

*Inria*

# Federated Learning with Packet Losses

A. Rodio, G. Neglia, F. Busacca, S. Mangione, S. Palazzo, F. Restuccia, I. Tinnirello

Tampa, 19 November 2023



Università  
di Catania

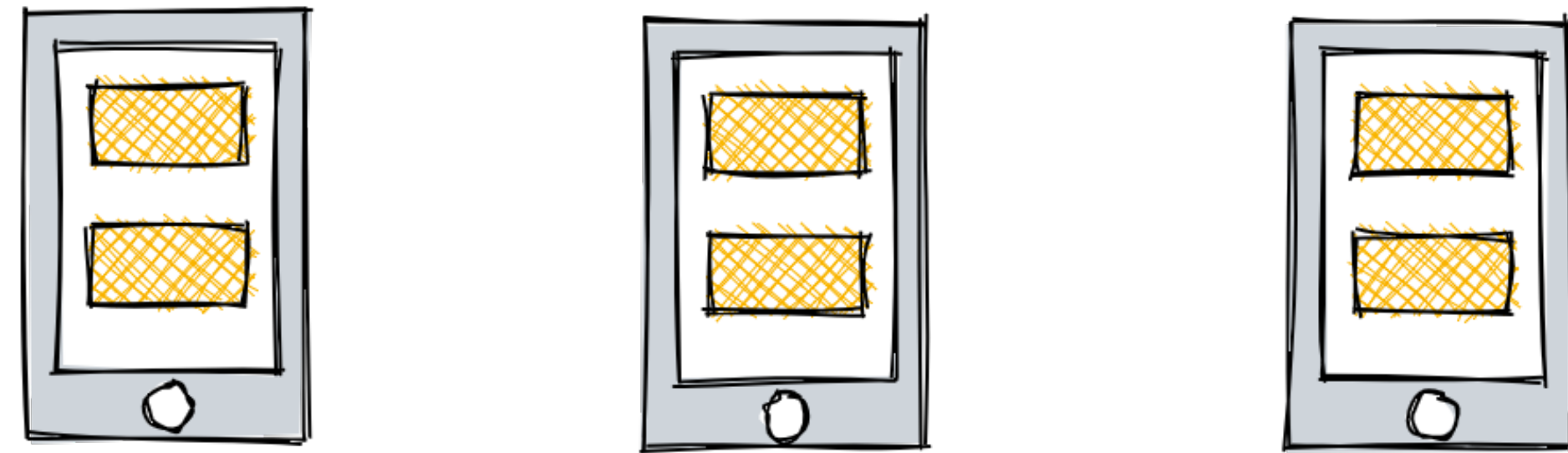


Università  
degli Studi  
di Palermo



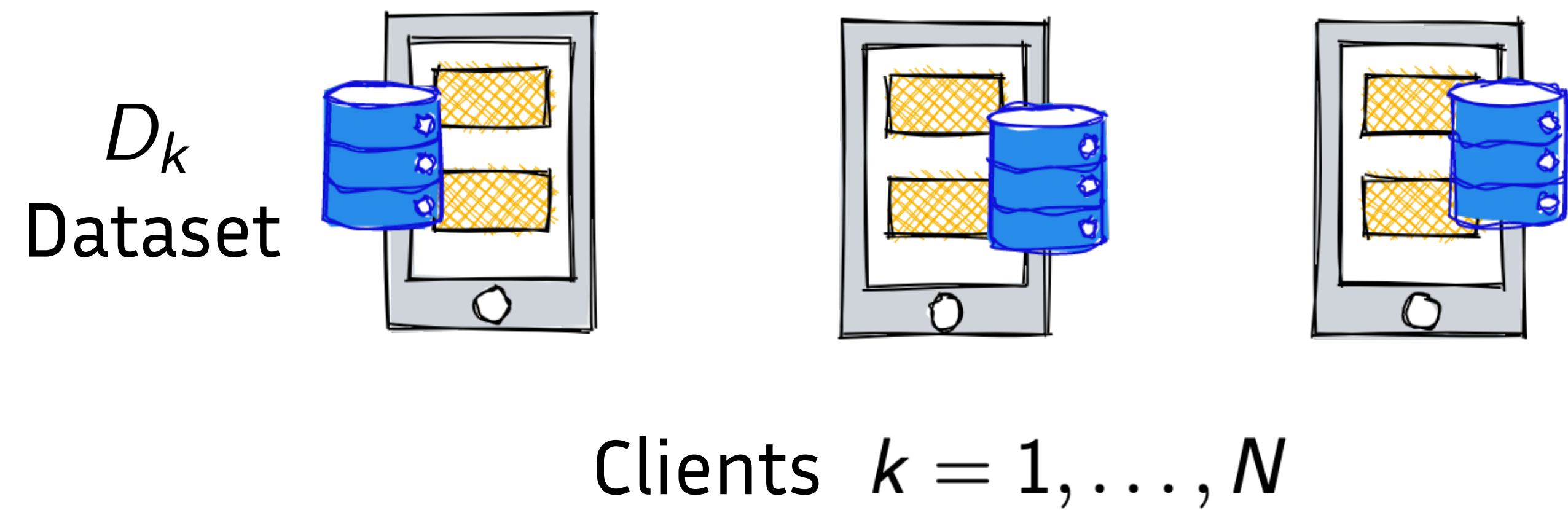
Northeastern  
University

# Centralized Learning

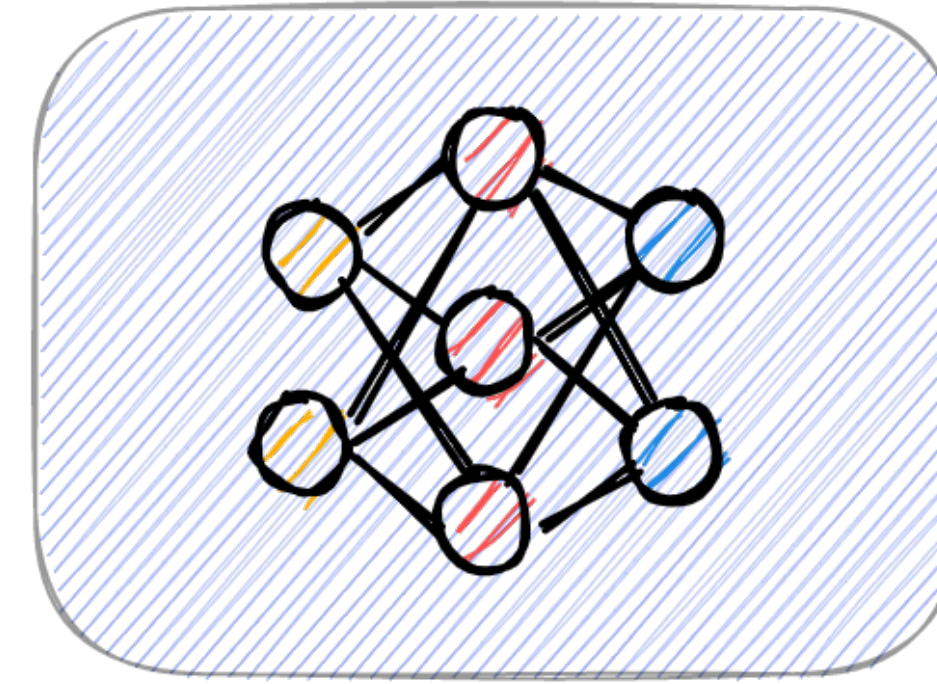


Clients  $k = 1, \dots, N$

# Centralized Learning

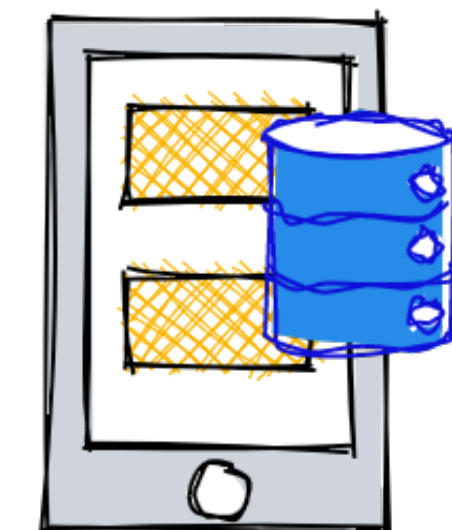
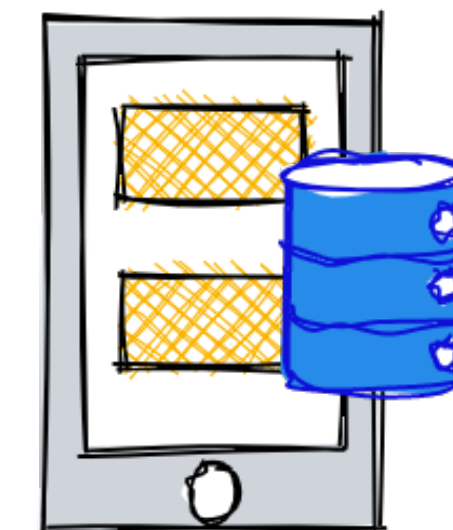
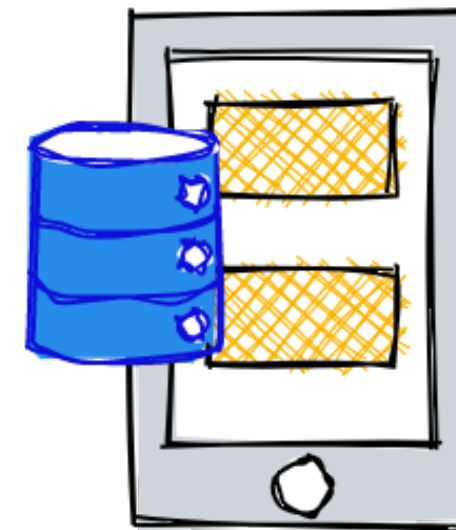


# Centralized Learning



Global model  
 $\mathbf{w} \in \mathbb{R}^d$

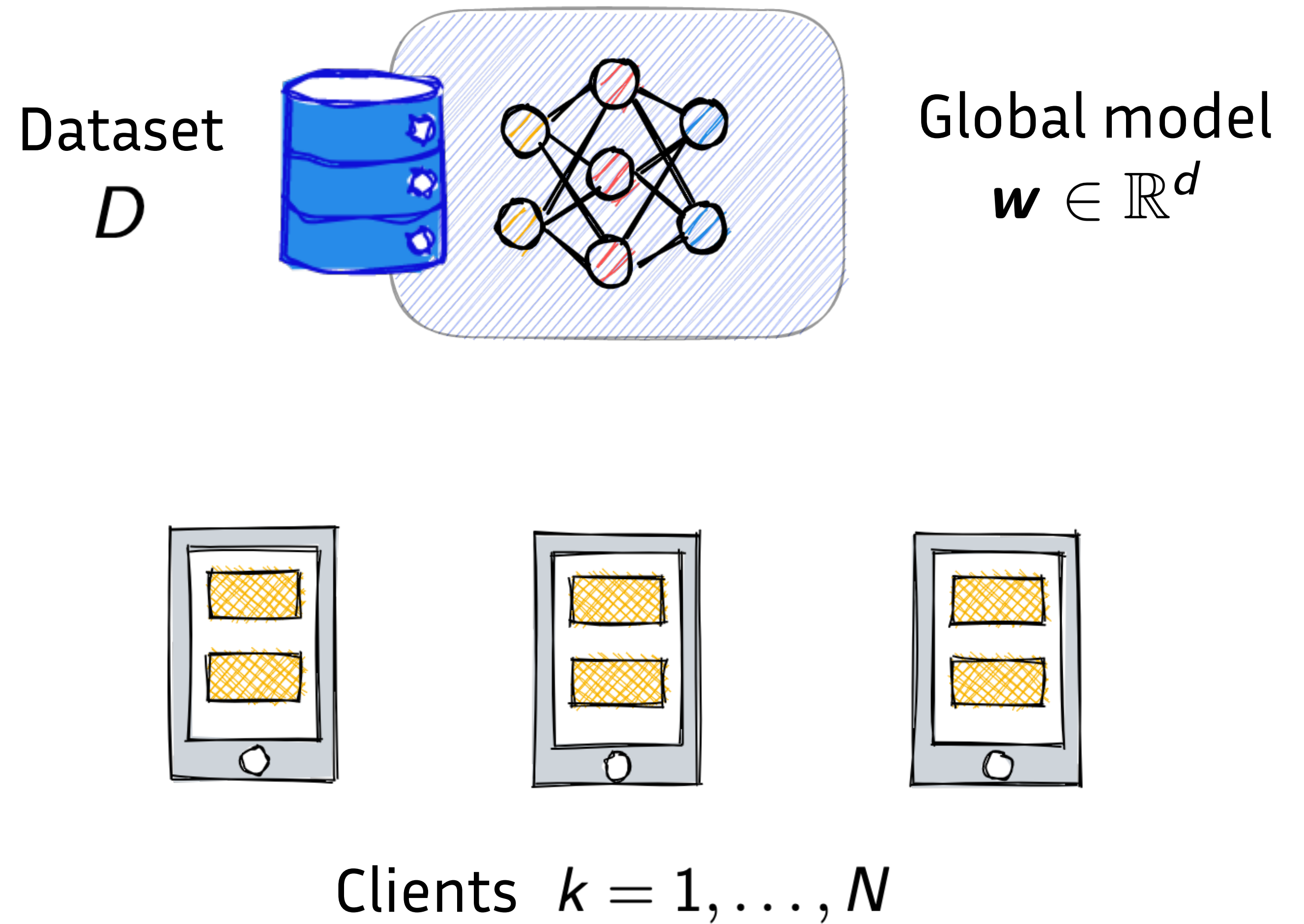
$D_k$   
Dataset



Clients  $k = 1, \dots, N$



# Centralized Learning

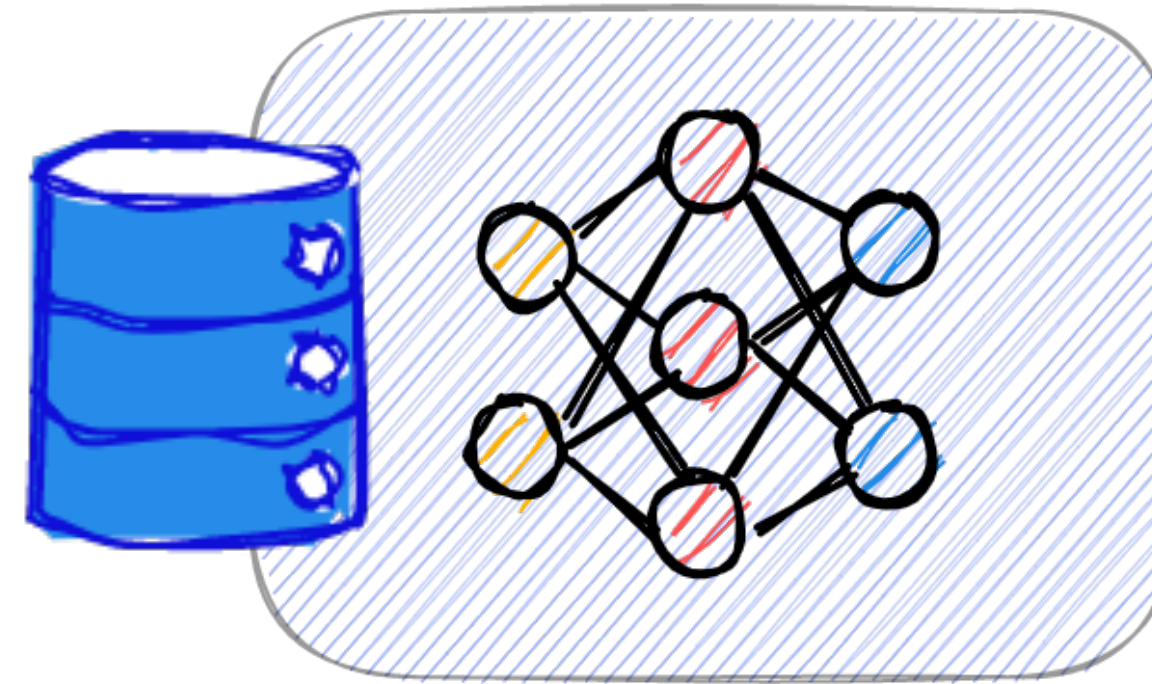


# Centralized Learning

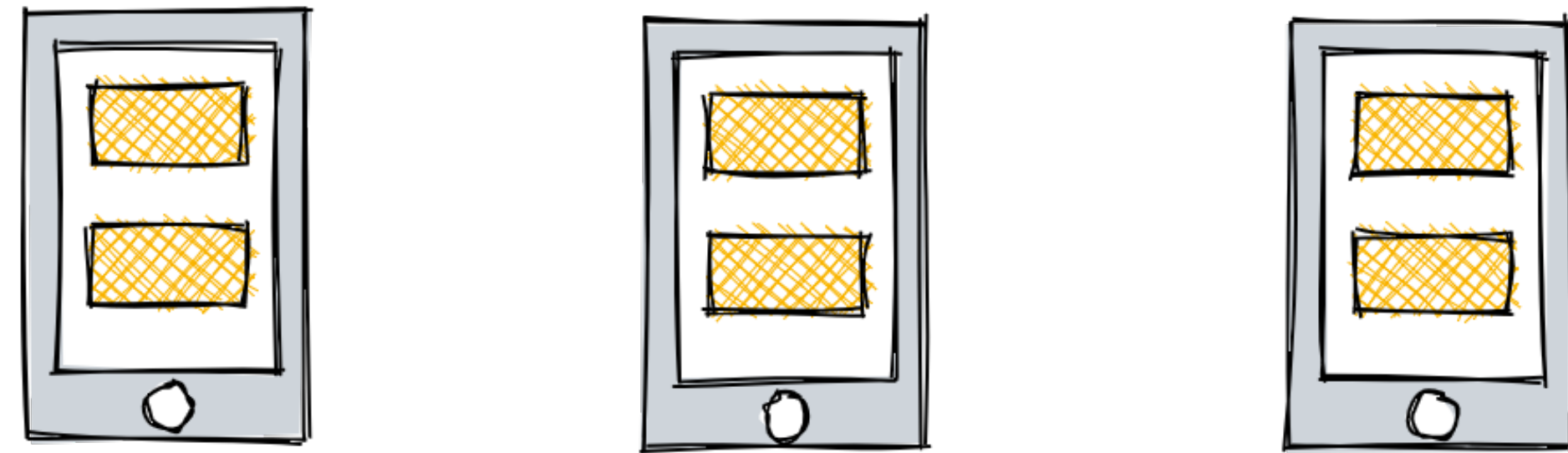
Solve the optimization problem

$$\min_{\mathbf{w}} \frac{1}{|D|} \sum_{d \in D} \ell(\mathbf{w}, d)$$

Dataset  
 $D$



Global model  
 $\mathbf{w} \in \mathbb{R}^d$



Clients  $k = 1, \dots, N$



# Centralized Learning

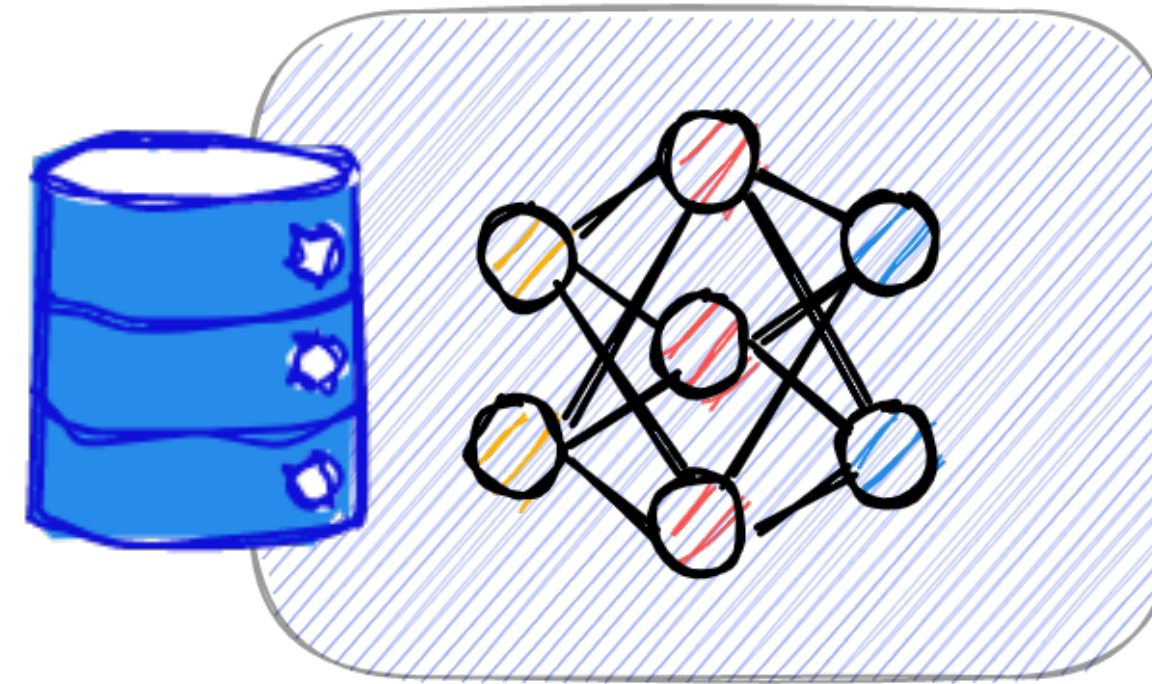
Solve the optimization problem

$$\min_{\mathbf{w}} \frac{1}{|D|} \sum_{d \in D} \ell(\mathbf{w}, d)$$

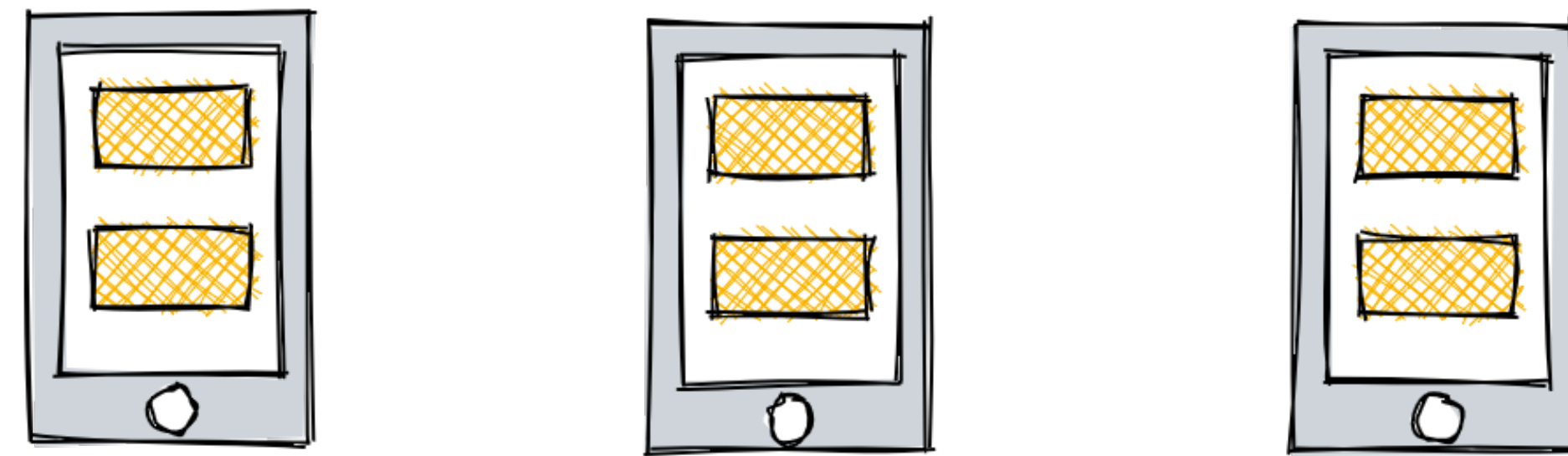
Data transfer

1. Communication cost
2. Privacy

Dataset  
 $D$



Global model  
 $\mathbf{w} \in \mathbb{R}^d$



Clients  $k = 1, \dots, N$

# Federated Learning

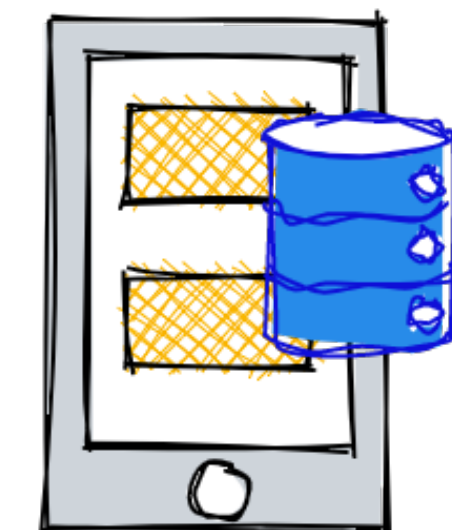
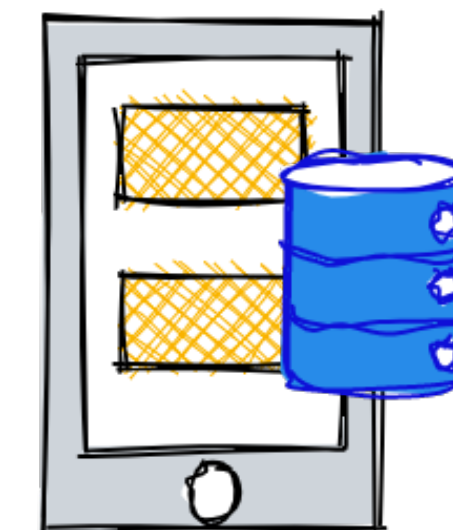
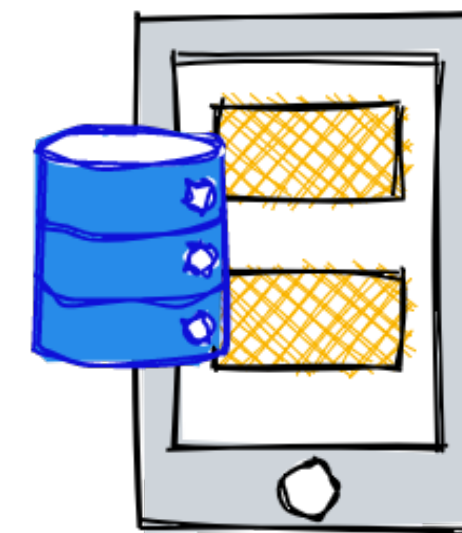
Solve the optimization problem

$$\min_{\mathbf{w}} F(\mathbf{w}) = \frac{1}{N} \sum_{k=1}^N F_k(\mathbf{w})$$

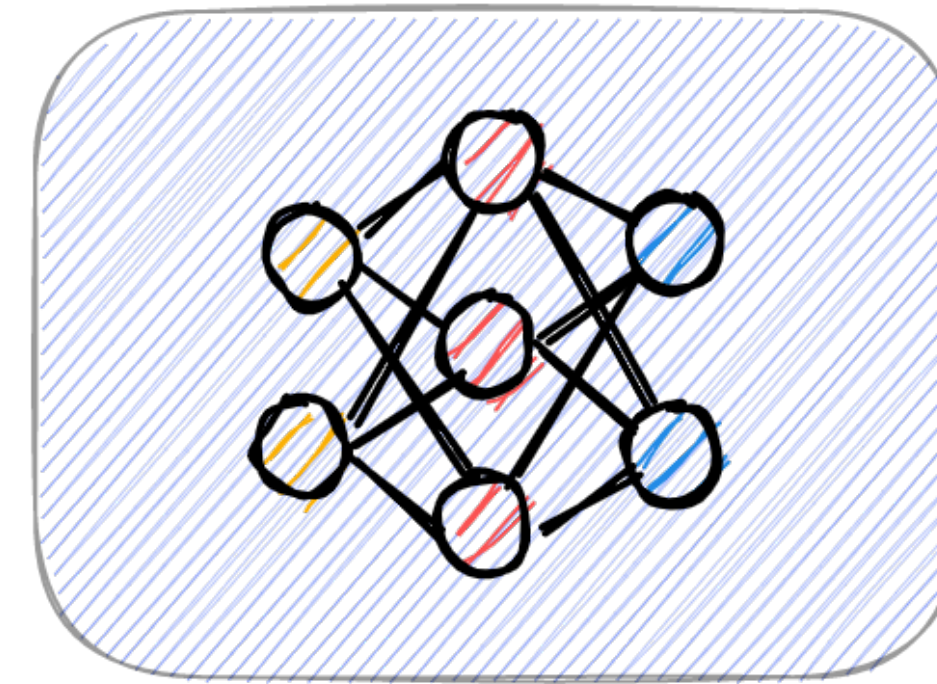
where

$$F_k(\mathbf{w}) = \frac{1}{|D_k|} \sum_{d_k \in D} \ell(\mathbf{w}, d_k)$$

$D_k$   
Dataset



Clients  $k = 1, \dots, N$



Global model

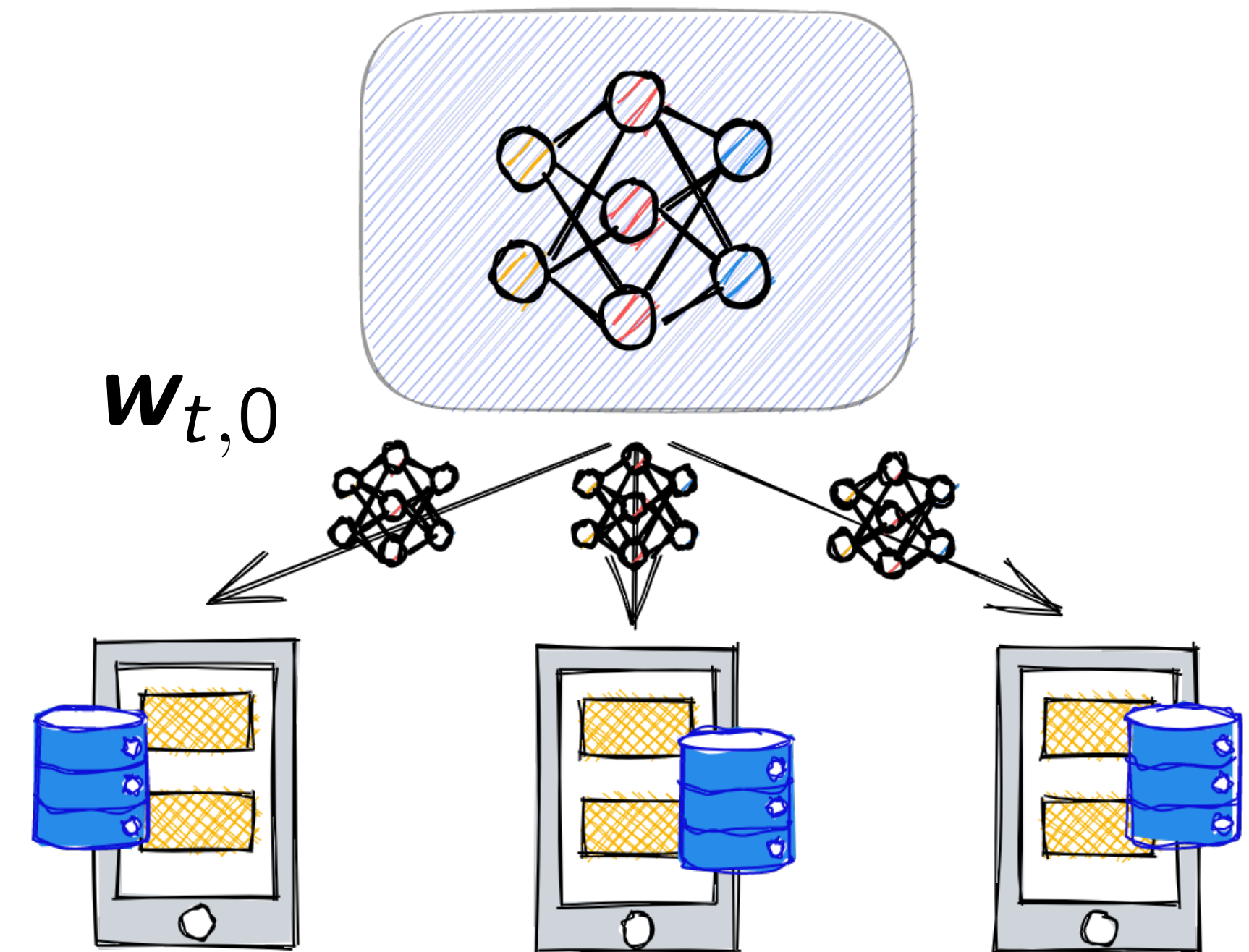
$$\mathbf{w} \in \mathbb{R}^d$$



# Federated Learning

**for**  $t \in \{1, \dots, T\}$  **do:**

(1) Server broadcasts the initial model



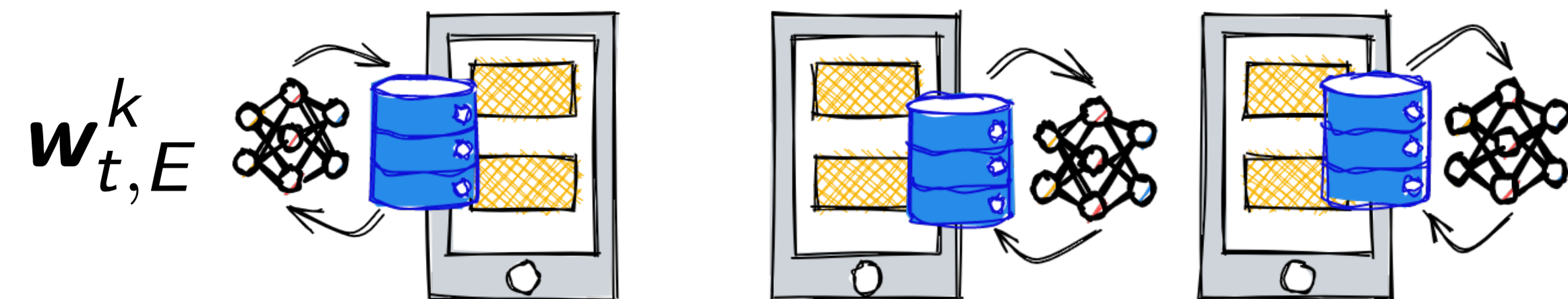
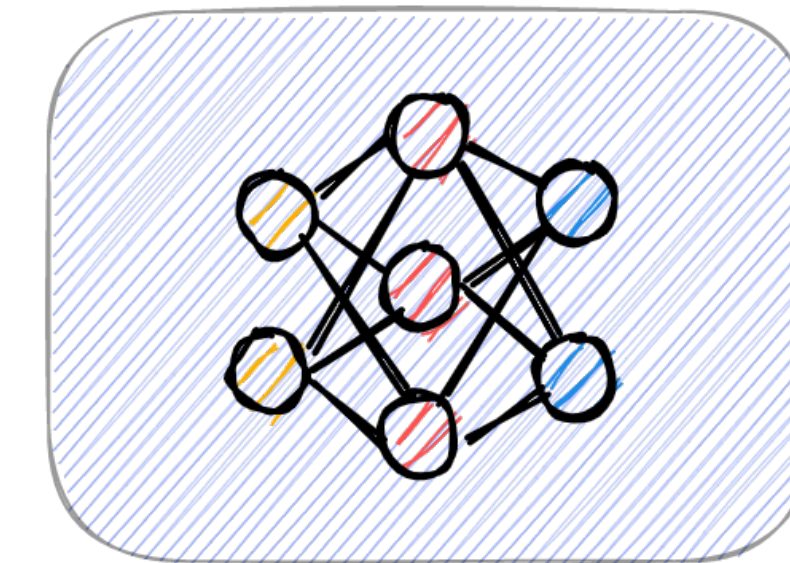
# Federated Learning

**for**  $t \in \{1, \dots, T\}$  **do:**

(2) Each client updates its local model

**for**  $j = 0, \dots, E - 1$  **do :**

$$\mathbf{w}_{t,j+1}^k = \mathbf{w}_{t,j}^k - \eta_t \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k)$$



# Federated Learning

for  $t \in \{1, \dots, T\}$  do:

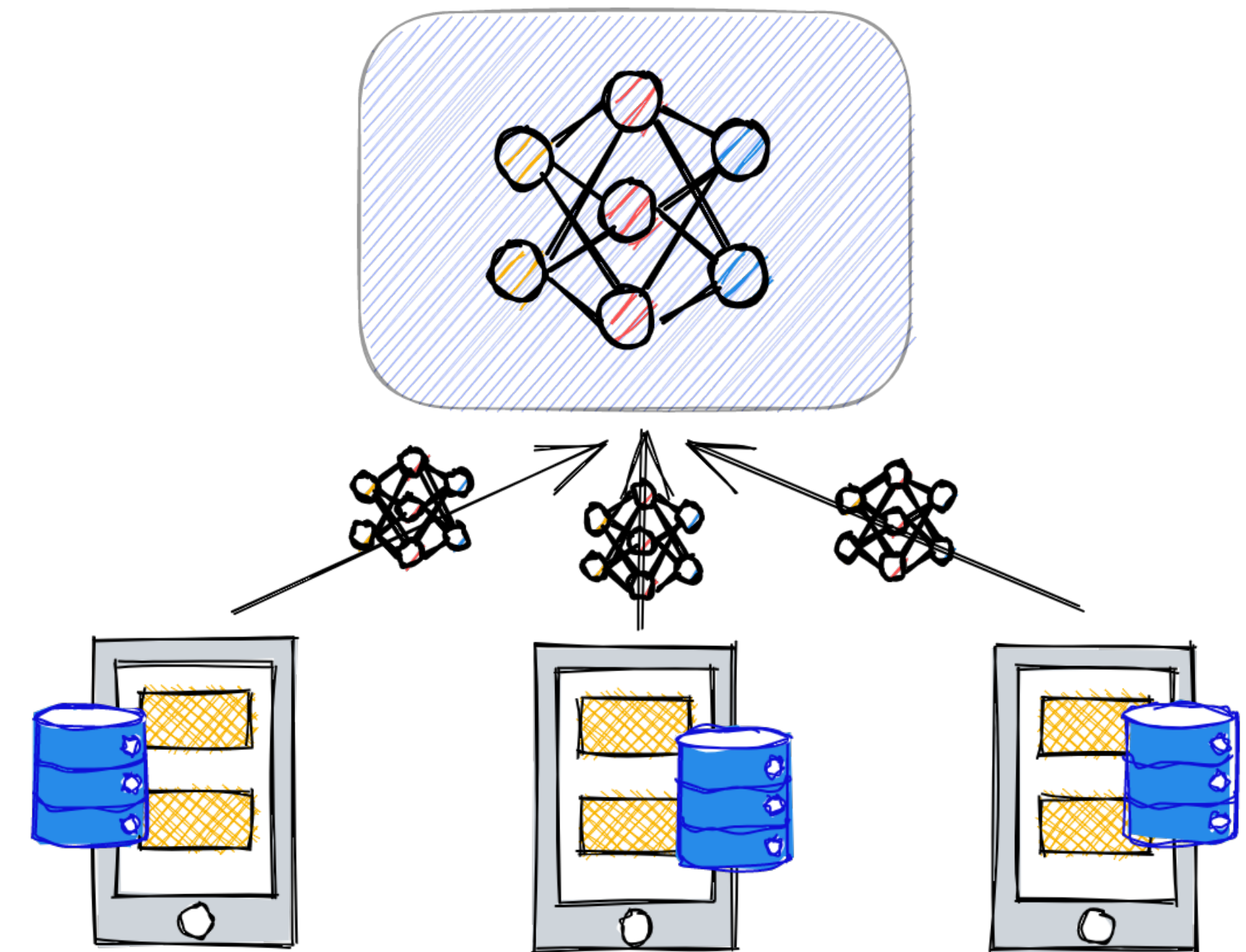
(3) Each client transmits

$$\mathbf{w}_{t,E}^k = \mathbf{w}_t - \underbrace{\eta_t \sum_{j=0}^{E-1} \nabla F_k(\mathbf{w}_{t,j}^k, \mathcal{B}_{t,j}^k)}_{\Delta_t^k}$$

OR

Send model  $\uparrow$

Send pseudo-gradient  $\uparrow$





# Federated Learning

for  $t \in \{1, \dots, T\}$  do:

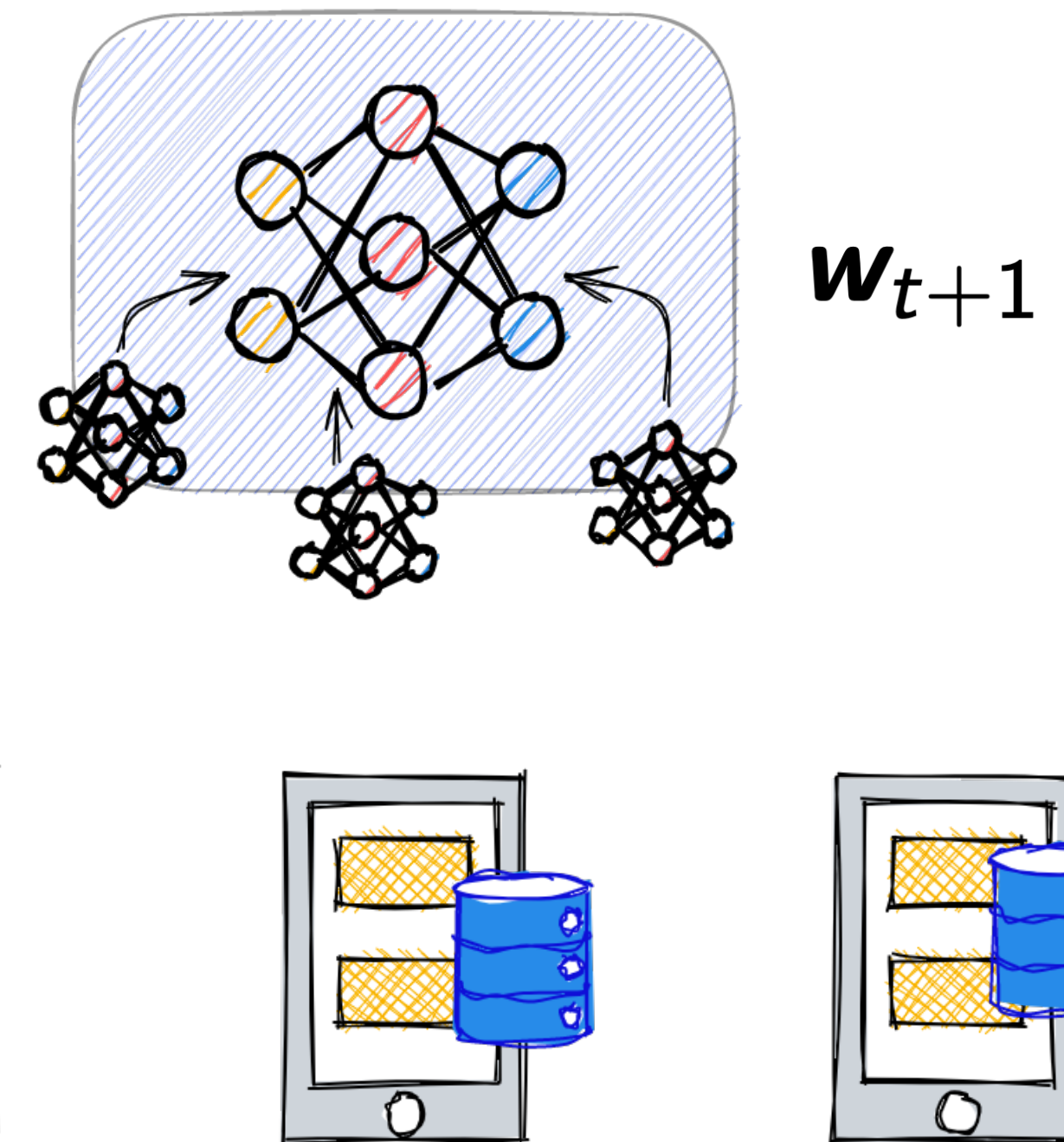
(4) Server aggregates

$$\mathbf{w}_{t+1}^{\text{DMA}} = \frac{1}{N} \sum_{k=1}^N \mathbf{w}_{t,E}^k$$

↑  
**Direct Model  
Aggregation (DMA)**

$$\mathbf{w}_{t+1}^{\text{PGA}} = \mathbf{w}_t + \frac{1}{N} \sum_{k=1}^N \Delta_t^k$$

OR ↑  
**Pseudo-Gradient  
Aggregation (PGA)**



# Federated Learning

for  $t \in \{1, \dots, T\}$  do:

(4) Server aggregates

$$\mathbf{w}_{t+1}^{\text{DMA}} = \frac{1}{N} \sum_{k=1}^N \mathbf{w}_{t,E}^k$$

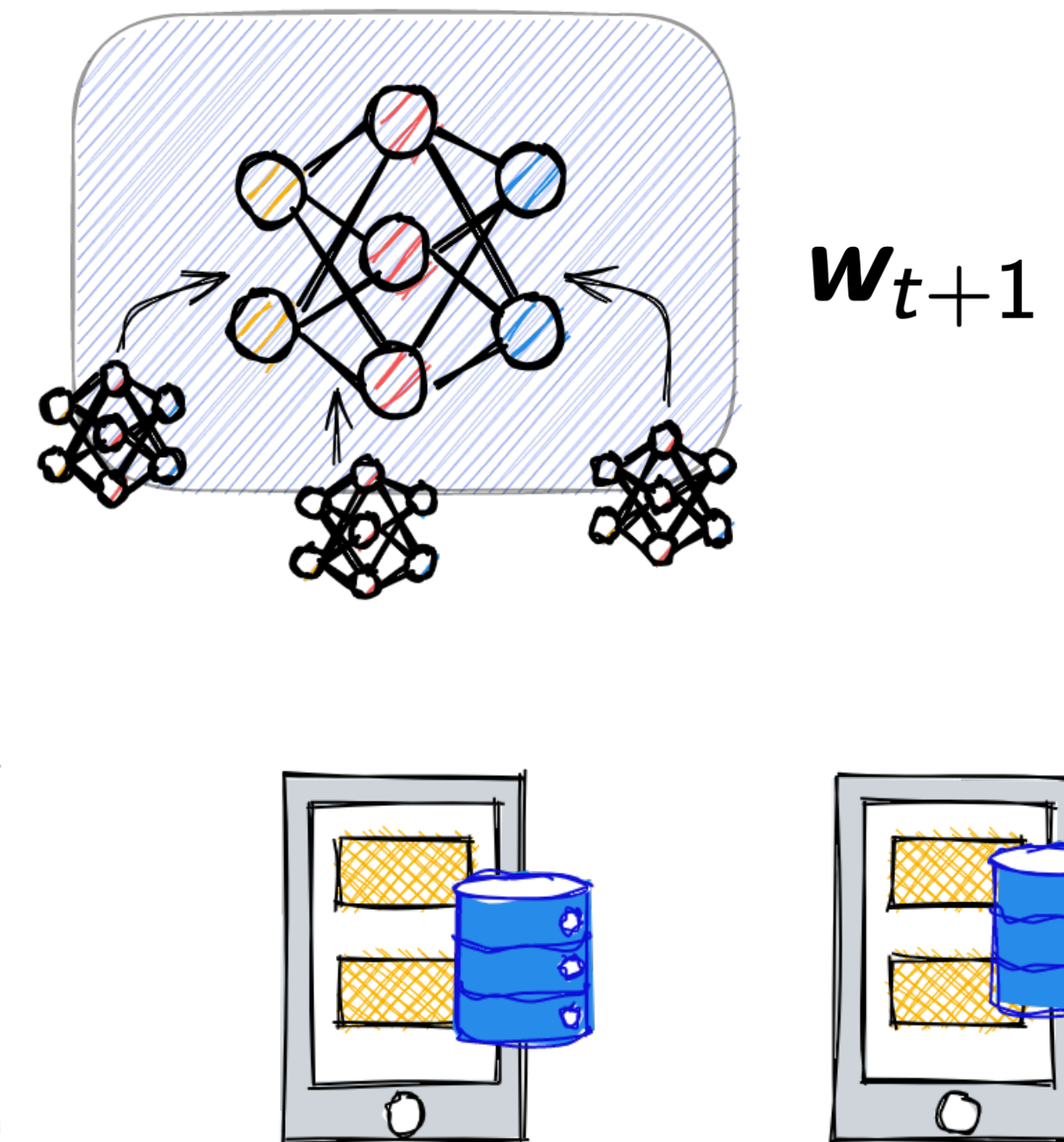
↑  
**Direct Model  
Aggregation (DMA)**

OR

$$\mathbf{w}_{t+1}^{\text{PGA}} = \mathbf{w}_t + \frac{1}{N} \sum_{k=1}^N \Delta_t^k$$

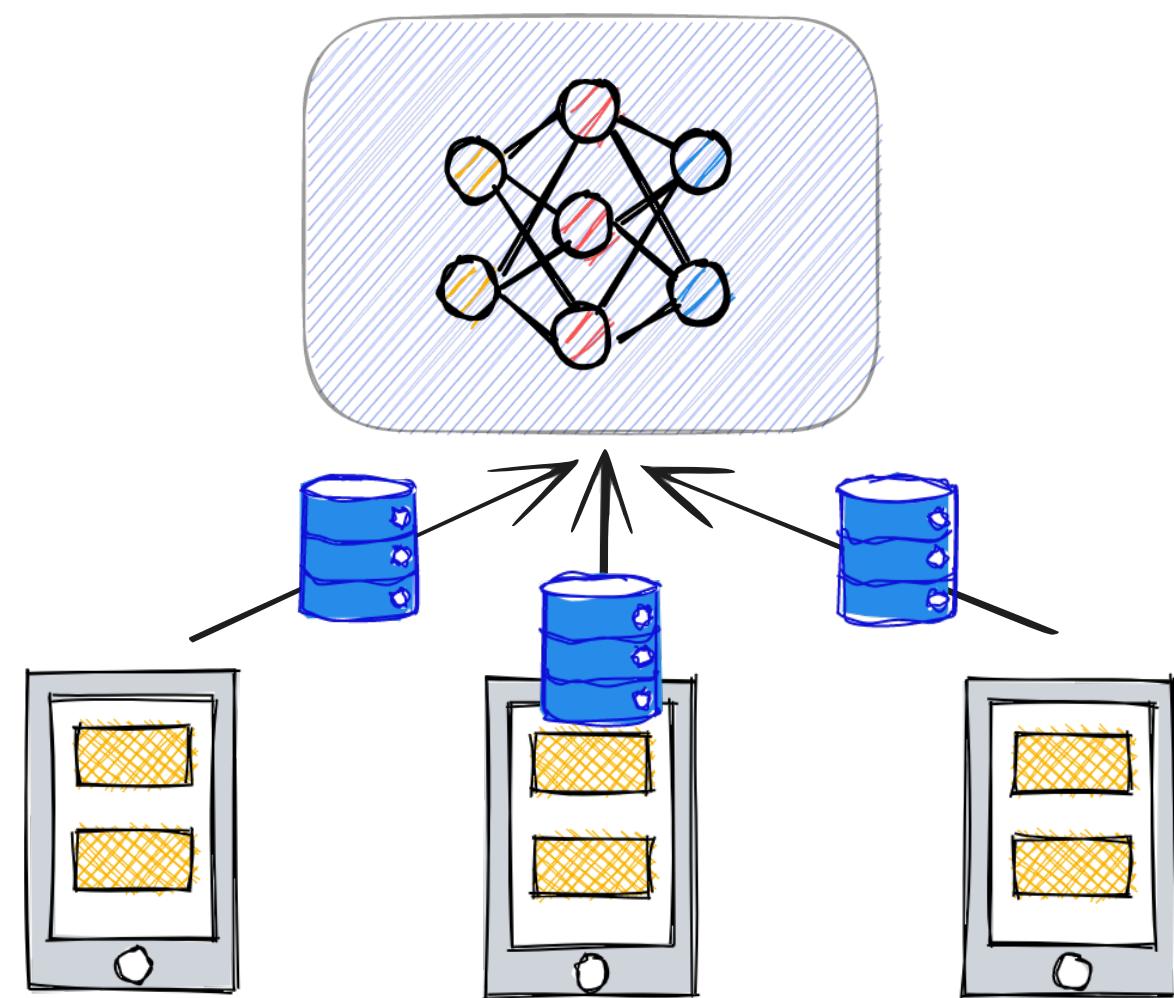
↑  
**Pseudo-Gradient  
Aggregation (PGA)**

(In lossless scenarios: DMA = PGA)



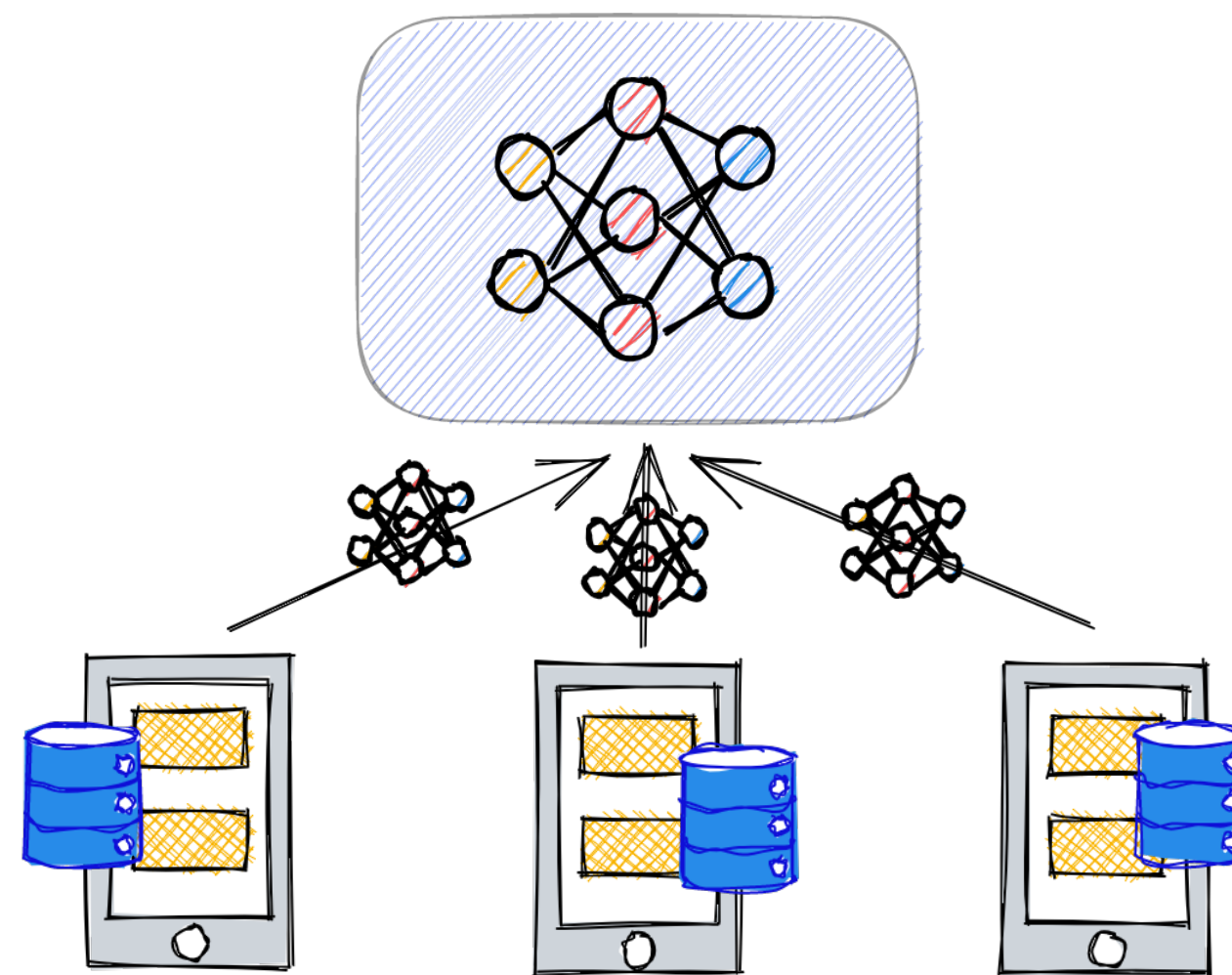
# Centralized vs Federated

## Centralized



Share data

## Federated



Share models / gradients

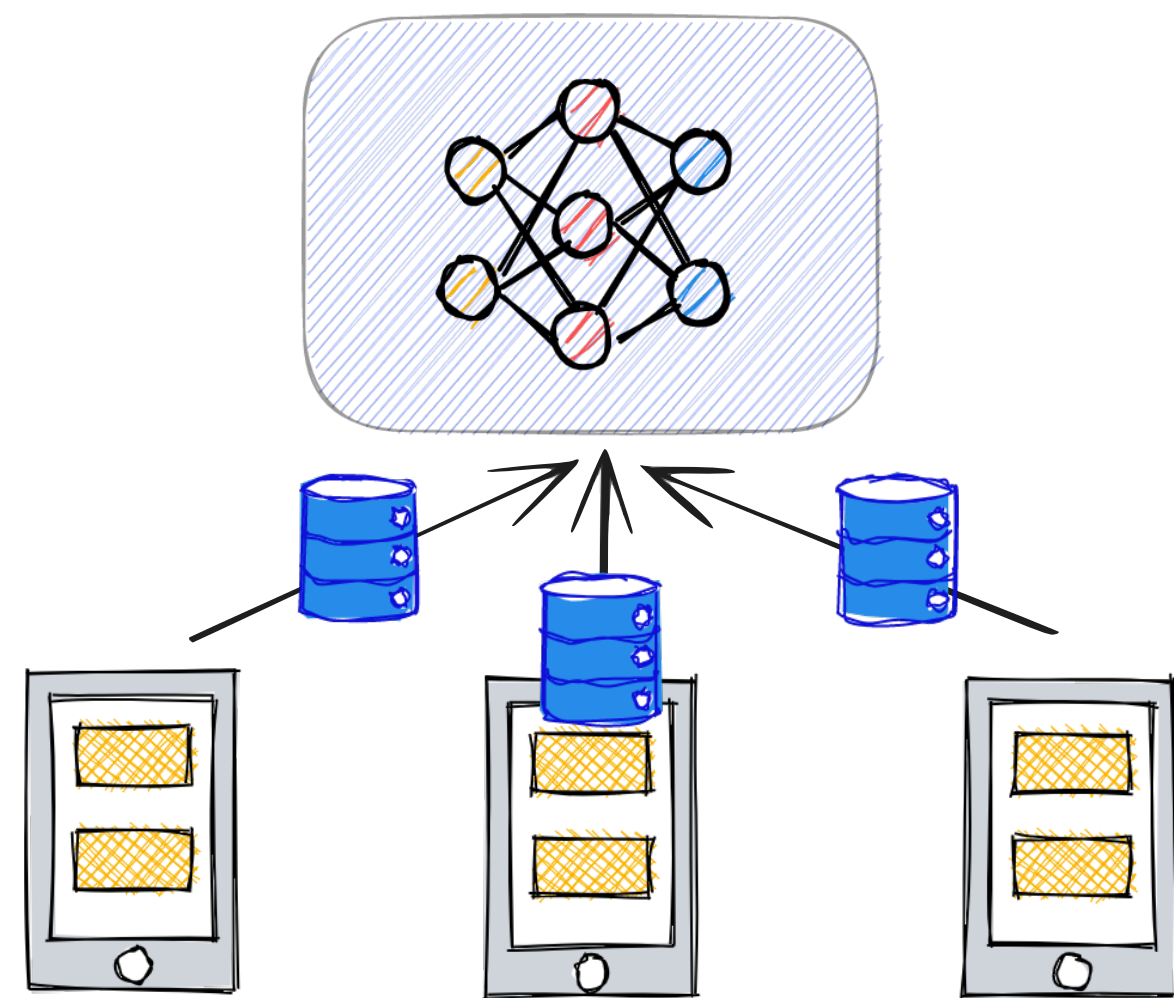
Pros:

1. Communication
2. Privacy



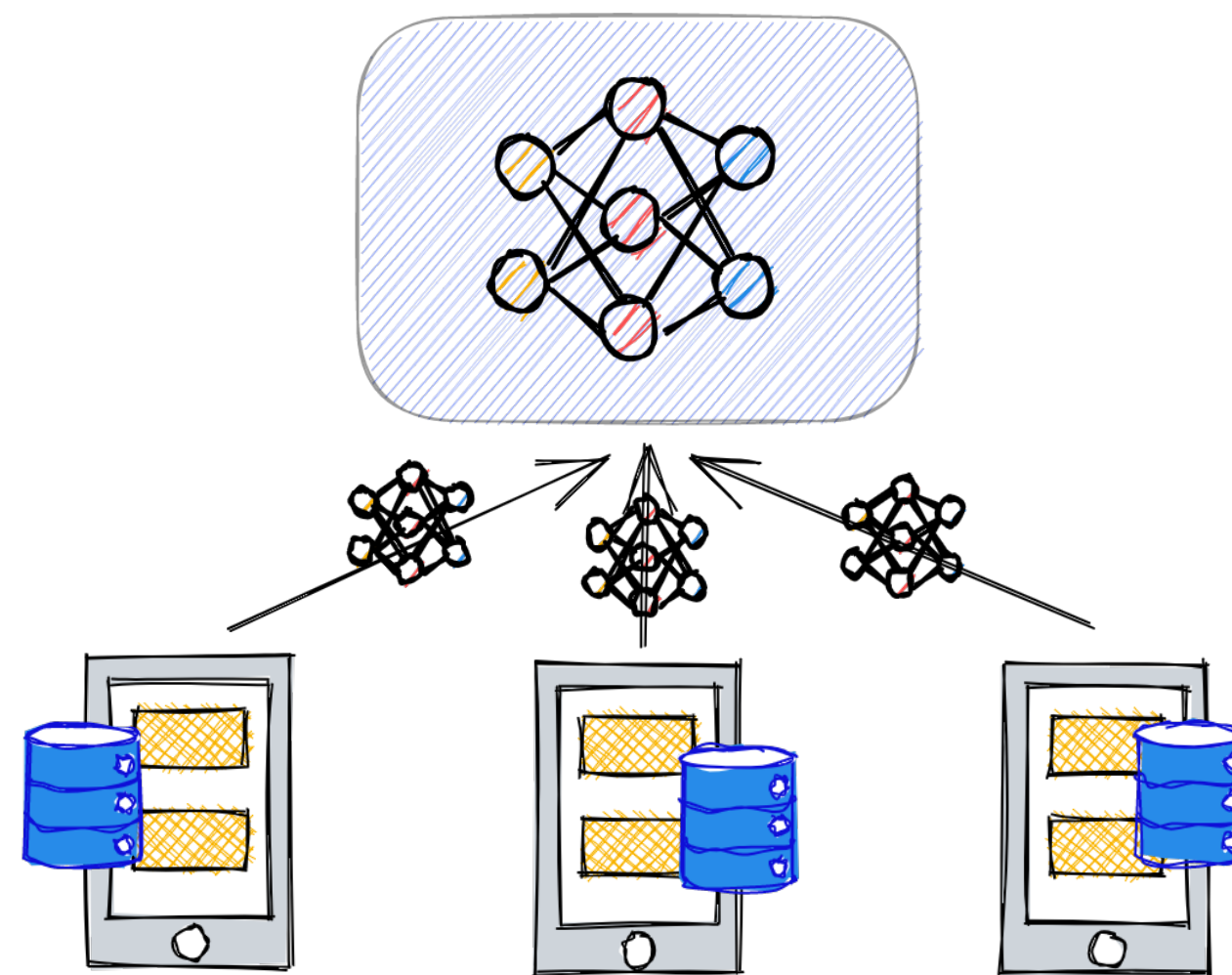
# Centralized vs Federated

## Centralized



Share data

## Federated



Share models / gradients

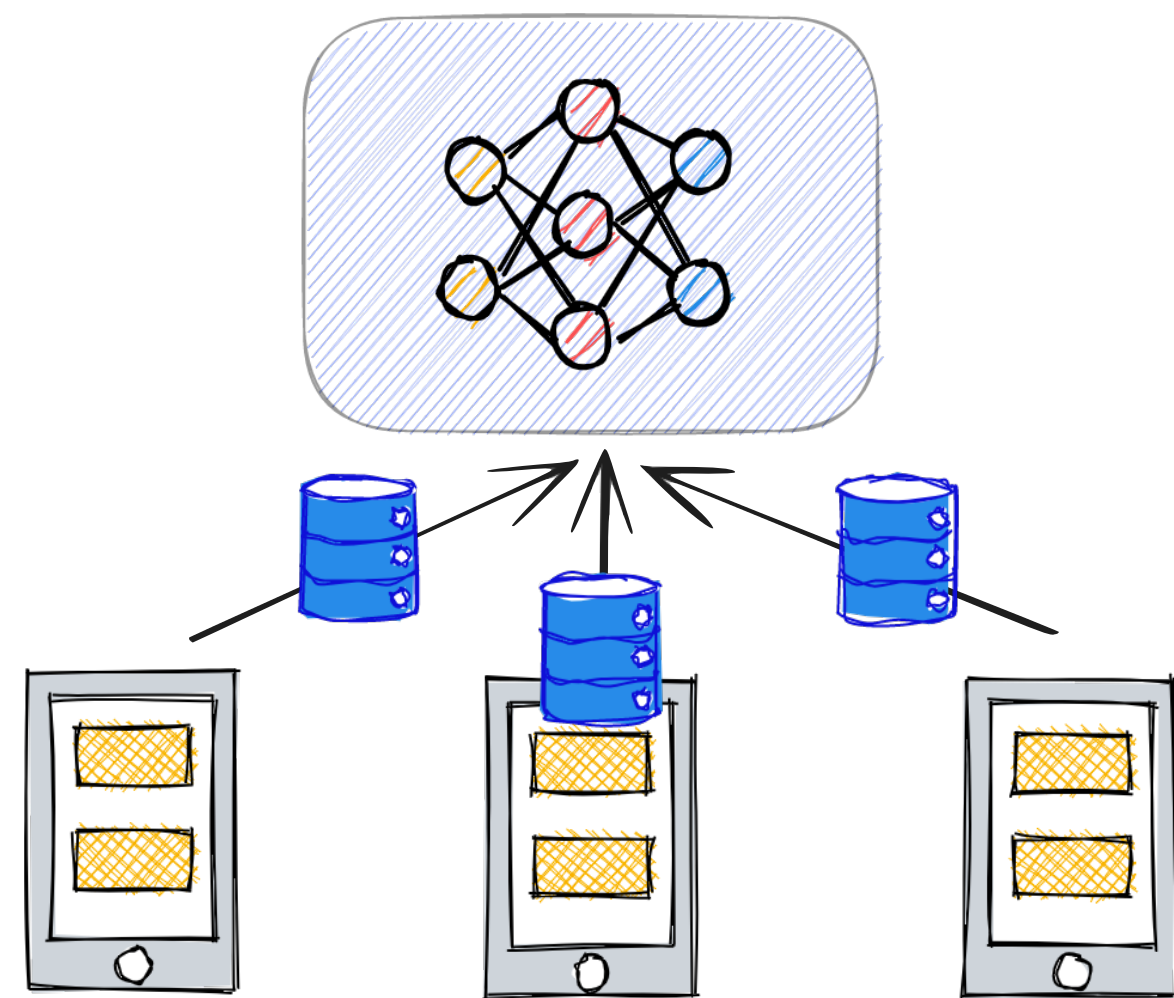
Pros:

1. Communication

2. Privacy

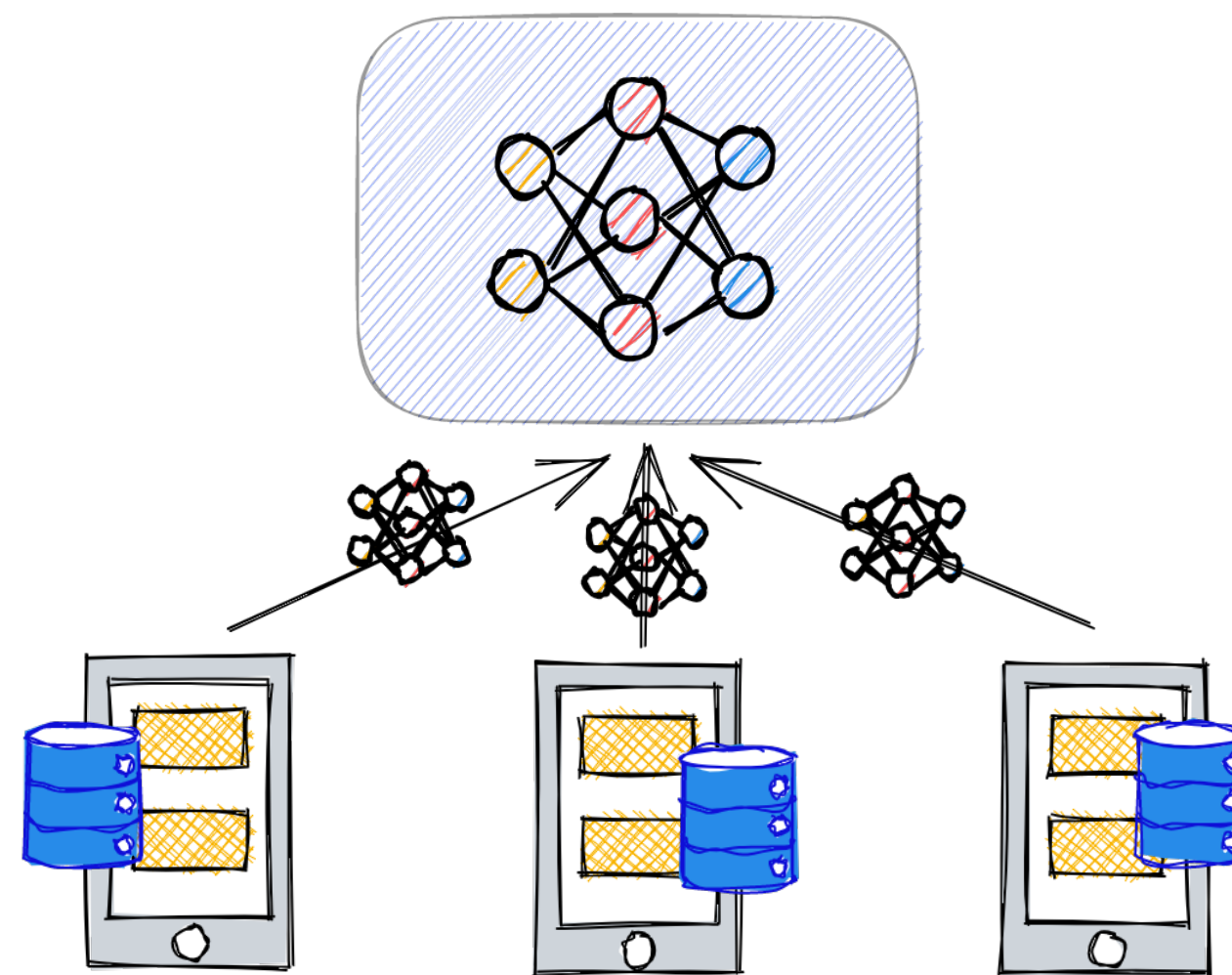
# Centralized vs Federated

## Centralized



Share data

## Federated



Share models / gradients

Pros:

1. Communication
2. Privacy

- **Common assumption:** clients are always available or uniform participation

# Lossy Communication Channels

## Previous works: loss mitigation

- Automatic Repeat Request (ARQ)
- Forward Error Correction (FEQ)

## Our motivations

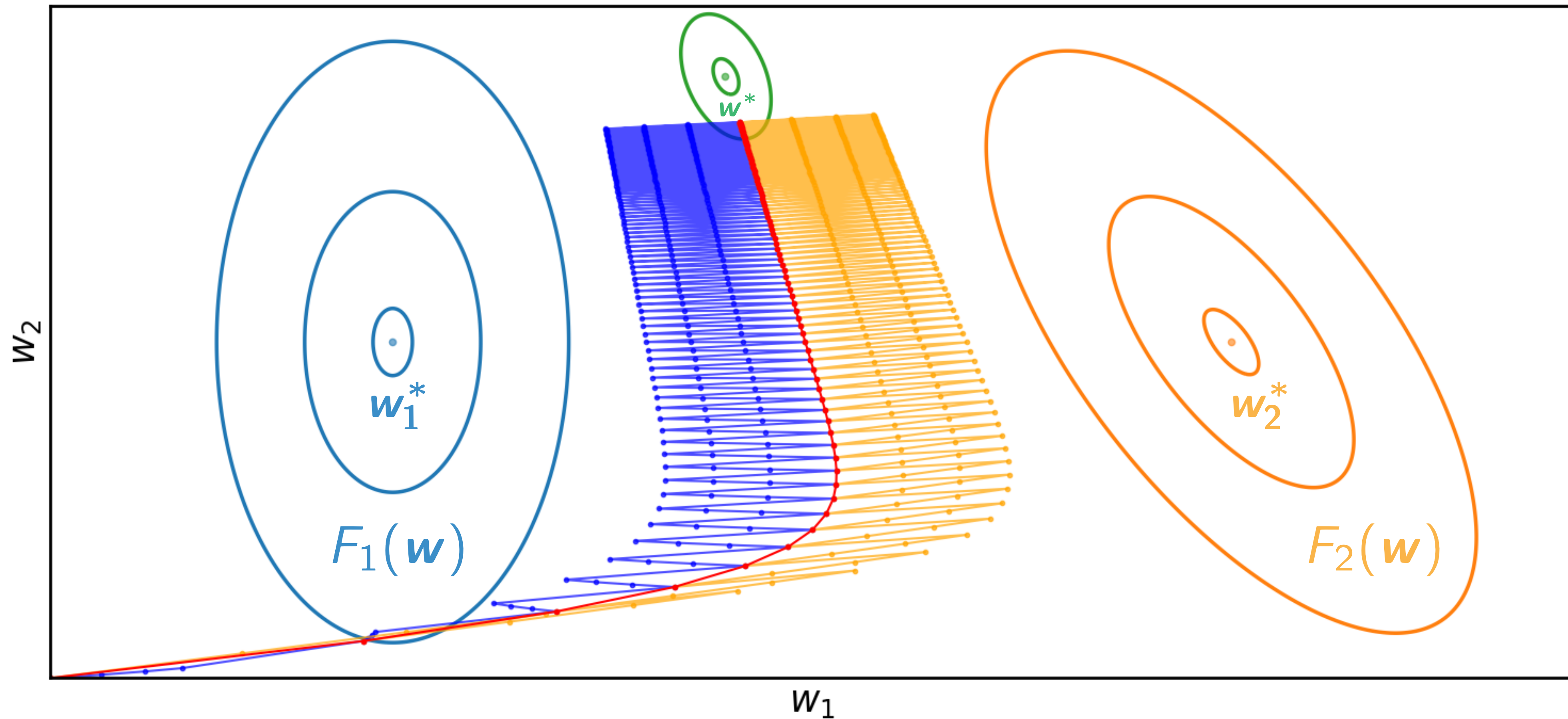
- Inevitable packet losses (e.g., retransmission failure)
- Larger training time and resource costs
- Robustness of gradient methods against limited errors

Can FL algorithms achieve optimal convergence despite packet losses?



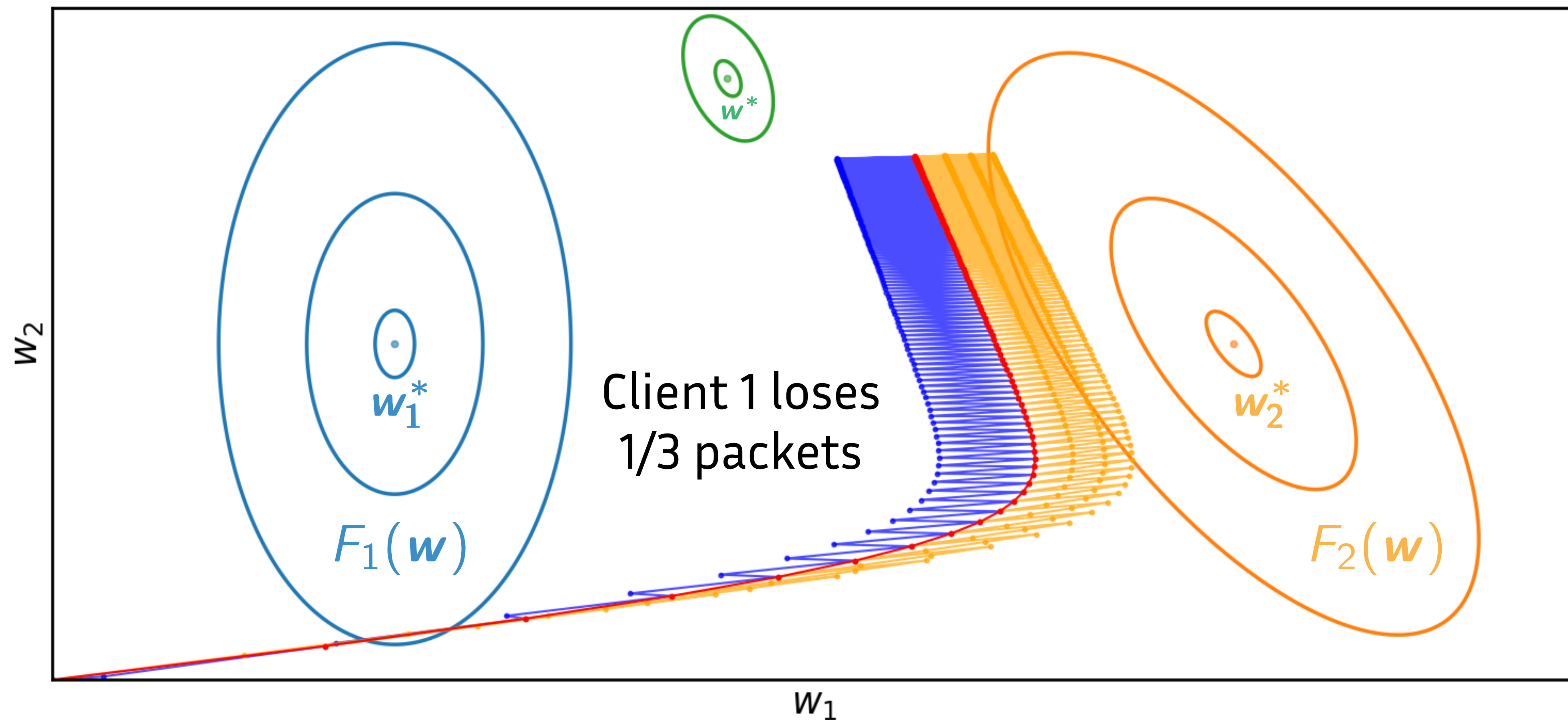
# Lossless Scenario

$$F(\mathbf{w}) = \frac{1}{2}F_1(\mathbf{w}) + \frac{1}{2}F_2(\mathbf{w})$$



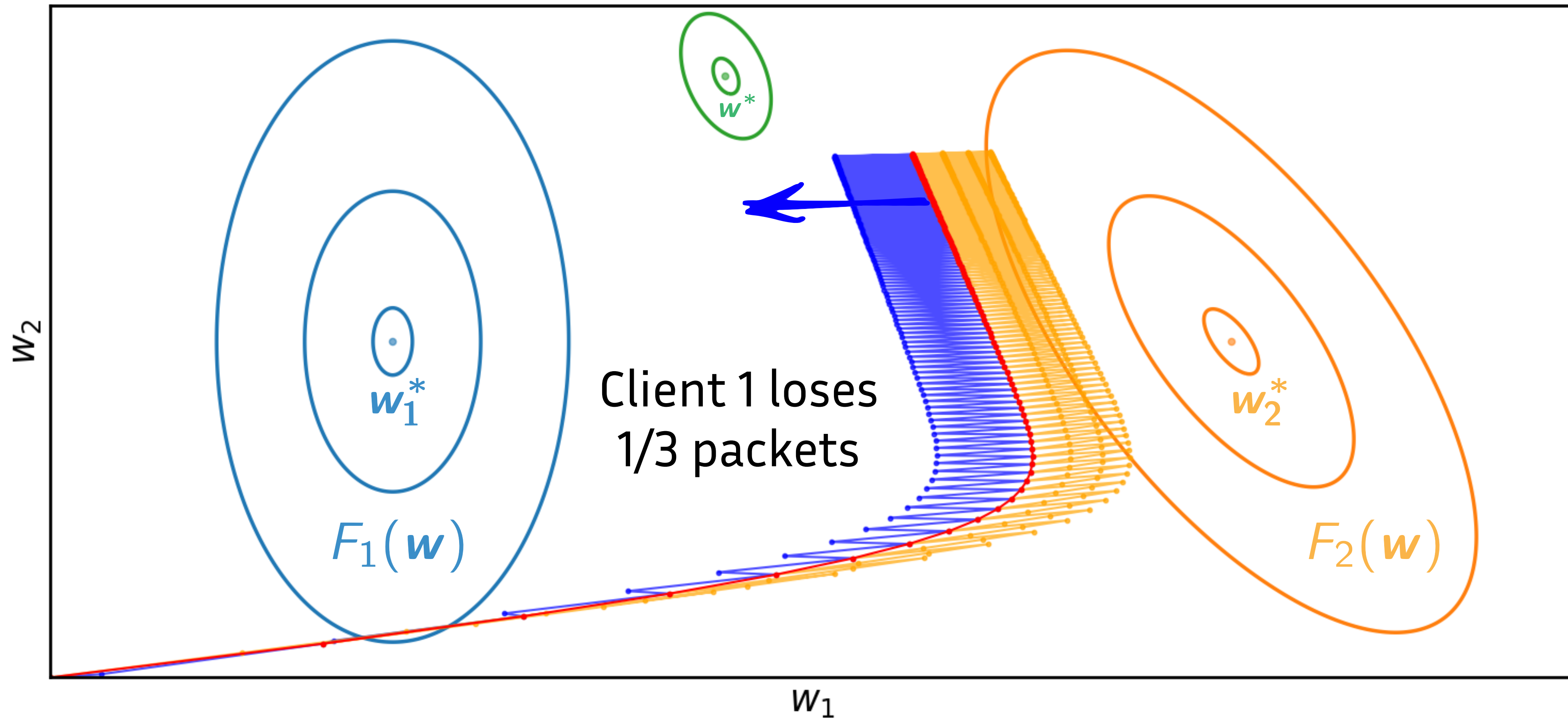
# Lossy Scenario

$$F(\mathbf{w}) = \frac{1}{2}F_1(\mathbf{w}) + \frac{1}{2}F_2(\mathbf{w})$$



# Lossy Scenario

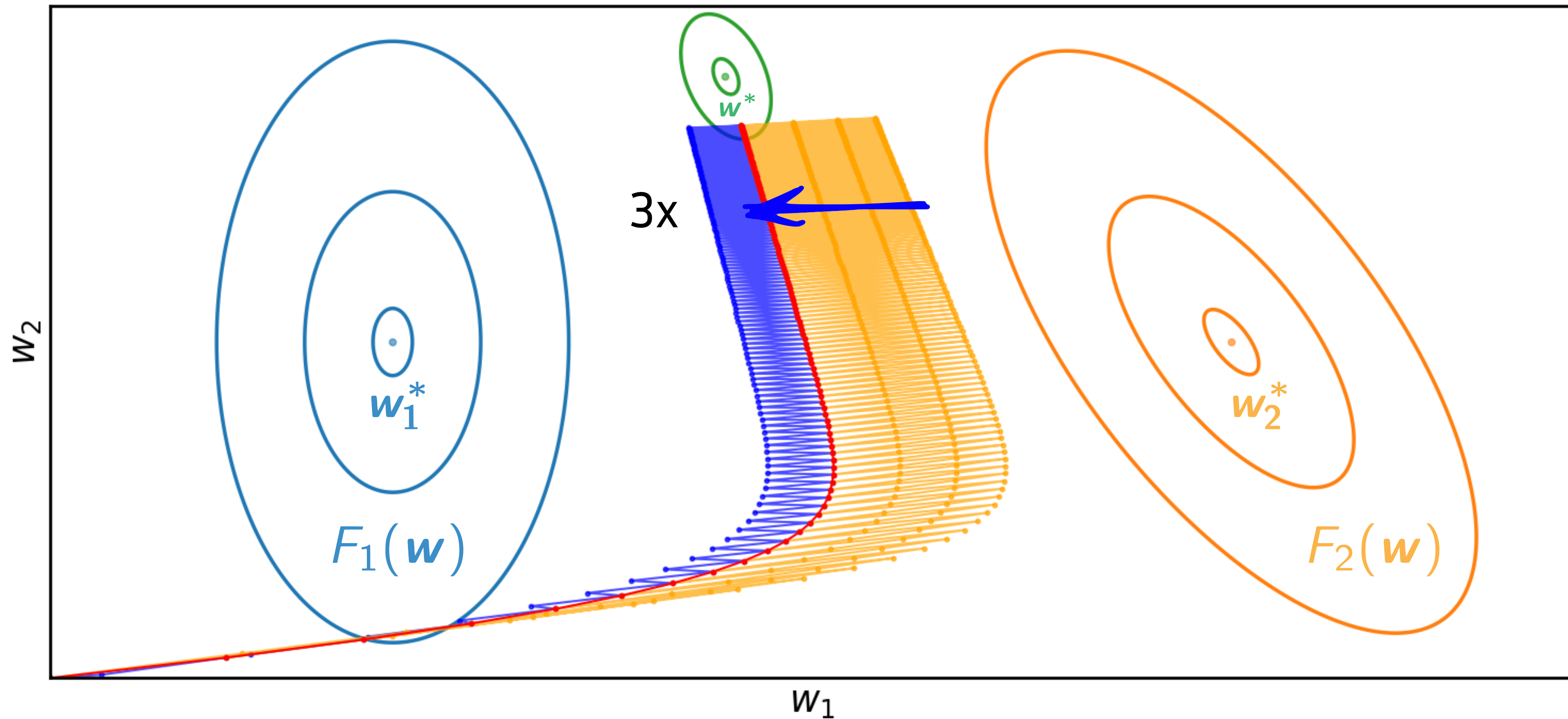
$$F(\mathbf{w}) = \frac{1}{2}F_1(\mathbf{w}) + \frac{1}{2}F_2(\mathbf{w})$$





# Lossy Scenario

$$F(\mathbf{w}) = \frac{1}{2}F_1(\mathbf{w}) + \frac{1}{2}F_2(\mathbf{w})$$



Yes, if 1) Aggregate Pseudo-Gradients  
2) Compensate for Packet Losses

# Aggregation for lossy channels

## Direct Model Aggregation (DMA)

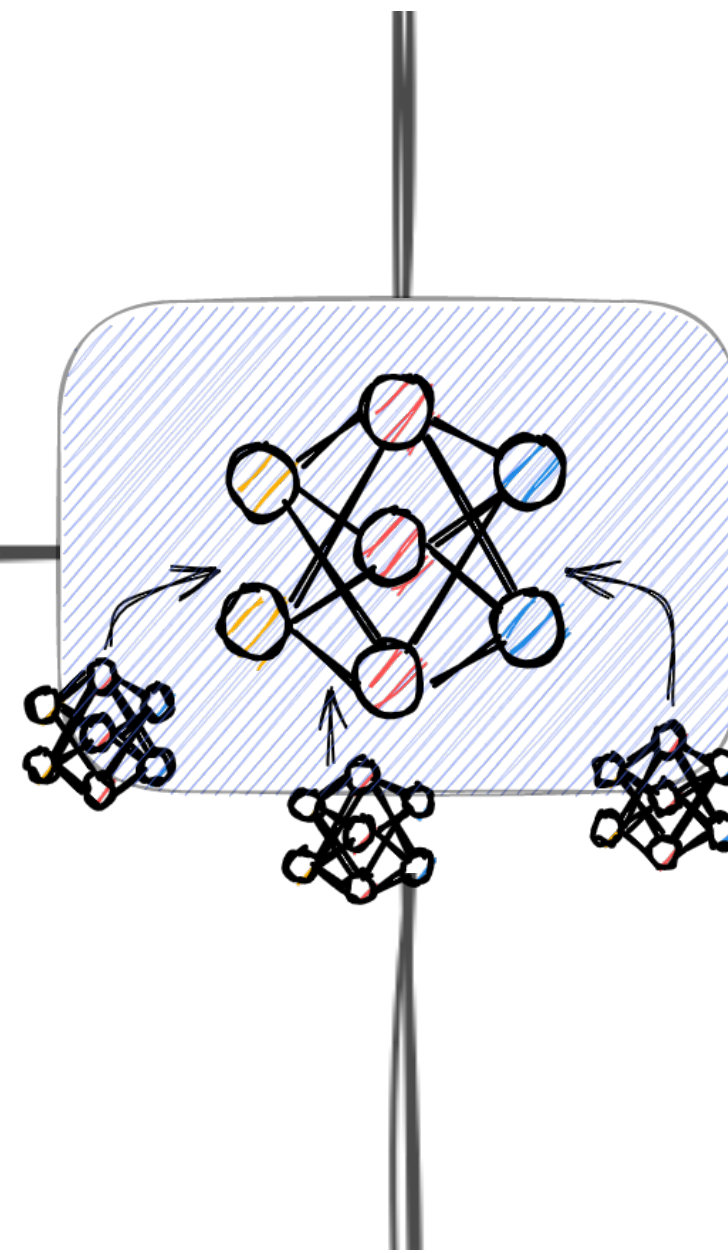
$$\mathbf{w}_{t+1}^{\text{DMA-PL}} = \frac{1}{|\mathcal{P}_t|} \sum_{k \in \mathcal{P}_t} \mathbf{w}_{t,E}^k$$

## Pseudo-Gradient Aggregation (PGA)

$$\mathbf{w}_{t+1}^{\text{PGA-PL}} = \mathbf{w}_t + \frac{1}{|\mathcal{P}_t|} \sum_{k \in \mathcal{P}_t} \Delta_t^k$$

## Unbiased DMA

$$\mathbf{w}_{t+1}^{\text{UDMA-PL}} = \frac{1}{N} \sum_{k \in \mathcal{P}_t} \frac{\mathbf{w}_{t,E}^k}{1 - p_k}$$



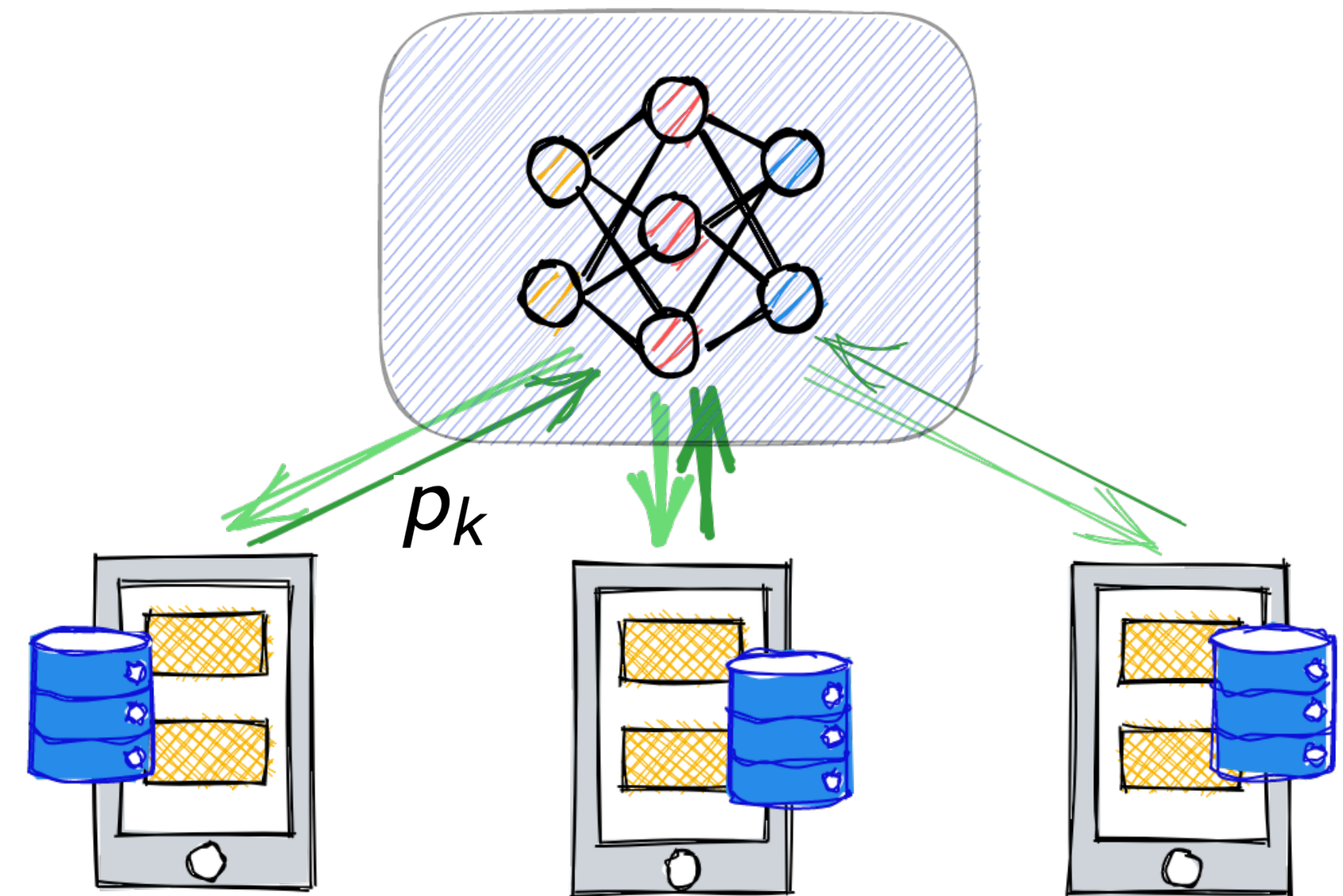
## Unbiased PGA (Ours)

$$\mathbf{w}_{t+1}^{\text{UPGA-PL}} = \mathbf{w}_t + \frac{1}{N} \sum_{k \in \mathcal{P}_t} \frac{\Delta_t^k}{1 - p_k}$$

- 1) Aggregate Pseudo-Gradients
- 2) Compensate for Packet Losses

# Assumptions to model lossy channels

- Loss probabilities  $p_k$  differ among clients
- Independent losses among clients
- For each client, IID losses over time
- Asymmetric channels (downlink/uplink)
- If ARQ or FEQ,  $p_k$  is the residual probability





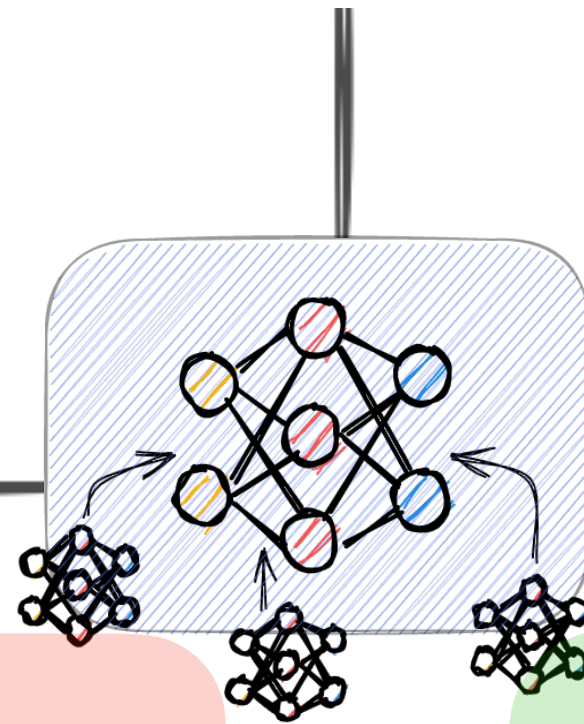
# Convergence Analysis

## Direct Model Aggregation

$$\mathbf{w}_{t+1}^{\text{DMA-PL}} = \frac{1}{|\mathcal{P}_t|} \sum_{k \in \mathcal{P}_t} \mathbf{w}_{t,E}^k$$

## Unbiased Pseudo-Gradient Aggregation

$$\mathbf{w}_{t+1}^{\text{UPGA-PL}} = \mathbf{w}_t + \frac{1}{N} \sum_{k \in \mathcal{P}_t} \frac{\Delta_t^k}{1 - p_k} \quad \text{(OURS)}$$



$$\mathbb{E}[F(\mathbf{w}_{t+1}^{\text{DMA-PL}})] - F^* \leq$$

$$\underbrace{A^t (F(\mathbf{w}_1) - F^*)}_{\text{vanishing term for small statistical heterogeneity}} + \underbrace{\frac{2\zeta_1}{L} \frac{1}{N} \sum_{k=1}^N p_k \frac{1 - A^t}{1 - A}}_{\text{non-vanishing error due to stat. het. and packet loss}}$$

vanishing term for small statistical heterogeneity

non-vanishing error due to stat. het. and packet loss

$$\mathbb{E}[F(\mathbf{w}_{t+1}^{\text{UPGA-PL}})] - F^* \leq$$

$$\underbrace{\frac{\kappa}{8\kappa + t} \left( \frac{2EC}{\mu} + 4L \|\mathbf{w}_1 - \mathbf{w}^*\|^2 \right)}_{\text{asymptotically vanishing term}}$$

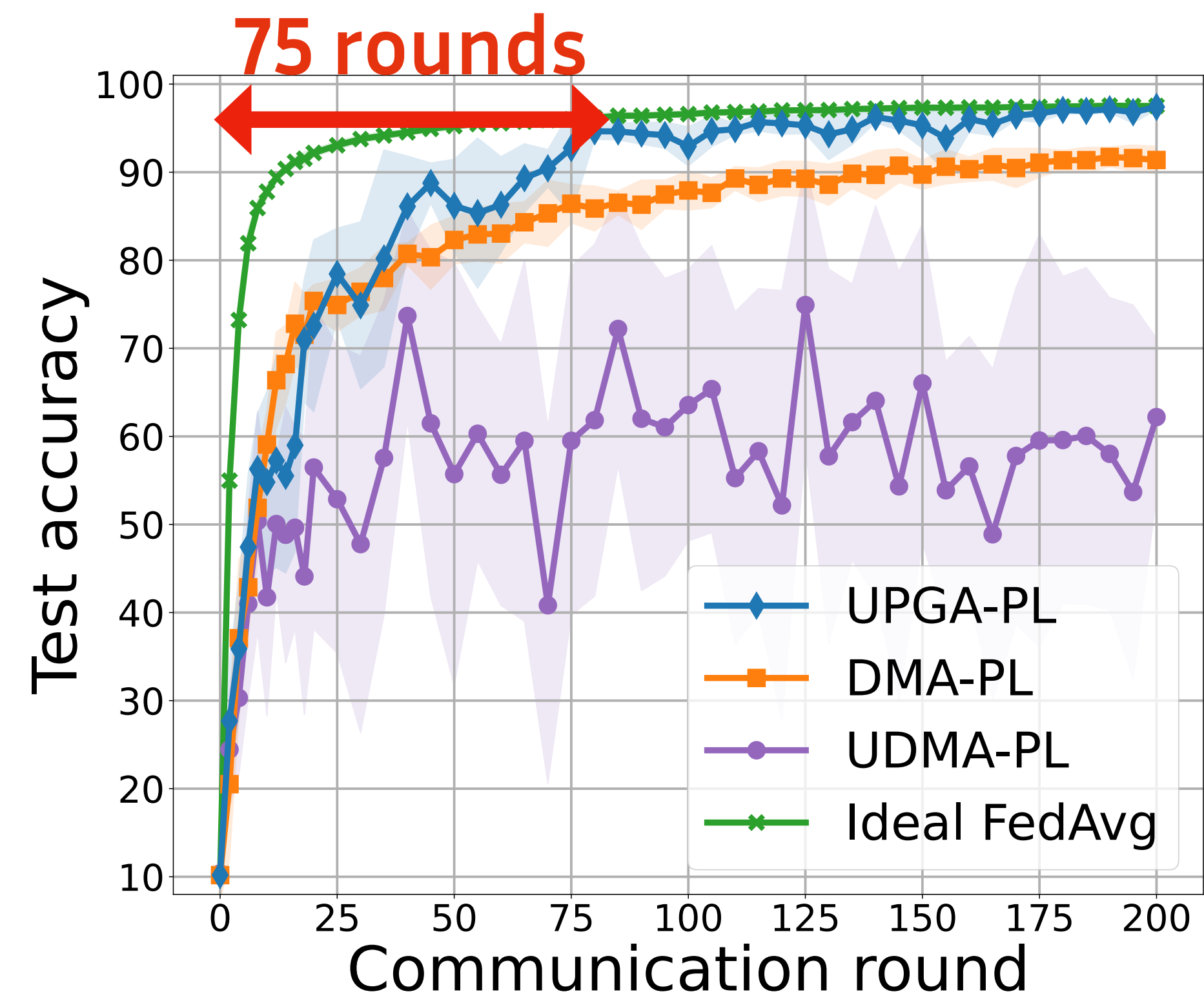
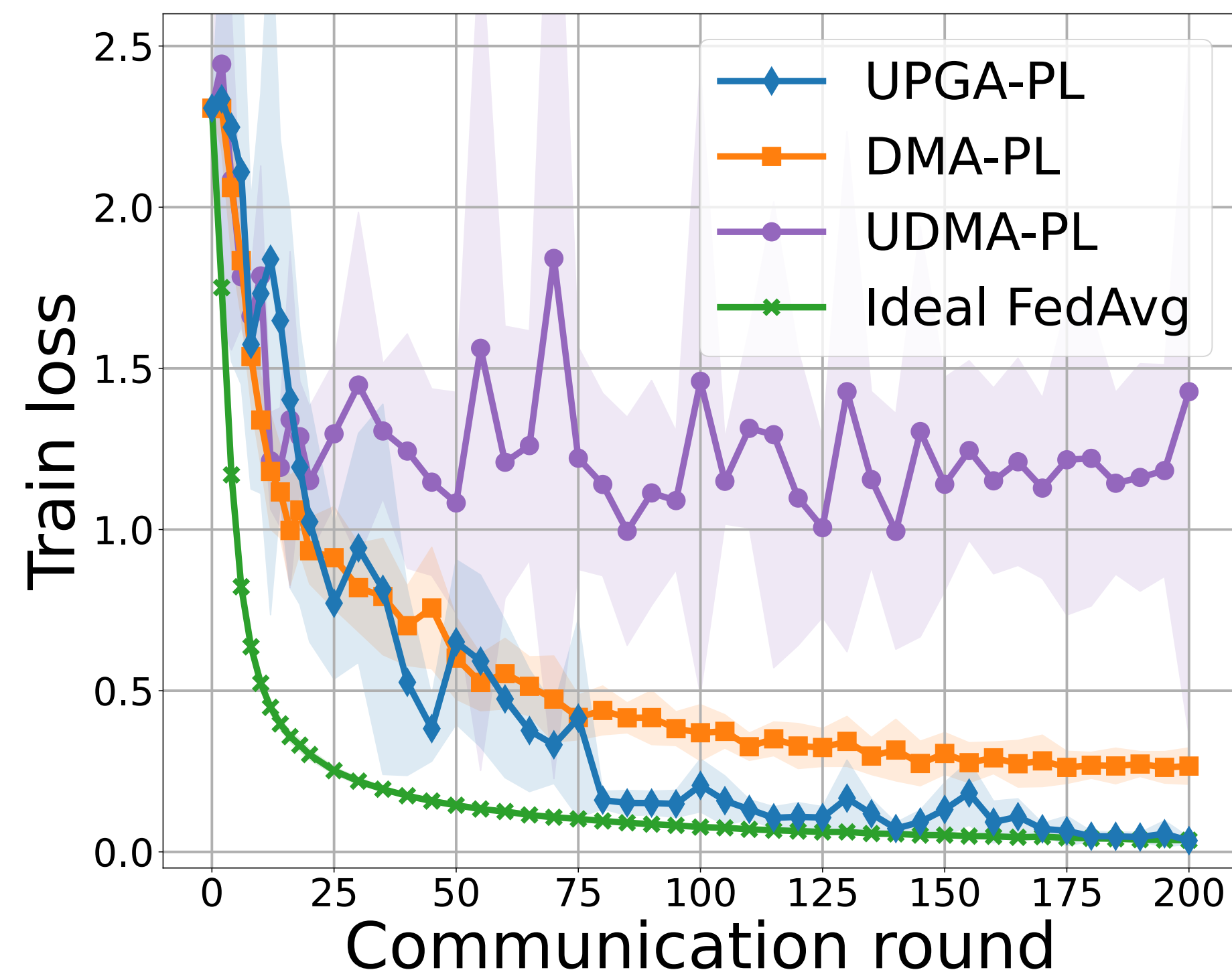
$$C := \frac{1}{N^2} \sum_{k=1}^N \sigma_k^2 + 2(E-1)^2 G^2 + 6L\Gamma + \underbrace{\frac{EG^2}{N^2} \sum_{k=1}^N \frac{p_k}{1 - p_k}}_{\text{effect of packet loss}}$$

A joint learning and communications framework for federated learning over wireless networks. Chen, Mingzhe, et al. IEEE Transactions on Wireless Communications, 2021.

UPGA-PL converges to the optimal model 😊

# Experimental Evaluation

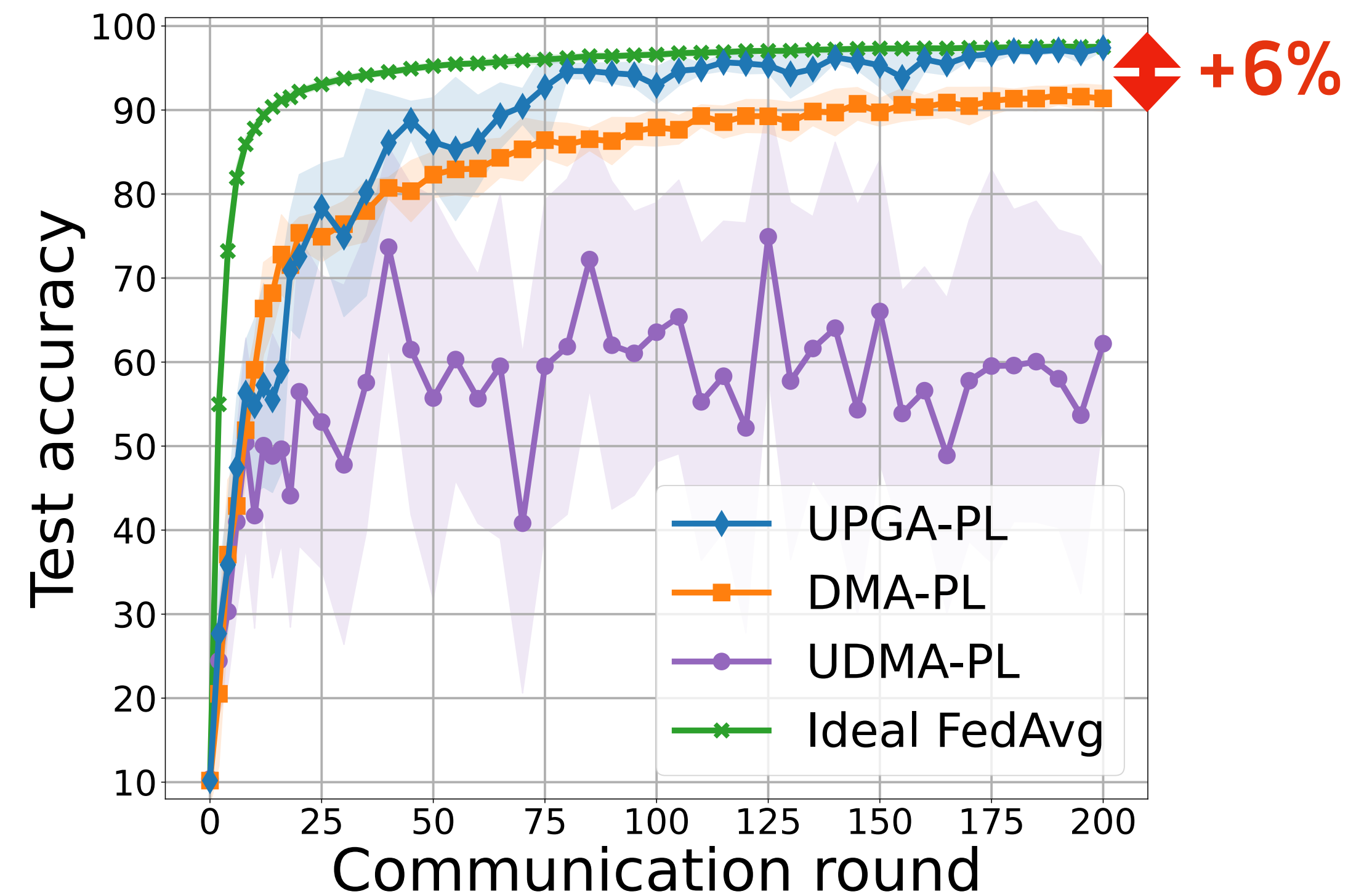
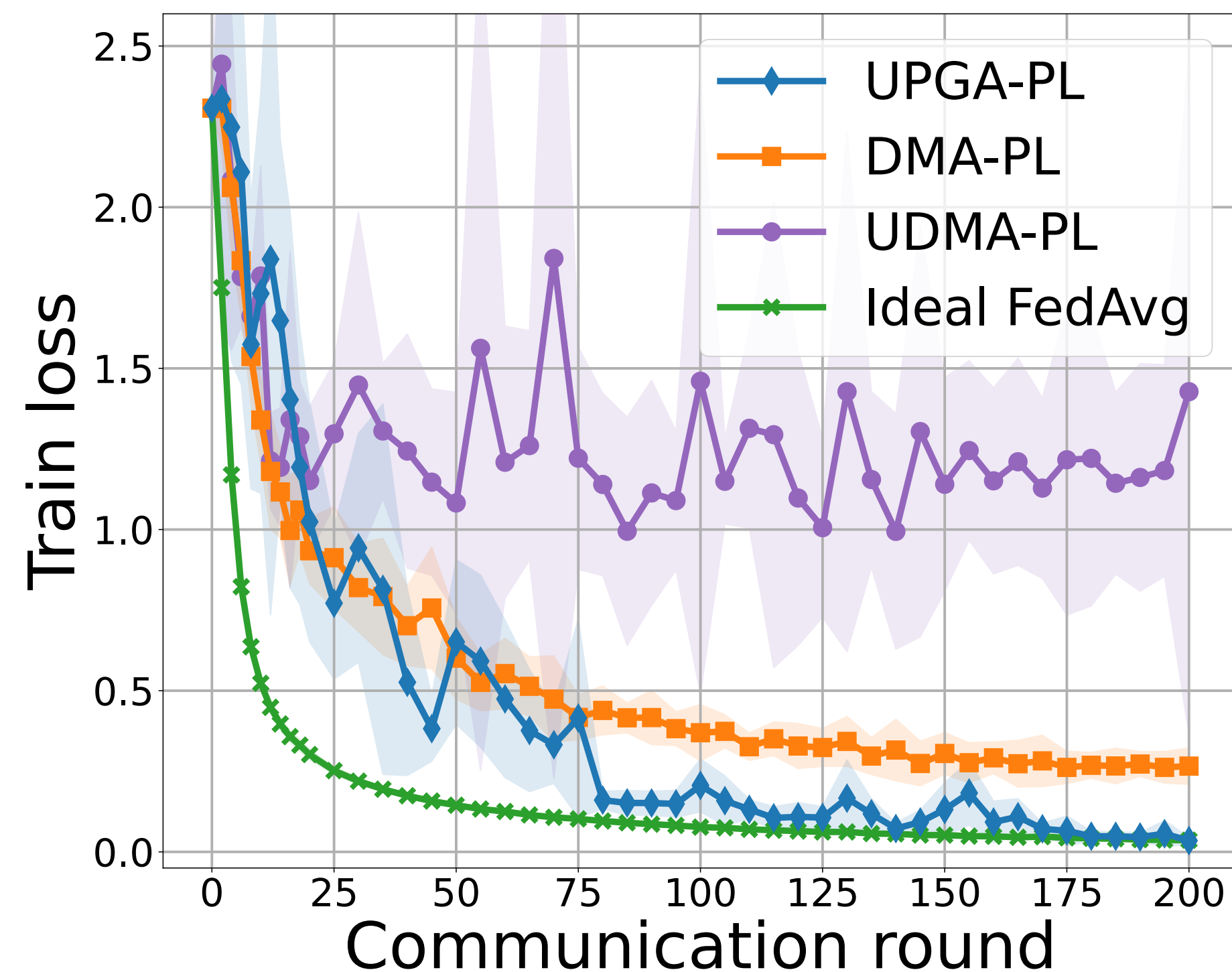
N = 10 clients equally split in two groups, with  $p_1 = \frac{1}{10}$ ,  $p_2 = \frac{9}{10}$ , MNIST dataset, CNN



UPGA-PL matches lossless performance in < 100 rounds

# Experimental Evaluation

N = 10 clients equally split in two groups, with  $p_1 = \frac{1}{10}$ ,  $p_2 = \frac{9}{10}$ , MNIST dataset, CNN

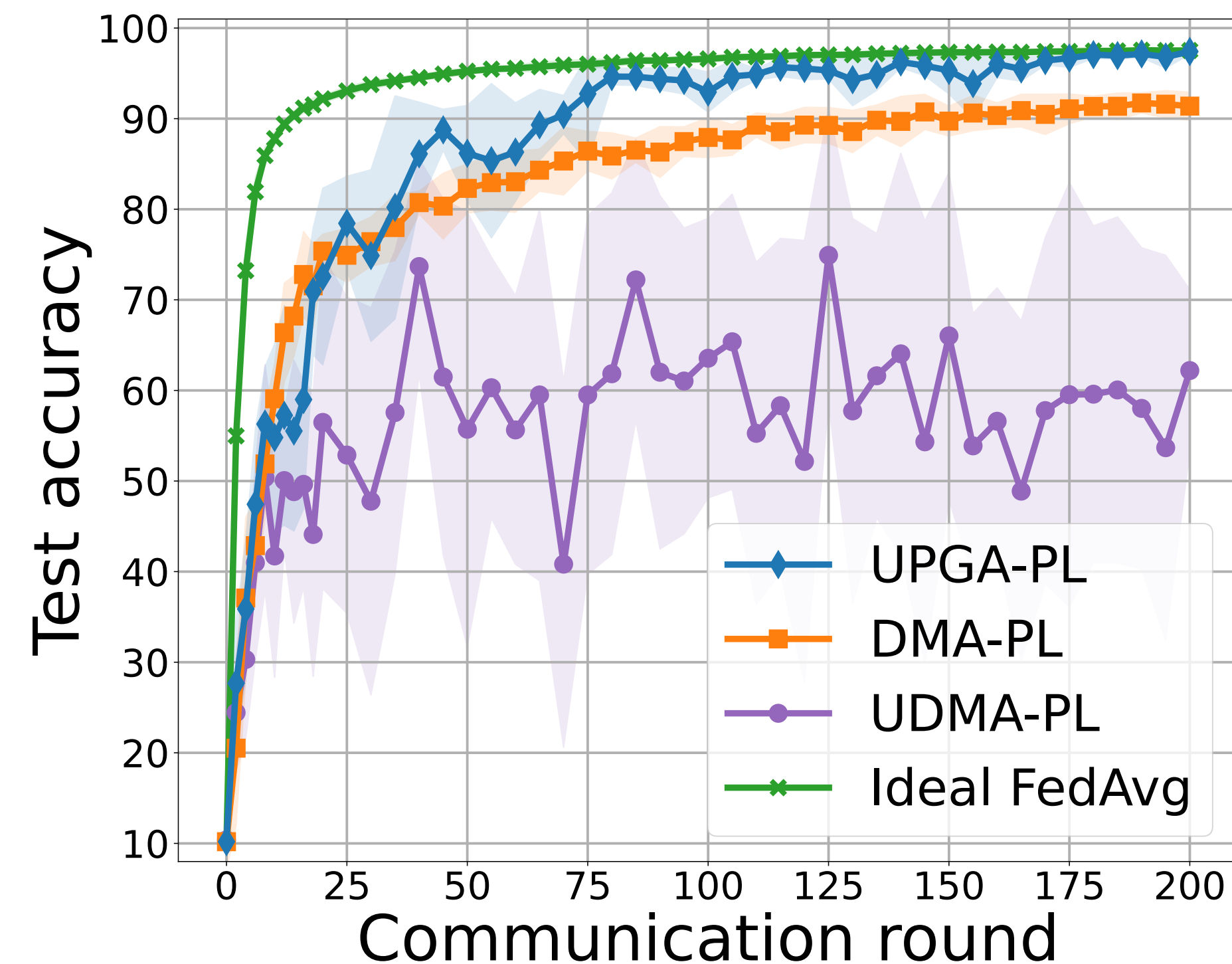
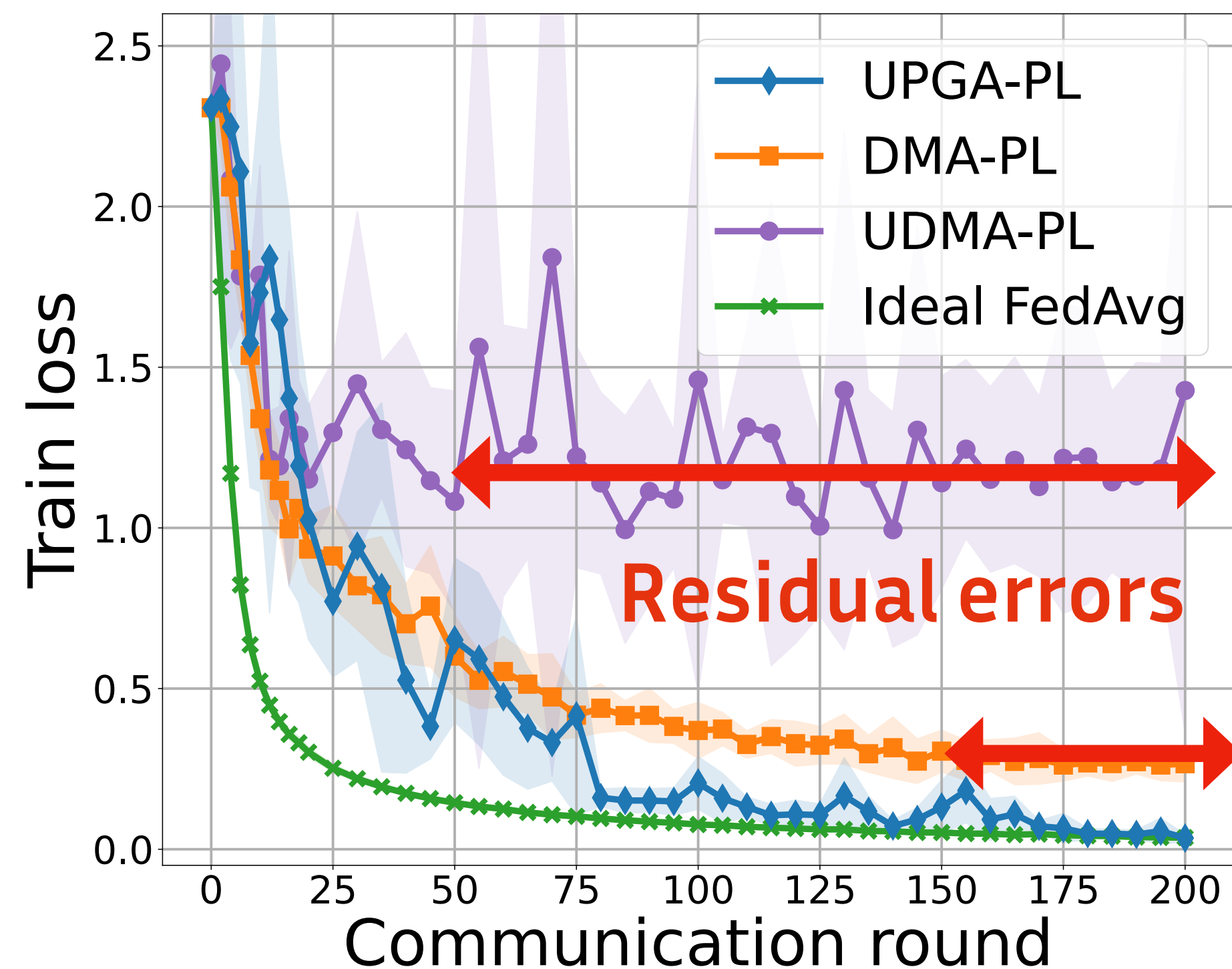


UPGA-PL improves MNIST performance by 6% over SOTA



# Experimental Evaluation

N = 10 clients equally split in two groups, with  $p_1 = \frac{1}{10}$ ,  $p_2 = \frac{9}{10}$ , MNIST dataset, CNN



DMA-PL and UDMA-PL exhibit non-vanishing errors

# Conclusions

- UPGA-PL has optimal convergence under asymmetric lossy channels
- UPGA-PL outperforms SOTA by filtering out losses
- UPGA-PL approaches ideal lossless channels with slightly slower convergence

Thank you for your attention!

Web



Code

