# Federated Learning Over Noisy Channel

Xizixiang Wei

[C] Xizixiang Wei, Cong Shen, Federated Learning in the Presence of Communication Errors, in Proc. IEEE International Conference on Communications (ICC), June 2021.

[J] Xizixiang Wei and Cong Shen, Federated Learning over Noisy Channels: Convergence Analysis and Design Examples, IEEE Transactions on Cognitive Communications and Networking, vol. 8, no. 2, pp. 1253-1268, June 2022.

# Background: Federated Learning

**1** Machine learning enables powerful applications

**2** Massive real-word data are generated at edge devices

**3** We want to keep sensitive data at edge devices

**Federated Learning (FL)**

- A distributed machine learning paradigm
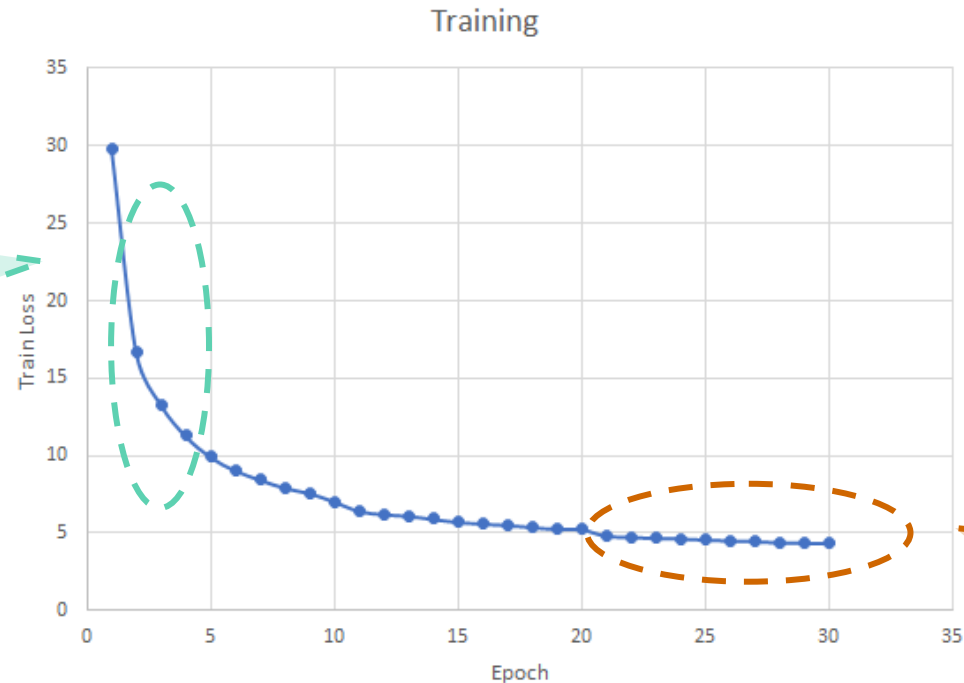- Obtain a global model while keeping data locally

# Motivation

Main difference from traditional communications
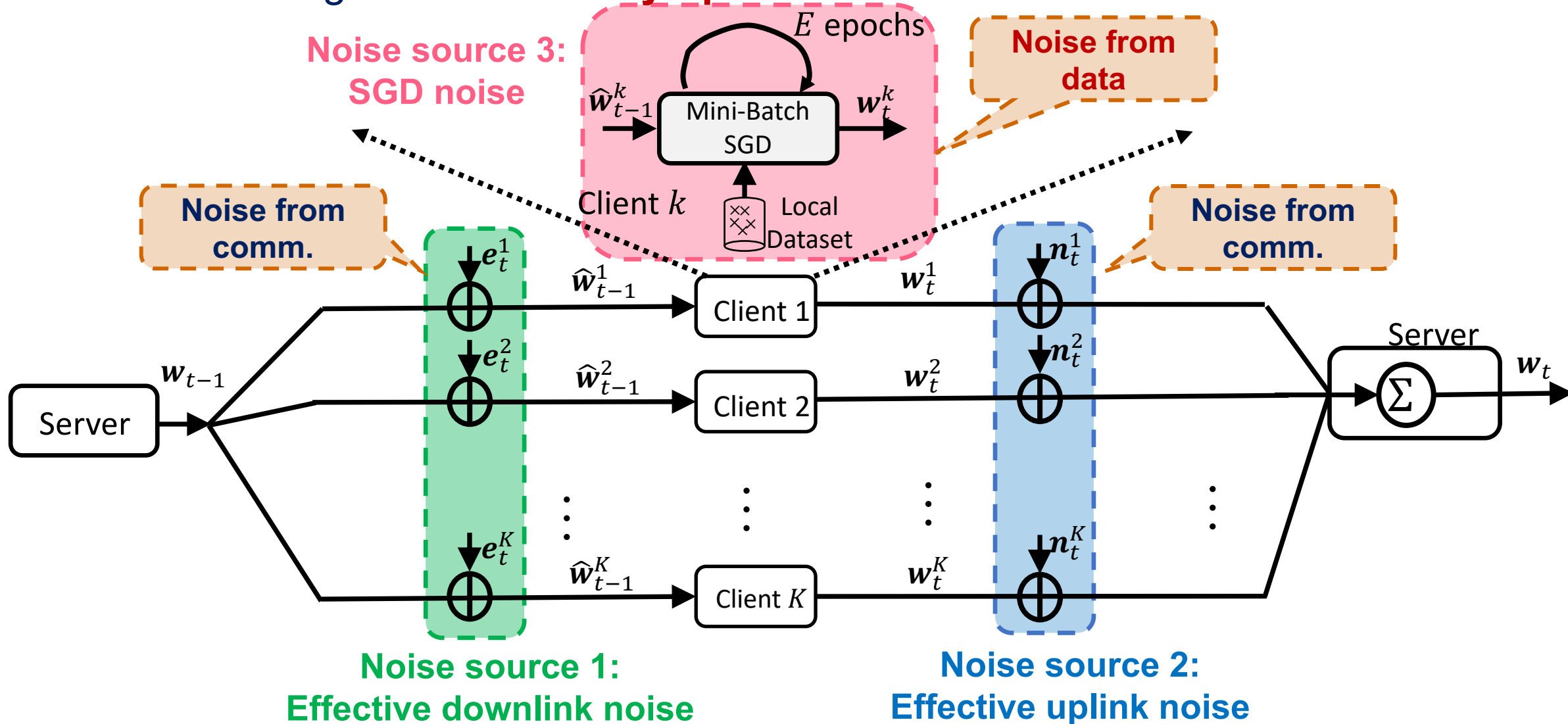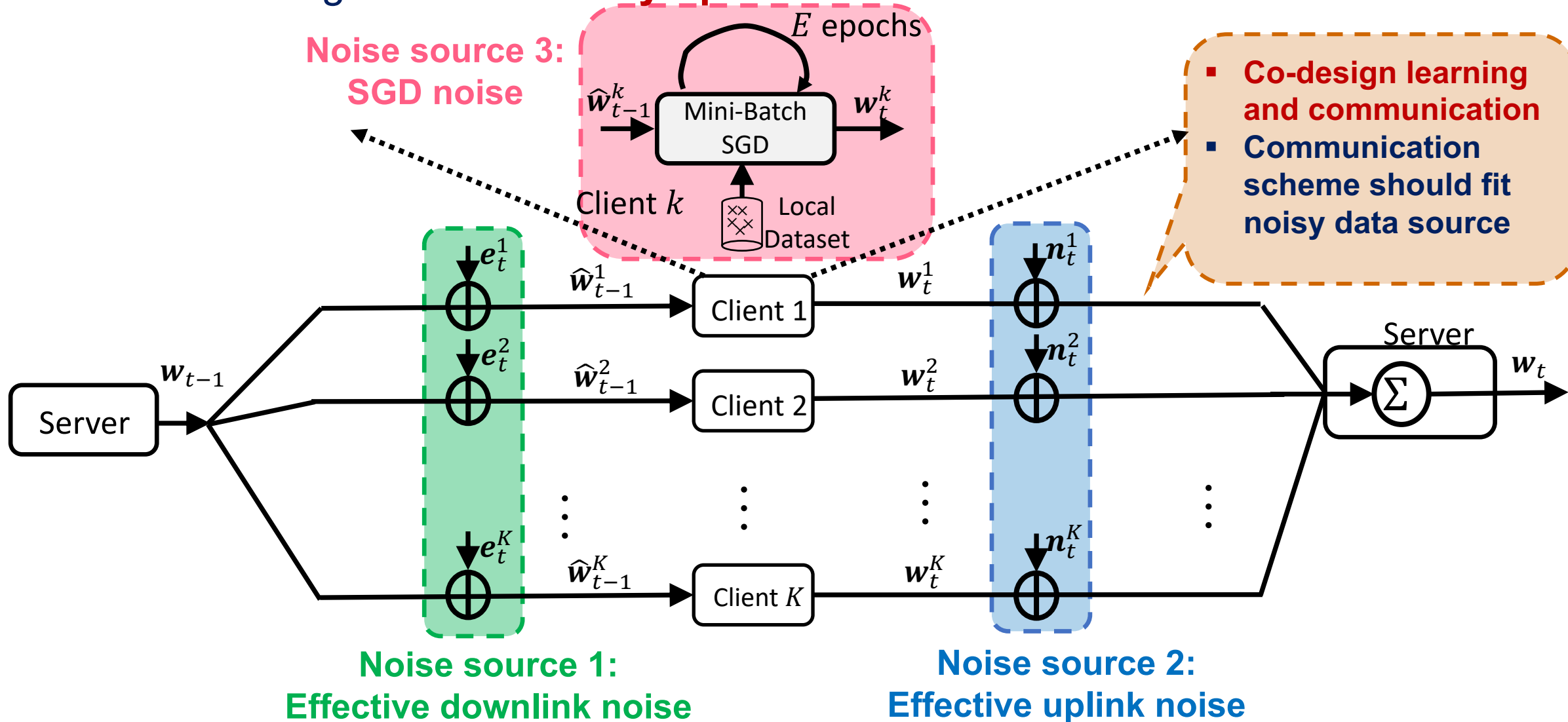- **Different** quality-of-service (QoS) requirements over time

# System model

Federated learning over **both noisy uplink and downlink** channels

# System model

Federated learning over **both noisy uplink and downlink** channels



Noise source 3: SGD noise

$E$ epochs

$\hat{\boldsymbol{w}}_{t-1}^k$ → Mini-Batch SGD → $\boldsymbol{w}_t^k$

Client $k$    Local Dataset

- **Co-design learning and communication**
- **Communication scheme should fit noisy data source**

Server — $\boldsymbol{w}_{t-1}$

$\boldsymbol{e}_t^1$    $\hat{\boldsymbol{w}}_{t-1}^1$    Client 1    $\boldsymbol{w}_t^1$    $\boldsymbol{n}_t^1$

$\boldsymbol{e}_t^2$    $\hat{\boldsymbol{w}}_{t-1}^2$    Client 2    $\boldsymbol{w}_t^2$    $\boldsymbol{n}_t^2$

$\boldsymbol{e}_t^K$    $\hat{\boldsymbol{w}}_{t-1}^K$    Client $K$    $\boldsymbol{w}_t^K$    $\boldsymbol{n}_t^K$

Server $\Sigma$ — $\boldsymbol{w}_t$

Noise source 1: Effective downlink noise

Noise source 2: Effective uplink noise

# SGD noise

Gradient descent (GD)

Loss function: $L_K(\boldsymbol{w})$



$\boldsymbol{w}_T^k$

$\nabla L_K(\boldsymbol{w}_t)$
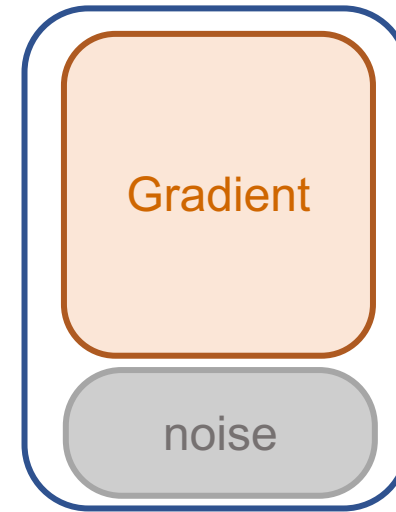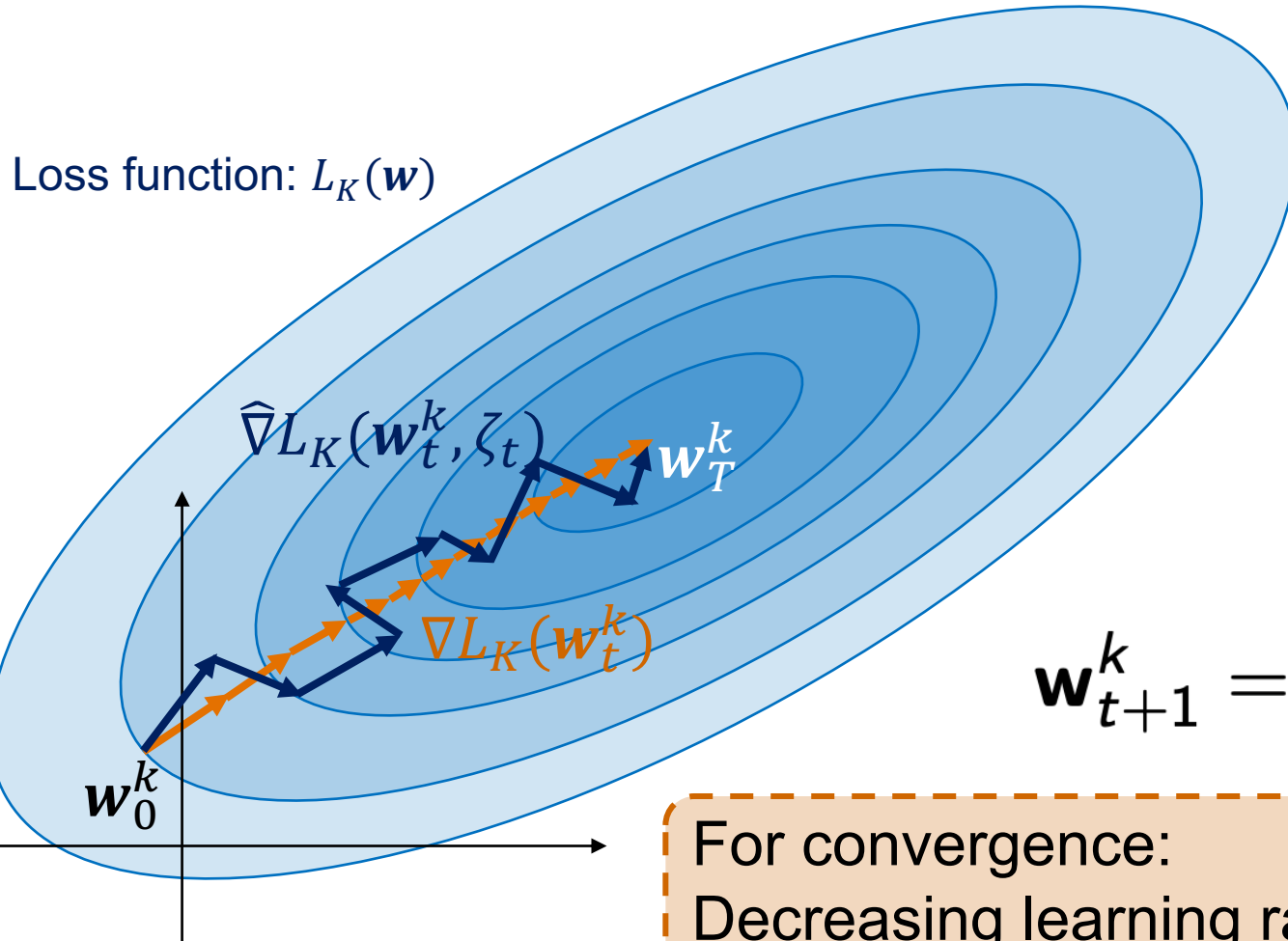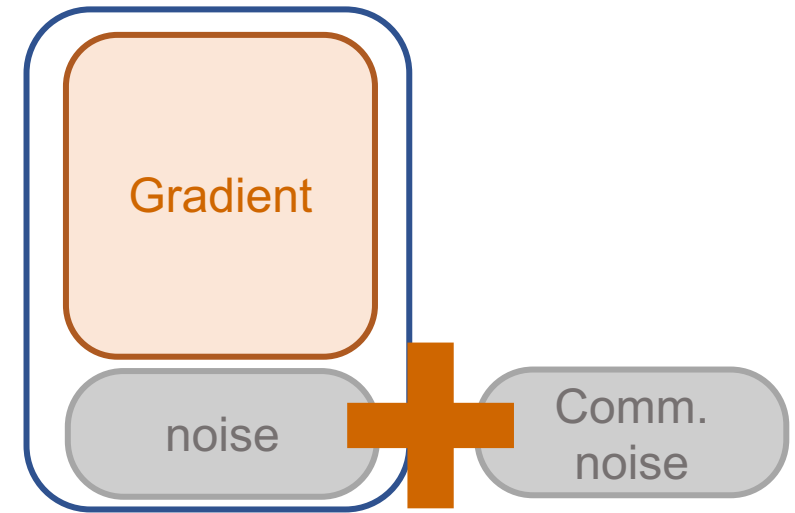
$\boldsymbol{w}_0^k$

# SGD noise

Gradient descent (GD)
Stochastic gradient descent (SGD)

$$\nabla L(\mathbf{w}_t^k) = \mathbb{E}[\hat{\nabla} L(\mathbf{w}_t^k, \zeta_t)]$$

stochastic gradient = gradient + noise

Loss function: $L_K(\boldsymbol{w})$



$\hat{\nabla} L_K(\boldsymbol{w}_t^k, \zeta_t)$

$\boldsymbol{w}_T^k$

$\nabla L_K(\boldsymbol{w}_t^k)$

$\boldsymbol{w}_0^k$

Gradient

noise

$$\mathbf{w}_{t+1}^k = \mathbf{w}_t^k - \boxed{\eta_t} \hat{\nabla} L(\mathbf{w}_t^k, \zeta_t)$$

For convergence:
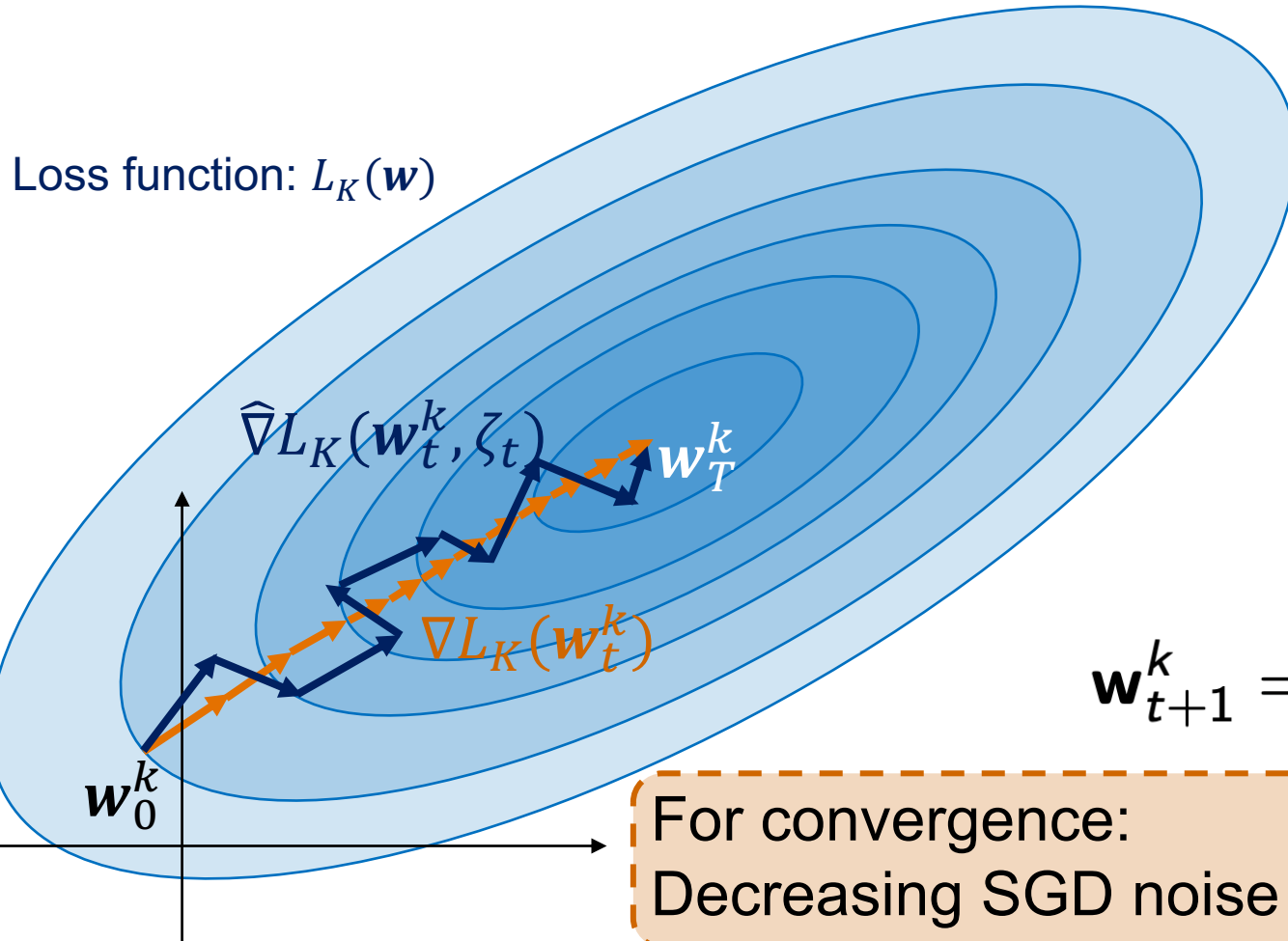Decreasing learning rate → Decreasing effective SGD noise

# SGD noise

Gradient descent (GD)
Stochastic gradient descent (SGD)

$$\nabla L(\mathbf{w}_t^k) = \mathbb{E}[\hat{\nabla} L(\mathbf{w}_t^k, \zeta_t)]$$

stochastic gradient = gradient + noise

Loss function: $L_K(\boldsymbol{w})$



$\hat{\nabla} L_K(\boldsymbol{w}_t^k, \zeta_t)$
$\boldsymbol{w}_T^k$
$\nabla L_K(\boldsymbol{w}_t^k)$
$\boldsymbol{w}_0^k$

Gradient

noise + Comm. noise

$$\mathbf{w}_{t+1}^k = \mathbf{w}_t^k - \eta_t \hat{\nabla} L(\mathbf{w}_t^k, \zeta_t) + \boxed{\mathbf{n}_{t+1}^k}$$

For convergence:
Decreasing SGD noise + **Decreasing effective comm. noise**

# Convergence over noisy channel

For $L$-smooth, $\mu$-strongly convex and bounded-gradient loss function

- Effective SNR control policy

**Effective DL noise power**          **Effective UL noise power**

$$\sigma_t^2 \sim \mathcal{O}\left(\frac{1}{t^2}\right)$$ **+** $$\zeta_t^2 \sim \mathcal{O}\left(\frac{1}{t^2}\right)$$

FL tasks with non-IID datasets and partial/full clients participation converge at rate $\mathcal{O}\left(\frac{1}{t}\right)$.

Channel noise should not dominate the SGD noise.

# Model differential for UL

Model differential
$$x_t^k = w_t^k - w_{t-1}$$

**Uplink comm.**

Global model recovery
$$w_t = w_{t-1} + \frac{1}{k}\sum_{k=1}^{K} x_t^k = \frac{1}{k}\sum_{k=1}^{K} w_t^k$$

# Model differential for UL

Model differential
$$x_t^k = w_t^k - w_{t-1}$$
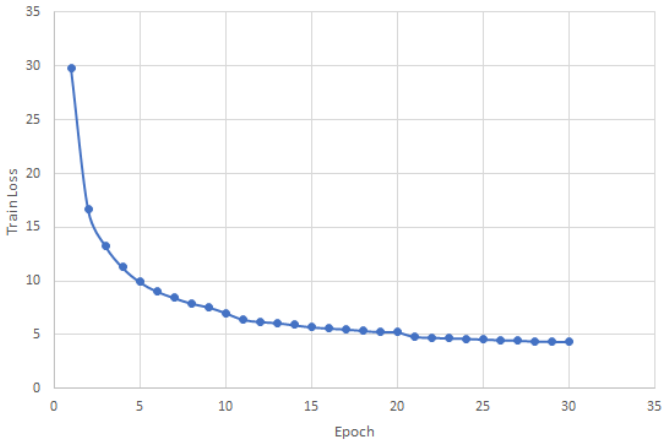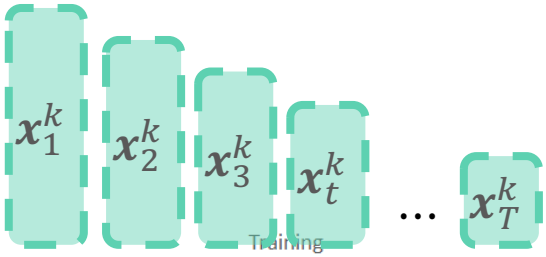
**Uplink comm.**

Global model recovery
$$w_t = w_{t-1} + \frac{1}{k}\sum_{k=1}^{K} x_t^k = \frac{1}{k}\sum_{k=1}^{K} w_t^k$$

$x_1^k$  $x_2^k$  $x_3^k$  $x_t^k$  ...  $x_T^k$
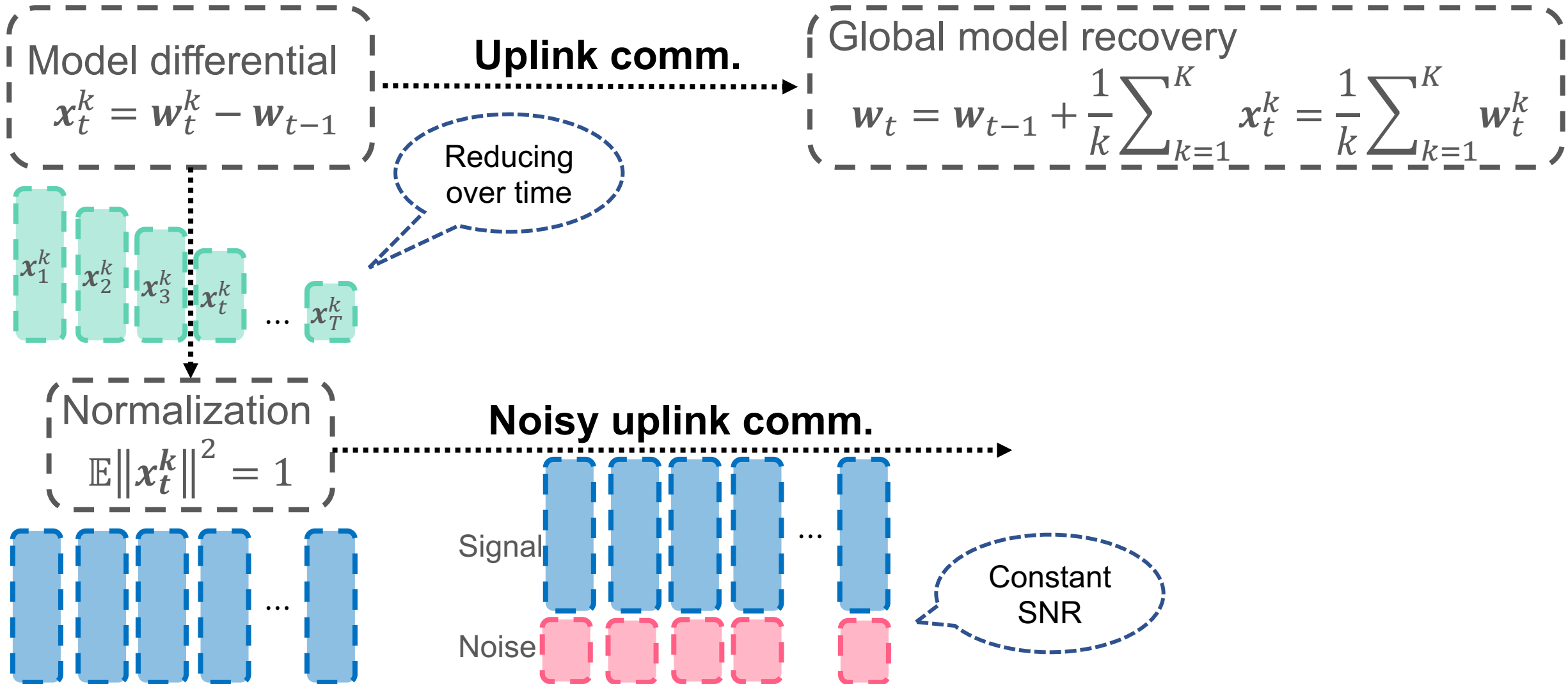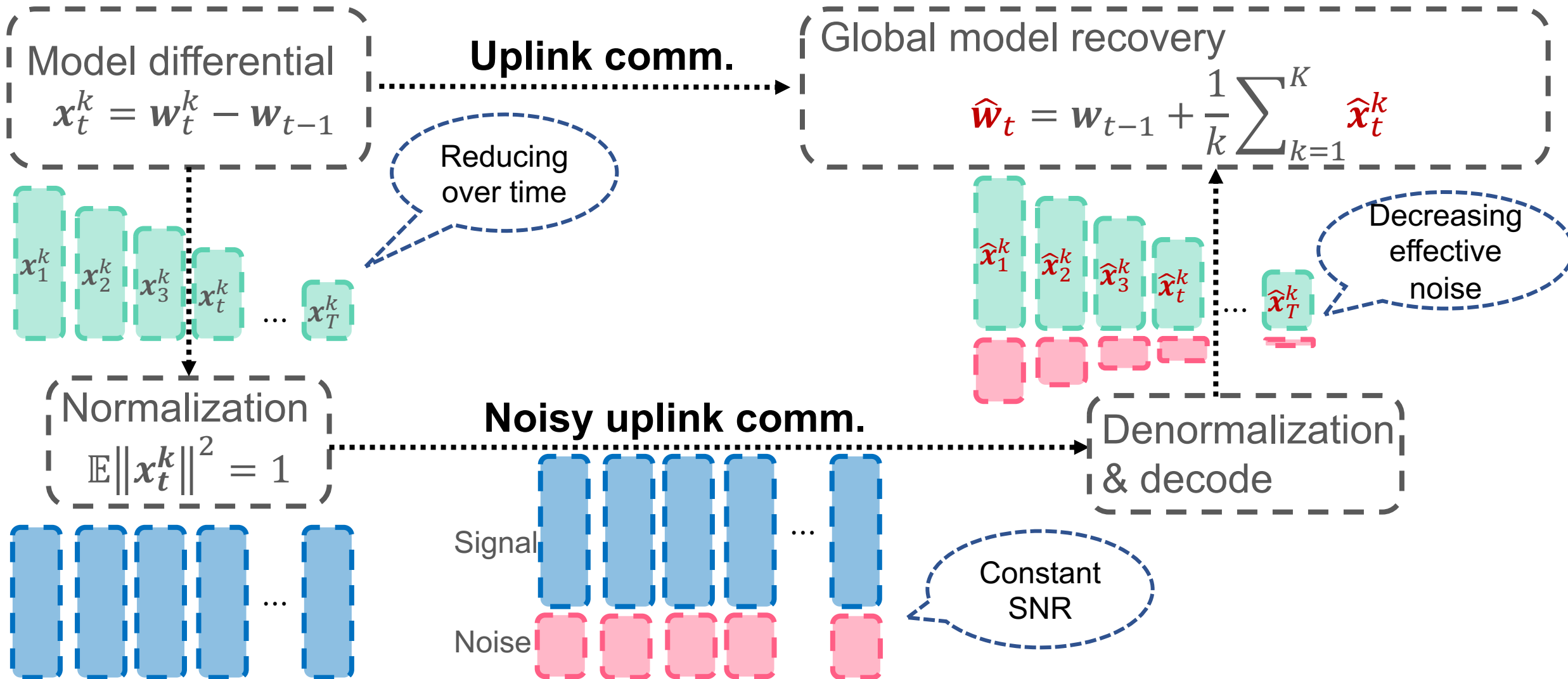
Training

Reducing over time

# Model differential for UL

# Model differential for UL

# Model differential for Uplink

# Convergence over noisy channel

For $L$-smooth, $\mu$-strongly convex and bounded-gradient loss function

- Effective SNR control policy for **uplink model differential**

**Effective DL noise power**

$$\sigma_t^2 \sim \mathcal{O}\left(\frac{1}{t^2}\right)$$

**+**

**Effective UL noise power**

$$\zeta_t^2 \sim \mathcal{O}(1)$$

FL tasks with non-**IID datasets** and partial clients participation converge at rate $\mathcal{O}\left(\frac{1}{t}\right)$.

We cannot adopt model differential for downlink due to partial participation.

# Experiment

- Noise free (ideal)

Under same budget

- Equal power allocation
- $\mathcal{O}(t^2)$-increased power allocation



CIFAR-10 IID

CIFAR-10 non-IID