

# Random Orthogonalization Design

Xizixiang Wei

[C] Xizixiang Wei, Cong Shen, Jing Yang and H. Vincent Poor, Random Orthogonalization for Federated Learning in Massive MIMO Systems, in Proc. IEEE International Conference on Communications (ICC), May 2022.

[J] Xizixiang Wei, Cong Shen, Jing Yang and H. Vincent Poor, Random Orthogonalization for Federated Learning in Massive MIMO Systems, IEEE Transactions on Wireless Communications, 2023.

# Background: Federated Learning

**1**

Machine learning enables powerful applications

**2**

Massive real-world data are generated at edge devices

**3**

We want to keep sensitive data at edge devices

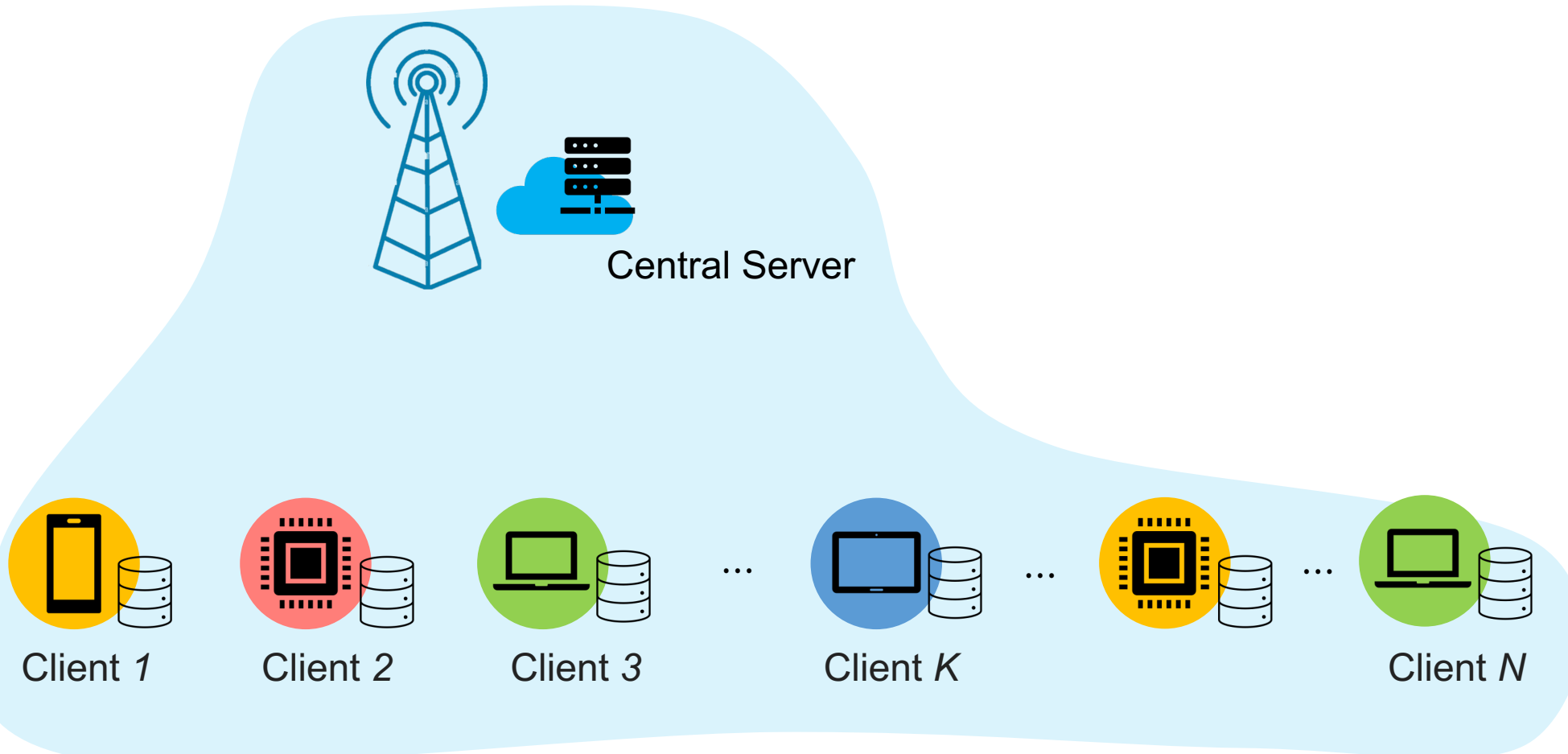


## Federated Learning (FL)

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

# FedAvg

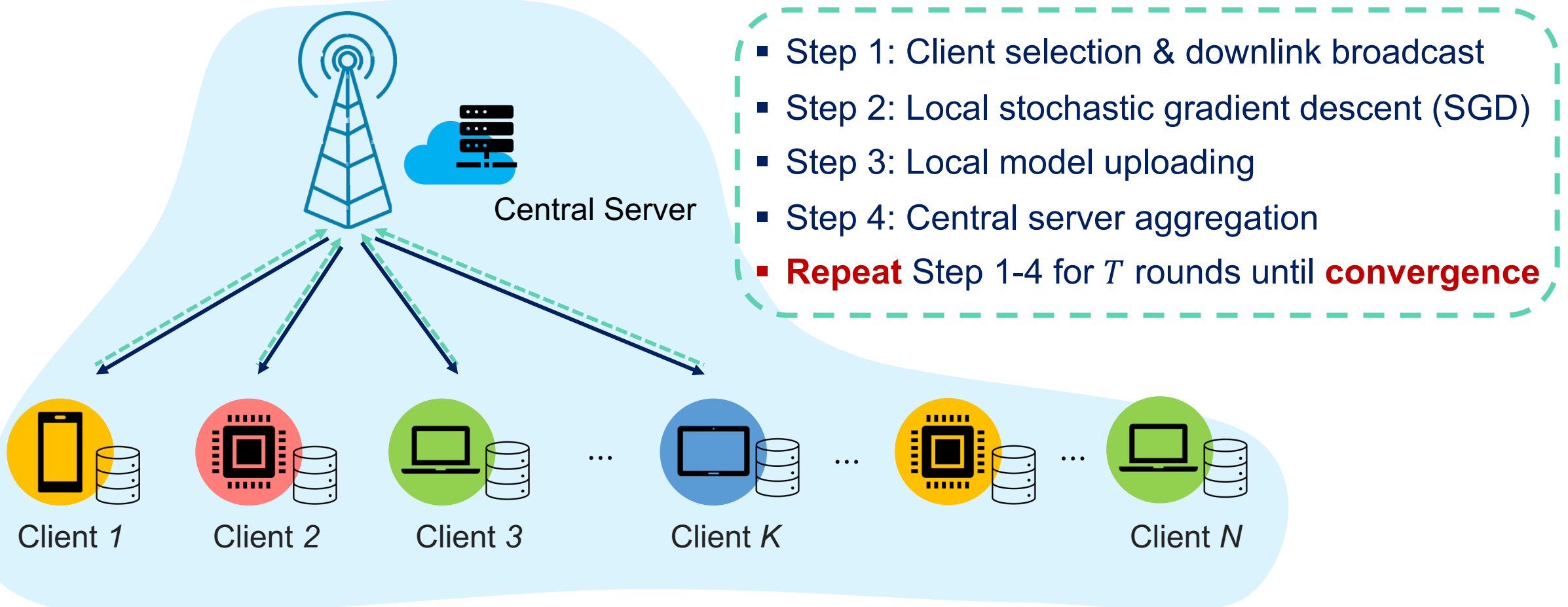
Goal: central server obtain a **global model** trained by **local data** at total  $N$  clients.



Ref: McMahan, Brendan, et al. "Communication-efficient learning of deep networks from decentralized data." *Artificial intelligence and statistics*. PMLR, 2017.

# FedAvg

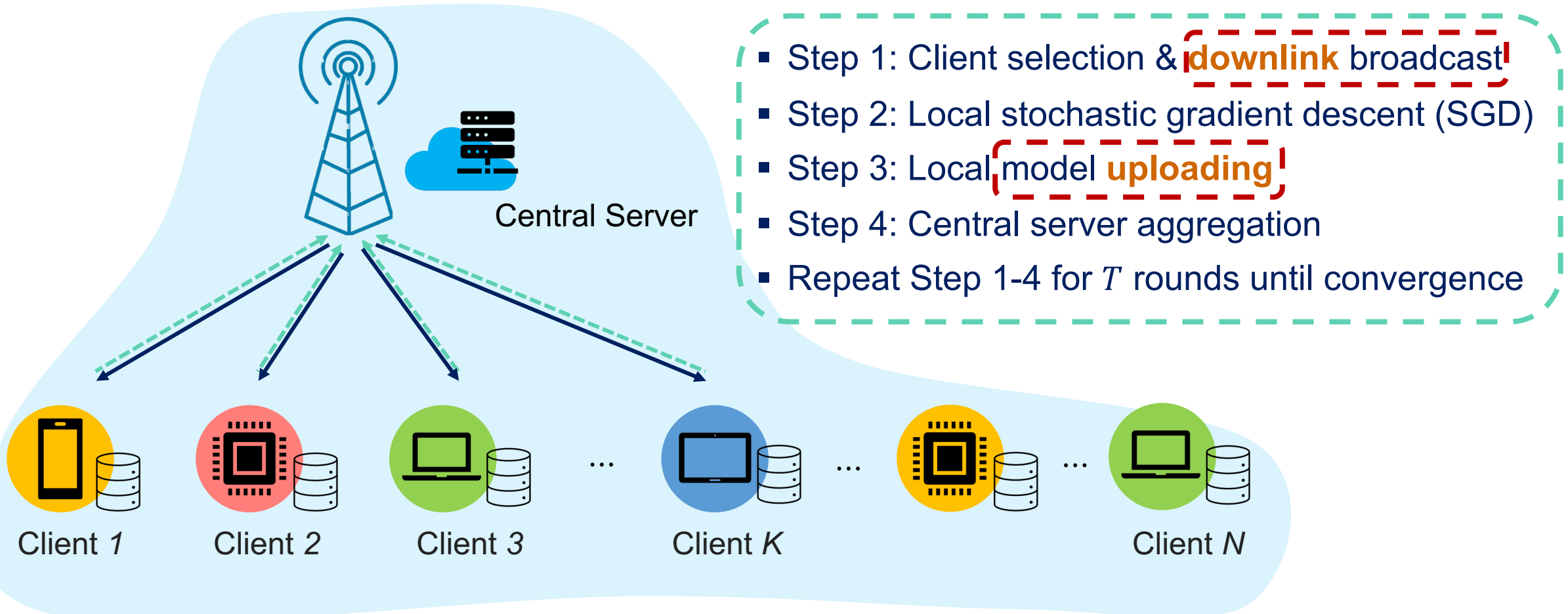
FedAvg: a composition of **multiple learning rounds**  
Each learning round contains **4 steps**.



Ref: McMahan, Brendan, et al. "Communication-efficient learning of deep networks from decentralized data." *Artificial intelligence and statistics*. PMLR, 2017.

# FedAvg

Communication is the **bottleneck** of federated learning.

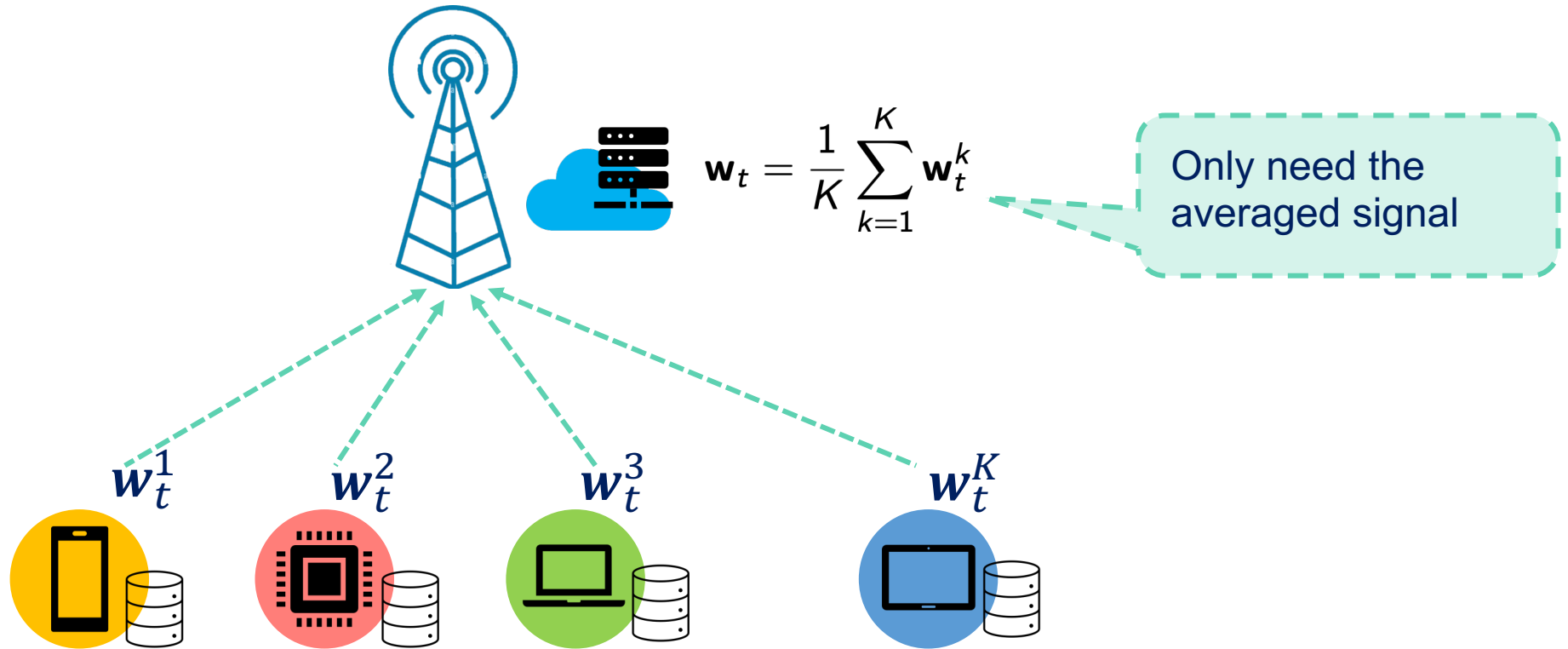


Ref: McMahan, Brendan, et al. "Communication-efficient learning of deep networks from decentralized data." *Artificial intelligence and statistics*. PMLR, 2017.

# Motivation

Main difference from traditional communications

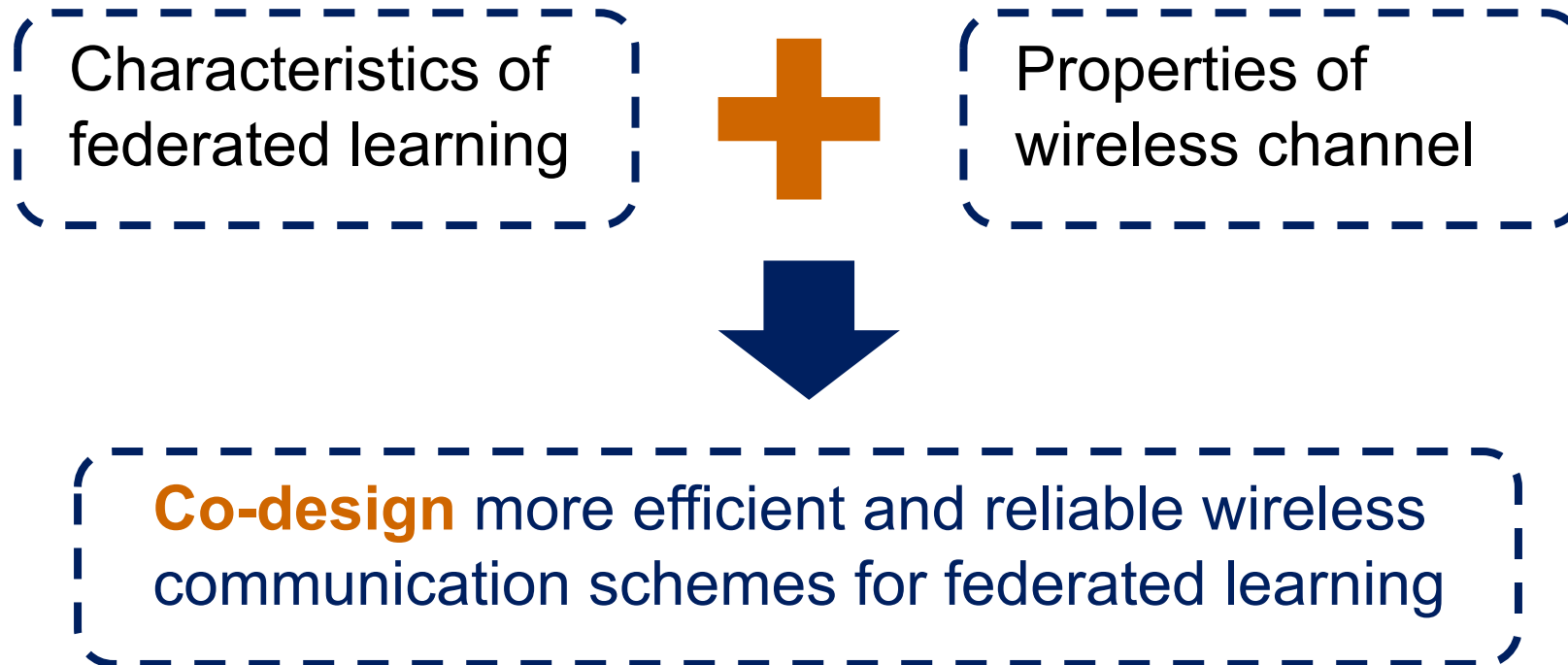
- Central server does **not** need to decode **individual model** in uplink comm.



# Motivation

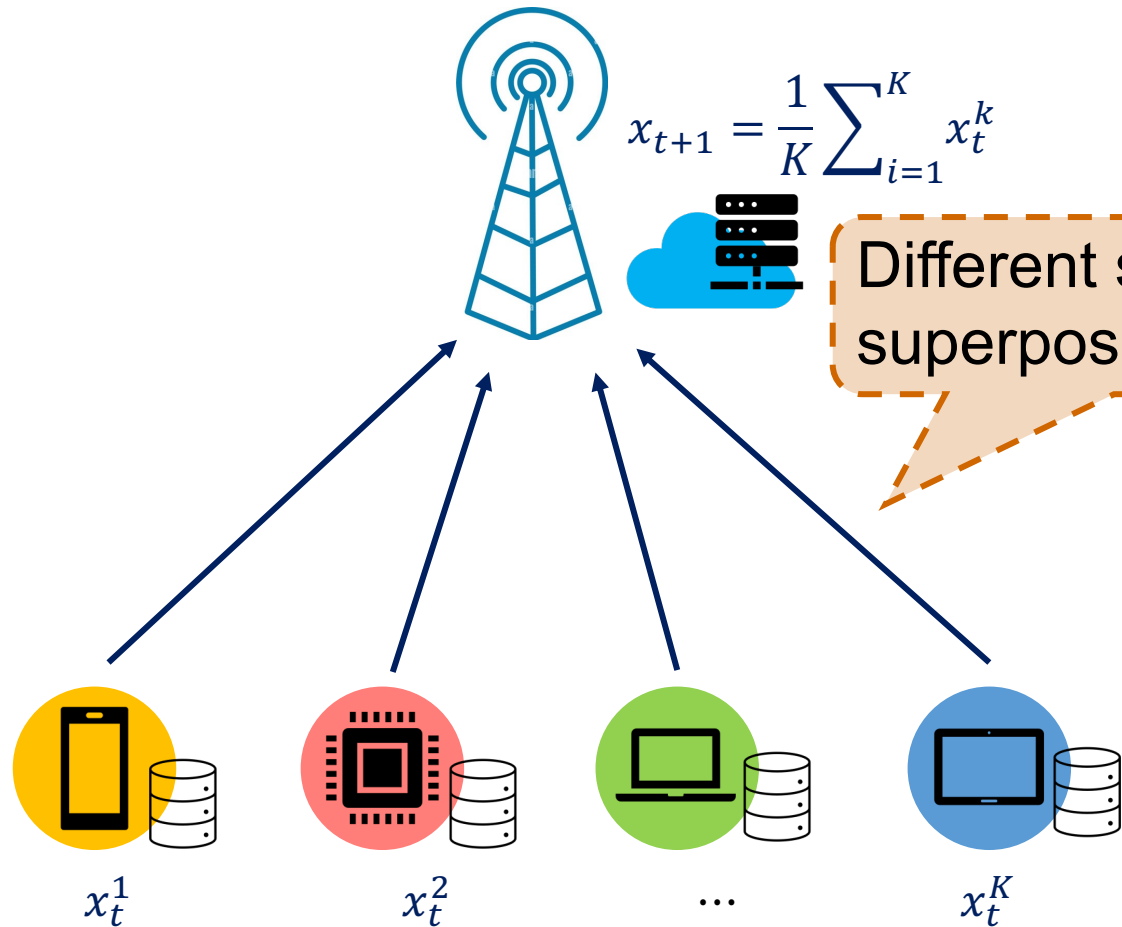
Main difference from traditional communications

- Central server does **not** need to decode **individual model** in uplink comm.



# Uplink of FL: scaling challenges

Over-the-Air Computation (**AirComp**) is a promising solution.



- All clients can be scheduled at the same time-frequency resource

$$x_{t+1} = y = \sum_{i=1}^K x_t^k$$

- In real wireless channel

$$\hat{x}_{t+1} = y = \sum_{i=1}^K h_k x_t^k + n$$

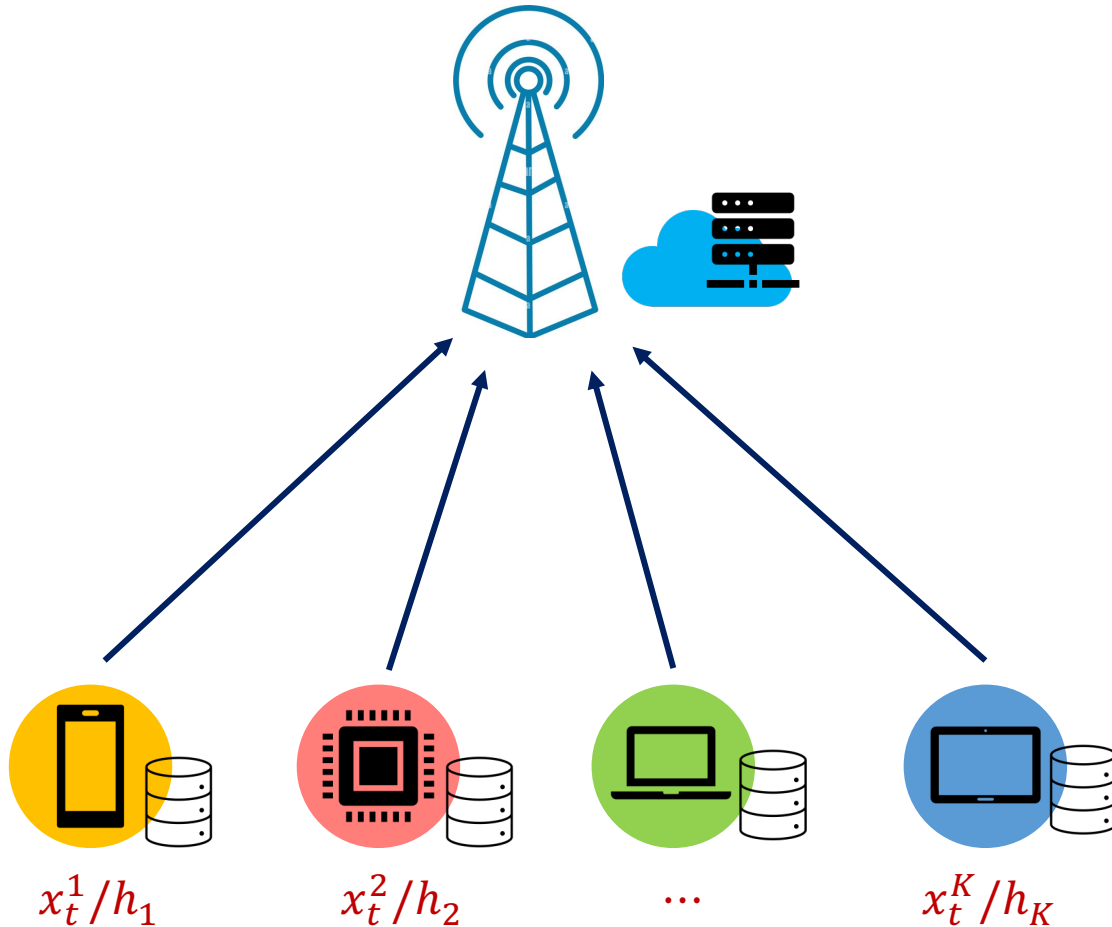
Biased estimation

Ref: Guangxu Zhu, Yong Wang, and Kaibin Huang. "Broadband analog aggregation for low-latency federated edge learning." IEEE Transactions on Wireless Communications, 2019.



# Uplink of FL: AirComp

Heuristic channel inversion



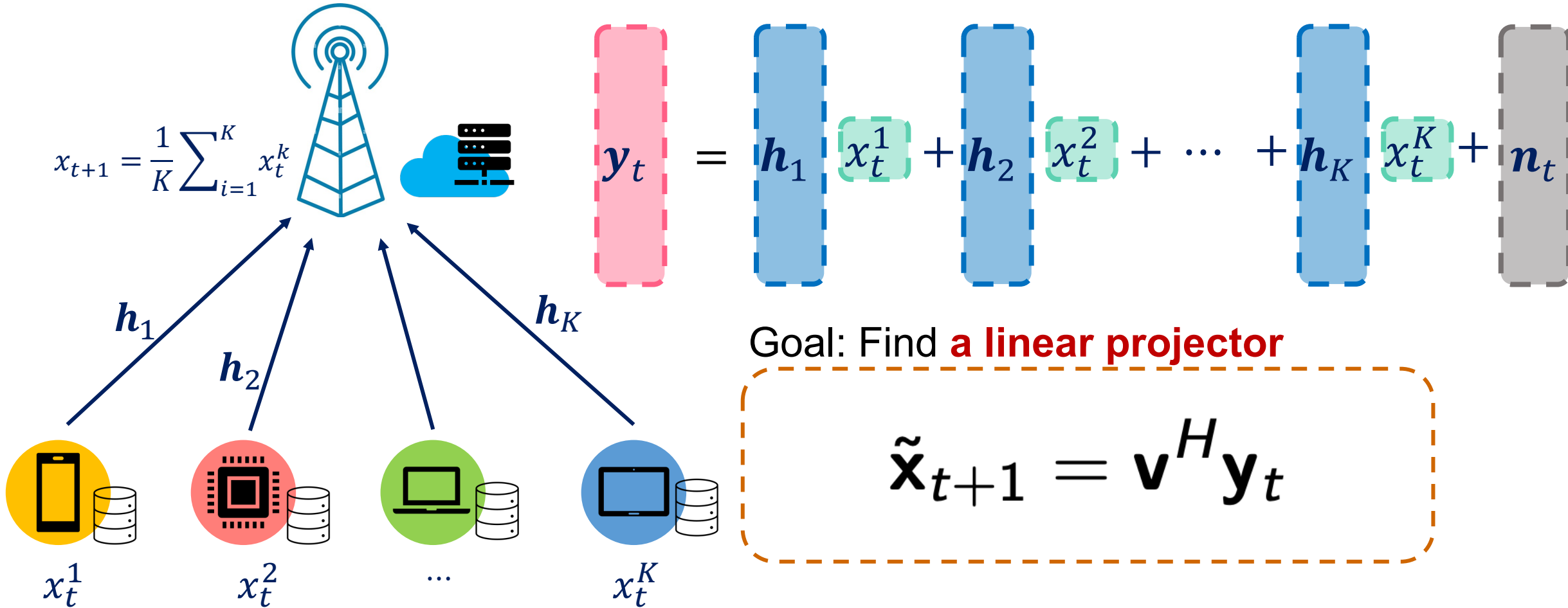
$$\hat{x}_{t+1} = y = \sum_{i=1}^K h_k \frac{x_t^k}{h_k} + n = \sum_{i=1}^K x_t^k + n$$

- Require channel state information at transmitters (CSIT)
- Increasing dynamics of signal
- Performance will “blow up” when deep fading

Ref: Guangxu Zhu, Yong Wang, and Kaibin Huang. "Broadband analog aggregation for low-latency federated edge learning." IEEE Transactions on Wireless Communications, 2019.

# Uplink of FL: AirComp + MIMO

Solution: Using **high-dimension**  $h_k \in \mathbb{C}^{M \times 1}$  provided by **massive MIMO**



# Channel Hardening and Favorable Propagation

IID Rayleigh fading channel model  $\mathbf{h}_k \sim CN(0, \frac{1}{M} \mathbf{I})$

Channel hardening

$$\mathbf{h}_k^H \mathbf{h}_k \rightarrow 1, \text{ as } M \rightarrow \infty.$$

Favorable propagation

$$\mathbf{h}_k^H \mathbf{h}_j \rightarrow 0, \text{ as } M \rightarrow \infty, \forall k \neq j.$$

Massive MIMO  $\xrightarrow{M \rightarrow \infty}$  Random Orthogonalization

Linear projector: sum channel

$$\mathbf{v} = \mathbf{h}_s = \sum_{k=1}^K \mathbf{h}_k$$

# Random Orthogonalization

Linear projection  $\mathbf{h}_s^H \mathbf{y}$ : an unbiased estimation

$$\tilde{x}_i = \mathbf{h}_s^H \mathbf{y}_i = \sum_{k \in [K]} \mathbf{h}_k^H \sum_{k \in [K]} \mathbf{h}_k x_{k,i} + \sum_{k \in [K]} \mathbf{h}_k^H \mathbf{n}_i$$

$$\stackrel{(a)}{=} \underbrace{\sum_{k \in [K]} \mathbf{h}_k^H \mathbf{h}_k x_{k,i}}_{\text{Signal}} + \underbrace{\sum_{k \in [K]} \sum_{j \in [K], j \neq k} \mathbf{h}_k^H \mathbf{h}_j x_{j,i}}_{\text{Interference}} + \underbrace{\sum_{k \in [K]} \mathbf{h}_k^H \mathbf{n}_i}_{\text{Noise}}$$

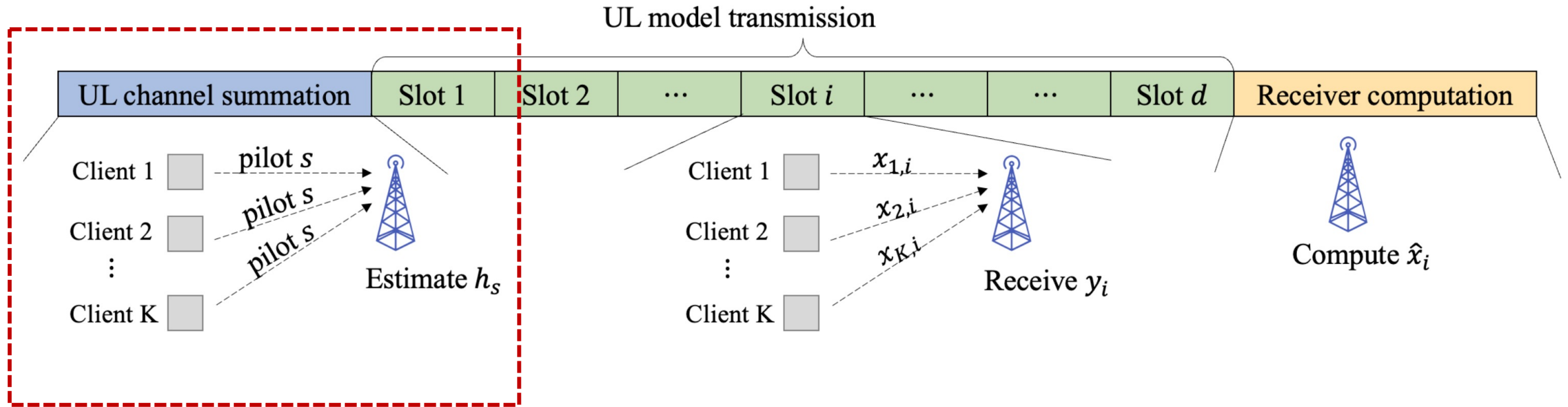
$$\stackrel{(b)}{\approx} \sum_{k \in [K]} x_{k,i}, \quad \forall i = 1, \dots, d.$$

$\mathbf{h}_s^H \mathbf{y}$  is an unbiased estimator of sum signal

$$\mathbf{h}_k^H \mathbf{h}_k \rightarrow 1, \text{ as } M \rightarrow \infty.$$

$$\mathbf{h}_k^H \mathbf{h}_j \rightarrow 0, \text{ as } M \rightarrow \infty,$$

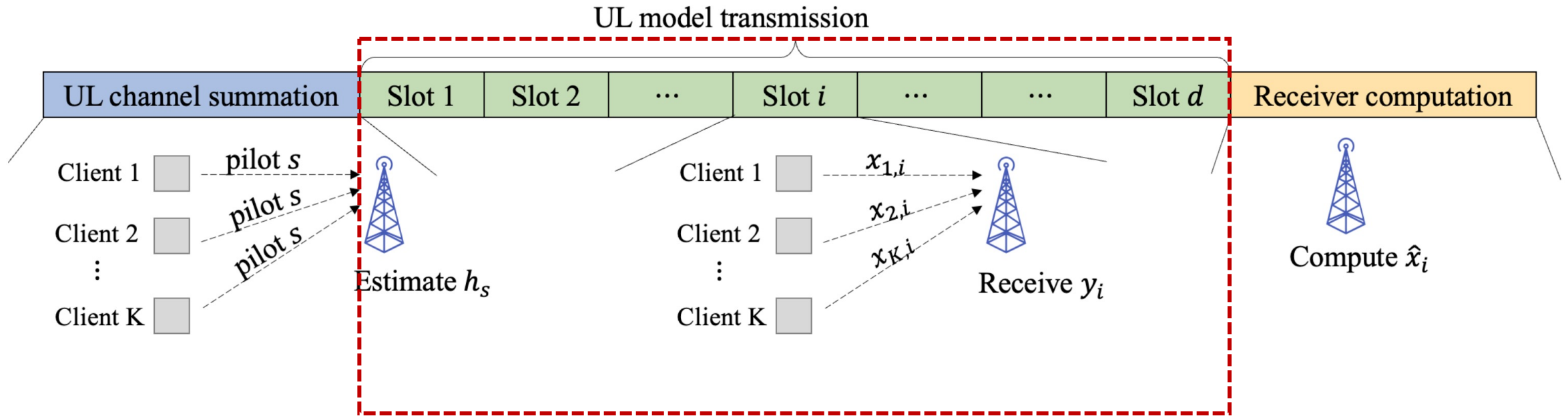
# RO: UL Design



$$\mathbf{y}_s = \sum_{k=1}^K \mathbf{h}_k s + \mathbf{n}_s \xrightarrow{\text{Estimate}} \mathbf{h}_s = \sum_{k=1}^K \mathbf{h}_k$$

- **Partial** CSI at the receiver (**CSIR**)
- Low communication overhead

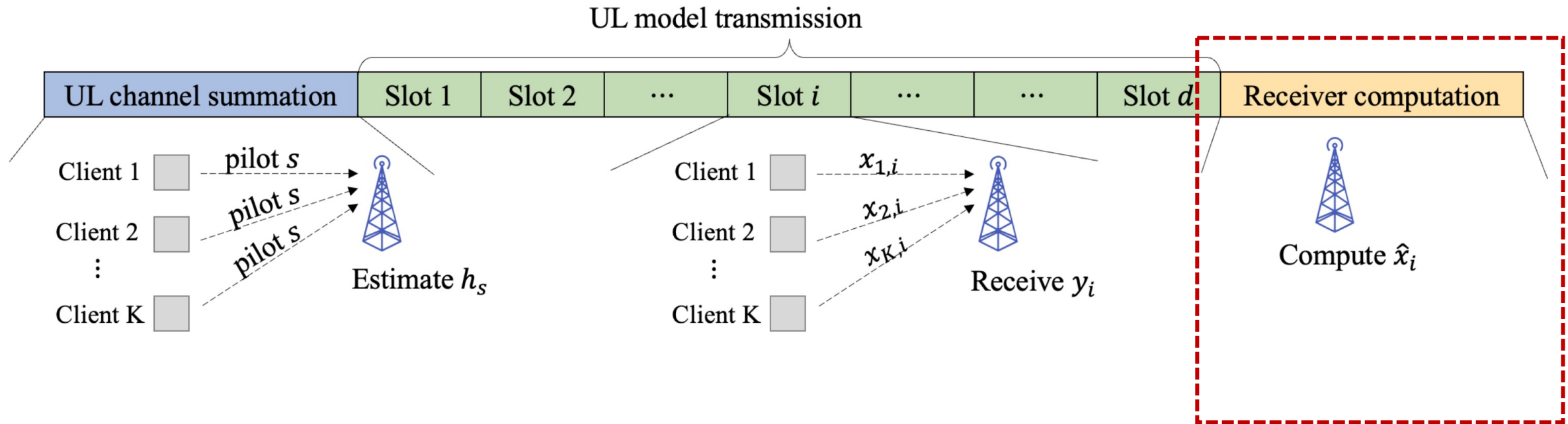
# RO: UL Design



$$y_t = h_1 x_t^1 + h_2 x_t^2 + \dots + h_K x_t^K + n_t$$

▪ No CSIT required

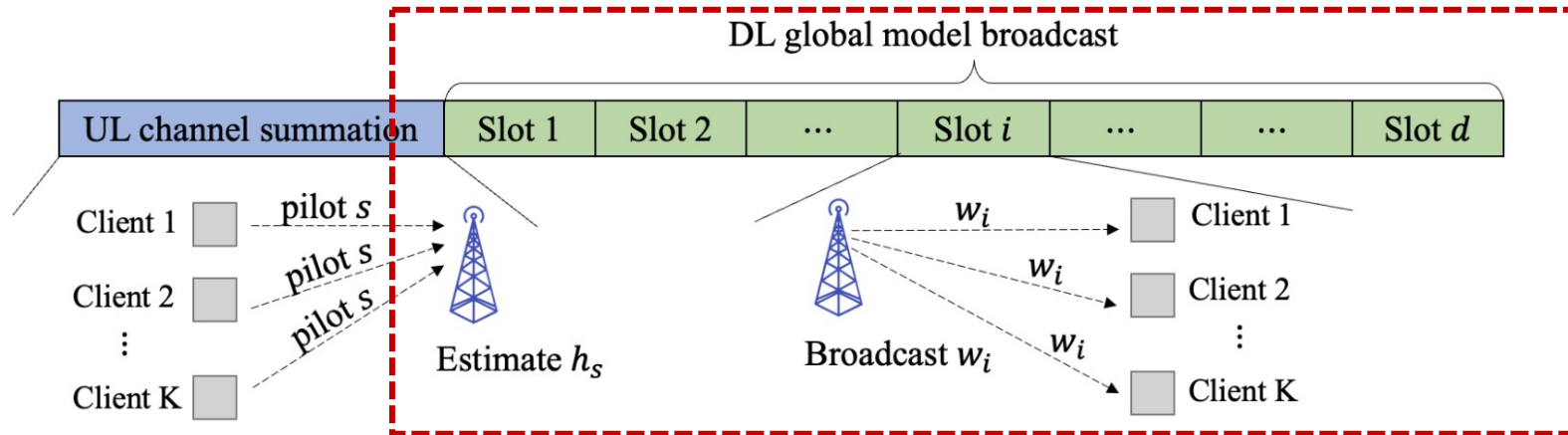
# RO: UL Design



$$\tilde{\mathbf{x}}_{t+1} = \mathbf{v}^H \mathbf{y}_t$$

- Low computational complexity

# RO: DL Design



Using  $h_s$  as the precoder: an efficient broadcast scheme

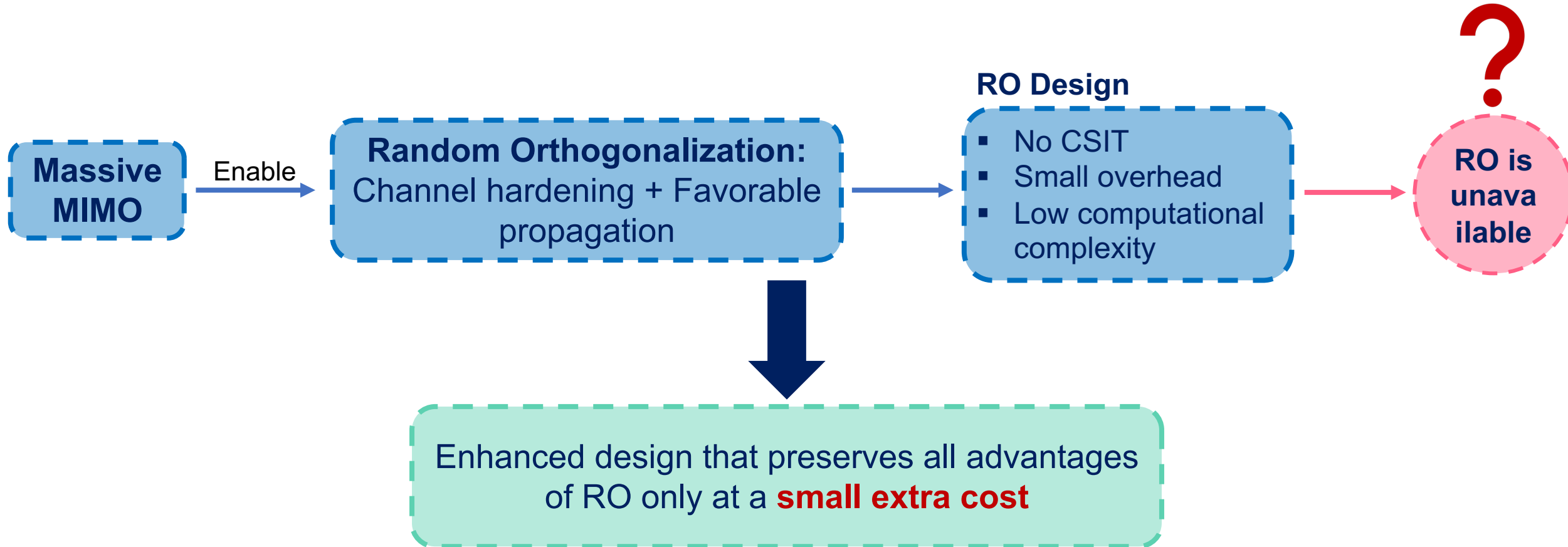
$$y_k = \mathbf{h}_k^H \mathbf{h}_s w_i + z_i^k \stackrel{(a)}{=} \underbrace{\mathbf{h}_k^H \mathbf{h}_k}_{\text{Signal}} w_i + \underbrace{\sum_{j \in [K], j \neq k} \mathbf{h}_k^H \mathbf{h}_j}_{\text{Interference}} w_i + \underbrace{z_i^k}_{\text{Noise}} \stackrel{(b)}{\approx} w_i \quad \forall i = 1, \dots, d.$$

$$\mathbf{h}_k^H \mathbf{h}_k \rightarrow 1, \text{ as } M \rightarrow \infty.$$

$$\mathbf{h}_k^H \mathbf{h}_j \rightarrow 0, \text{ as } M \rightarrow \infty,$$



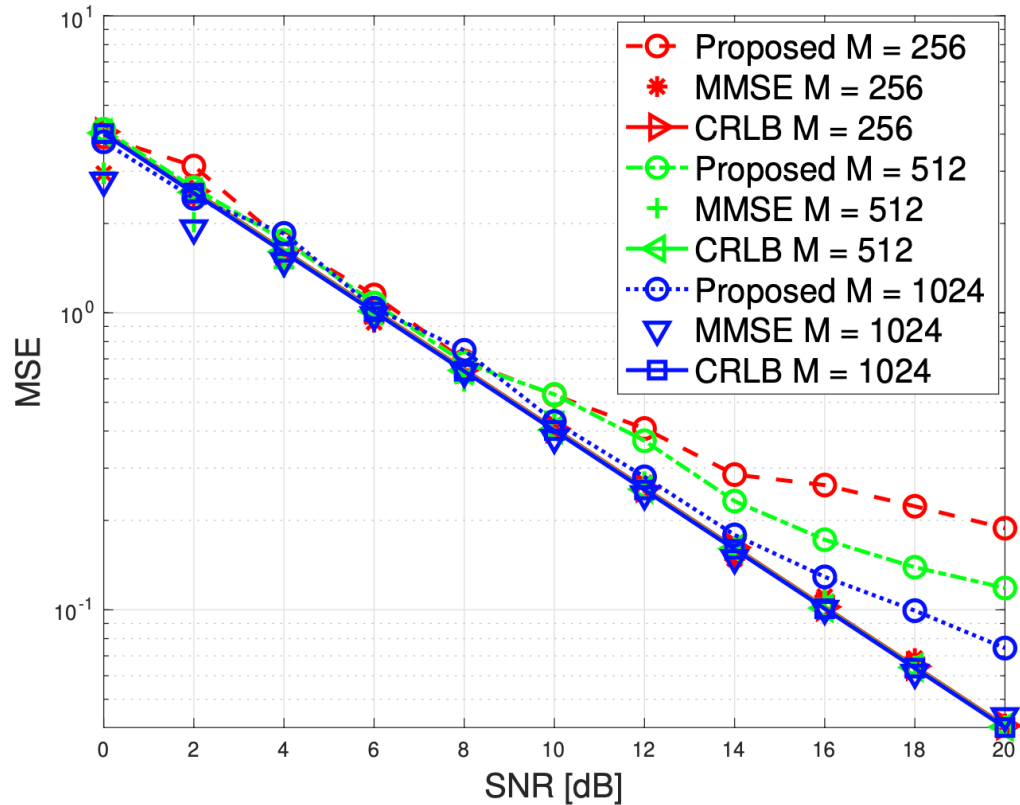
# Summary and Enhanced Design



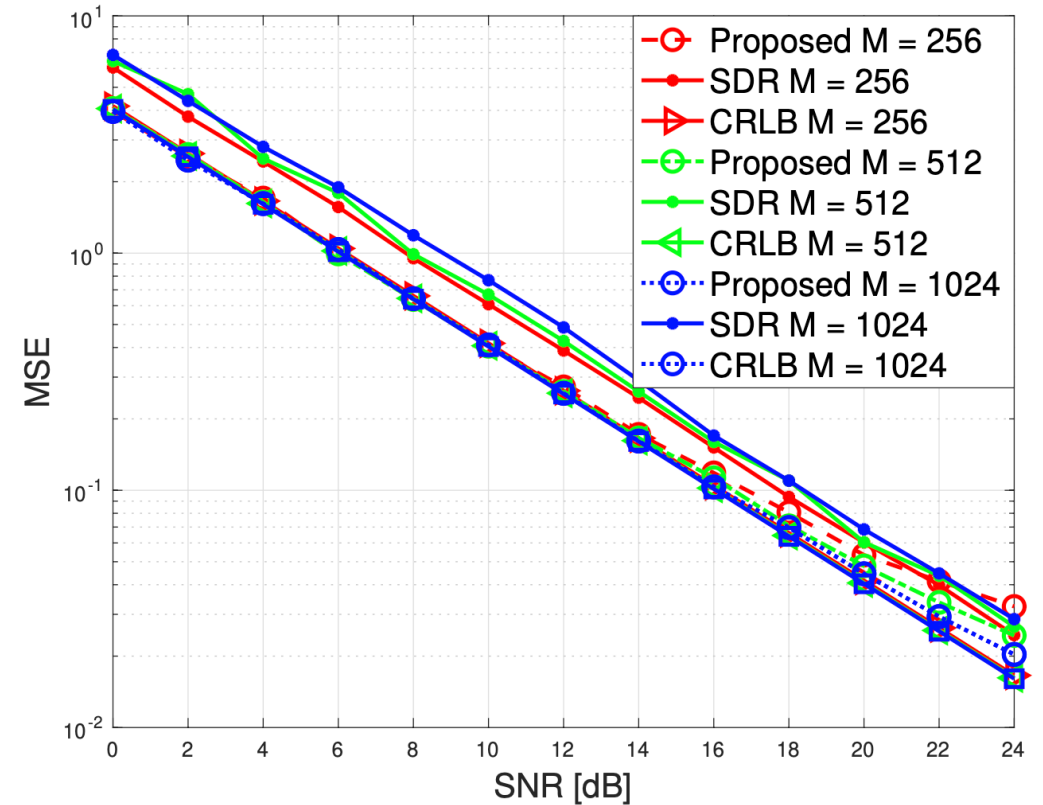
# Performance

- Communication performance

## Uplink

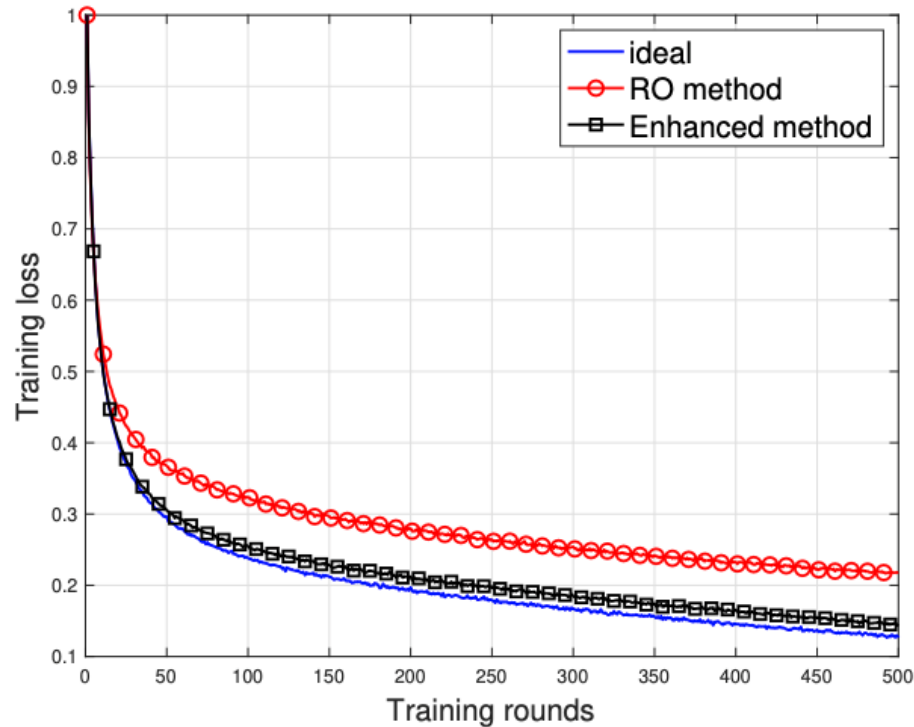


## Downlink

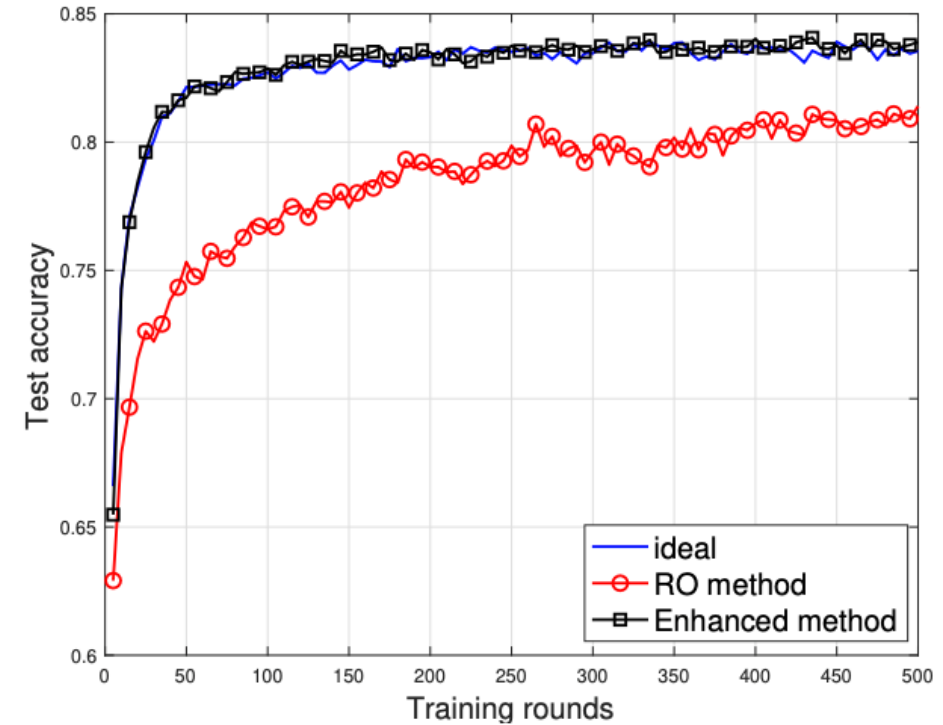


# Performance

- Learning performance



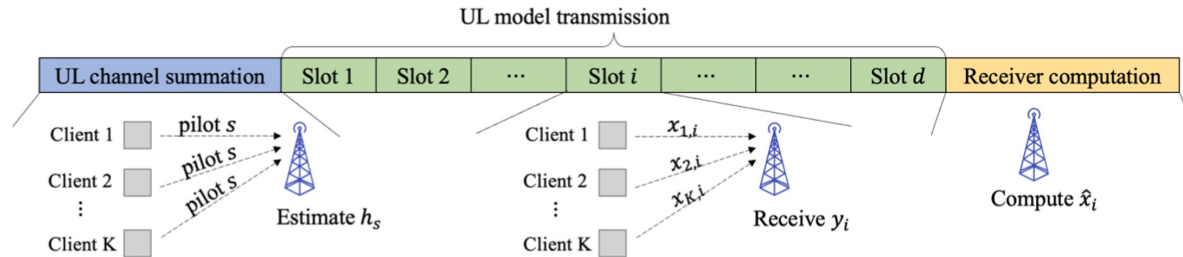
(e) CIFAR-10 uplink+downlink



(f) CIFAR-10 uplink+downlink

**Backup**

# Random Orthogonalization: Uplink Design



- **Step 1: Uplink channel summation**

All clients transmit a common pilot signal  $s$  synchronously. The received signal at the BS is

$$\mathbf{y}_s = \sum_{k \in [K]} \mathbf{h}_k s + \mathbf{n}_s,$$

so that the BS can estimate the summation channel  $\mathbf{h}_s \triangleq \sum_{k \in [K]} \mathbf{h}_k$

- **Step 2: Uplink model transmission**

All clients transmit model differential parameters to the BS in  $d$  shared time slots.

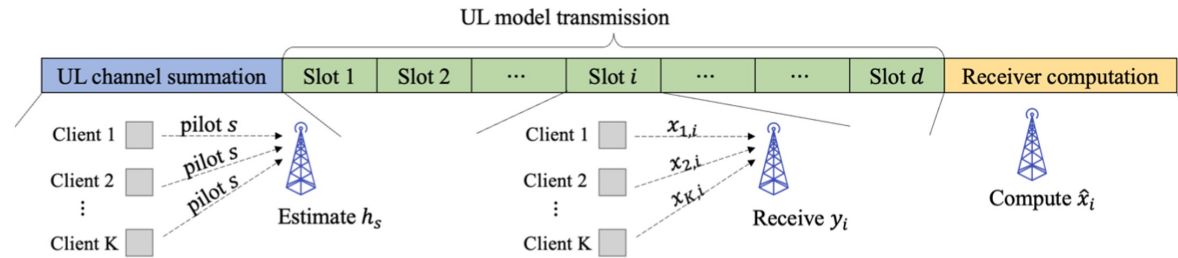
$$\mathbf{y}_i = \sum_{k \in [K]} \mathbf{h}_k x_{k,i} + \mathbf{n}_i, \quad \forall i = 1, \dots, d.$$

- **Step 3: Receiver computation**

The BS estimates aggregated parameter via a simple **linear projection** operation:

$$\tilde{x}_i = \mathbf{h}_s^H \mathbf{y}_i = \sum_{k \in [K]} \mathbf{h}_k^H \sum_{k \in [K]} \mathbf{h}_k x_{k,i} + \sum_{k \in [K]} \mathbf{h}_k^H \mathbf{n}_i$$

# Random Orthogonalization: Uplink Design



- Linear projection: **an unbiased estimation**

$$\tilde{x}_i = \mathbf{h}_s^H \mathbf{y}_i = \sum_{k \in [K]} \mathbf{h}_k^H \sum_{k \in [K]} \mathbf{h}_k x_{k,i} + \sum_{k \in [K]} \mathbf{h}_k^H \mathbf{n}_i$$

$$\stackrel{(a)}{=} \underbrace{\sum_{k \in [K]} \mathbf{h}_k^H \mathbf{h}_k x_{k,i}}_{\text{Signal}} + \underbrace{\sum_{k \in [K]} \sum_{j \in [K], j \neq k} \mathbf{h}_k^H \mathbf{h}_j x_{j,i}}_{\text{Interference } \bullet} + \underbrace{\sum_{k \in [K]} \mathbf{h}_k^H \mathbf{n}_i}_{\text{Noise}}$$

$$\stackrel{(b)}{\approx} \sum_{k \in [K]} x_{k,i}, \quad \forall i = 1, \dots, d.$$

## Advantages:

Only require partial CSIT

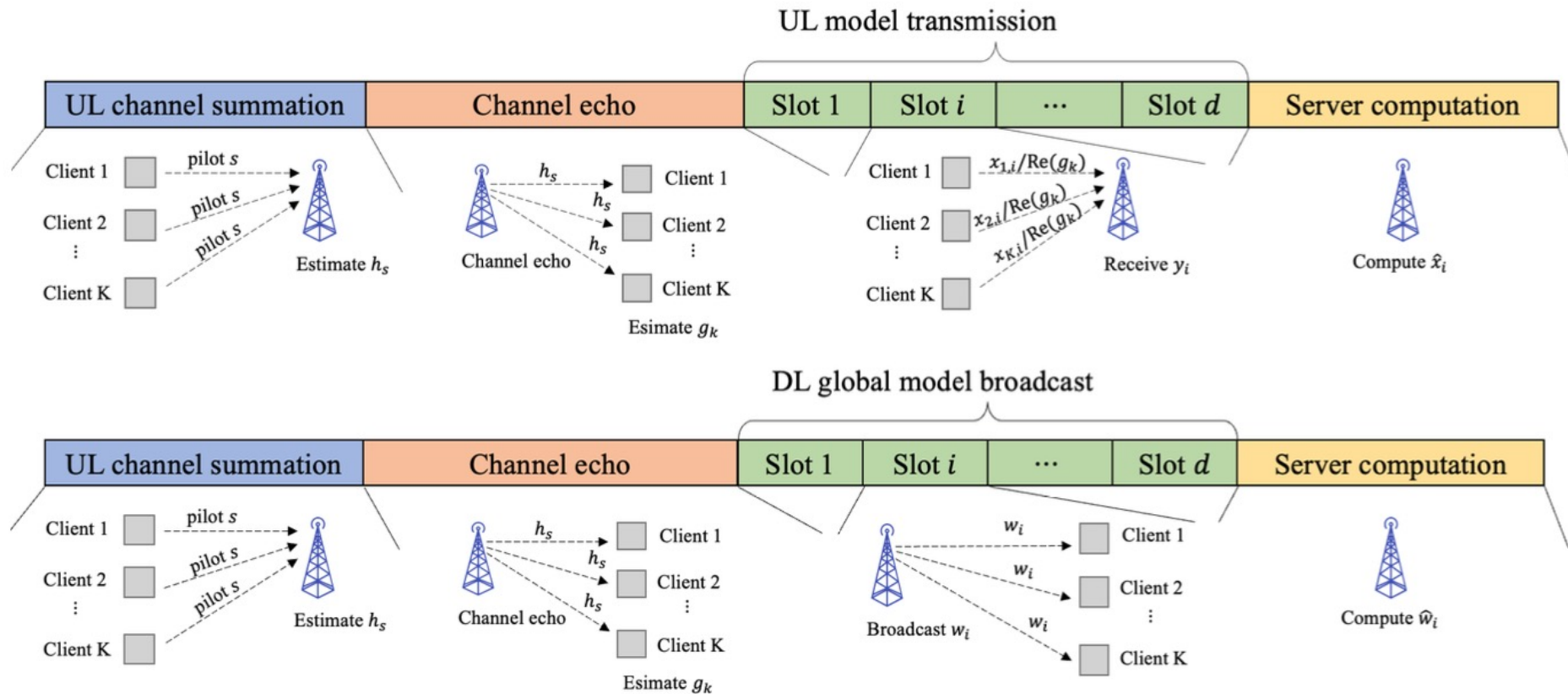
Extremely low complexity

No individual parameter decoded

$$\mathbf{h}_k^H \mathbf{h}_k \rightarrow 1, \text{ as } M \rightarrow \infty.$$

$$\mathbf{h}_k^H \mathbf{h}_j \rightarrow 0, \text{ as } M \rightarrow \infty,$$

# Enhanced Design



**Channel echo:**  $g_k = \mathbf{h}_k^H \mathbf{h}_s$

# Convergence Analysis

**Assumption 1.  $L$ -smooth:**  $\forall \mathbf{v}$  and  $\mathbf{w}$ ,  $\|f_k(\mathbf{v}) - f_k(\mathbf{w})\| \leq L \|\mathbf{v} - \mathbf{w}\|$ ;

**Assumption 2.  $\mu$ -strongly convex:**  $\forall \mathbf{v}$  and  $\mathbf{w}$ ,  $\langle f_k(\mathbf{v}) - f_k(\mathbf{w}), \mathbf{v} - \mathbf{w} \rangle \geq \mu \|\mathbf{v} - \mathbf{w}\|^2$ ;

**Assumption 3. Bounded variance for unbiased mini-batch SGD:**  $\forall k \in [N]$ ,

$$\mathbb{E}[\nabla \tilde{f}_k(\mathbf{w})] = \nabla f_k(\mathbf{w}) \quad \text{and} \quad \mathbb{E} \left\| \nabla f_k(\mathbf{w}) - \nabla \tilde{f}_k(\mathbf{w}) \right\|^2 \leq H_k^2;$$

**Assumption 4. Uniformly bounded gradient:**  $\forall k \in [N]$ ,  $\mathbb{E} \left\| \nabla \tilde{f}_k(\mathbf{w}) \right\|^2 \leq H^2$  for all mini-batch data.

Preserve  $O\left(\frac{1}{T}\right)$   
convergence rate  
of SGD

**Theorem 1 (Convergence for random orthogonalization).**

With Assumptions 1-4, for some  $\gamma \geq 0$ , if we select the learning rate as  $\eta_t = \frac{2}{\mu(t+\gamma)}$ , we have

$$\mathbb{E}[f(\mathbf{w}_t)] - f^* \leq \frac{L}{2(t+\gamma)} \left[ \frac{4B}{\mu^2} + (1+\gamma) \|\mathbf{w}_0 - \mathbf{w}^*\|^2 \right], \quad (14)$$

for any  $t \geq 1$ , where

$$B \triangleq \left[ 1 + \frac{K}{M} + \frac{1}{\text{SNR}} \right] \frac{H^2}{K}. \quad (15)$$