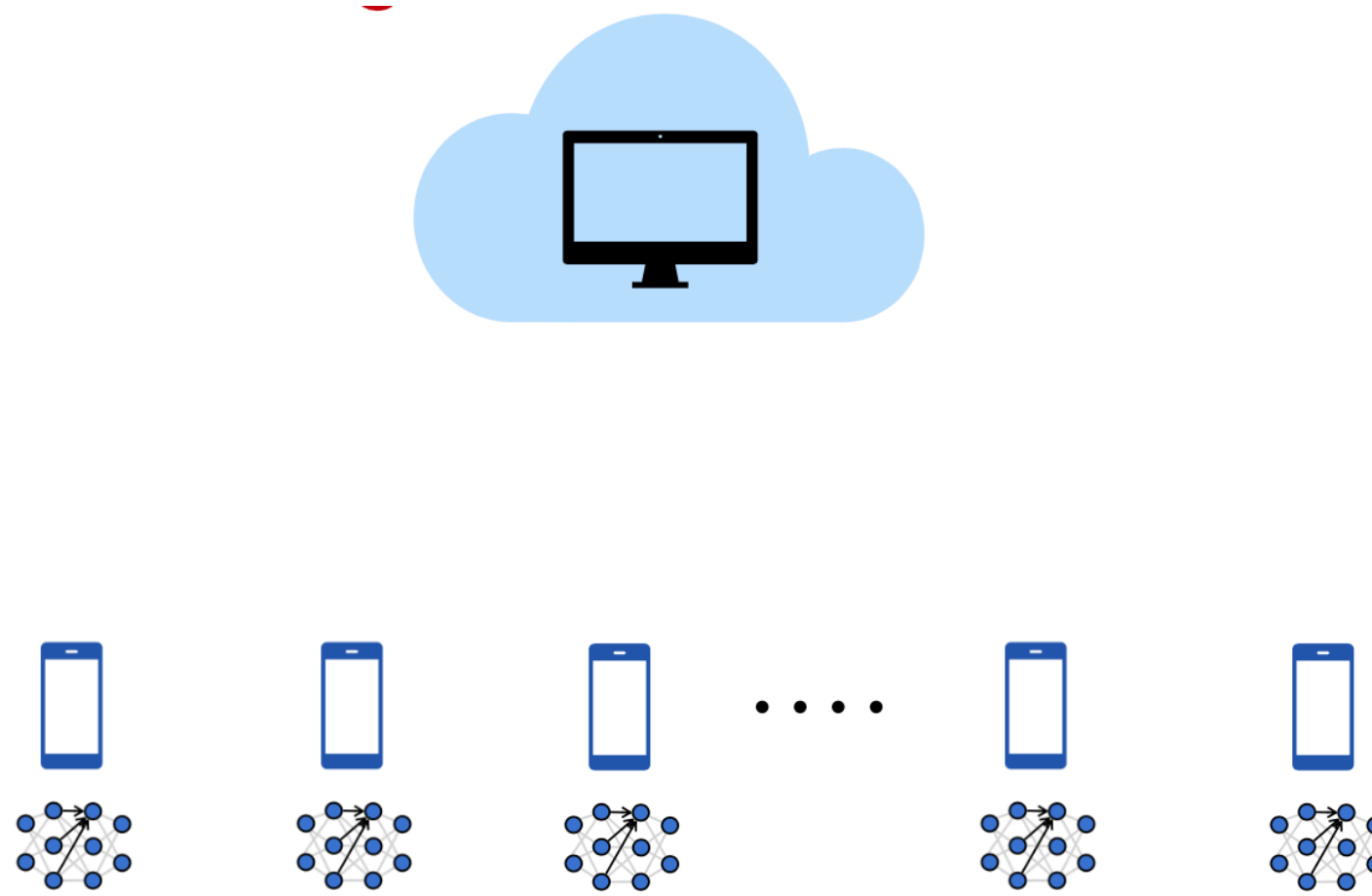
# Federated Learning



# Federated Learning



# Federated Average

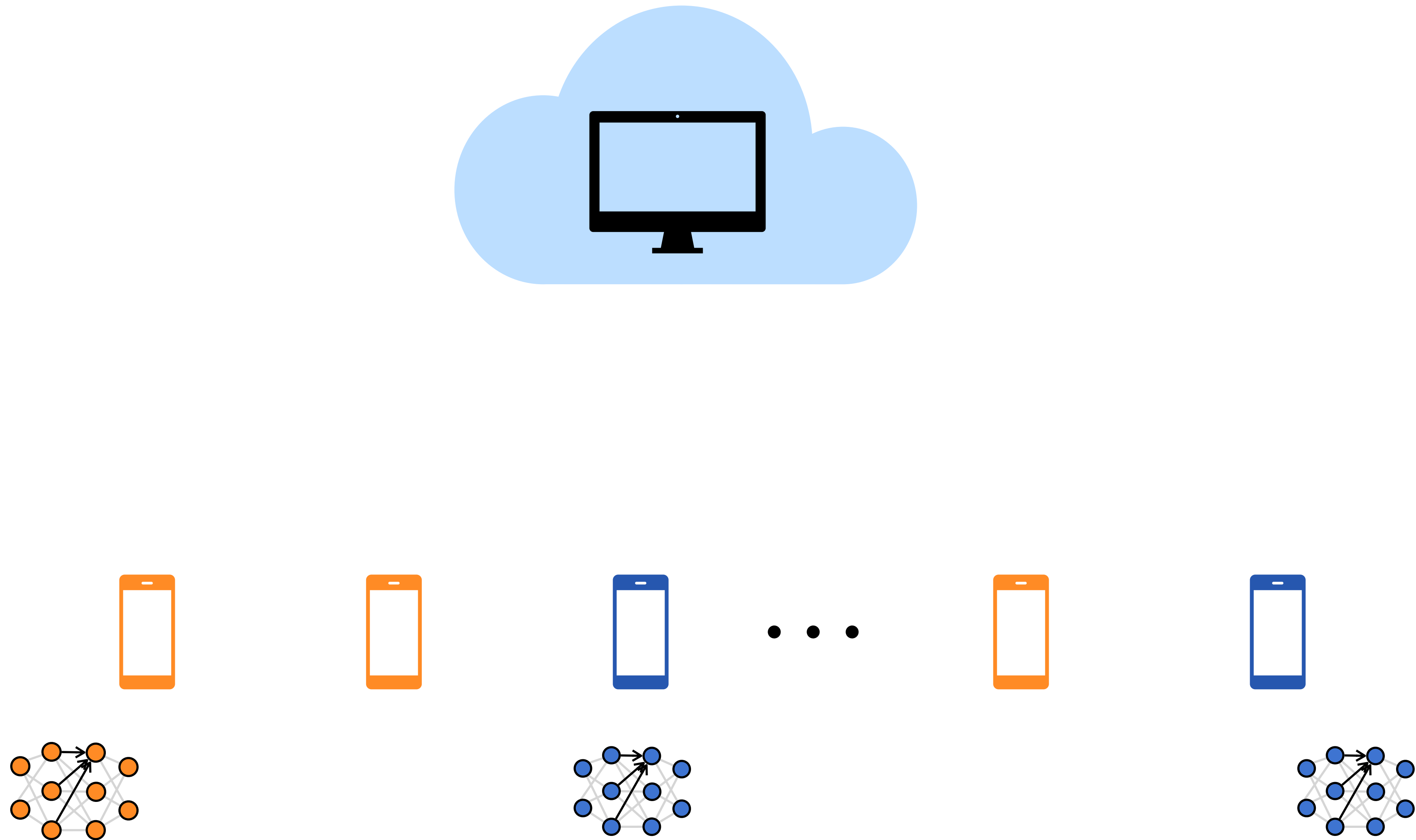


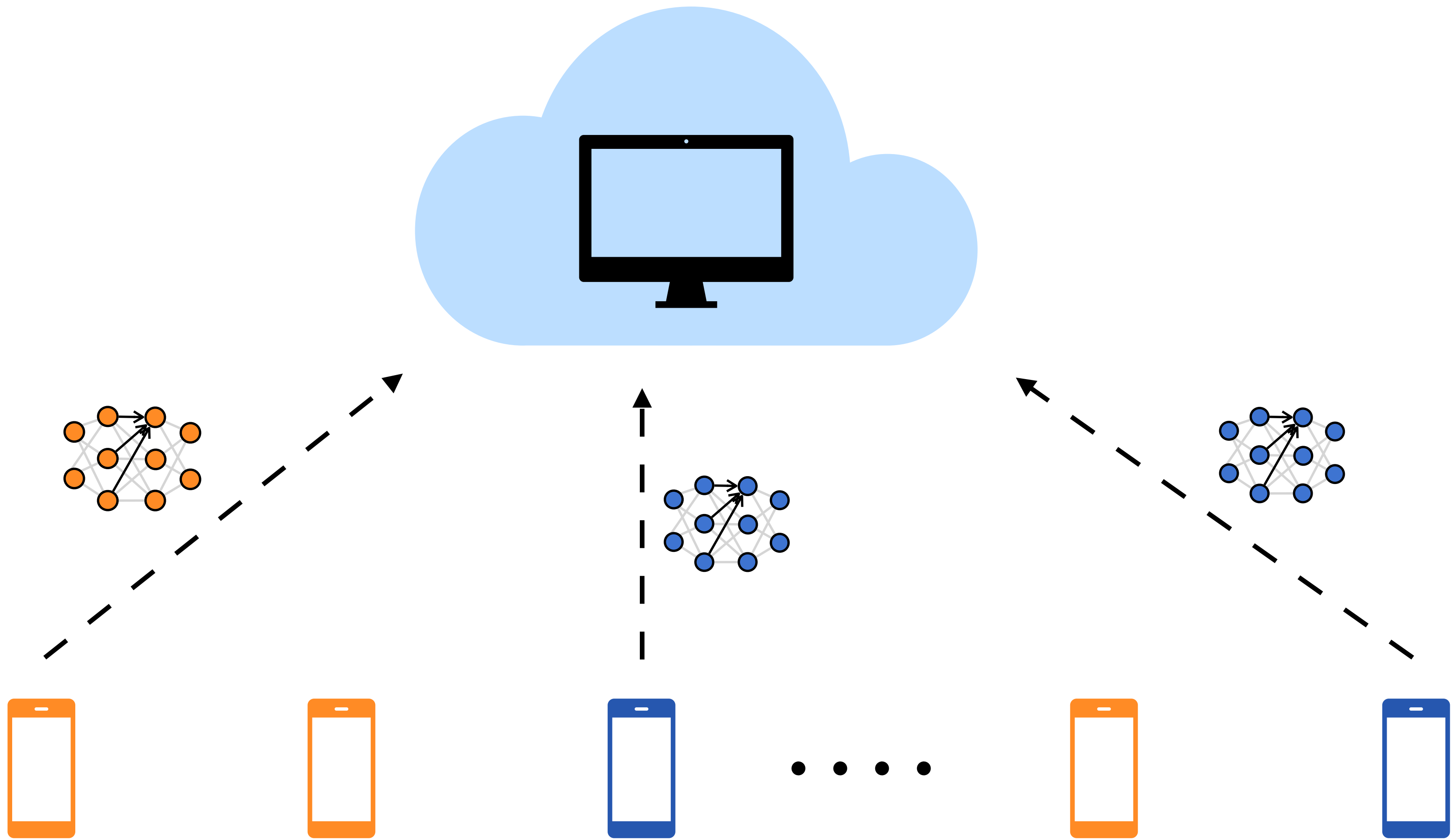
**FedAvg [McMahan et al 16']:**

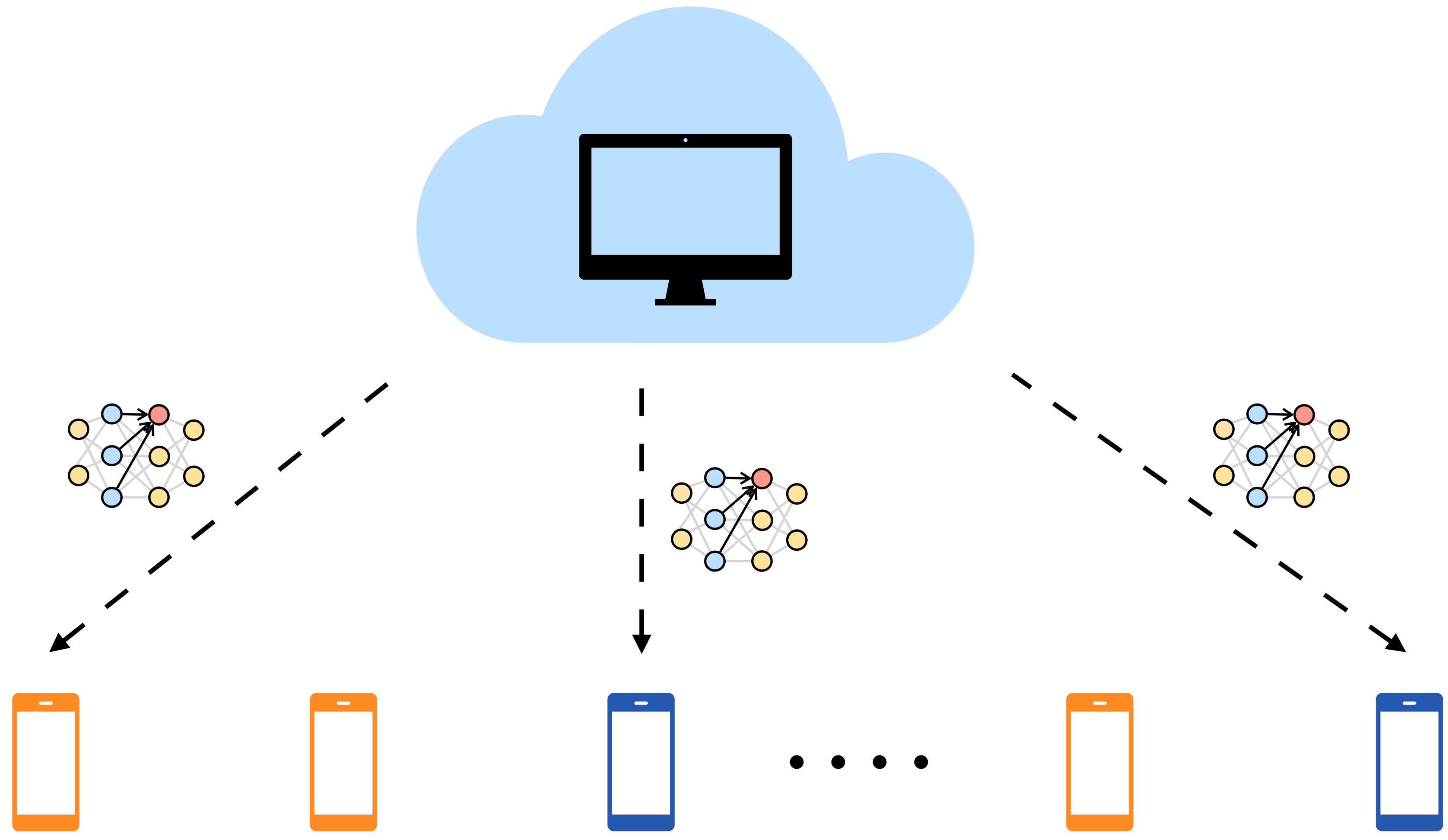
$$\begin{cases} d\theta_t^i = -\nabla L_i(\bar{\theta}_t) + \text{Noise} , & i = 1, \dots, N. \\ \bar{\theta}_t = \frac{1}{N} \sum_{i=1}^N \theta_t^i \end{cases}$$



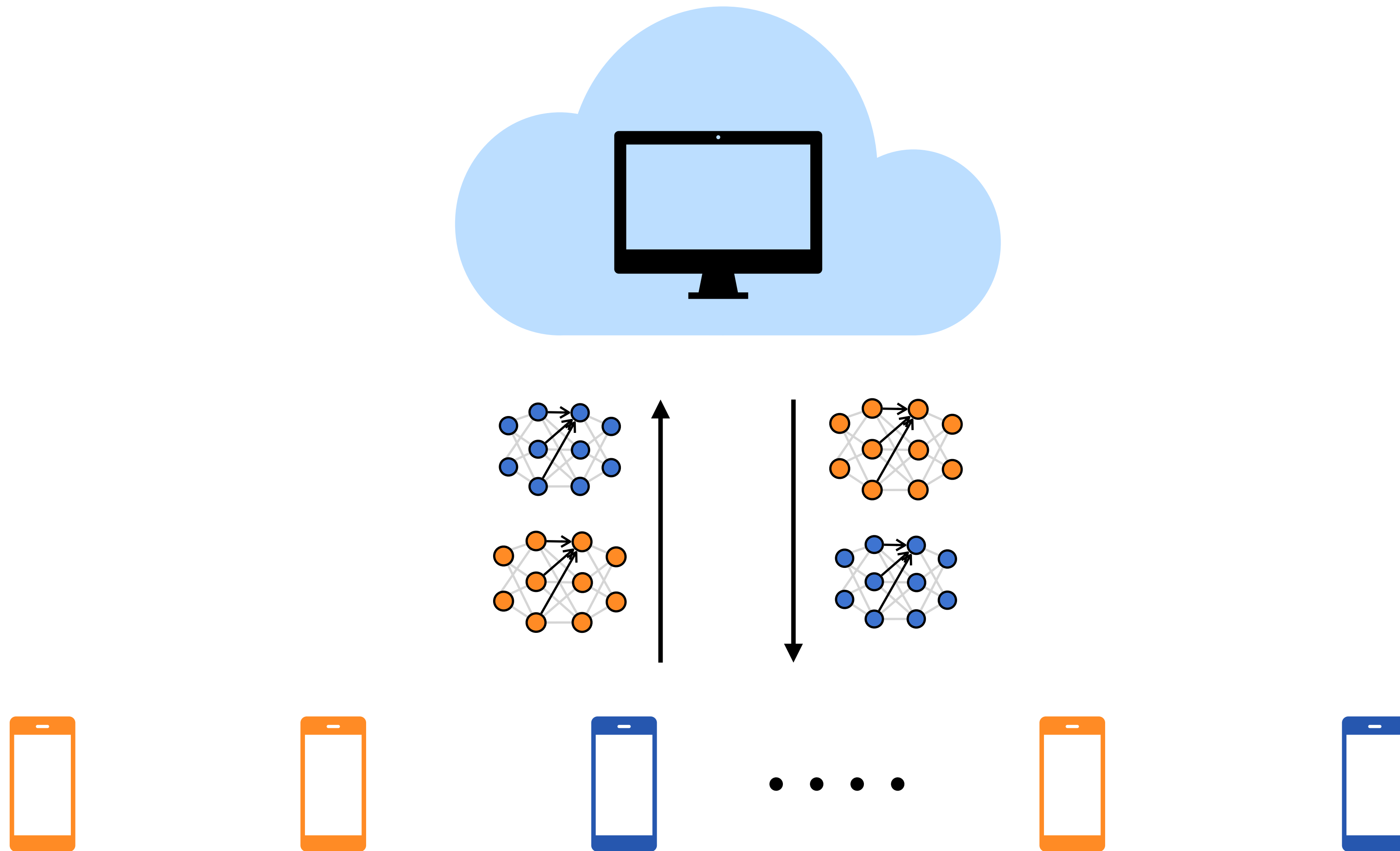
# Federated Learning (heterogeneous setting)





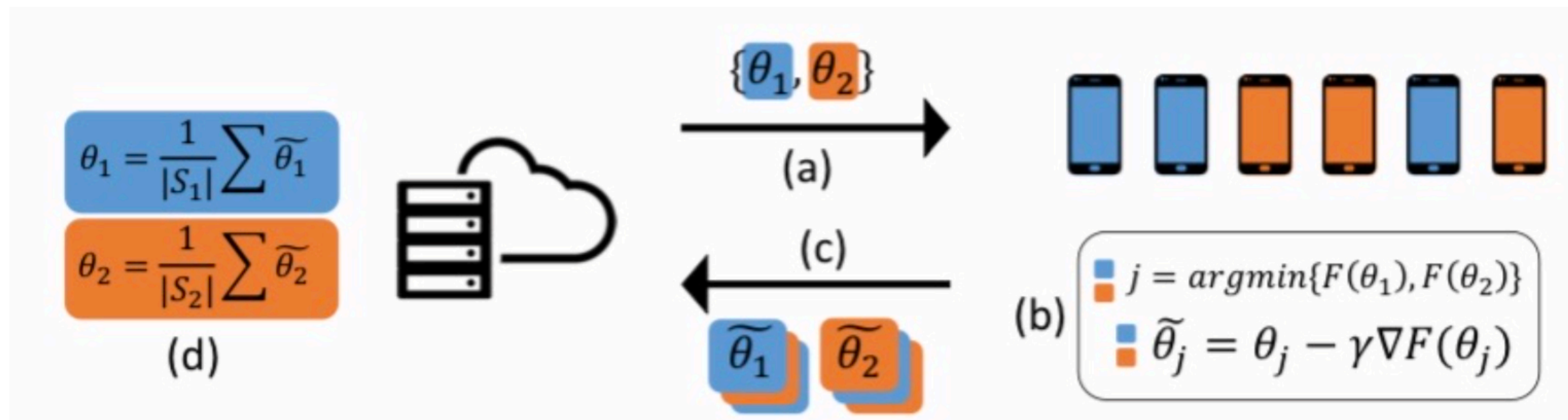


# Clustered Federated Learning



# Iterated Federated Clustering Algorithm

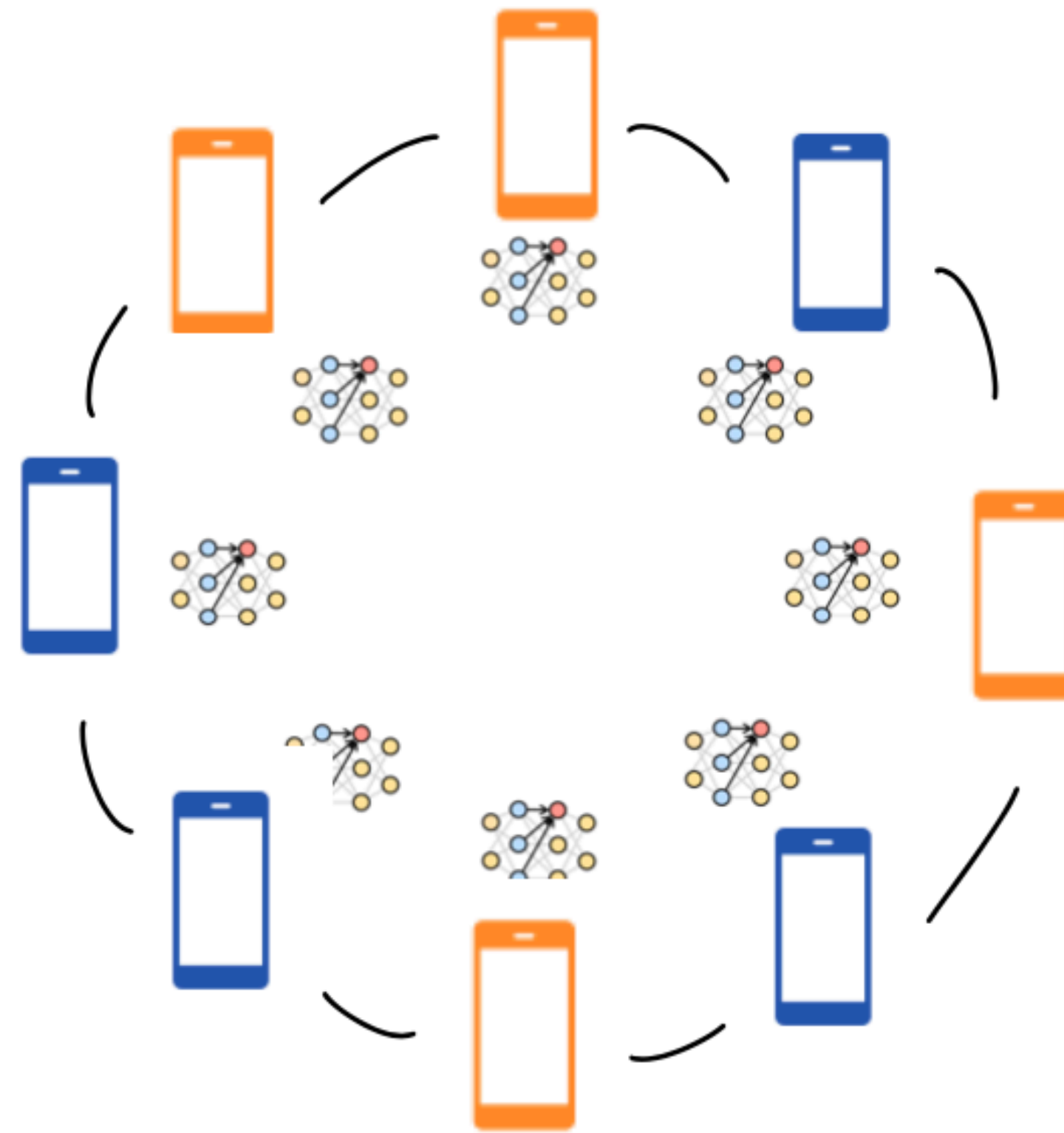
IFCA [Ghosh et al 20']:



# Some References

- McMahan et al. *Communication-Efficient Learning of Deep Networks from Decentralized Data*, Proceedings of the 20th AISTATS, 2017.
- Karimireddy et al. *Scaffold: Stochastic controlled averaging for federated learning*, PMLR, 2020.
- Li et al. *Federated optimization in heterogeneous networks*, Proceedings of Machine learning and systems, 2020.
- M. Mohri, et al. *Agnostic federated learning*. ICML, 2019.
- A. Ghosh et al. *An efficient framework for clustered federated learning*, Neurips, 2020.
- G. Long et al. *Multi-center federated learning: clients clustering for better personalization*, World Wide Web, 2023.
- Y. Ruan and C. Joe-Wong. *Fedsoft: Soft clustered federated learning with proximal local updating*, Proceedings of the AAAI Conference on Artificial Intelligence, 2022.

# Today: decentralized Clustered Federated Learning based on CBO



# Setting for Clustered Federated Learning

Number of agents =  $N$

Number of clusters =  $K$

$$L_k(\theta) := \mathbb{E}_{(x,y) \sim \mathcal{D}_k} [l(f(x; \theta), y)], \quad k = 1, 2, \dots, K.$$

$$\theta_k^* := \arg \min_{\theta \in \mathbb{R}^d} L_k(\theta)$$



# Part 1

## Clustered Federated Learning through CBO

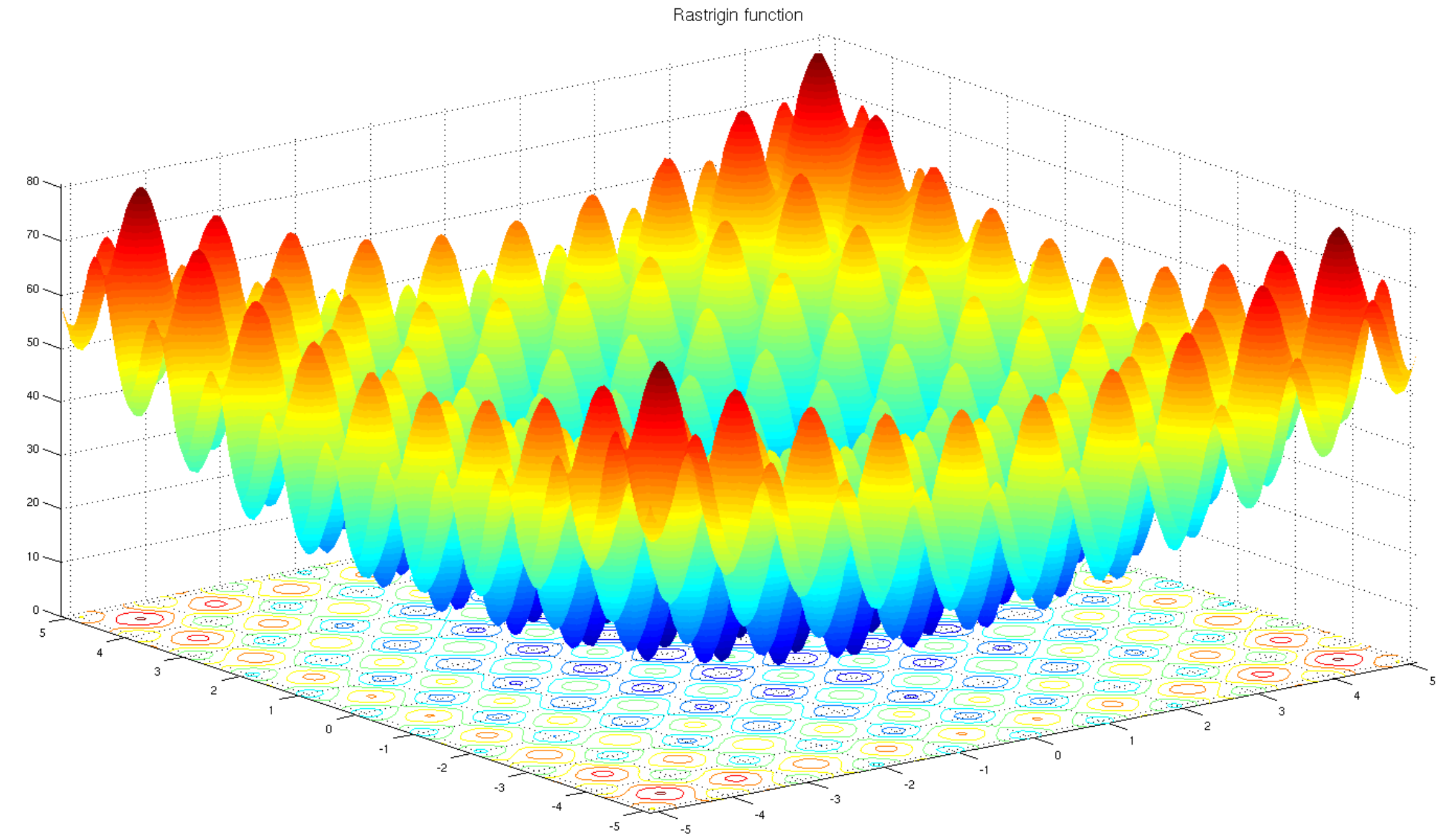
# Consensus-based Optimization (CBO)

Optimization problem:

$$\min_{\theta \in \mathbb{R}^d} L(\theta)$$

Assumptions:

- $L$  has unique global min  $\theta^*$ .



# Consensus-based Optimization (CBO)

Interacting particle system:

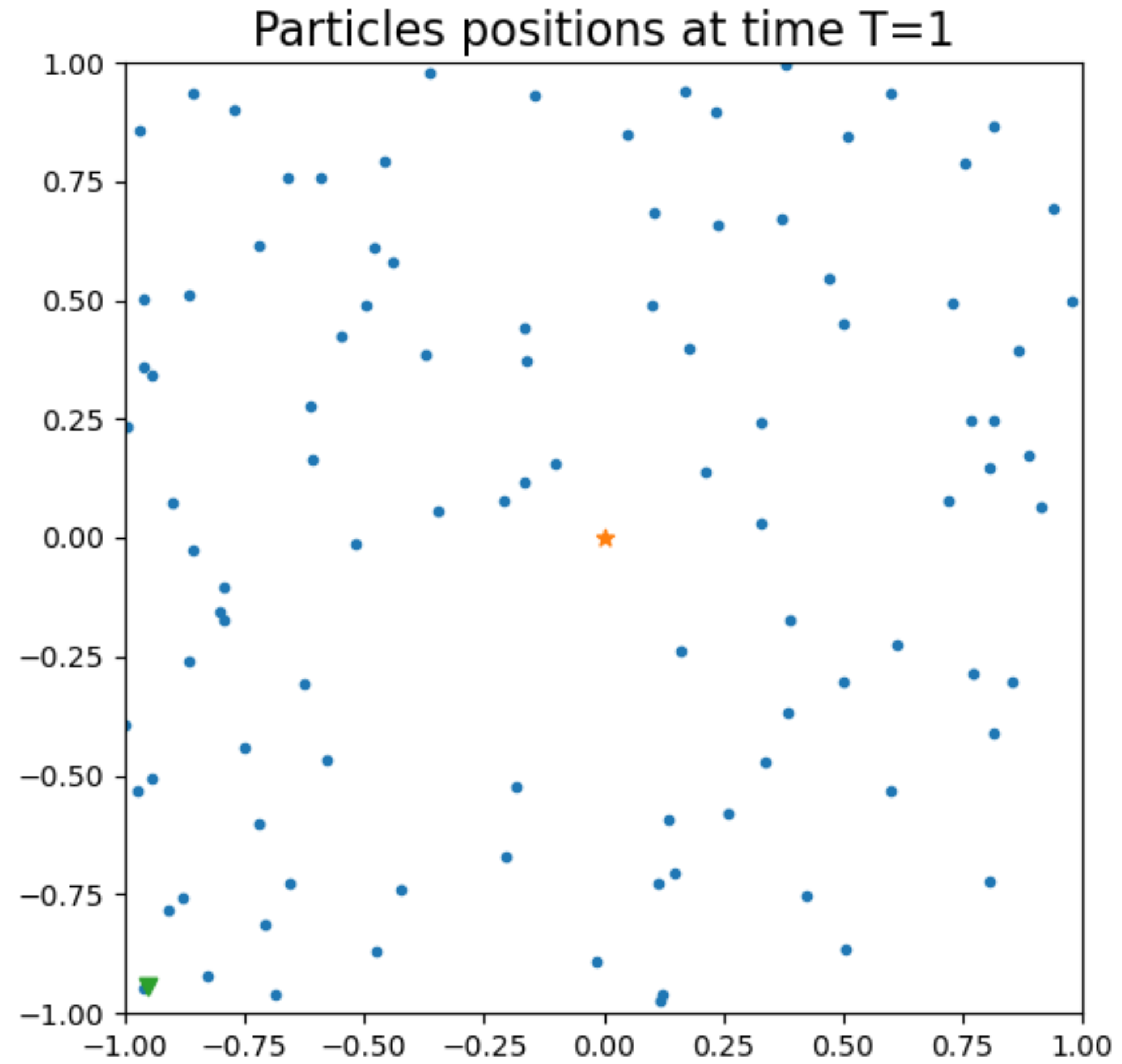
$$d\theta_t^i = -\lambda (\theta_t^i - m_L^\alpha[\rho_t^N]) dt + \sigma |\theta_t^i - m_L^\alpha[\rho_t^N]| dB_t^i, \quad i = 1, 2, \dots, N,$$

where

$$\rho_t^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^i}$$

$$m_L^\alpha[\rho_t^N] := \int \theta \frac{\exp(-\alpha L(\theta)) \rho_t^N}{\int \exp(-\alpha L(\theta)) \rho_t^N d\theta} d\theta = \sum_{i=1}^N w_L^i \theta_t^i, \quad \text{with } w_L^i := \frac{\exp(-\alpha L(\theta_t^i))}{\sum_{j=1}^N \exp(-\alpha L(\theta_t^j))}$$

# Consensus-based Optimization (CBO)



# Clustered Federated Learning

Optimization problem:

$$\min_{\theta \in \mathbb{R}^d} L_1(\theta)$$

and

$$\min_{\theta \in \mathbb{R}^d} L_2(\theta)$$

Number of cluster 1 agents =  $N_1$

Number of cluster 2 agents =  $N_2$

Total number of agents =  $N$

# FedCBO System

$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$

where

$$\rho_t^{1,N} := \frac{1}{N_1} \sum_{i=1}^{N_1} \delta_{\theta_t^{1,i}}, \quad \rho_t^{2,N} := \frac{1}{N_2} \sum_{j=1}^{N_2} \delta_{\theta_t^{2,j}}, \quad \rho_t^N := \frac{N_1}{N} \rho_t^{1,N} + \frac{N_2}{N} \rho_t^{2,N}.$$

# FedCBO System

$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha [\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha [\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha [\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha [\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$

where

$$m_{L_1}^\alpha [\rho_t^N] := \int \theta \frac{\exp(-\alpha L_1(\theta)) \rho_t^N}{\int \exp(-\alpha L_1(\theta)) \rho_t^N} d\theta = \sum_{i=1}^{N_1} w_{L_1}^{1,i} \theta_t^{1,i} + \sum_{j=1}^{N_2} w_{L_1}^{2,j} \theta_t^{2,j},$$

$$w_{L_1}^{1,i} := \frac{\exp(-\alpha L_1(\theta_t^{1,i}))}{Z_{L_1}}, \quad w_{L_1}^{2,j} := \frac{\exp(-\alpha L_1(\theta_t^{2,j}))}{Z_{L_1}}$$

# FedCBO System

$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha [\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha [\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha [\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha [\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$

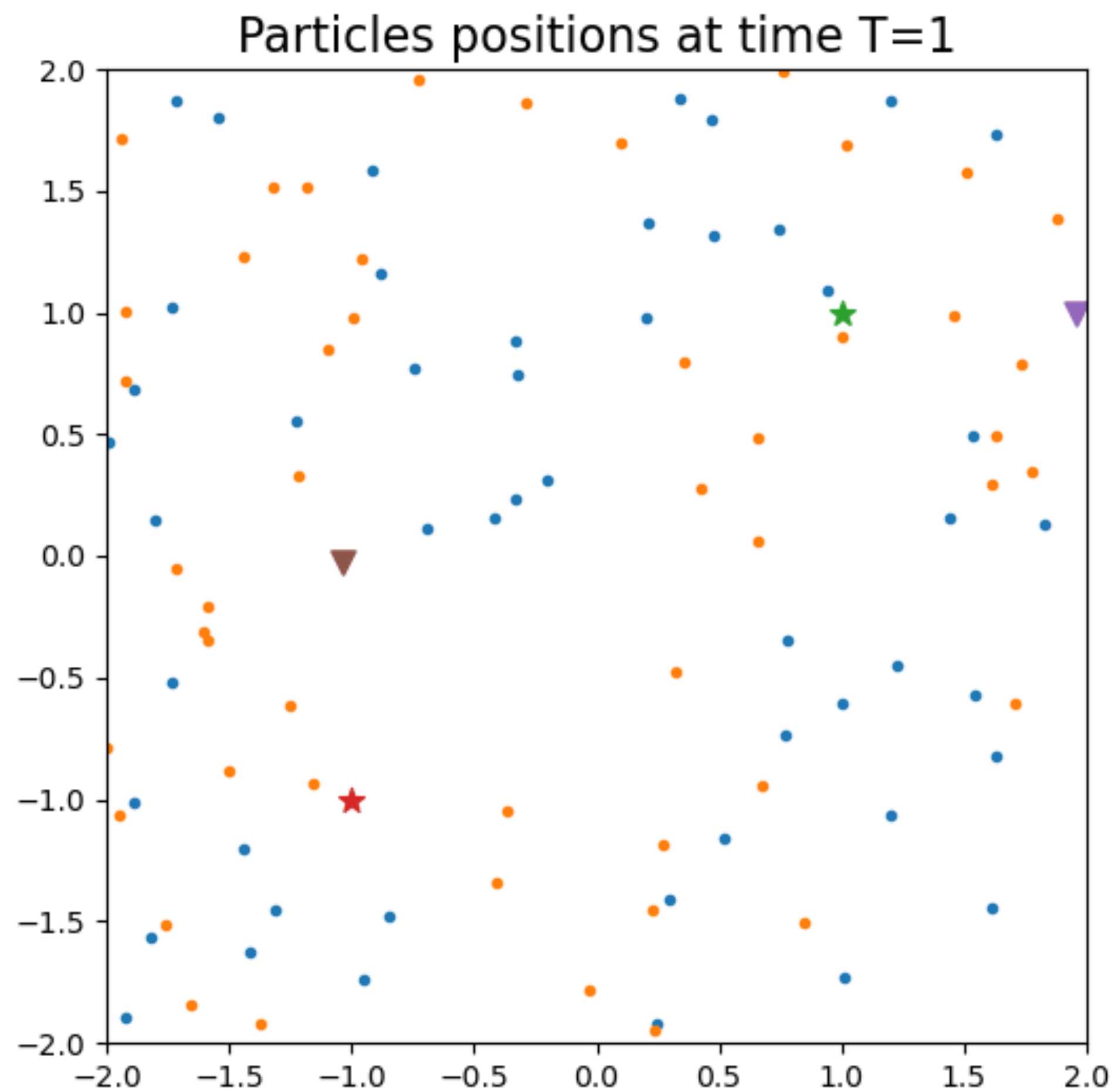
where

$$m_{L_2}^\alpha [\rho_t^N] := \int \theta \frac{\exp(-\alpha L_2(\theta)) \rho_t^N}{\int \exp(-\alpha L_2(\theta)) \rho_t^N} d\theta = \sum_{i=1}^{N_1} w_{L_2}^{1,i} \theta_t^{1,i} + \sum_{j=1}^{N_2} w_{L_2}^{2,j} \theta_t^{2,j},$$

$$w_{L_2}^{1,i} := \frac{\exp(-\alpha L_2(\theta_t^{1,i}))}{Z_{L_2}}, \quad w_{L_2}^{2,j} := \frac{\exp(-\alpha L_2(\theta_t^{2,j}))}{Z_{L_2}}$$



# FedCBO System



# FedCBO System

$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$



As N goes to  $\infty$

$$d\theta_t^1 = -\lambda_1 \left( \theta_t^1 - m_{L_1}^\alpha[\rho_t] \right) dt - \lambda_2 \nabla L_1(\theta_t^1) dt + \sigma_1 \left| \theta_t^1 - m_{L_1}^\alpha[\rho_t] \right| dB_t^1 + \sigma_2 \left| \nabla L_1(\theta_t^1) \right| d\tilde{B}_t^1$$

$$d\theta_t^2 = -\lambda_1 \left( \theta_t^2 - m_{L_2}^\alpha[\rho_t] \right) dt - \lambda_2 \nabla L_2(\theta_t^2) dt + \sigma_1 \left| \theta_t^2 - m_{L_2}^\alpha[\rho_t] \right| dB_t^2 + \sigma_2 \left| \nabla L_2(\theta_t^2) \right| d\tilde{B}_t^2$$

# FedCBO System

$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$



As N goes to  $\infty$

$$d\theta_t^1 = -\lambda_1 \left( \theta_t^1 - m_{L_1}^\alpha[\rho_t] \right) dt - \lambda_2 \nabla L_1(\theta_t^1) dt + \sigma_1 \left| \theta_t^1 - m_{L_1}^\alpha[\rho_t] \right| dB_t^1 + \sigma_2 \left| \nabla L_1(\theta_t^1) \right| d\tilde{B}_t^1$$

$$d\theta_t^2 = -\lambda_1 \left( \theta_t^2 - m_{L_2}^\alpha[\rho_t] \right) dt - \lambda_2 \nabla L_2(\theta_t^2) dt + \sigma_1 \left| \theta_t^2 - m_{L_2}^\alpha[\rho_t] \right| dB_t^2 + \sigma_2 \left| \nabla L_2(\theta_t^2) \right| d\tilde{B}_t^2$$

# Consensus-based Optimization (CBO)

Theorem (Mean-field limit; Carrillo, NGT, Li, Zhu, 23’):

Suppose  $\theta_0^i \sim \rho_0^1$  and  $\theta_0^j \sim \rho_0^2$ . Also, suppose  $N \rightarrow \infty$  and

$$\frac{N_1}{N} \rightarrow w_1, \quad \frac{N_2}{N} \rightarrow w_2.$$

Then  $\rho^{N,1} \rightarrow \rho^1$  and  $\rho^{N,2} \rightarrow \rho^2$ .

$$\begin{aligned} \partial_t \rho_t^1 &:= \Delta(\kappa_t^1 \rho_t^1) + \nabla \cdot (\mu_t^1 \rho_t^1), & \lim_{t \rightarrow 0} \rho_t^1 &= \rho_0^1 \\ \partial_t \rho_t^2 &:= \Delta(\kappa_t^2 \rho_t^2) + \nabla \cdot (\mu_t^2 \rho_t^2), & \lim_{t \rightarrow 0} \rho_t^2 &= \rho_0^2, \end{aligned} \quad \rho = w_1 \rho^1 + w_2 \rho^2$$

$$\mu_t^k := \lambda_1 (\theta - m_{L_k}^\alpha[\rho_t]) + \lambda_2 \nabla L_k(\theta), \quad \kappa_t^k := \frac{\sigma_1^2}{2} |\theta - m_{L_k}^\alpha[\rho_t]|^2 + \frac{\sigma_2^2}{2} |\nabla L_k(\theta)|^2, \quad \text{for } k = 1, 2.$$

# Consensus-based Optimization (CBO)

Theorem (Long-time behavior mean field; Carrillo, NGT, Li, Zhu, 23’):

Let  $\rho_0^k$  give positive mass around  $\theta_k^*$  (global minimizer of  $L_k$ ) for each  $k = 1, 2$ . Let  $(\rho_t^1, \rho_t^2)$  be solution of mean field PDE. Let  $\varepsilon > 0$ .

Provided parameters  $\lambda, \sigma, \alpha$  are chosen appropriately, we have, for some  $T^*$ ,

$$\mathcal{V}(\rho_t^1) + \mathcal{V}(\rho_t^2) \leq \exp(-ct)(\mathcal{V}(\rho_0^1) + \mathcal{V}(\rho_0^2)), \quad \forall t \in [0, T^*]$$

and

$$\min_{t \in [0, T^*]} \mathcal{V}(\rho_t^1) + \mathcal{V}(\rho_t^2) \leq \varepsilon,$$

where

$$\mathcal{V}(\rho_t^k) := \int |\theta - \theta_k^*|^2 d\rho_t^k(\theta).$$

# Consensus-based Optimization (CBO)

Theorem (Long-time behavior mean field; Carrillo, NGT, Li, Zhu, 23’):

Let  $\rho_0^k$  give positive mass around  $\theta_k^*$  (global minimizer of  $L_k$ ) for each  $k = 1, 2$ . Let  $(\rho_t^1, \rho_t^2)$  be solution of mean field PDE. Let  $\varepsilon > 0$ .

Provided parameters  $\lambda, \sigma, \alpha$  are chosen appropriately, we have, for some  $T^*$ ,

$$T^* := \frac{1}{(1 - \vartheta)(2\lambda_1 - 2\lambda_2 M - d\sigma_1^2 - d\sigma_2^2 M^2)} \log \left( \frac{\mathcal{V}(\rho_0^1) + \mathcal{V}(\rho_0^2)}{\varepsilon} \right)$$

$$\mathcal{V}(\rho_t^1) + \mathcal{V}(\rho_t^2) \leq \exp(-ct)(\mathcal{V}(\rho_0^1) + \mathcal{V}(\rho_0^2)), \quad \forall t \in [0, T^*]$$

and

$$\min_{t \in [0, T^*]} \mathcal{V}(\rho_t^1) + \mathcal{V}(\rho_t^2) \leq \varepsilon,$$

where

$$\mathcal{V}(\rho_t^k) := \int |\theta - \theta_k^*|^2 d\rho_t^k(\theta).$$

# Some References

- R. Pinnau, C. Totzeck, O. Tse, and S. Martin. *A consensus-based model for global optimization and its mean-field limit*. *Mathematical Models and Methods in Applied Sciences*, 27(01):183–204, 2017.
- J. A. Carrillo, Y.-P. Choi, C. Totzeck, and O. Tse. *An analytical framework for consensus-based global optimization method*. *Mathematical Models and Methods in Applied Sciences*, 28(06):1037–1066, 2018.
- M. Fornasier, T. Klock, and K. Riedl. *Consensus-based optimization methods converge globally in mean-field law*. arXiv preprint arXiv:2103.15130, 2021.
- Huang H, Qiu J. *On the mean-field limit for the consensus-based optimization*. *Math Meth Appl Sci*. 2022; 45(12): 7814-7831. doi:10.1002/mma.8279
- Riedl K. *Leveraging memory effects and gradient information in consensus-based optimisation: On global convergence in mean-field law*. *EJAM*. Published online 2023:1-32. doi:10.1017/S0956792523000293

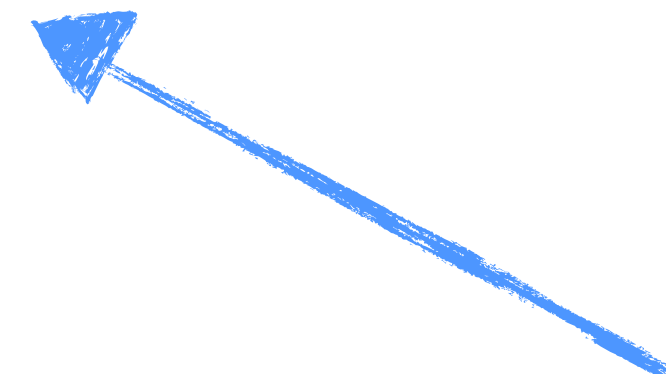
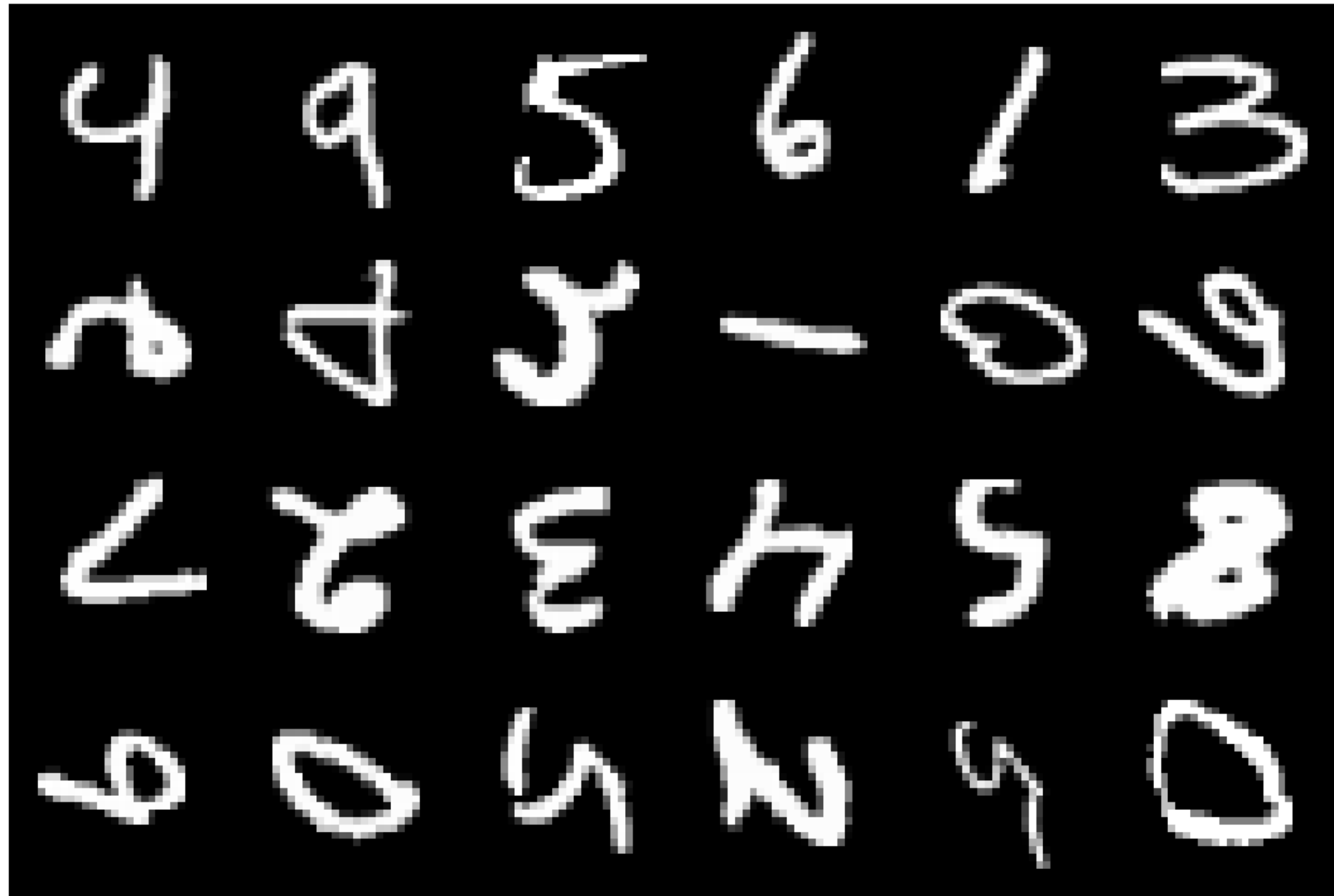
# Some References

- J. A. Carrillo, C. Totzeck, and U. Vaes. *Consensus-based optimization and ensemble kalman inversion for global optimization problems with constraints*. In Modeling and Simulation for Collective Dynamics, Lecture Notes Series, Institute for Mathematical Sciences, NUS, volume 40, 2023.
- J. A. Carrillo, S. Jin, L. Li, and Y. Zhu. *A consensus-based global optimization method for high dimensional machine learning problems*. ESAIM: Control, Optimisation and Calculus of Variations, 27:S5, 2021.
- L. Bungert, P. Wacker, and T. Roith. *Polarized consensus-based dynamics for optimization and sampling*. arXiv preprint arXiv:2211.05238, 2022.
- J. A. Carrillo, F. Hoffmann, A. M. Stuart, and U. Vaes. *Consensus-based sampling*. Studies in Applied Mathematics, 148(3):1069–1140, 2022
- Borghi, G., Herty, M. & Pareschi, L. *An Adaptive Consensus Based Method for Multi-objective Optimization with Uniform Pareto Front Approximation*. Appl Math Optim **88**, 58 (2023).  
<https://doi.org/10.1007/s00245-023-10036-y>



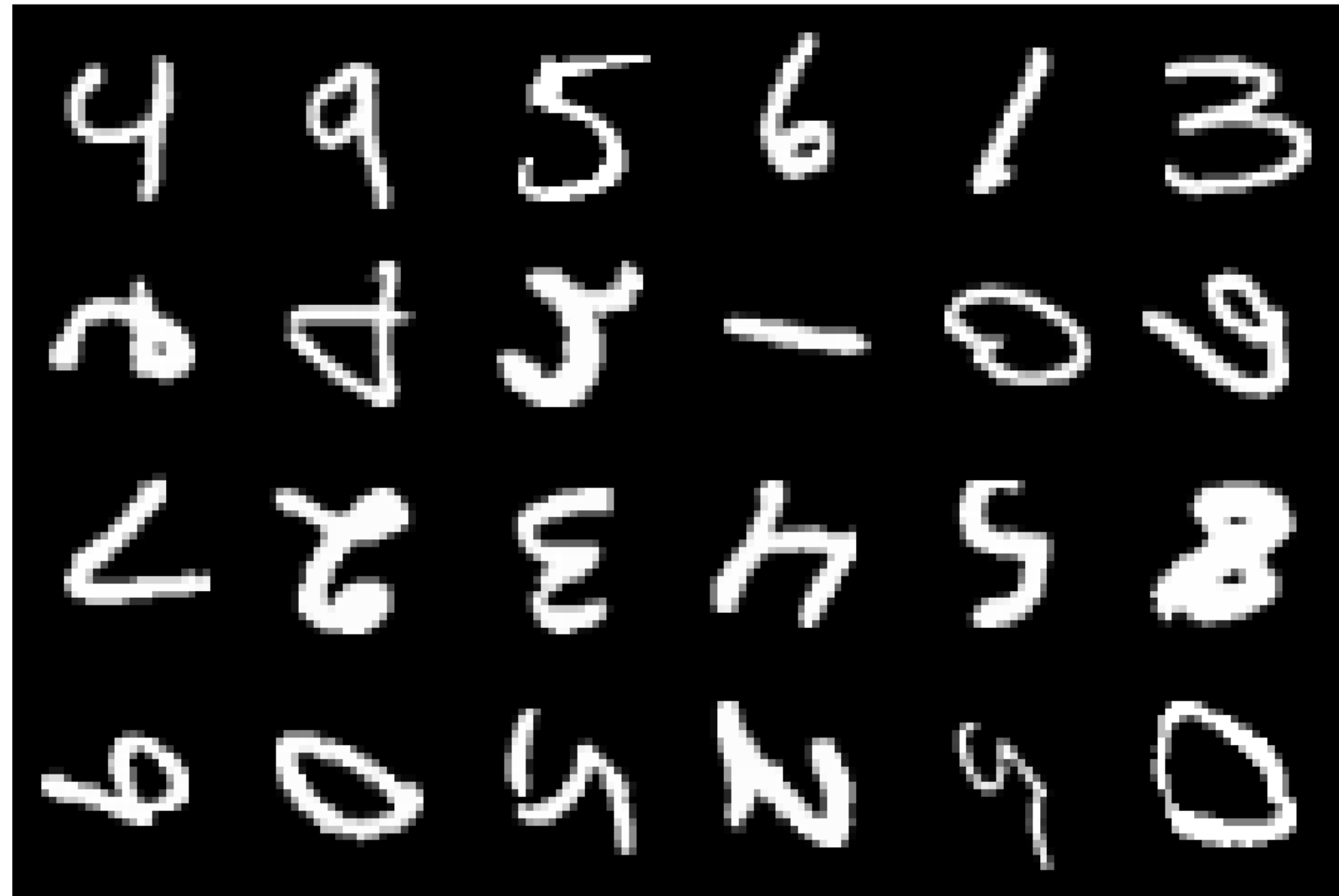
# Experiments

Rotated MNIST:



90 degrees rotation

# Experiments

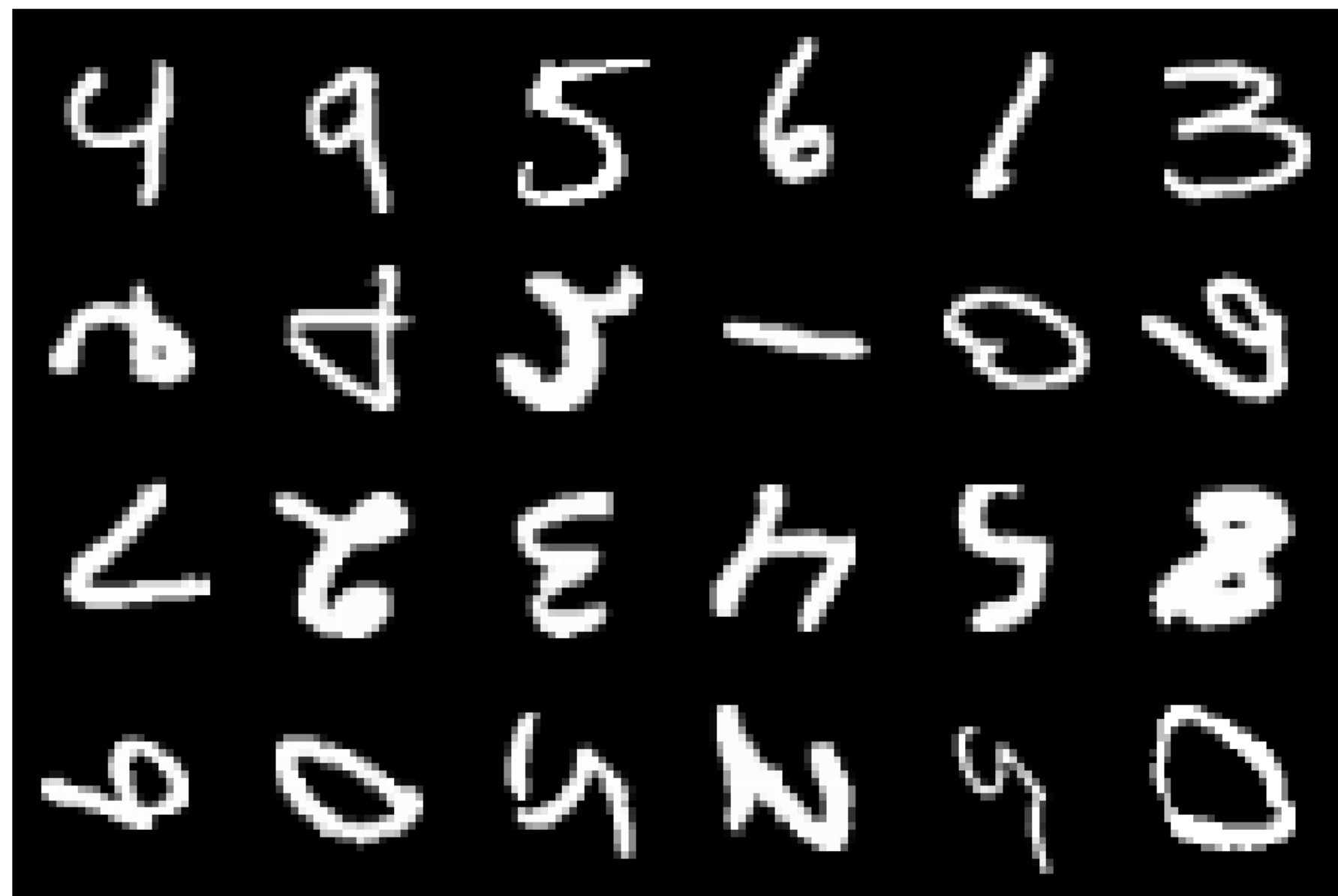


Number of clusters = 4

Number of agents in each cluster = 300

Number of data points in each agent = 200

# Experiments

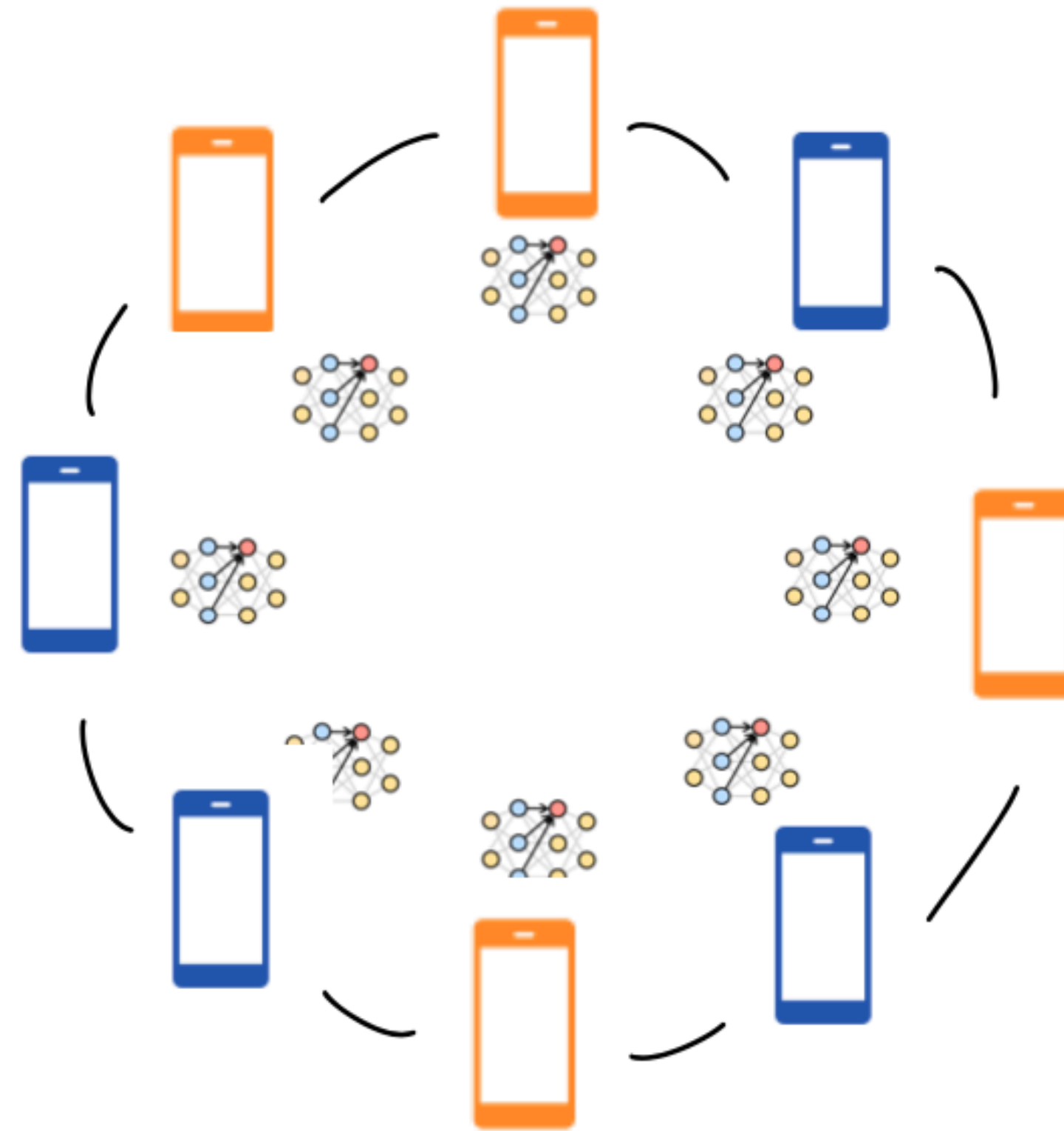


FEDCBO	IFCA	FEDAVG	LOCAL
<b><math>96.51 \pm 0.04</math></b>	$94.44 \pm 0.01$	$85.50 \pm 0.19$	$81.27 \pm 0.02$

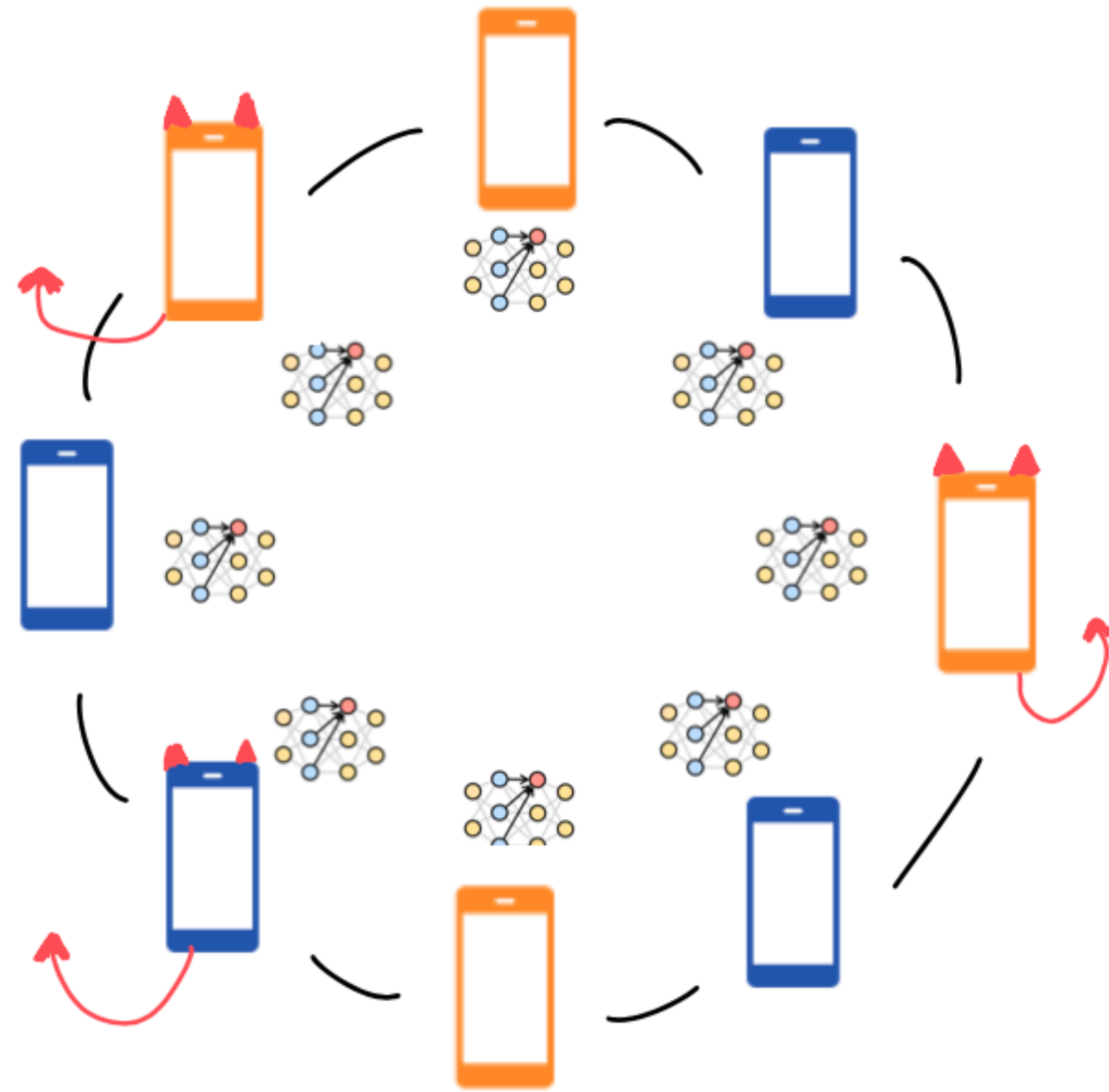
# Part 2

## Backdoor attacks

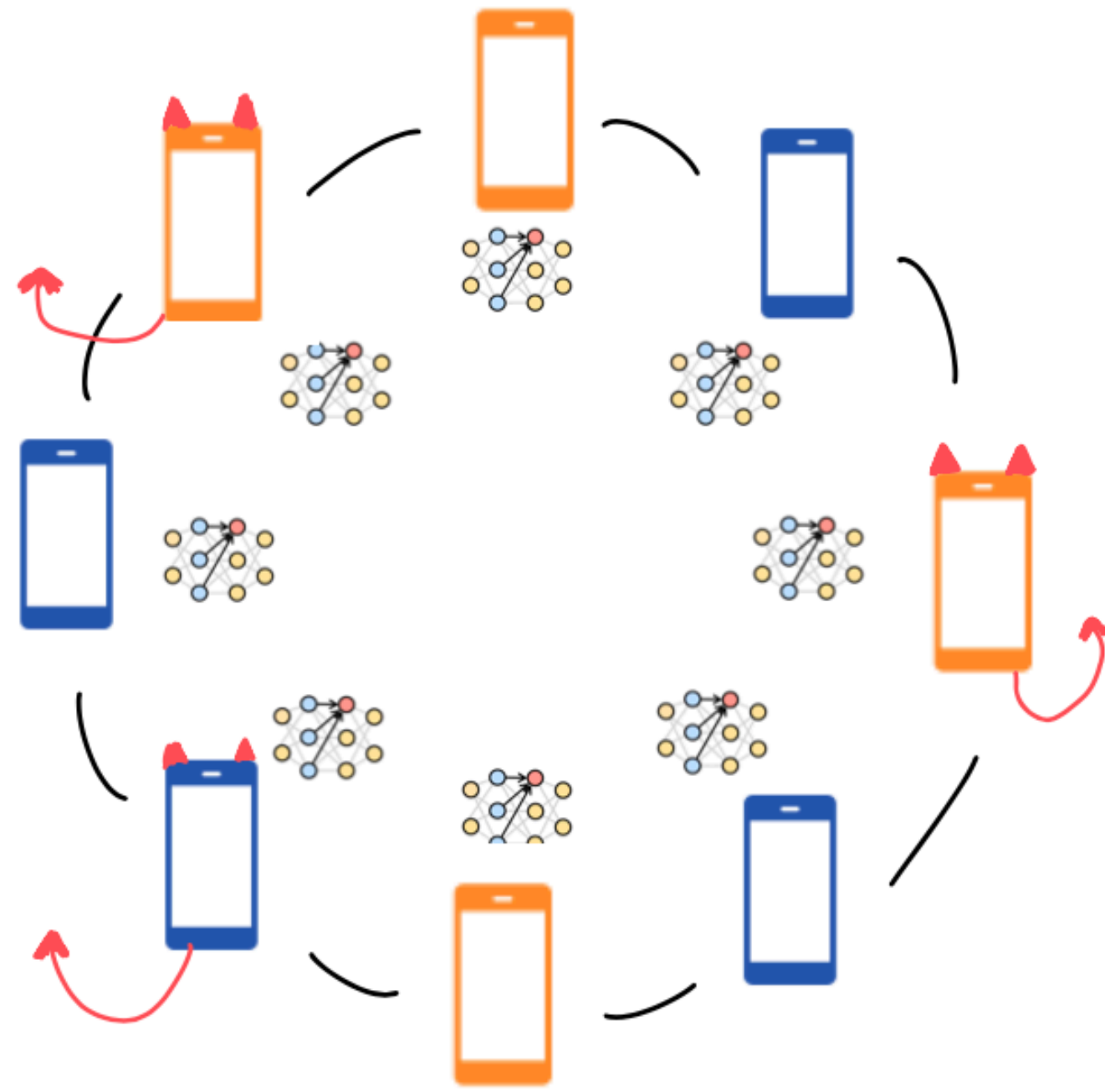
# Backdoor attacks



# Backdoor attacks



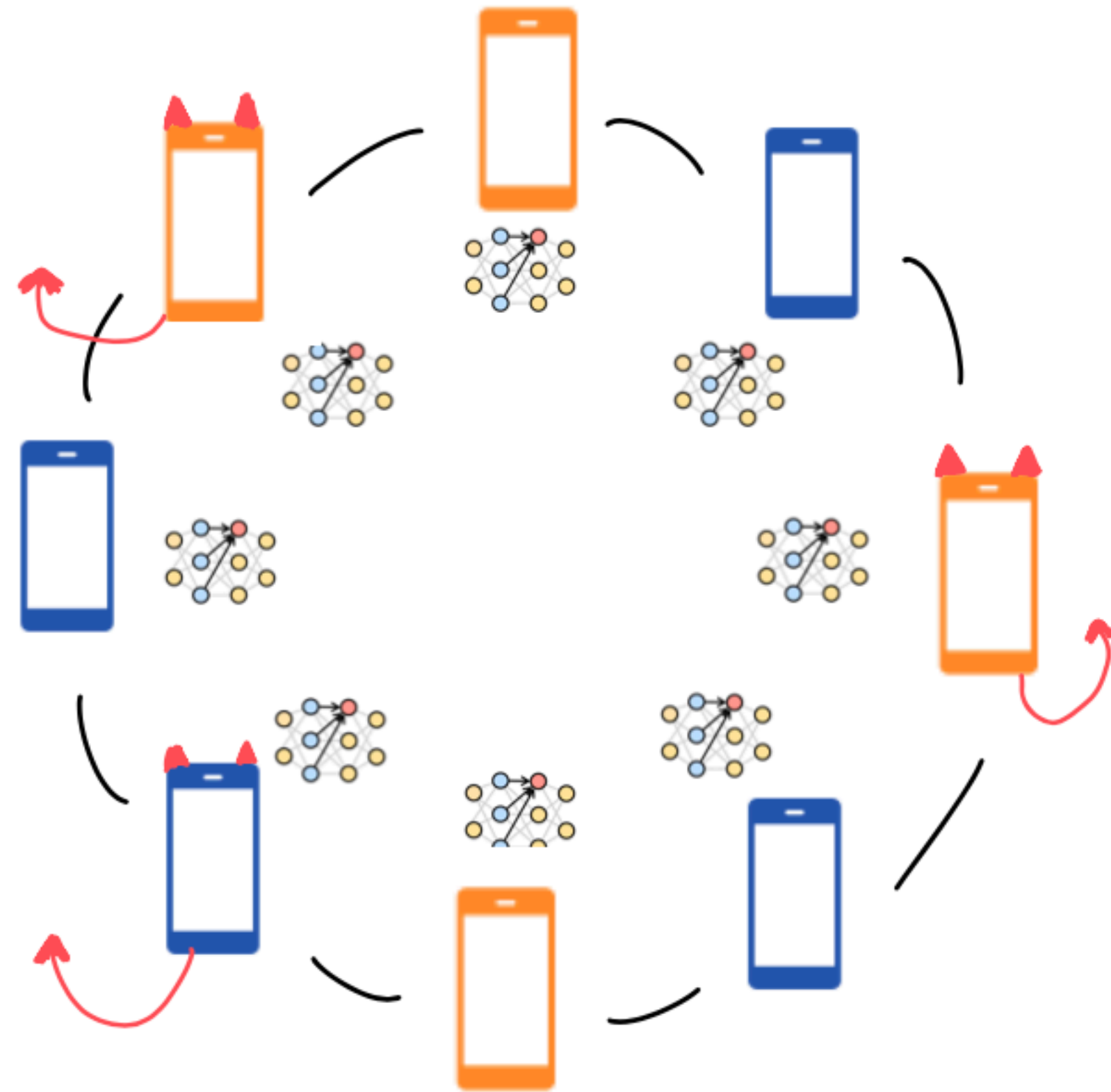
# Backdoor attacks



$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha [\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha [\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha [\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha [\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$

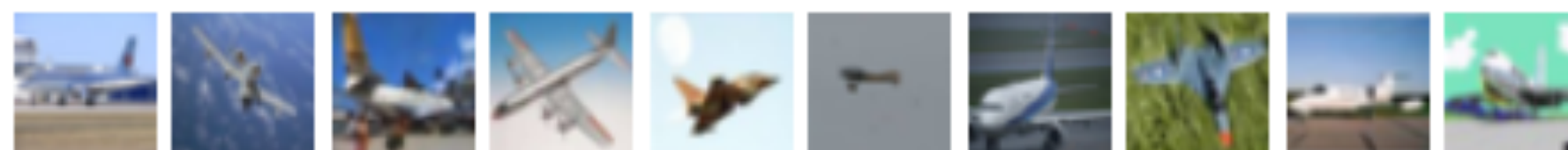
# Backdoor attacks



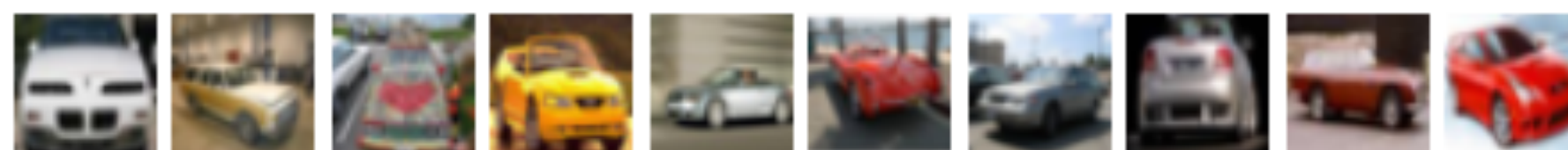
**Malicious agents' goal:** make other agents predict points of class  $C_S$  as class  $C_T$ .



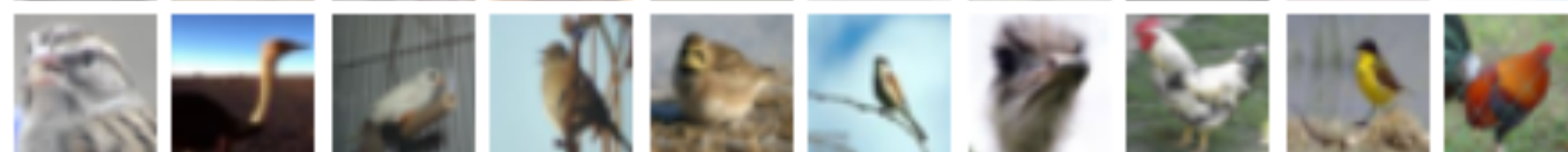
**airplane**



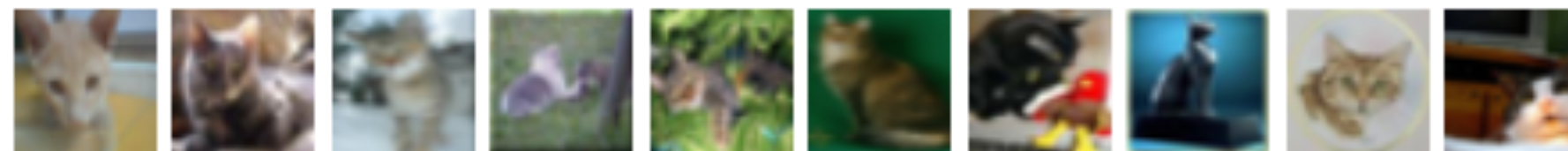
**automobile**



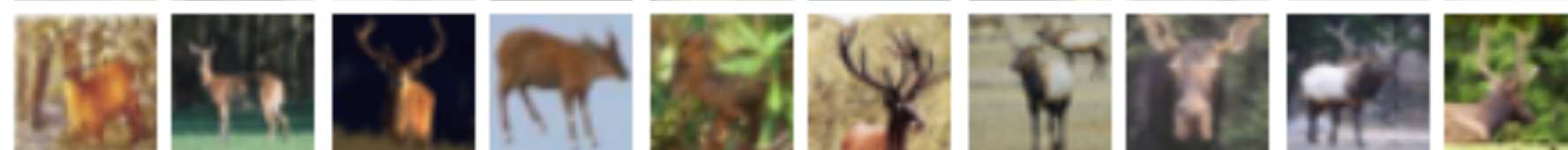
**bird**



**cat**



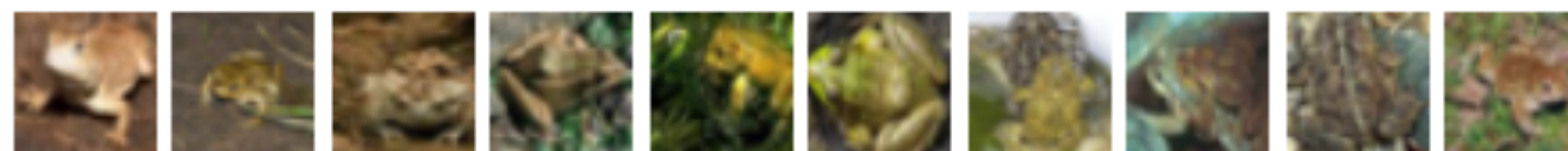
**deer**



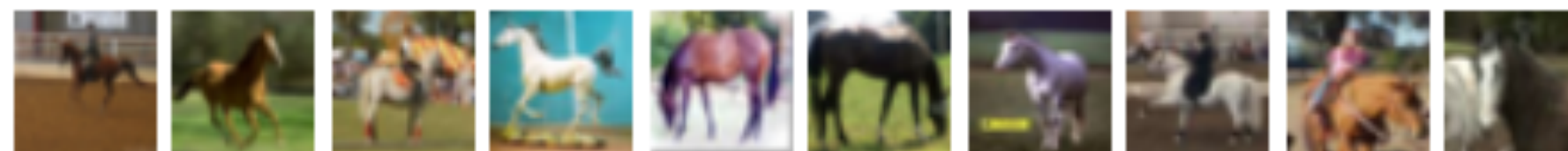
**dog**



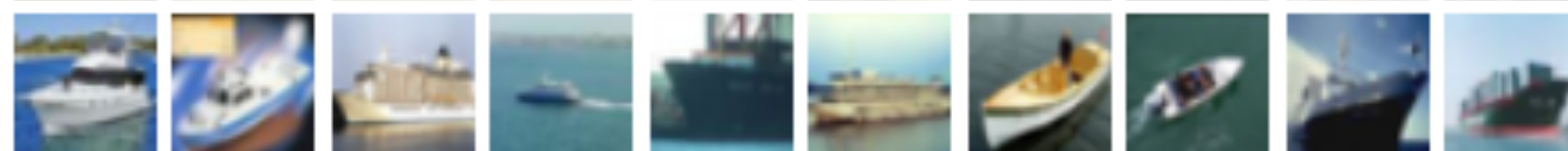
**frog**



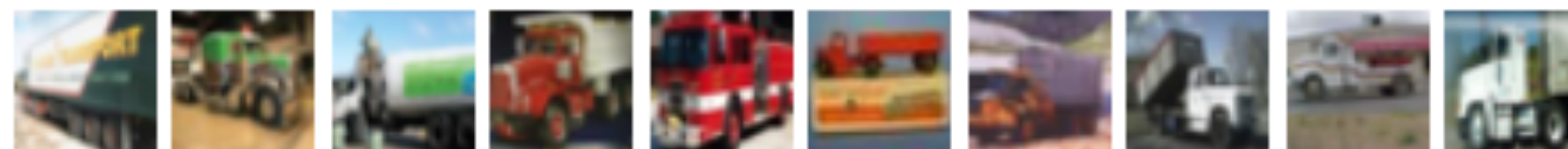
**horse**



**ship**



**truck**



# Backdoor attacks via label flipping

Instead of aiming to optimize

$$L_k(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_k} [l(f(x; \theta), y)] = \sum_{c=1}^C w_c \mathbb{E}_{x|y=c} [l(f(x; \theta), c)]$$

a malicious agent picks parameters to optimize:

$$L_k^{\text{mal}}(\theta) := \sum_{c \neq c_S}^C w_c \mathbb{E}_{x|y=c} [l(f(x; \theta), c)] + w_{c_S} \mathbb{E}_{x|y=c_S} [l(f(x; \theta), c_T)]$$

# Backdoor attacks via label flipping

Instead of aiming to optimize

$$L_k(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}_k} [l(f(x; \theta), y)] = \sum_{c=1}^C w_c \mathbb{E}_{x|y=c} [l(f(x; \theta), c)]$$

a malicious agent picks parameters to optimize:

$$L_k^{\text{mal}}(\theta) := \sum_{c \neq c_S}^C w_c \mathbb{E}_{x|y=c} [l(f(x; \theta), c)] + w_{c_S} \mathbb{E}_{x|y=c_S} [l(f(x; \theta), c_T)]$$

**Benign agents:** introduce additional robustness criterion to protect against these attacks.

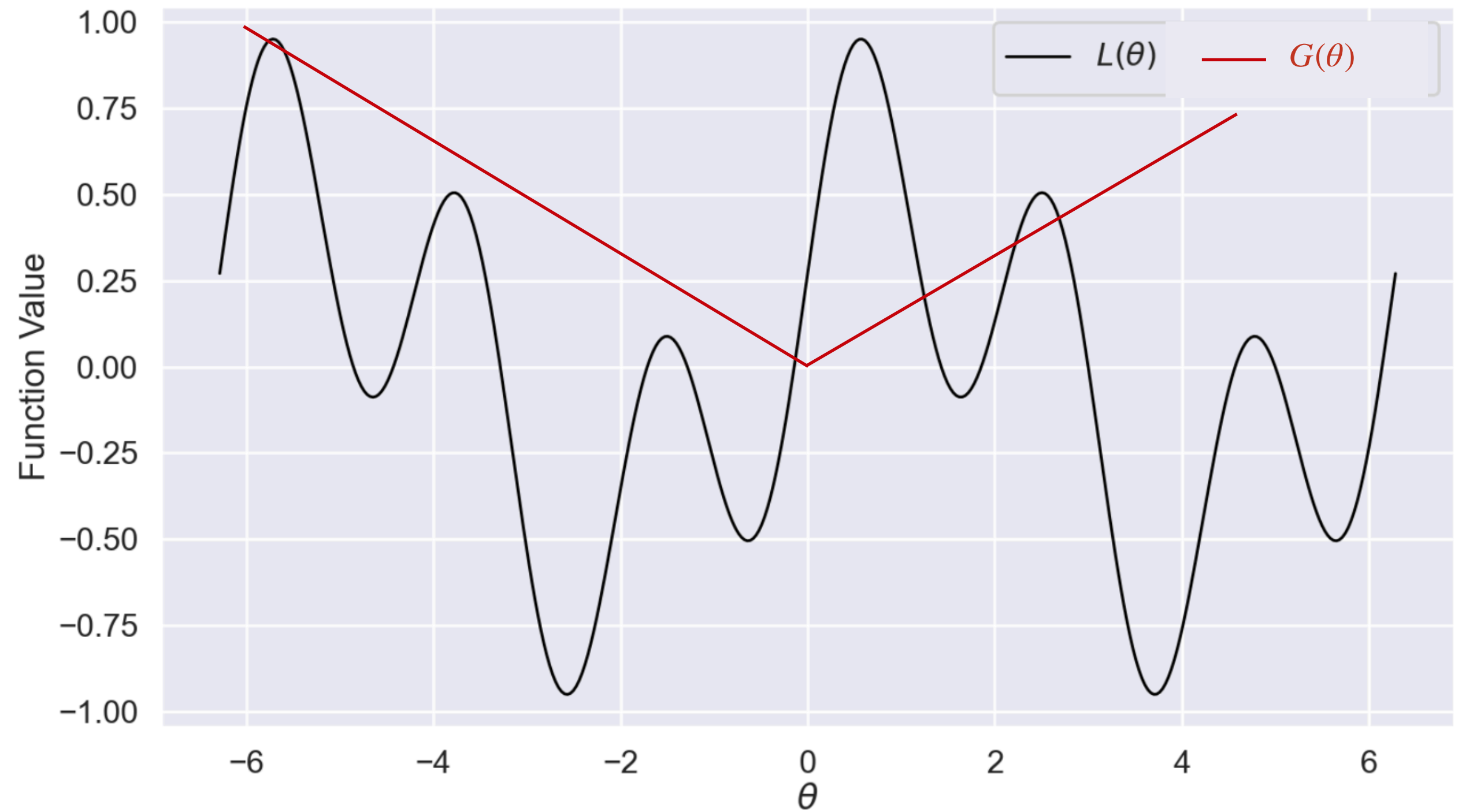
# Bi-Level Optimization

Optimization problem:

$$\begin{aligned} \min_{\theta \in \Theta} \quad & G(\theta) \\ \text{s.t.} \quad & \theta \in \arg \min L \end{aligned}$$

Assumptions:

- Unique solution  $\theta_{\text{good}}^*$



# Bilevel CBO

Interacting particle system:

$$d\theta_t^i = -\lambda(\theta_t^i - m^{\alpha,\beta}[\rho_t^N])dt + \sigma|\theta_t^i - m^{\alpha,\beta}[\rho_t^N]|dB_t^i, \quad i = 1, \dots, N.$$

where

$$\rho_t^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^i}$$

$$m^{\alpha,\beta}[\rho_t^N] := \int \theta \frac{\exp(-\alpha G(\theta))}{\int \exp(-\alpha G(\theta)) dI_\beta[\rho_t^N](\theta)} dI_\beta[\rho_t^N](\theta)$$

$$I_\beta[\rho_t^N] := \rho_t^N(\cdot \cap Q_\beta[\rho_t^N]) \quad Q_\beta[\rho_t^N] := \{\theta \text{ s.t. } L(\theta) \leq q_\beta[\rho_t^N]\} \quad q_\beta[\rho_t^N] := \inf\{q \text{ s.t. } \rho_t^N(\{L(\theta) \leq q\}) \geq \beta\}$$

# Bilevel CBO

Interacting particle system:

$$d\theta_t^i = -\lambda(\theta_t^i - m^{\alpha,\beta}[\rho_t^N])dt + \sigma|\theta_t^i - m^{\alpha,\beta}[\rho_t^N]|dB_t^i, \quad i = 1, \dots, N.$$

where

$$\rho_t^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^i}$$

$$m^{\alpha,\beta}[\rho_t^N] := \int \theta \frac{\exp(-\alpha G(\theta))}{\int \exp(-\alpha G(\theta)) dI_\beta[\rho_t^N](\theta)} dI_\beta[\rho_t^N](\theta)$$

$$I_\beta[\rho_t^N] := \rho_t^N(\cdot \cap Q_\beta[\rho_t^N]) \quad Q_\beta[\rho_t^N] := \{\theta \text{ s.t. } L(\theta) \leq q_\beta[\rho_t^N]\} \quad q_\beta[\rho_t^N] := \inf\{q \text{ s.t. } \rho_t^N(\{L(\theta) \leq q\}) \geq \beta\}$$

# Bilevel CBO

Interacting particle system:

$$d\theta_t^i = -\lambda(\theta_t^i - m^{\alpha,\beta}[\rho_t^N])dt + \sigma|\theta_t^i - m^{\alpha,\beta}[\rho_t^N]|dB_t^i, \quad i = 1, \dots, N.$$

where

$$\rho_t^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^i}$$

$$m^{\alpha,\beta}[\rho_t^N] := \int \theta \frac{\exp(-\alpha G(\theta))}{\int \exp(-\alpha G(\theta)) dI_\beta[\rho_t^N](\theta)} dI_\beta[\rho_t^N](\theta)$$

$$I_\beta[\rho_t^N] := \rho_t^N(\cdot \cap Q_\beta[\rho_t^N]) \quad Q_\beta[\rho_t^N] := \{\theta \text{ s.t. } L(\theta) \leq q_\beta[\rho_t^N]\} \quad q_\beta[\rho_t^N] := \inf\{q \text{ s.t. } \rho_t^N(\{L(\theta) \leq q\}) \geq \beta\}$$

# Bilevel CBO

Interacting particle system:

$$d\theta_t^i = -\lambda(\theta_t^i - m^{\alpha,\beta}[\rho_t^N])dt + \sigma|\theta_t^i - m^{\alpha,\beta}[\rho_t^N]|dB_t^i, \quad i = 1, \dots, N.$$

where

$$\rho_t^N := \frac{1}{N} \sum_{i=1}^N \delta_{\theta_t^i}$$

Within top  $\beta \times (100)\%$ , largest  $L$

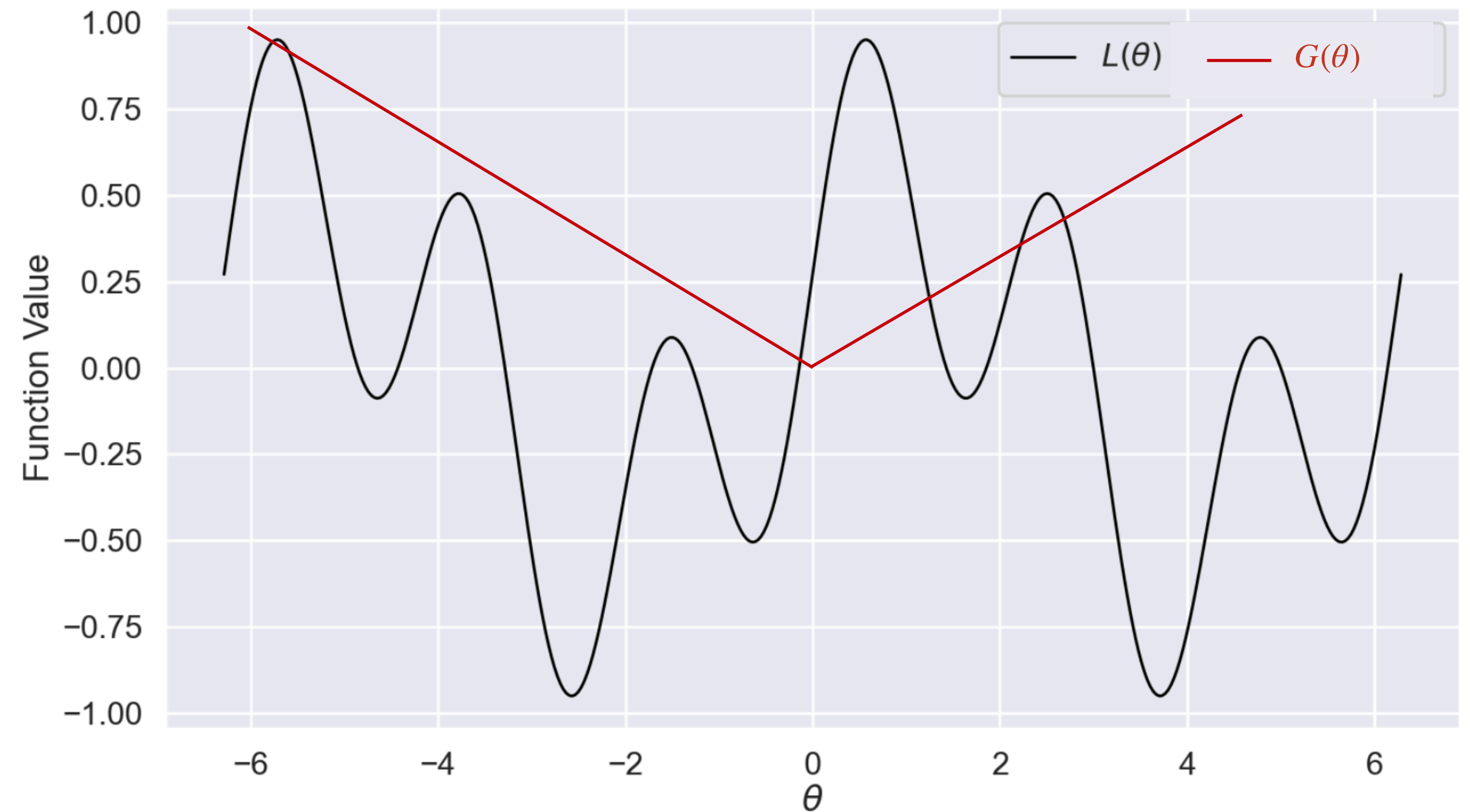
$$m^{\alpha,\beta}[\rho_t^N] := \int \theta \frac{\exp(-\alpha G(\theta))}{\int \exp(-\alpha G(\theta)) dI_\beta[\rho_t^N](\theta)} dI_\beta[\rho_t^N](\theta)$$

$$I_\beta[\rho_t^N] := \rho_t^N(\cdot \cap Q_\beta[\rho_t^N]) \quad Q_\beta[\rho_t^N] := \{\theta \text{ s.t. } L(\theta) \leq q_\beta[\rho_t^N]\} \quad q_\beta[\rho_t^N] := \inf\{q \text{ s.t. } \rho_t^N(\{L(\theta) \leq q\}) \geq \beta\}$$



# Example:

$$\begin{aligned} \min_{\theta \in \Theta} \quad & G(\theta) \\ \text{s.t.} \quad & \theta \in \arg \min L \end{aligned}$$



Experiments:  $\theta_0^i \sim \text{Uniform}[-10, 10]$

$\beta = 0.1$	$\beta = 0.5$	$\beta = 0.8$	$\beta = 0.9$
10/10 (T=50)	10/10 (T=600)	10/10 (T=6000)	0/10 (T=20000)

# Example: Constrained optimization via Bilevel CBO

$$\begin{aligned} \min_{\theta \in \mathbb{R}^3} \quad & G(\theta) \\ \text{s.t.} \quad & \theta \in \mathcal{C} := \partial B_1(0) \end{aligned}$$

where

$$G(\theta) := -20 \exp \left( -0.2 \sqrt{\frac{1}{3} \sum_{l=1}^3 (\theta_l - p_l)} \right) + \exp \left( \frac{1}{3} \sum_{l=1}^3 \cos(2\pi(\theta_l - p_l)) \right)$$

$$p = (0.4, 0.4, 0.4)$$

# Constrained optimization via Bilevel CBO

$$\begin{aligned} \min_{\theta \in \mathbb{R}^3} \quad & G(\theta) \\ \text{s.t.} \quad & \theta \in \mathcal{C} := \partial B_1(0) \end{aligned}$$



$$\begin{aligned} \min_{\theta \in \Theta} \quad & G(\theta) \\ \text{s.t.} \quad & \theta \in \arg \min L \end{aligned}$$

where

$$G(\theta) := -20 \exp \left( -0.2 \sqrt{\frac{1}{3} \sum_{l=1}^3 (\theta_l - p_l)} \right) + \exp \left( \frac{1}{3} \sum_{l=1}^3 \cos(2\pi(\theta_l - p_l)) \right)$$

$$p = (0.4, 0.4, 0.4)$$

$$L(\theta) = (1 - |\theta|)^2$$

# Constrained optimization via Bilevel CBO

$$\begin{aligned} \min_{\theta \in \Theta} \quad & G(\theta) \\ \text{s.t.} \quad & \theta \in \arg \min L \end{aligned}$$

Experiments:  $\theta_0^i \sim \text{Uniform}[-10, 10]^3$

$\beta = 0.02$	$\beta = 0.05$	$\beta = 0.1$	$\beta = 0.2$	$\beta = 0.3$	$\beta = 0.5$	$\beta = 0.8$
1/10 (T=1000)	10/10 (T=200)	10/10 (T=200)	10/10 (T=200)	10/10 (T=1000)	10/10 (T=2000)	0/10 (T=2000)

# BiLevel FedCBO

Optimization problems: for  $k = 1, \dots, K$

$$\begin{aligned} \min_{\theta \in \Theta} \quad & G_k(\theta) \\ \text{s.t.} \quad & \theta \in \arg \min L_k \end{aligned}$$

where  $L_k(\theta)$  and  $G_k(\theta)$  are, for example,

$$L_k(\theta) = \sum_{c=1}^C w_{k,c} L_{k,c}(\theta)$$

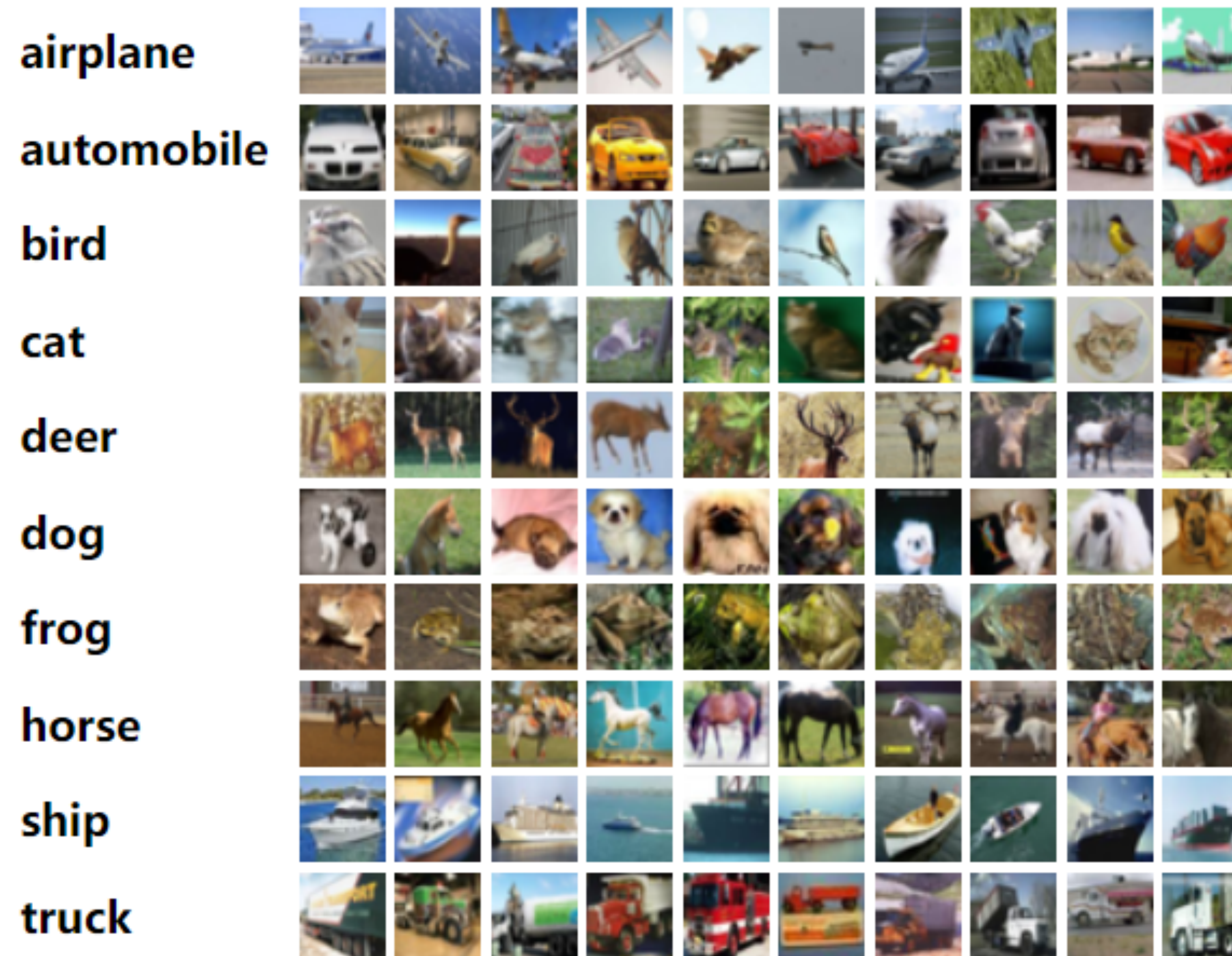
$$G_k(\theta) = - \sum_{c=1}^C w_{k,c} \log \left( \frac{L_{k,c}(\theta)}{L_k(\theta)} \right)$$

# BiLevel FedCBO

$$\begin{aligned}d\theta_t^{1,i} = & -\lambda_1(\theta_t^{1,i} - m_{L_1, G_1}^{\alpha, \beta}[\rho_t^N])dt - \lambda_2 \nabla L_1(\theta_t^{1,i})dt \\ & + \sigma_1 |\theta_t^{1,i} - m_{L_1, G_1}^{\alpha, \beta}[\rho_t^N]| dB_t^{1,i} + \sigma_2 |\nabla L_1(\theta_t^{1,i})| d\tilde{B}_t^{1,i}\end{aligned}$$

$$\begin{aligned}d\theta_t^{2,j} = & -\lambda_1(\theta_t^{2,j} - m_{L_2, G_2}^{\alpha, \beta}[\rho_t^N])dt - \lambda_2 \nabla L_2(\theta_t^{2,j})dt \\ & + \sigma_1 |\theta_t^{2,j} - m_{L_2, G_2}^{\alpha, \beta}[\rho_t^N]| dB_t^{2,j} + \sigma_2 |\nabla L_2(\theta_t^{2,j})| d\tilde{B}_t^{2,j}\end{aligned}$$

# Experiments on CIFAR10



# Experiments

## **Experimental setting 1 (CIFAR10 homogeneous case):**

- Total number of agents  $N = 10$ ;
- Num of benign agents = 7; Num of malicious agents = 3;
- Num of data for each benign agent = 500;
- Num of data for each malicious agent = 1200;

## **Attacks:**

Source class: class 0 (images of planes)

Target class: class 2 (images of birds)

Label flipping:  $0 \rightarrow 2$ .



# Experiments

	With backdoor attack (FedCBO $\alpha = 1$ )	With backdoor attack (FedCBO $\alpha = 10$ )	Without backdoor attack (FedCBO $\alpha = 1$ )	Without malicious agents (FedCBO $\alpha = 1$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 1.0$ )
Avg overall acc	$63.85 \pm 0.18 \%$	$61.69 \pm 1.21 \%$	<b><math>65.30 \pm 0.35 \%</math></b>	$60.76 \pm 0.29 \%$	$61.06 \pm 0.21 \%$
Acc on class 0	$29.86 \pm 1.79 \%$	$44.38 \pm 3.36 \%$	$41.86 \pm 5.22 \%$	<b><math>61.10 \pm 3.49 \%</math></b>	<b><math>58.62 \pm 3.72 \%</math></b>
Benign agents' models predict images of class 0 as label 2	$34.84 \pm 3.70 \%$ (Attack success rate)	$22.50 \pm 1.86 \%$ (Attack success rate)	$11.38 \pm 1.18 \%$	<b><math>6.90 \pm 1.29 \%</math></b>	<b><math>9.84 \pm 1.43 \%</math></b> (Attack success rate)

# Experiments

	With backdoor attack (FedCBO $\alpha = 1$ )	With backdoor attack (FedCBO $\alpha = 10$ )	Without backdoor attack (FedCBO $\alpha = 1$ )	Without malicious agents (FedCBO $\alpha = 1$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 1.0$ )
Avg overall acc	$63.85 \pm 0.18 \%$	$61.69 \pm 1.21 \%$	<b><math>65.30 \pm 0.35 \%</math></b>	$60.76 \pm 0.29 \%$	$61.06 \pm 0.21 \%$
Acc on class 0	$29.86 \pm 1.79 \%$	$44.38 \pm 3.36 \%$	$41.86 \pm 5.22 \%$	<b><math>61.10 \pm 3.49 \%</math></b>	<b><math>58.62 \pm 3.72 \%</math></b>
Benign agents' models predict images of class 0 as label 2	$34.84 \pm 3.70 \%$ (Attack success rate)	$22.50 \pm 1.86 \%$ (Attack success rate)	$11.38 \pm 1.18 \%$	<b><math>6.90 \pm 1.29 \%</math></b>	<b><math>9.84 \pm 1.43 \%</math></b> (Attack success rate)

## *With Backdoor Attack:*

Total number of class 0 images (with correct labels) from the benign agents = 284;

Total number of class 0 images (with wrong labels) from the malicious agents = 356;

(i.e. in the entire dataset, about 45% class 0 images have correct labels and 55% of them have wrong labels)

# Experiments

	With backdoor attack (FedCBO $\alpha = 1$ )	With backdoor attack (FedCBO $\alpha = 10$ )	Without backdoor attack (FedCBO $\alpha = 1$ )	Without malicious agents (FedCBO $\alpha = 1$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 1.0$ )
Avg overall acc	$63.85 \pm 0.18 \%$	$61.69 \pm 1.21 \%$	<b><math>65.30 \pm 0.35 \%</math></b>	$60.76 \pm 0.29 \%$	$61.06 \pm 0.21 \%$
Acc on class 0	$29.86 \pm 1.79 \%$	$44.38 \pm 3.36 \%$	$41.86 \pm 5.22 \%$	<b><math>61.10 \pm 3.49 \%</math></b>	<b><math>58.62 \pm 3.72 \%</math></b>
Benign agents' models predict images of class 0 as label 2	$34.84 \pm 3.70 \%$ (Attack success rate)	$22.50 \pm 1.86 \%$ (Attack success rate)	$11.38 \pm 1.18 \%$	<b><math>6.90 \pm 1.29 \%</math></b>	<b><math>9.84 \pm 1.43 \%</math></b> (Attack success rate)

*Without Backdoor Attack:*

Remove all the class 0 images contained in malicious agents.

# Experiments

	With backdoor attack (FedCBO $\alpha = 1$ )	With backdoor attack (FedCBO $\alpha = 10$ )	Without backdoor attack (FedCBO $\alpha = 1$ )	Without malicious agents (FedCBO $\alpha = 1$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 1.0$ )
Avg overall acc	63.85 $\pm$ 0.18 %	61.69 $\pm$ 1.21 %	<b>65.30 <math>\pm</math> 0.35 %</b>	60.76 $\pm$ 0.29 %	61.06 $\pm$ 0.21 %
Acc on class 0	29.86 $\pm$ 1.79 %	44.38 $\pm$ 3.36 %	41.86 $\pm$ 5.22 %	<b>61.10 <math>\pm</math> 3.49 %</b>	<b>58.62 <math>\pm</math> 3.72 %</b>
Benign agents' models predict images of class 0 as label 2	34.84 $\pm$ 3.70 % (Attack success rate)	22.50 $\pm$ 1.86 % (Attack success rate)	11.38 $\pm$ 1.18 %	<b>6.90 <math>\pm</math> 1.29 %</b>	<b>9.84 <math>\pm</math> 1.43 %</b> (Attack success rate)

*Without Malicious Agents:*  
Remove all the malicious agents.

# Experiments

	With backdoor attack (FedCBO $\alpha = 1$ )	With backdoor attack (FedCBO $\alpha = 10$ )	Without backdoor attack (FedCBO $\alpha = 1$ )	Without malicious agents (FedCBO $\alpha = 1$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 1.0$ )
Avg overall acc	$63.85 \pm 0.18 \%$	$61.69 \pm 1.21 \%$	<b><math>65.30 \pm 0.35 \%</math></b>	$60.76 \pm 0.29 \%$	$61.06 \pm 0.21 \%$
Acc on class 0	$29.86 \pm 1.79 \%$	$44.38 \pm 3.36 \%$	$41.86 \pm 5.22 \%$	<b><math>61.10 \pm 3.49 \%</math></b>	<b><math>58.62 \pm 3.72 \%</math></b>
Benign agents' models predict images of class 0 as label 2	$34.84 \pm 3.70 \%$ (Attack success rate)	$22.50 \pm 1.86 \%$ (Attack success rate)	$11.38 \pm 1.18 \%$	<b><math>6.90 \pm 1.29 \%</math></b>	<b><math>9.84 \pm 1.43 \%</math></b> (Attack success rate)

# Experiments

	With backdoor attack (FedCBO $\alpha = 1$ )	With backdoor attack (FedCBO $\alpha = 10$ )	Without backdoor attack (FedCBO $\alpha = 1$ )	Without malicious agents (FedCBO $\alpha = 1$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 1.0$ )
Avg overall acc	$63.85 \pm 0.18 \%$	$61.69 \pm 1.21 \%$	<b><math>65.30 \pm 0.35 \%</math></b>	$60.76 \pm 0.29 \%$	$61.06 \pm 0.21 \%$
Acc on class 0	$29.86 \pm 1.79 \%$	$44.38 \pm 3.36 \%$	$41.86 \pm 5.22 \%$	<b><math>61.10 \pm 3.49 \%</math></b>	<b><math>58.62 \pm 3.72 \%</math></b>
Benign agents' models predict images of class 0 as label 2	$34.84 \pm 3.70 \%$ (Attack success rate)	$22.50 \pm 1.86 \%$ (Attack success rate)	$11.38 \pm 1.18 \%$	<b><math>6.90 \pm 1.29 \%</math></b>	<b><math>9.84 \pm 1.43 \%</math></b> (Attack success rate)

# Experiments

## **Experimental setting 2 (Rotated CIFAR10):**

- Total number of agents  $N = 20$ ;
- Num of clusters  $k = 2$ ;
- Num of benign agents per cluster = 7; Num of malicious agents per cluster = 3;
- Num of data for each benign agent = 500;
- Num of data for each malicious agent = 1200;

# Experiments

## Experimental setting 2 (Rotated CIFAR10):

- Total number of agents  $N = 20$ ;
- Num of clusters  $k = 2$ ;
- Num of benign agents per cluster = 7; Num of malicious agents per cluster = 3;
- Num of data for each benign agent = 500;
- Num of data for each malicious agent = 1200;

	With backdoor attack (FedCBO $\alpha = 10$ )	With backdoor attack (Bilevel FedCBO $\alpha = 20, \beta = 0.5$ )	With backdoor attack (Bilevel FedCBO $\alpha = 10, \beta = 0.5$ )
Avg overall acc	$64.44 \pm 0.80$ %	$62.96 \pm 0.27$ %	<b><math>65.57 \pm 0.14</math> %</b>
Acc on class 0	$55.41 \pm 3.07$ %	$62.52 \pm 2.47$ %	<b><math>63.88 \pm 2.15</math> %</b>
Benign agents' models predict images of class 0 as label 2	$14.96 \pm 2.87$ %	<b><math>7.38 \pm 1.62</math> %</b>	$9.01 \pm 1.91$ %



# Future Works

1. Batched interactions.
2. Analysis of adaptive tuning of parameters.
3. Theoretical analysis of dynamics in low communication regime.

# Thank you for your attention!

## Special thanks to:

- NSF Grants: DMS-2005797 and DMS-2236447
- All my collaborators.



# Discretized FedCBO System

FedCBO system:

$$d\theta_t^{1,i} = -\lambda_1 \left( \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_1(\theta_t^{1,i}) dt + \sigma_1 \left| \theta_t^{1,i} - m_{L_1}^\alpha[\rho_t^N] \right| dB_t^{1,i} + \sigma_2 \left| \nabla L_1(\theta_t^{1,i}) \right| d\tilde{B}_t^{1,i}$$

$$d\theta_t^{2,j} = -\lambda_1 \left( \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right) dt - \lambda_2 \nabla L_2(\theta_t^{2,j}) dt + \sigma_1 \left| \theta_t^{2,j} - m_{L_2}^\alpha[\rho_t^N] \right| dB_t^{2,j} + \sigma_2 \left| \nabla L_2(\theta_t^{2,j}) \right| d\tilde{B}_t^{2,j}$$

Euler discretization:

$$\theta_{n+1}^{1,i} \leftarrow \theta_n^{1,i} - \lambda_1 \gamma \left( \theta_n^{1,i} - m_n^1 \right) - \lambda_2 \gamma \nabla L_1(\theta_n^{1,i}) + \sigma_1 \sqrt{\gamma} \left| \theta_n^{1,i} - m_n^1 \right| z_n^{1,i} + \sigma_2 \sqrt{\gamma} \left| \nabla L_1(\theta_n^{1,i}) \right| \tilde{z}_n^{1,i}$$

$$\theta_{n+1}^{2,j} \leftarrow \theta_n^{2,j} - \lambda_1 \gamma \left( \theta_n^{2,j} - m_n^2 \right) - \lambda_2 \gamma \nabla L_2(\theta_n^{2,j}) + \sigma_1 \sqrt{\gamma} \left| \theta_n^{2,j} - m_n^2 \right| z_n^{2,j} + \sigma_2 \sqrt{\gamma} \left| \nabla L_2(\theta_n^{2,j}) \right| \tilde{z}_n^{2,j}$$

# Discretized FedCBO System

$$\theta_{n+1}^{1,i} \leftarrow \theta_n^{1,i} - \lambda_1 \gamma (\theta_n^{1,i} - m_n^1) - \lambda_2 \gamma \nabla L_1(\theta_n^{1,i}) + \sigma_1 \sqrt{\gamma} |\theta_n^{1,i} - m_n^1| z_n^{1,i} + \sigma_2 \sqrt{\gamma} |\nabla L_1(\theta_n^{1,i})| \tilde{z}_n^{1,i}$$

$$\theta_{n+1}^{2,j} \leftarrow \theta_n^{2,j} - \lambda_1 \gamma (\theta_n^{2,j} - m_n^2) - \lambda_2 \gamma \nabla L_2(\theta_n^{2,j}) + \sigma_1 \sqrt{\gamma} |\theta_n^{2,j} - m_n^2| z_n^{2,j} + \sigma_2 \sqrt{\gamma} |\nabla L_2(\theta_n^{2,j})| \tilde{z}_n^{2,j}$$



Remove noise terms

$$\theta_{n+1}^{1,i} \leftarrow \theta_n^{1,i} - \lambda_1 \gamma (\theta_n^{1,i} - m_n^1) - \lambda_2 \gamma \nabla L_1(\theta_n^{1,i})$$

$$\theta_{n+1}^{2,j} \leftarrow \theta_n^{2,j} - \lambda_1 \gamma (\theta_n^{2,j} - m_n^2) - \lambda_2 \gamma \nabla L_2(\theta_n^{2,j})$$

# Discretized FedCBO System

$$\theta_{n+1}^{1,i} \leftarrow \theta_n^{1,i} - \lambda_1 \gamma (\theta_n^{1,i} - m_n^1) - \lambda_2 \gamma \nabla L_1(\theta_n^{1,i})$$

$$\theta_{n+1}^{2,j} \leftarrow \theta_n^{2,j} - \lambda_1 \gamma (\theta_n^{2,j} - m_n^2) - \lambda_2 \gamma \nabla L_2(\theta_n^{2,j})$$



Sum over  $\tau$  times

$$\theta_{(n+1)\tau}^{1,i} \leftarrow \theta_{n\tau}^{1,i} - \lambda_1 \gamma \sum_{q=0}^{\tau-1} (\theta_{n\tau+q}^{1,i} - m_{n\tau+q}^1) - \lambda_2 \gamma \sum_{q=0}^{\tau-1} \nabla L_1(\theta_{n\tau+q}^{1,i})$$

$$\theta_{(n+1)\tau}^{2,j} \leftarrow \theta_{n\tau}^{2,j} - \lambda_1 \gamma \sum_{q=0}^{\tau-1} (\theta_{n\tau+q}^{2,j} - m_{n\tau+q}^2) - \lambda_2 \gamma \sum_{q=0}^{\tau-1} \nabla L_2(\theta_{n\tau+q}^{2,j})$$

# Discretized FedCBO System

$$\theta_{(n+1)\tau}^{1,i} \leftarrow \theta_{n\tau}^{1,i} - \lambda_1 \gamma \sum_{q=0}^{\tau-1} \left( \theta_{n\tau+q}^{1,i} - m_{n\tau+q}^1 \right) - \lambda_2 \gamma \sum_{q=0}^{\tau-1} \nabla L_1(\theta_{n\tau+q}^{1,i})$$

$$\theta_{(n+1)\tau}^{2,j} \leftarrow \theta_{n\tau}^{2,j} - \lambda_1 \gamma \sum_{q=0}^{\tau-1} \left( \theta_{n\tau+q}^{2,j} - m_{n\tau+q}^1 \right) - \lambda_2 \gamma \sum_{q=0}^{\tau-1} \nabla L_2(\theta_{n\tau+q}^{2,j})$$

# Splitting Scheme

Step 1:

$$\widehat{\theta}_{n\tau}^{1,i} \leftarrow \theta_{n\tau}^{1,i}, \quad \widehat{\theta}_{n\tau}^{2,j} \leftarrow \theta_{n\tau}^{2,j}$$

Step 2:

$$\widehat{\theta}_{n\tau+q+1}^{1,i} \leftarrow \widehat{\theta}_{n\tau+q}^{1,i} - \lambda_2 \gamma \nabla L_1(\widehat{\theta}_{n\tau+q}^{1,i}), \quad \widehat{\theta}_{n\tau+q+1}^{2,j} \leftarrow \widehat{\theta}_{n\tau+q}^{2,j} - \lambda_2 \gamma \nabla L_2(\widehat{\theta}_{n\tau+q}^{2,j}) \quad \text{for } q = 0, \dots, \tau - 1.$$

Step 3:

$$\theta_{(n+1)\tau}^{1,i} \leftarrow \widehat{\theta}_{(n+1)\tau}^{1,i} - \lambda_1 \gamma \left( \widehat{\theta}_{(n+1)\tau}^{1,i} - m_{(n+1)\tau}^1 \right), \quad \theta_{(n+1)\tau}^{2,j} \leftarrow \widehat{\theta}_{(n+1)\tau}^{2,j} - \lambda_1 \gamma \left( \widehat{\theta}_{(n+1)\tau}^{2,j} - m_{(n+1)\tau}^2 \right)$$

# FedCBO Algorithm

---

## Algorithm 1 FedCBO

---

**Input:** Initialized model  $\theta_0^j \in \mathbb{R}^d, j \in [N]$ ; Number of iterations  $T$ ; Number of local gradient steps  $\tau$ ; Number of models downloaded  $M$ ; CBO system hyperparameters  $\lambda_1, \lambda_2, \alpha$ ; Discretization step size  $\gamma$ ; Initialized sampling likelihood  $P_0 \in \mathbb{R}^{N \times (N-1)}$ ;

- 1: **for**  $n = 0, \dots, T - 1$  **do**
- 2:    $G_n \leftarrow$  random subset of agents (participating devices);
- 3:   **LocalUpdate**( $\theta_n^j, \tau, \lambda_2, \gamma$ ) for  $j \in G_n$ ;
- 4:   **LocalAggregation**(agent  $j$ ) for  $j \in G_n$ ;
- 5: **end for**

**Output:**  $\theta_T^j$  for  $j \in [N]$ .

**LocalUpdate**( $\hat{\theta}_0, \tau, \lambda_2, \gamma$ ) at  $j$ -th agent

- 6: **for**  $q = 0, \dots, \tau - 1$  **do**
  - 7:   (stochastic) gradient descent  $\hat{\theta}_{q+1} \leftarrow \hat{\theta}_q - \lambda_2 \gamma \nabla L_j(\hat{\theta}_q)$ ;
  - 8: **end for**
  - 9: **return**  $\hat{\theta}_\tau$ ;
-



# FedCBO Algorithm

---

**Algorithm 2** LocalAggregation(agent  $j$ )

---

**Input:** Agent  $j$ 's model  $\theta_n^j \in \mathbb{R}^d$ ; Participating devices at  $n$  iteration  $G_n$ ; Sampling likelihood  $P_n^j \in \mathbb{R}^{N-1}$ ;  
CBO system hyperparameters  $\lambda_1, \alpha$ ; Discretization step size  $\gamma$ ; Random sample proportion  $\varepsilon \in (0, 1)$ ;  
Number of models downloaded  $M$ ;

1:  $A_n \leftarrow \varepsilon\text{-greedySampling}(P_n^j, G_n, M)$ ;

2: Agent  $j$  downloads models  $\theta_n^i$  for  $i \in A_n$ ;

3: Evaluate models  $\theta_n^i$  on agent  $j$ 's data set respectively and denote the corresponding loss as  $L_j^i$ ;

4: Calculate consensus point  $m_j$  by

$$(19) \quad m_j \leftarrow \frac{1}{\sum_{i \in A_n} \mu_j^i} \sum_{i \in A_n} \theta_n^i \mu_j^i, \quad \text{with } \mu_j^i = \exp(-\alpha L_j^i)$$

5: Update agent  $j$ 's model by

$$(20) \quad \theta_{n+1}^j \leftarrow \theta_n^j - \lambda_1 \gamma (\theta_n^j - m_j),$$

6: Update sampling likelihood  $P_n^j$  by

$$(21) \quad P_{n+1}^{j,i} \leftarrow P_n^{j,i} + (L_j^j - L_j^i), \quad \text{for } i \in A_n$$

**Output:**  $\theta_{n+1}^j, P_{n+1}^j$

$\varepsilon\text{-greedySampling}(P_n^j, G_n, M)$

7: Randomly sample  $\varepsilon * M$  number of agents from  $G_n$ , denoted as  $A_n^1$ ;

8: Select  $(1 - \varepsilon) * M$  numbers of agents in  $G_n \setminus A_n^1$  with top value  $P_j^{j,i}, i \in G_n \setminus A_n^1$ , denoted as  $A_n^2$ ;

9: **return**  $A_n = A_n^1 \cup A_n^2$

---