



Diffusion Federated Dataset



¹Seok-Ju Hahn and ²Junghye Lee

¹Argonne National Laboratory

²Seoul National University

TL; DR

- Synthetic Data from Federated Inference of Pre-trained Diffusion Models

TL; DR

~~Learning~~

- Synthetic Data from Federated Inference of Pre-trained Diffusion Models
 - No parameter exchange
 - ✓ Other signals needs to be exchanged

TL; DR

- Synthetic Data from Federated Inference of ~~Learning~~ Pre-trained Diffusion Models

TL; DR

~~Learning~~

- Synthetic Data from Federated Inference of **Pre-trained** Diffusion Models
 - No local **update** is required during communication
 - ✓ Local **inference** is required instead

TL; DR

~~Learning~~

- **Synthetic Data** from Federated Inference of **Pre-trained** Diffusion Models
 - Synthetic data in FL be from **a mixture** of local distributions

TL; DR

~~Learning~~

- **Synthetic Data** from Federated Inference of **Pre-trained** Diffusion Models
 - Synthetic data in FL be from **a mixture** of local distributions

$$p^*(\boldsymbol{x}) = \sum_{i=1}^K w_i p_i(\boldsymbol{x})$$

TL; DR

- **Synthetic Data** from Federated Inference of **Pre-trained** Diffusion Models
 - Synthetic data in FL be from **a mixture** of local distributions

$$p^*(\mathbf{x}) = \sum_{i=1}^K w_i p_i(\mathbf{x})$$

mixing coefficient
(unknown)

TL; DR

- **Synthetic Data** from Federated Inference of **Pre-trained** Diffusion Models
 - Synthetic data in FL be from **a mixture** of local distributions

$$p^*(\mathbf{x}) = \sum_{i=1}^K \textcircled{w_i} p_i(\mathbf{x})$$

Usually set as
 $w_i \propto n_i$ or $w_i = \frac{1}{K}$

Scope

- Cross-silo FL setting

Scope

- Cross-silo FL setting
 - Where a moderate number of reliable, and known, and addressable clients participate^[1]

Scope

- Cross-silo FL setting
 - Where a moderate number of reliable, and known, and addressable clients participate^[1]
 - Where each client has sufficient computing power

Scope

- Cross-silo FL setting
 - Where a moderate number of reliable, and known, and addressable clients participate^[1]
 - Where each client has sufficient computing power
 - Where each client suffers from low diversity of data or low sample size^[2]

[1] Advances and Open Problems in Federated Learning (Kairouz et al., 2019)

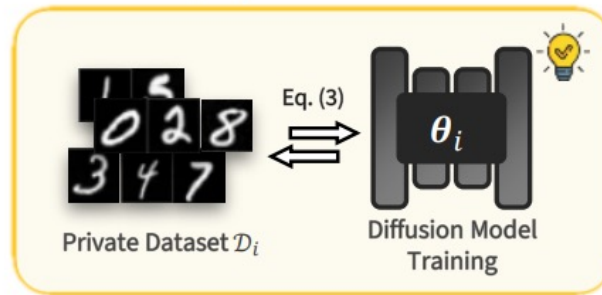
[2] Flamby: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings (du Terrail et al., 2022)

Overview

- Synthetic Data from Federated Inference of Pre-trained Diffusion Models

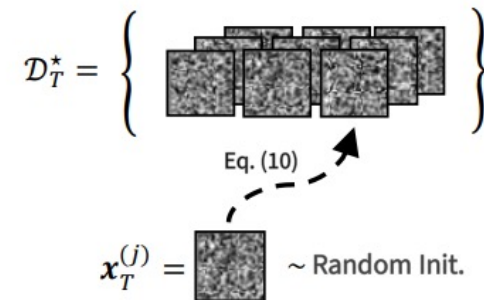
Ⓐ Preparation of diffusion models

Client $i = 1, \dots, K$



Ⓑ Initialization of synthetic data

Server



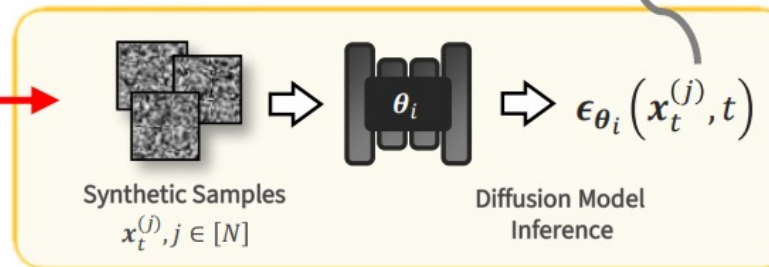
Ⓒ Iterative refinement via cooperative inference

Server repeats $t = T, \dots, 1$

$$x_t^{(j)} + \eta_t \nabla_{x_t^{(j)}} \log p^*(x_t^{(j)}) + \sqrt{2\eta_t} z_t \xrightarrow{\text{Eq. (7)}} x_{t-1}^{(j)}$$



Client $i = 1, \dots, K$

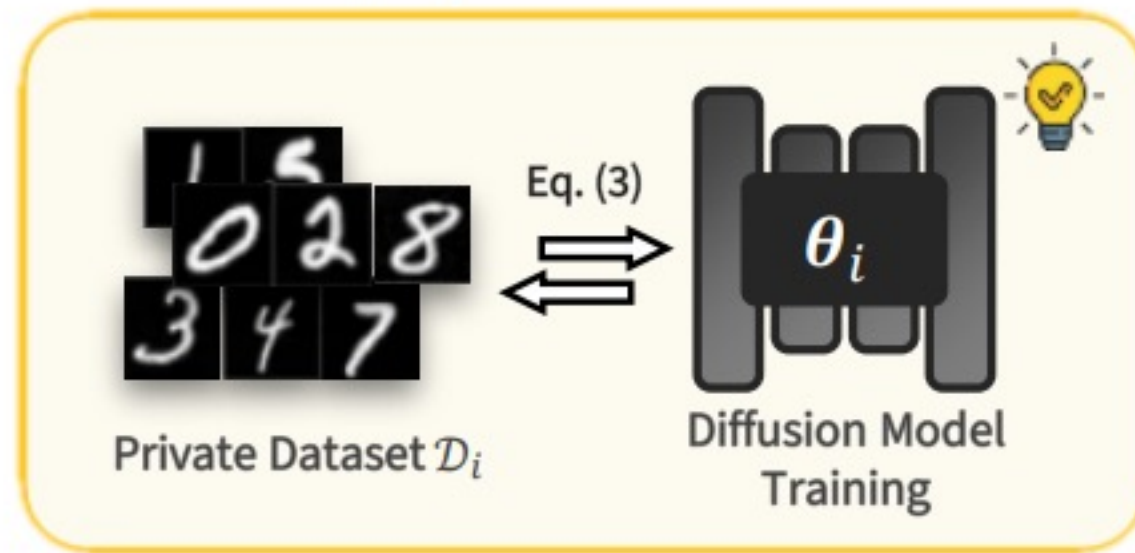


Breakdown

- A) Each client trains its own local diffusion model

① Preparation of diffusion models

Client $i = 1, \dots, K$



Breakdown

- B) Server randomly initializes synthetic dataset

⑤ Initialization of synthetic data

Server

$$\mathcal{D}_T^* = \left\{ \begin{array}{c} \text{[Stack of 5 noisy images]} \end{array} \right\}$$

Eq. (10)

$$\mathbf{x}_T^{(j)} = \begin{array}{c} \text{[Noisy image]} \end{array} \sim \text{Random Init.}$$

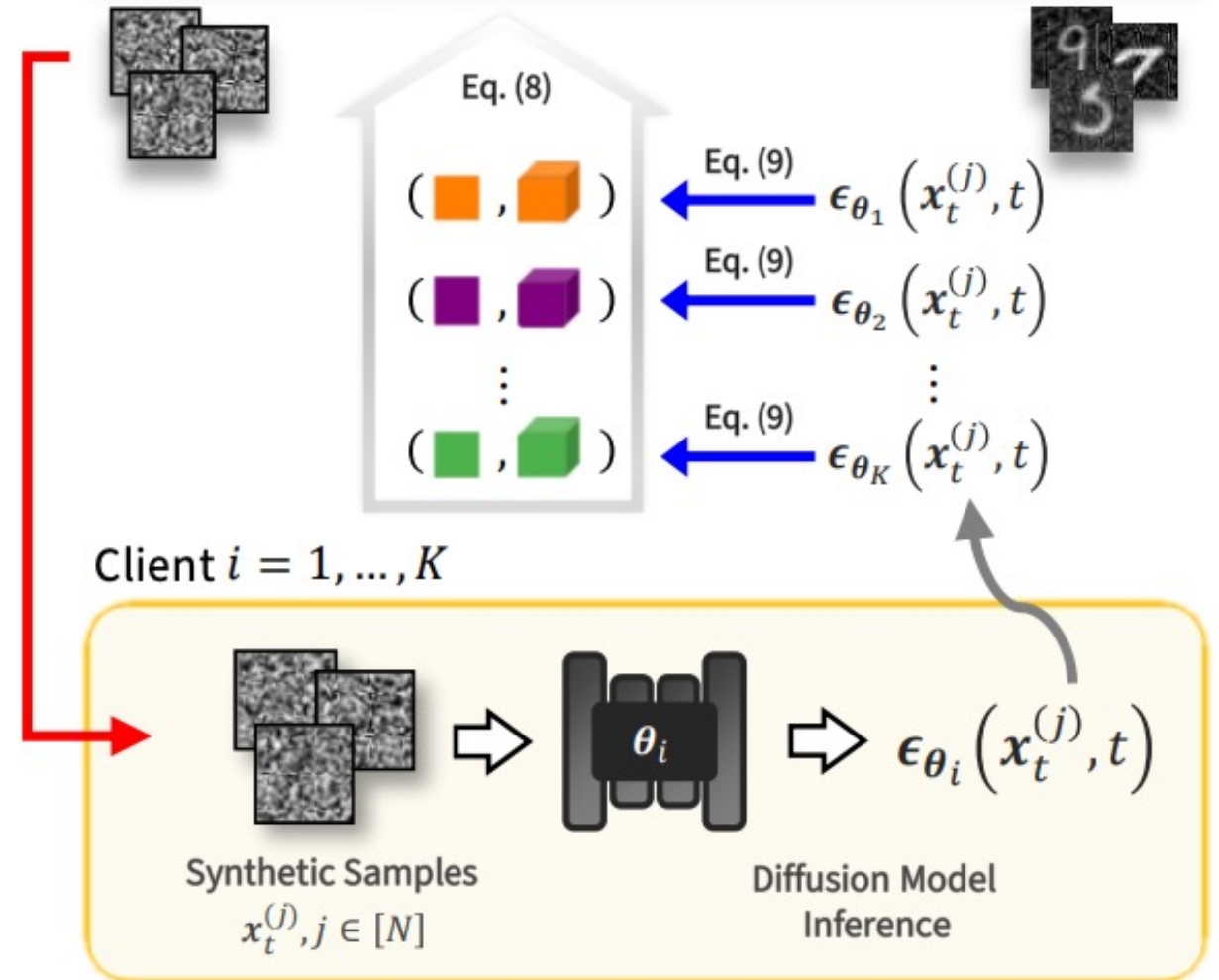
Breakdown

- C) Cooperative sampling of the synthetic dataset
- Each local diffusion model represents a local distribution

© Iterative refinement via cooperative inference

Server repeats $t = T, \dots, 1$

$$x_t^{(j)} + \eta_t \nabla_{x_t^{(j)}} \log p^*(x_t^{(j)}) + \sqrt{2\eta_t} z_t \xrightarrow{\text{Eq. (7)}} x_{t-1}^{(j)}$$



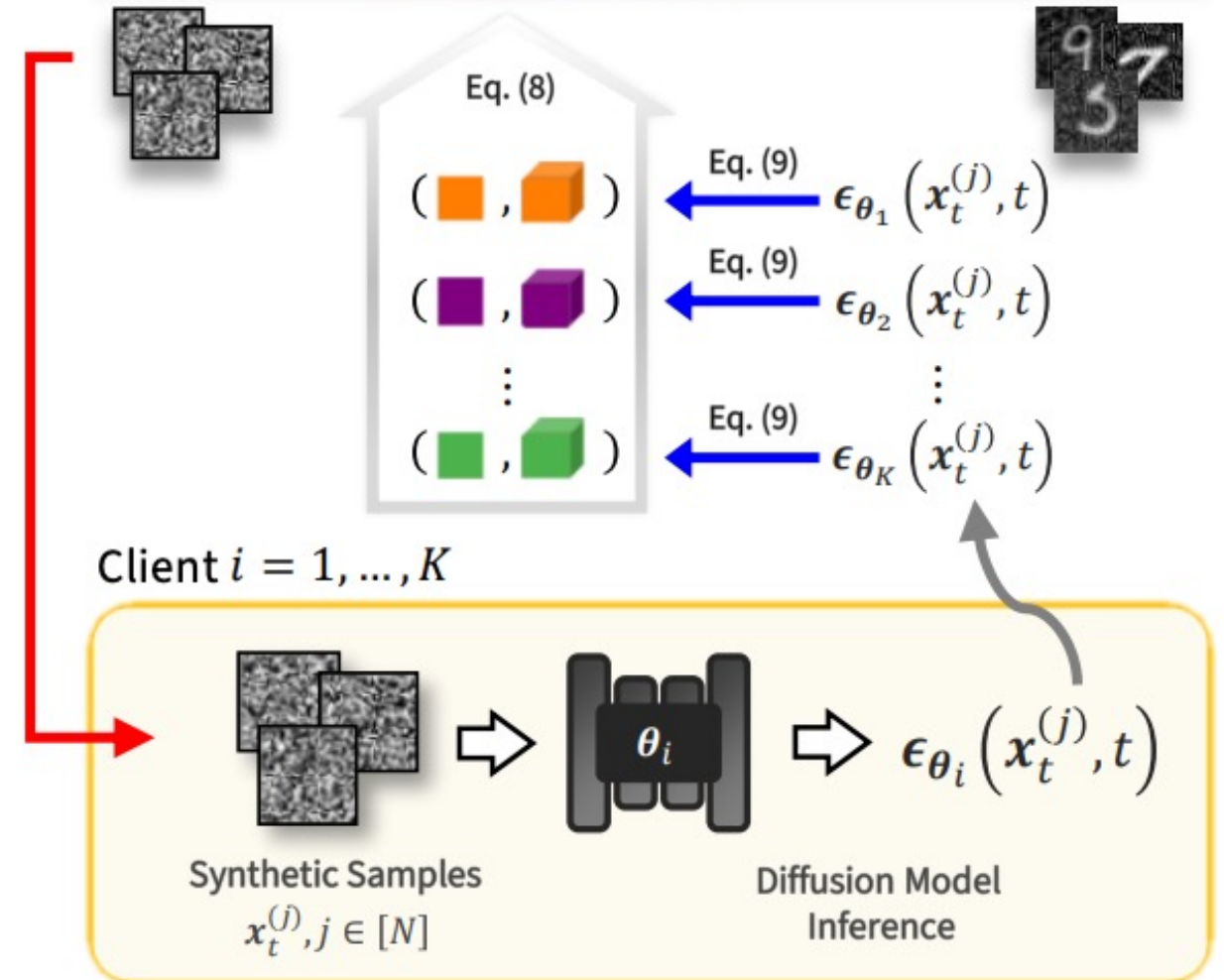
Breakdown

- Eq. (7) (Black box itself)
 - Unadjusted Langevin Algorithm (ULA^[3]) (first-order MCMC sampler)

© Iterative refinement via cooperative inference

Server repeats $t = T, \dots, 1$

$$x_t^{(j)} + \eta_t \nabla_{x_t^{(j)}} \log p^*(x_t^{(j)}) + \sqrt{2\eta_t} z_t \xrightarrow{\text{Eq. (7)}} x_{t-1}^{(j)}$$



Breakdown

- Eq. (8) (grad. of log-prob. w.r.t. \mathbf{x})

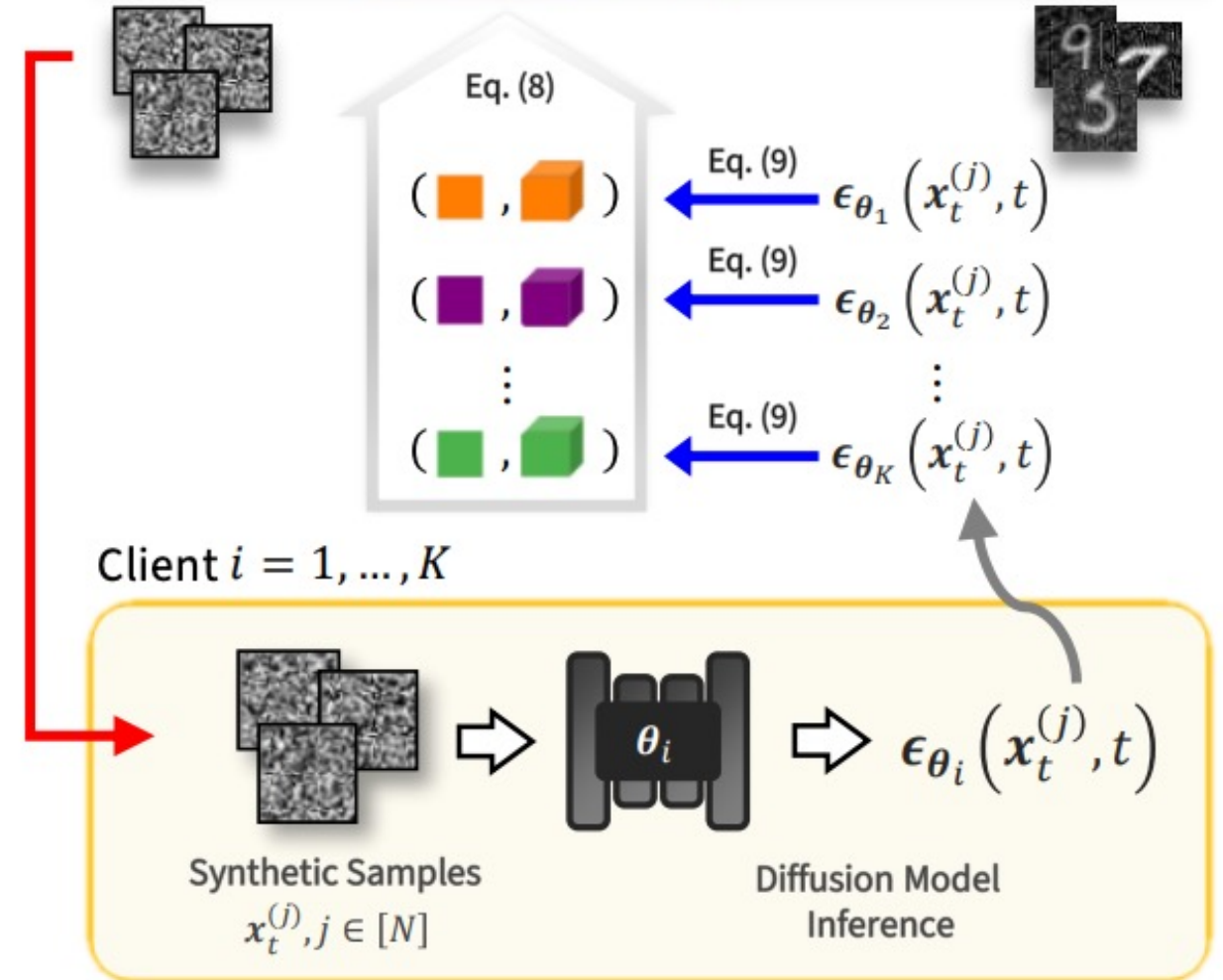
$$\nabla_{\mathbf{x}} \log p^*(\mathbf{x}) = \sum_{i=1}^K \tilde{w}_i \nabla_{\mathbf{x}} \log p_{\theta_i}(\mathbf{x})$$

$$\tilde{w}_i = \frac{w_i \exp(-\lambda f_{\theta_i}(\mathbf{x}))}{\sum_{j=1}^K w_j \exp(-\lambda f_{\theta_j}(\mathbf{x}))}$$

© Iterative refinement via cooperative inference

Server repeats $t = T, \dots, 1$

$$\mathbf{x}_t^{(j)} + \eta_t \nabla_{\mathbf{x}_t^{(j)}} \log p^*(\mathbf{x}_t^{(j)}) + \sqrt{2\eta_t} \mathbf{z}_t \xrightarrow{\text{Eq. (7)}} \mathbf{x}_{t-1}^{(j)}$$



Breakdown

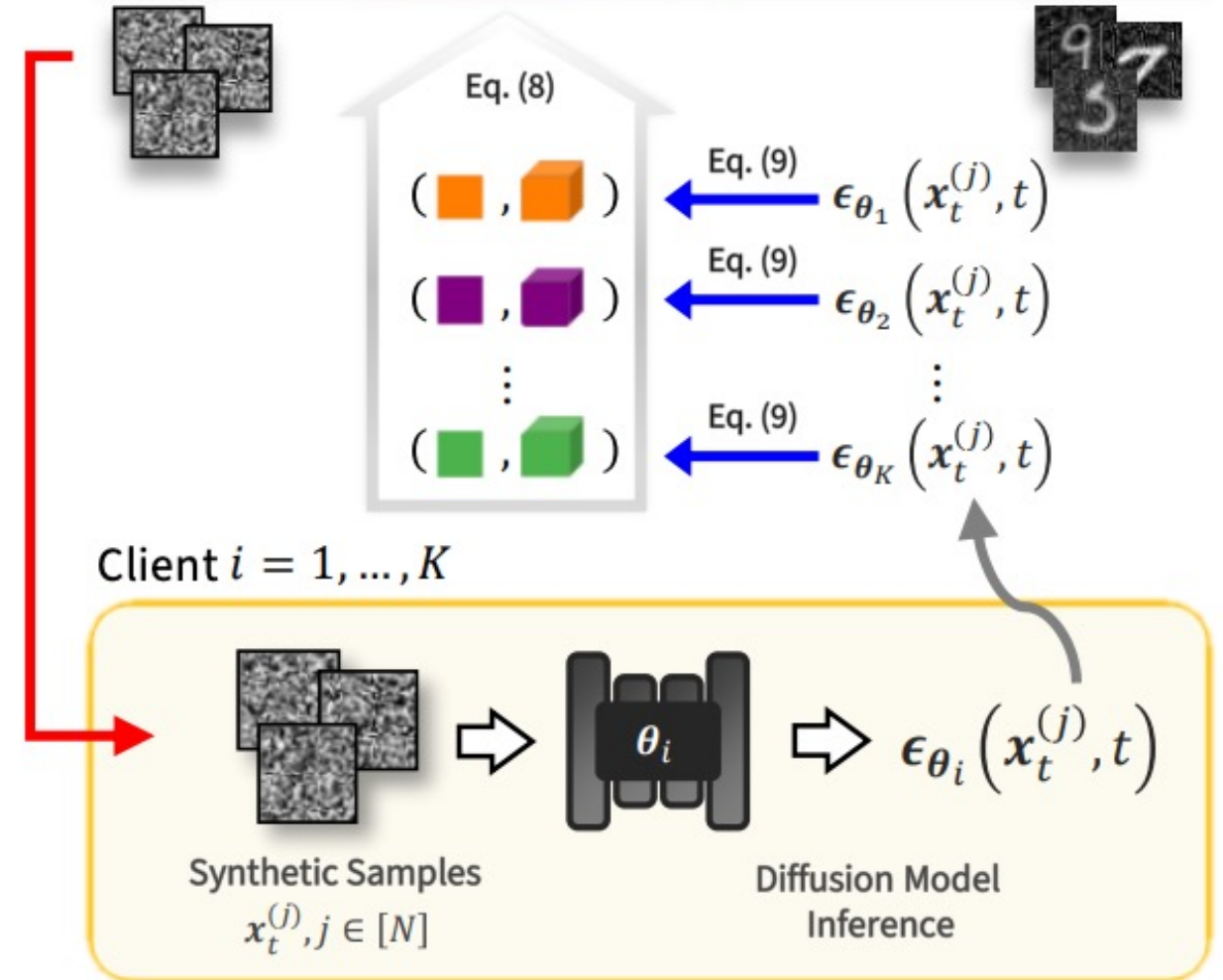
- Eq. (9) (output to **score** conversion)

$$\nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t, t) \approx -\frac{\lambda}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{x}_t, t)$$

© Iterative refinement via cooperative inference

Server repeats $t = T, \dots, 1$

$$\mathbf{x}_t^{(j)} + \eta_t \nabla_{\mathbf{x}_t^{(j)}} \log p^*(\mathbf{x}_t^{(j)}) + \sqrt{2\eta_t} \mathbf{z}_t \xrightarrow{\text{Eq. (7)}} \mathbf{x}_{t-1}^{(j)}$$



Wrap-up

- Data-centric federated synthetic data generation framework
- ... as a cooperative inference of pre-trained diffusion models
- Thanks to their connection to energy-based models (EBMs)
- Non-asymptotic convergence guarantee
- Easy plug-and-play of differential privacy guarantee
- Lighter in communication compared to parameter averaging