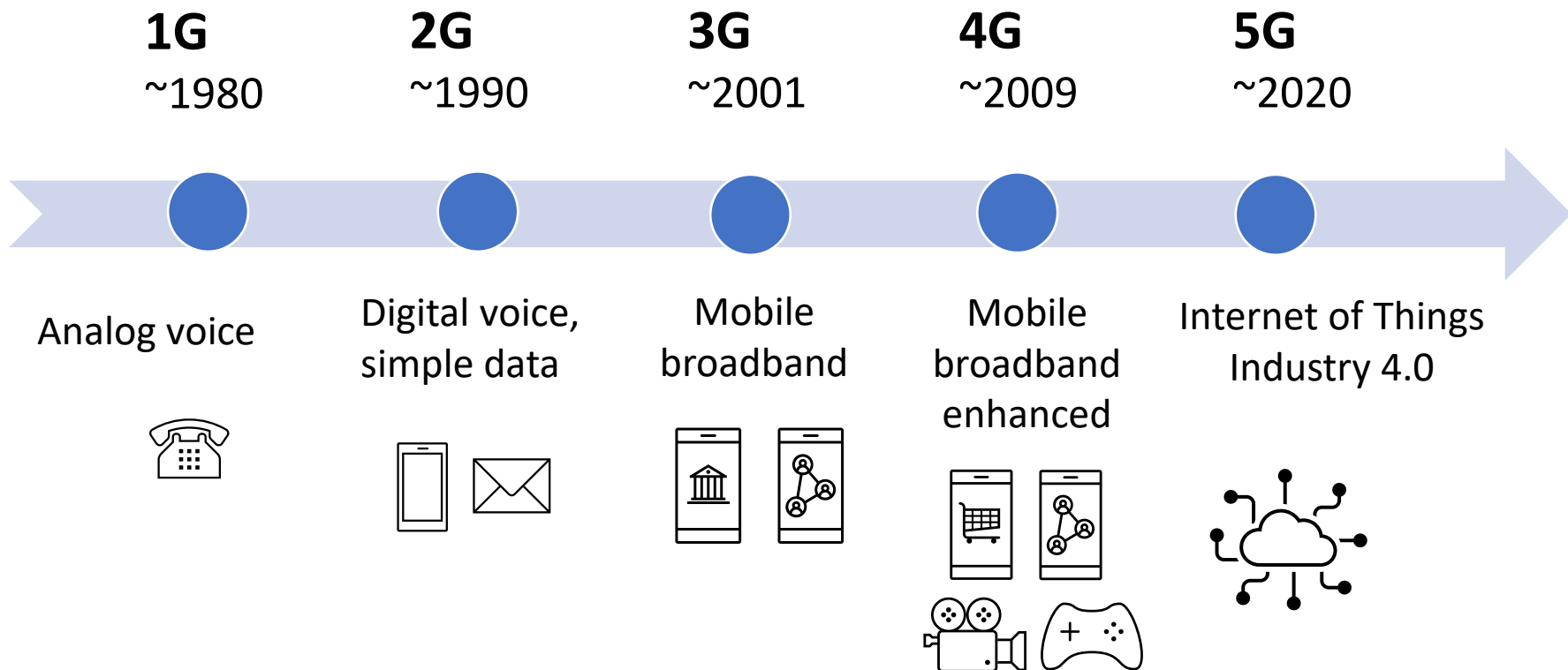


# Distributed Intelligence over Wireless Networks

IEEE SPS - EURASIP 6G Summer School  
2022-09-01

# Evolution of Mobile Standards



# Demand-Driven Wireless Technology



- Text messages
- Video streaming
- Social media
- Online gaming
- Remote surgery
- Status monitoring
- ...

- Data rate
- Connectivity
- Spectral Efficiency
- Energy efficiency
- Latency
- Reliability
- Timing
- ...

# Wireless Communications—Shannon's Legacy

- Shannon-Hartley theorem on channel capacity for memoryless point-to-point AWGN channel:

$$C = B \cdot \log_2 \left( 1 + \frac{S}{N} \right) \text{ bits per second}$$

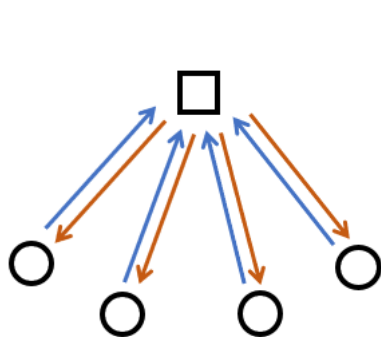
- Shannon's channel coding theorem
  - If transmission rate  $R < C$ , there exists a coding scheme that achieves arbitrarily low error probability
- Wireless communication system is usually designed to achieve the **maximum number of bits** that can be **reliably transmitted** under some resource constraints.

# Wireless Bottleneck in Distributed Systems

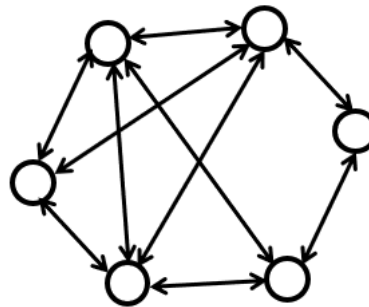
- Point to point



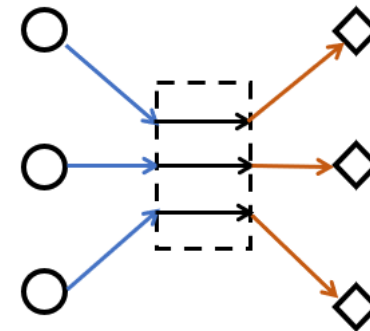
- Distributed systems



(a) Server-based



(b) Peer-to-peer

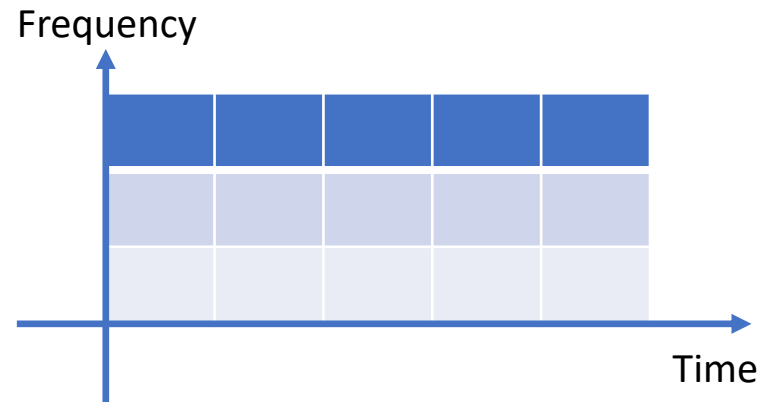


(c) Shared medium access

# Wireless Resource Limitations

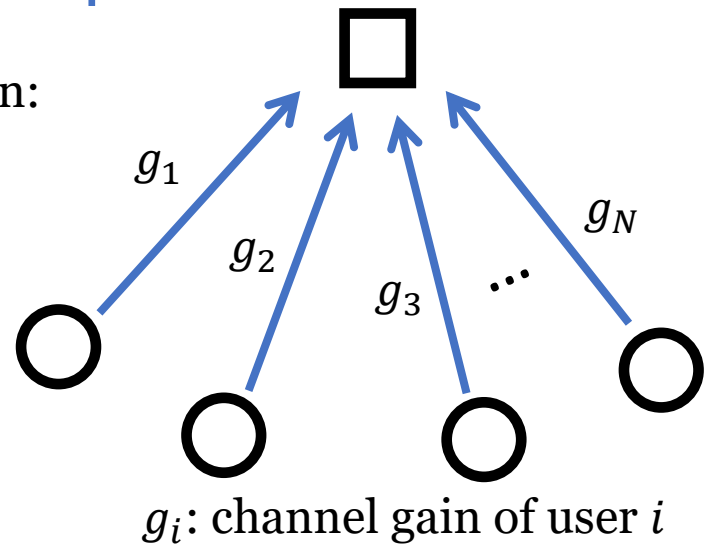
Wireless communication resources:

- Frequency
- Time
- Space



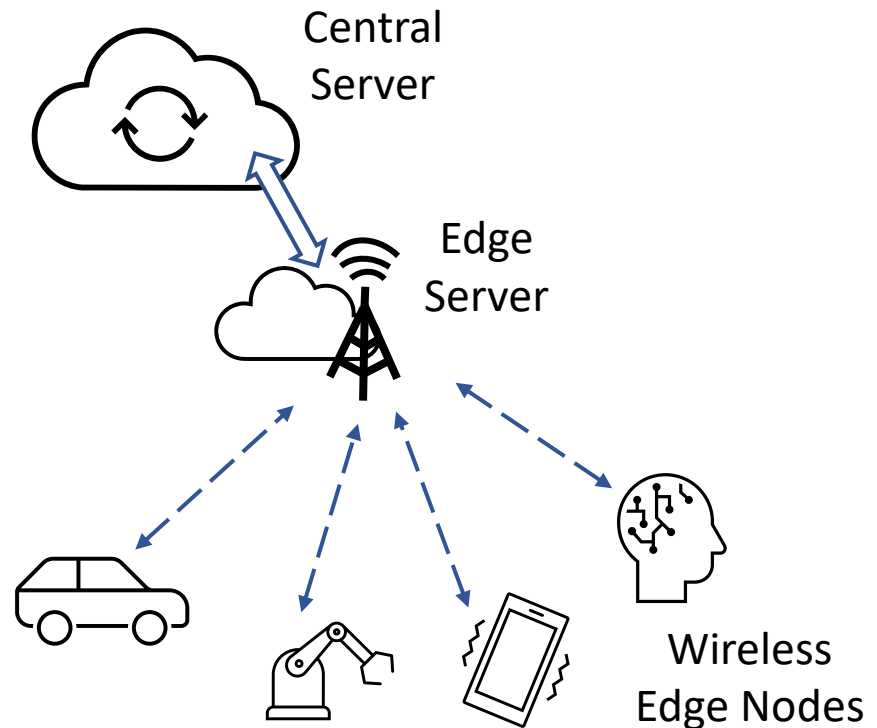
**Rate-oriented** design for resource allocation:

- Maximize  $\sum_{i=1}^N R_i$
- Maximize  $\min_i \{R_i\}$
- Others



# 6G: Era of Connected Intelligence

- For human perception
  - Messages carry specific meaning, must be correctly received
  - Higher rate implies higher QoE
- For machine perception
  - Error-free communication might not be needed
  - Performance depends on the tasks (training, inference, control...)
  - **Goal-oriented** comm. design



# Outline

- ❖ From centralized to distributed machine learning (ML)
- ❖ Optimization for ML
  - Stochastic gradient method
  - Distributed ML meets wireless
- ❖ Federated edge learning
  - Impact of resource allocation
  - Example on scheduling and aggregation design
- ❖ Model aggregation over distributed nodes
  - Over-the-Air (OtA) computation
  - OtA with multiple receivers
  - OtA FL with multiple antennas
  - Applications and challenges of OtA aggregation
- ❖ Summary

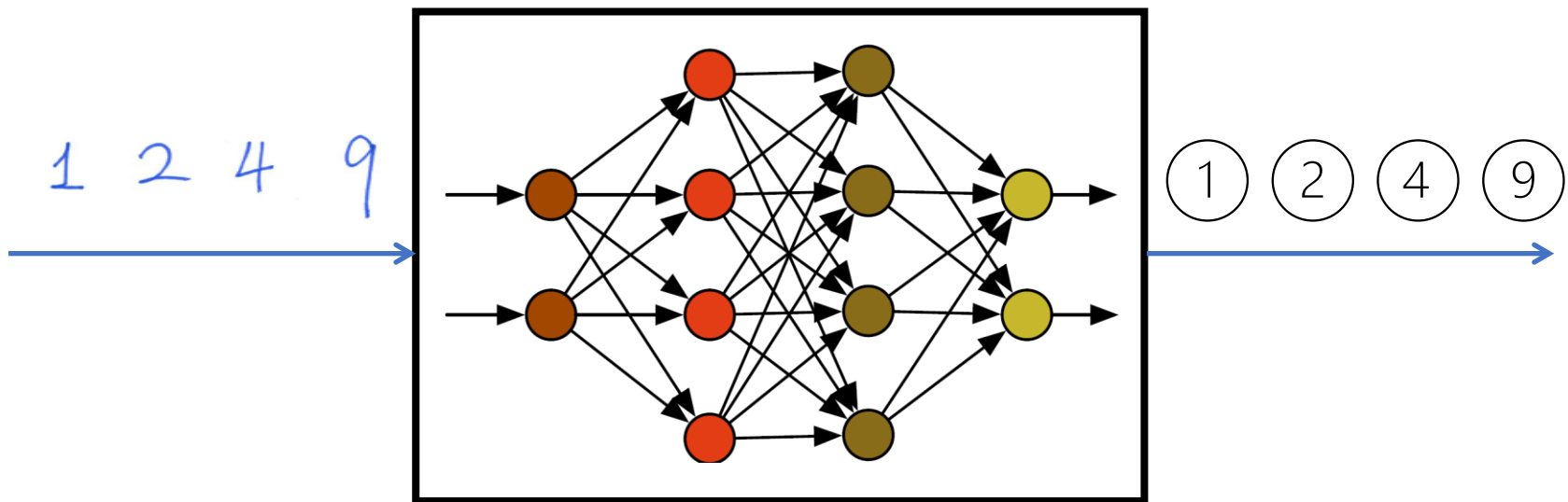


# Classes of ML Algorithms

- Supervised learning
  - Regression
  - Classification ← Focus of this lecture
- Unsupervised learning
  - Clustering
  - Anomaly detection
- Reinforcement learning

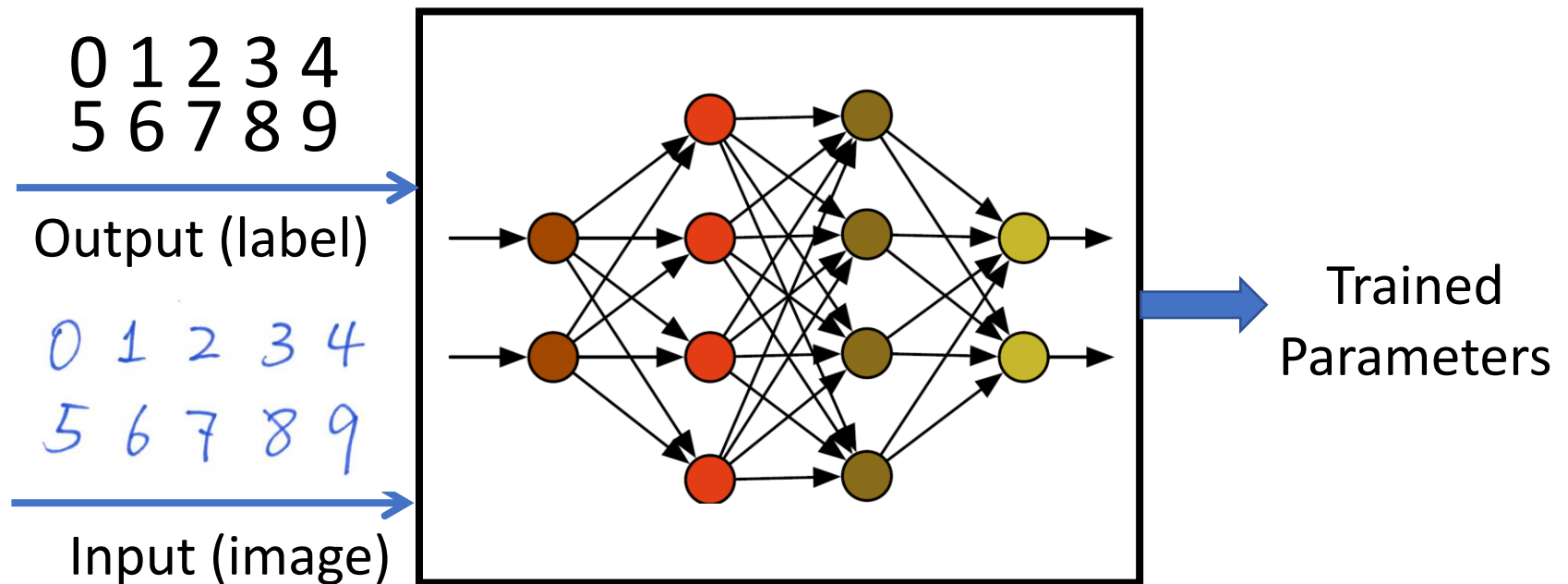
# Example on Supervised Learning

Hand-written digits recognition with deep neural networks

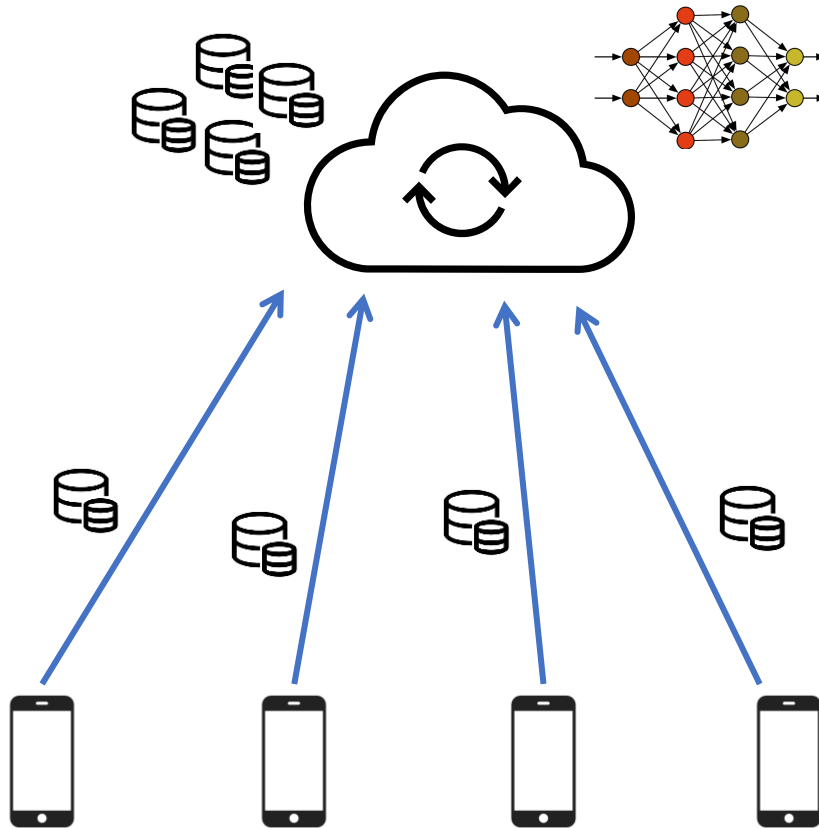


# Model Training

**Training** of deep learning algorithms



# Centralized Training



- Privacy issue
- Uploading cost
- Latency



Server

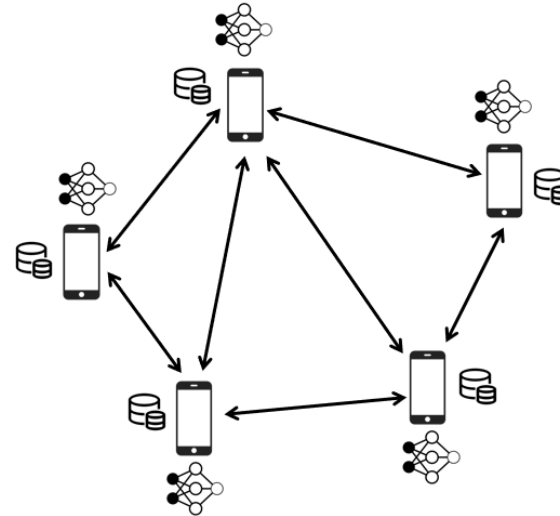
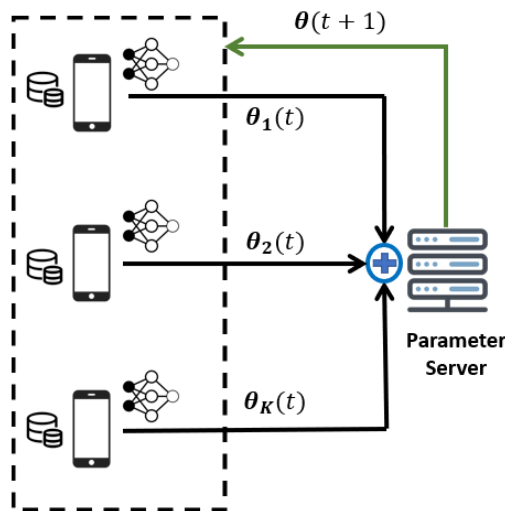


End user devices



Training data

# Collaborative Training with Decentralized Data



- A group of agents  $\mathcal{N} = \{1, \dots, N\}$  collaborate in training a common ML model, parameterized by  $\boldsymbol{\theta} \in \mathbb{R}^d$
- $\mathcal{S}_k$ : local training data at device  $k$ .  $\mathcal{S} = \cup_{k \in \mathcal{N}} \mathcal{S}_k$ : entire dataset in the system
- Optimization objective: find  $\boldsymbol{\theta}^*$  that minimizes the global loss function

$$F(\boldsymbol{\theta}) = \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} l(\boldsymbol{\theta}, x)$$

# Optimization for ML

- Expected risk

$$R(\boldsymbol{\theta}) = \mathbb{E}[f(\boldsymbol{\theta}; \xi)]$$

$\boldsymbol{\theta}$ : model parameters,  $\xi$ : randomness in data,  $f: \mathbb{R}^d \rightarrow \mathbb{R}$

- Empirical risk

$$F(\boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N f_i(\boldsymbol{\theta})$$

- Objective: minimize $_{\boldsymbol{\theta}}$   $F(\boldsymbol{\theta})$
- Stochastic Gradient (SG) method

$$\boldsymbol{\theta}(t+1) = \boldsymbol{\theta}(t) - \alpha_t g(t)$$

$g(t)$ : stochastic gradient vector with  $\mathbb{E}[g(t)] = \nabla F(\boldsymbol{\theta})$ . Example:

$$g(t) = \begin{cases} \nabla f(\boldsymbol{\theta}(t), \xi_t) \\ \frac{1}{n} \sum_{i=1}^n \nabla f(\boldsymbol{\theta}(t), \xi_{t,i}) \end{cases}$$

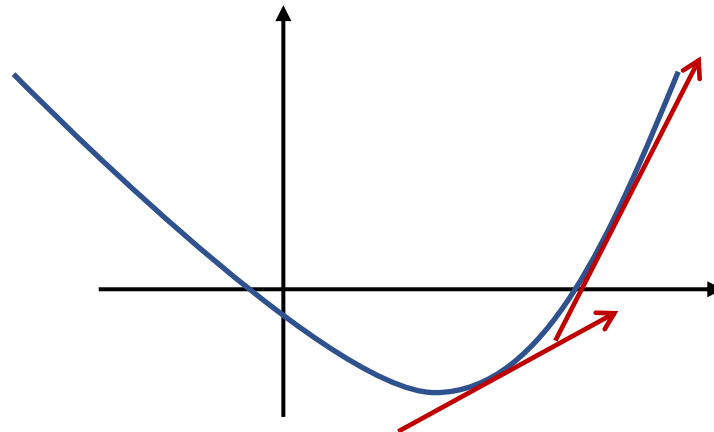
# Assumptions

## ➤ L-smoothness

$F$  is continuously differentiable, and the gradient function  $\nabla F(\boldsymbol{\theta})$  is Lipschitz continuous with constant  $L > 0$ , i.e.,

$$\|\nabla F(\boldsymbol{\theta}_1) - \nabla F(\boldsymbol{\theta}_2)\| \leq L\|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\|$$

Intuitive explanation: the gradient function does not change too quickly



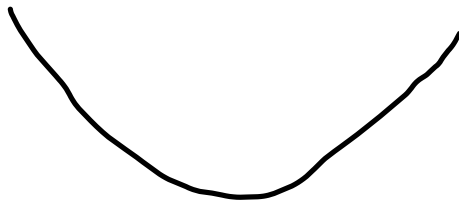
# Assumptions

## ➤ Strong convexity

$F$  is strongly convex, i.e., there exists a constant  $\mu > 0$  such that for any  $\theta_1, \theta_2 \in \mathbb{R}^d$ , we have

$$F(\theta_2) - F(\theta_1) \geq \nabla F(\theta_1)^T (\theta_2 - \theta_1) + \frac{\mu}{2} \|\theta_2 - \theta_1\|^2$$

Intuitive explanation: the objective function is not just convex, but also “curvy” enough



YES



NO



# Assumptions

## ➤ Bound on second-order moment

The second order moment of the stochastic gradient vector satisfies

$$\mathbb{E}_{\xi_t} [\|g(t)\|_2^2] \leq M + M_g \|\nabla F(\boldsymbol{\theta}(t))\|^2$$



$$\mathbb{E}_{\xi_t} [\|g(t) - \nabla F(\boldsymbol{\theta}(t))\|_2^2] \leq M + M_v \|\nabla F(\boldsymbol{\theta}(t))\|^2$$

Standard assumption on bounded variance

$$\mathbb{E} [\|g(t) - \nabla F(\boldsymbol{\theta}(t))\|^2] \leq \sigma^2$$

- not realistic unless parameter space is compact
- often not compatible with strong convexity assumption

# Convergence and Optimality Gap

Let  $\boldsymbol{\theta}^* = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} F(\boldsymbol{\theta})$  and  $F^* = F(\boldsymbol{\theta}^*)$ , under these assumptions, with fixed step size  $\alpha_t = \bar{\alpha}$ , and  $0 < \bar{\alpha} < \frac{1}{LM_G}$ , then

$$\mathbb{E}[F(\boldsymbol{\theta}(t)) - F^*] \leq \frac{\bar{\alpha}LM}{2\mu} + (1 - \bar{\alpha}\mu)^t (F(\boldsymbol{\theta}(0)) - F^*)$$
$$\xrightarrow{t \rightarrow \infty} \frac{\bar{\alpha}LM}{2\mu}$$

- With fixed step size, SGD converges to some neighborhood of  $\boldsymbol{\theta}^*$
- Less "noise" in stochastic gradient computation, smaller optimality gap

# Federated Learning: Collaborative Training over Decentralized Data

Federated Learning is one instance of distributed and parallel SGD with

- Massive and heterogeneous devices
- Non-IID and unbalanced data
- Limited communication

Two well-known problems:

- Straggler issue
- Communication constraints

# Federated Learning at the Wireless Edge

During communication round  $t$



Broadcast global model  $\theta(t)$



On-device training with local data ( $E$  iterations)

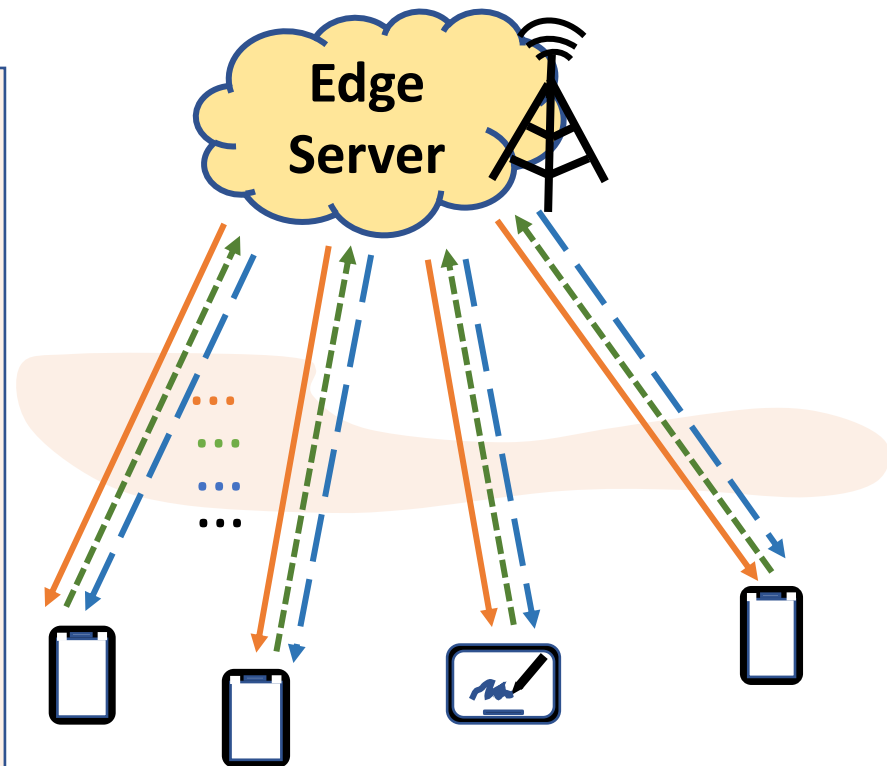
$$\theta_k(t, \tau + 1) = \theta_k(t, \tau) - \eta(t) \nabla F_k(\theta(t, \tau))$$

Upload model update  $\theta_k(t, E)$  to server



Broadcast new global model after weighted average of  $\theta_k(t, E)$

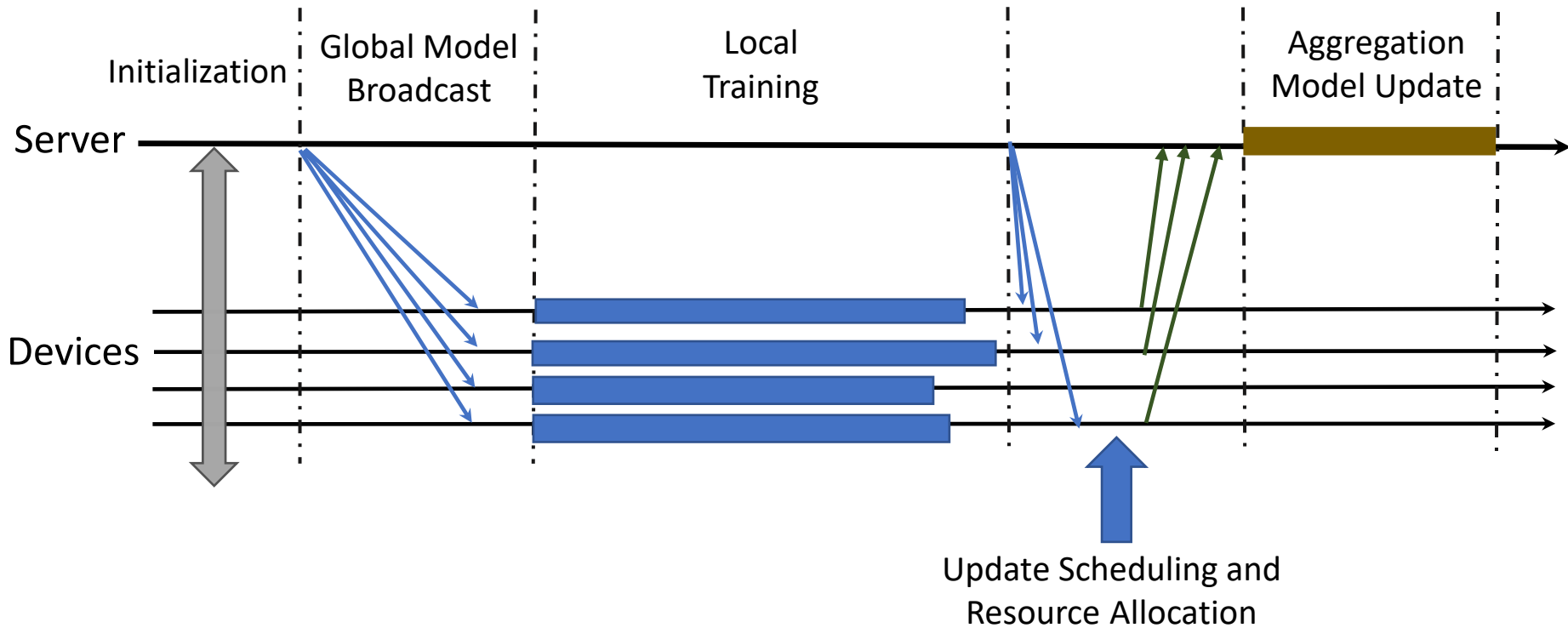
$$\theta(t + 1) = \sum_{\forall k} w_k \theta_k(t, E)$$



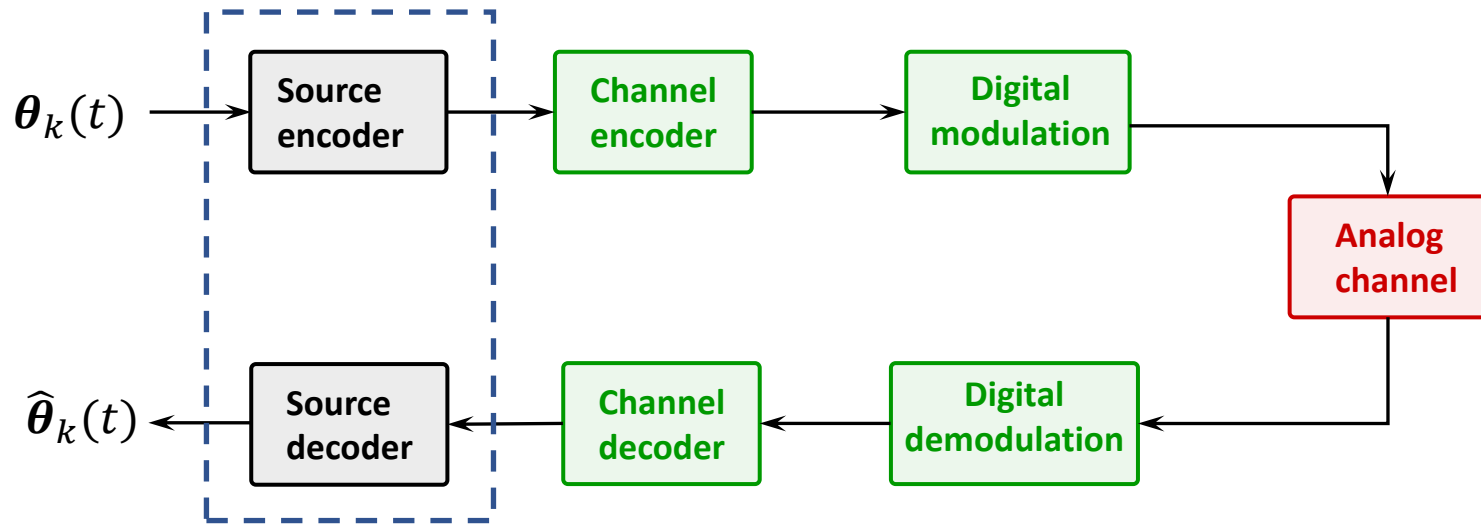
Downlink: model distribution

Uplink: update uploading (**most affected by resource constraints**)

# Training Procedure



# Federated Learning with Digital Transmission

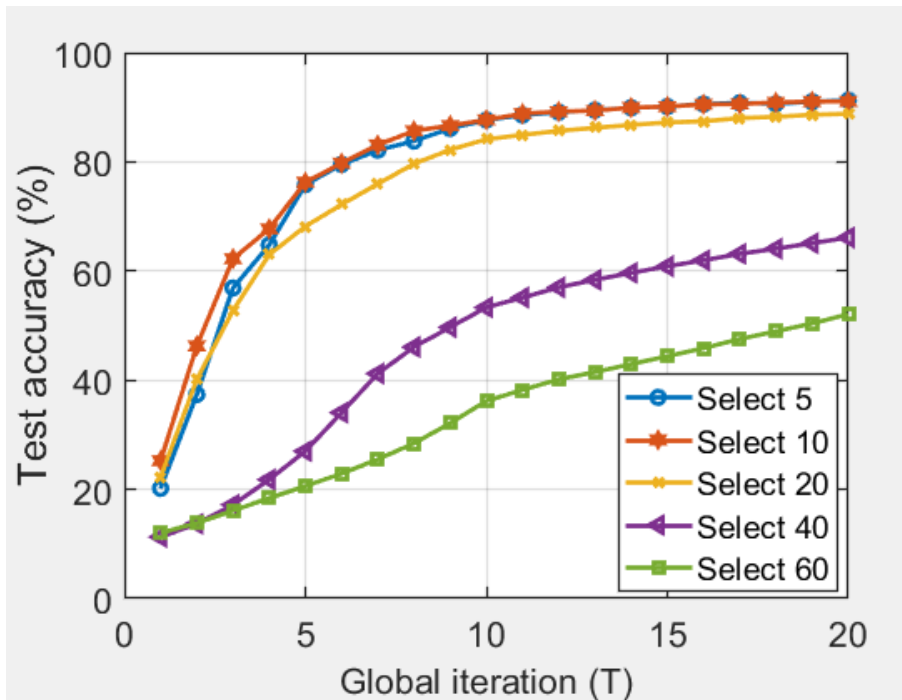


Example: sparsification + quantization

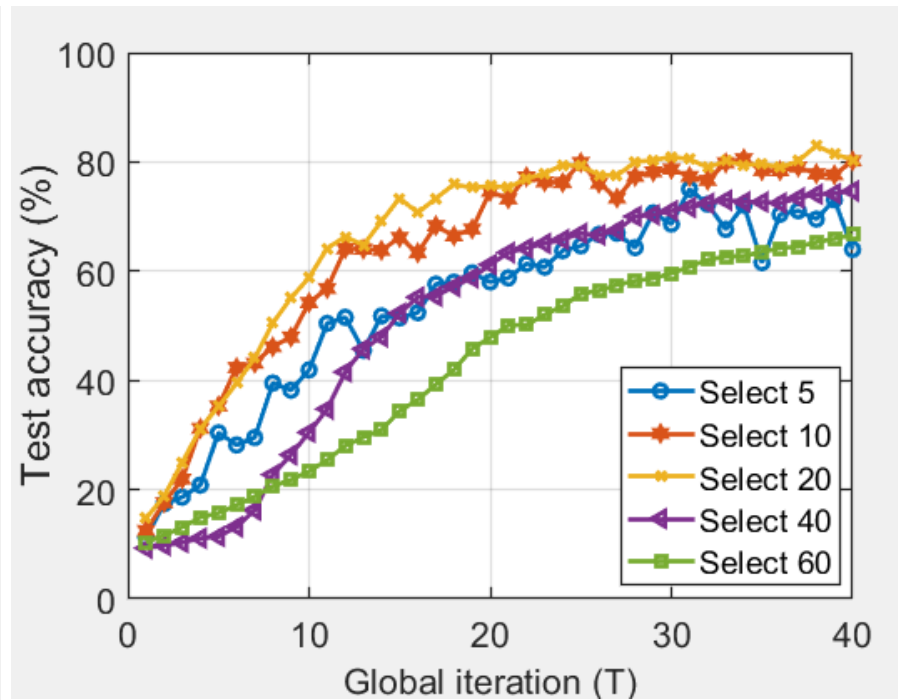
Resource limitation => { Limited number of scheduled users  
Limited number of information bits => Compression loss

# Example: Impact of Resource Allocation

With IID data



With non-IID data

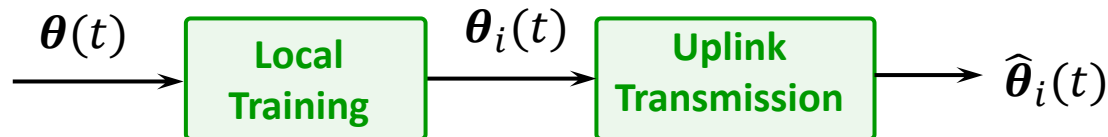


Ideally, schedule users with the most significant impact on the learning result.

# Distributed ML meets Wireless

During uplink transmission of model updates:

- Ideally, we want to obtain  $\boldsymbol{\theta}(\mathcal{K}) = \frac{1}{|\mathcal{K}|} \sum_{i \in \mathcal{K}} \boldsymbol{\theta}_i(t)$ , where  $\mathcal{K}$  is the set of all participating agents
- In practice, we get  $\hat{\boldsymbol{\theta}}(\mathcal{K}_s) = \frac{1}{|\mathcal{K}_s|} \sum_{i \in \mathcal{K}_s} \hat{\boldsymbol{\theta}}_i(t)$ , where  $\mathcal{K}_s$  is the set of scheduled agents





# Distributed ML meets Wireless

Mean Squared Error (MSE) of aggregated model

$$\text{MSE} = \underbrace{\mathbb{E} \left[ \|\hat{\boldsymbol{\theta}}(\mathcal{K}_s) - \mathbb{E}[\hat{\boldsymbol{\theta}}(\mathcal{K}_s)]\|^2 \right]}_{\text{variance}} + \underbrace{\|\mathbb{E}[\hat{\boldsymbol{\theta}}(\mathcal{K}_s)] - \boldsymbol{\theta}(\mathcal{K})\|^2}_{\text{bias}}$$

- Bias: affected by the number of scheduled users and their data representation
- Variance: affected by compression loss (depending on how many bits we can reliably transmit)

# Joint Data- and Channel-aware Scheduling

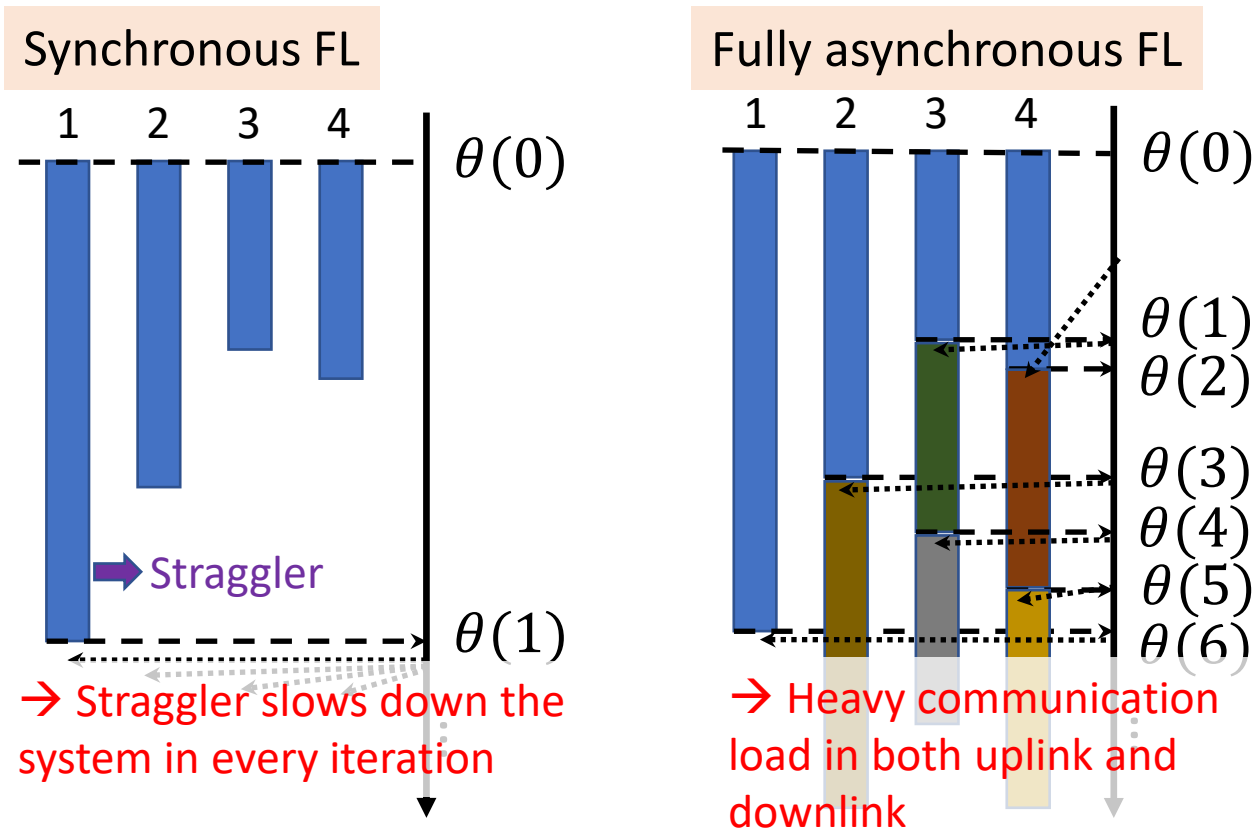
## Motivation:

- Schedule users with more homogeneous data representation to reduce bias
- Prioritize users with better channels to reduce variance

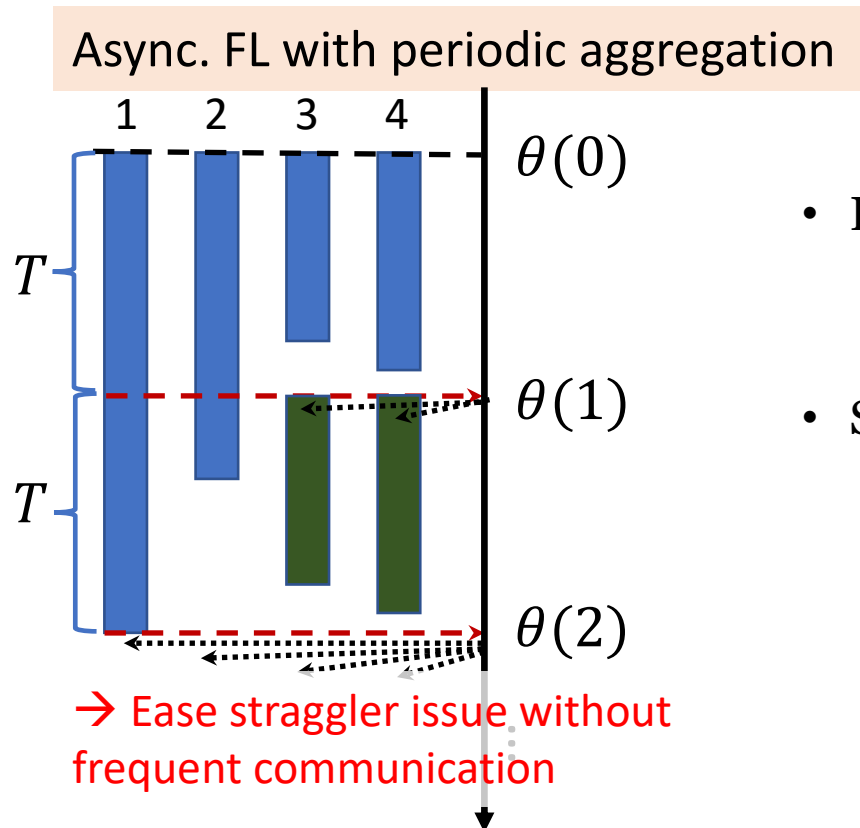
For supervised learning with labelled data, define  $\mathbf{b}_k = [b_k^1, \dots, b_k^L]$  as the label distribution of local data set at user  $k$ .

- First, select  $\Pi'(t) \in \mathcal{K}$  users with largest channel gain
- Then, select a subset  $\Pi(t)$  that minimizes  $\sum_{l=1}^L |\sum_{k \in \Pi(t)} (b_k^l - \bar{b})|^2$

# Synchronous vs. Asynchronous Training



# Asynchronous FL with Periodic Aggregation



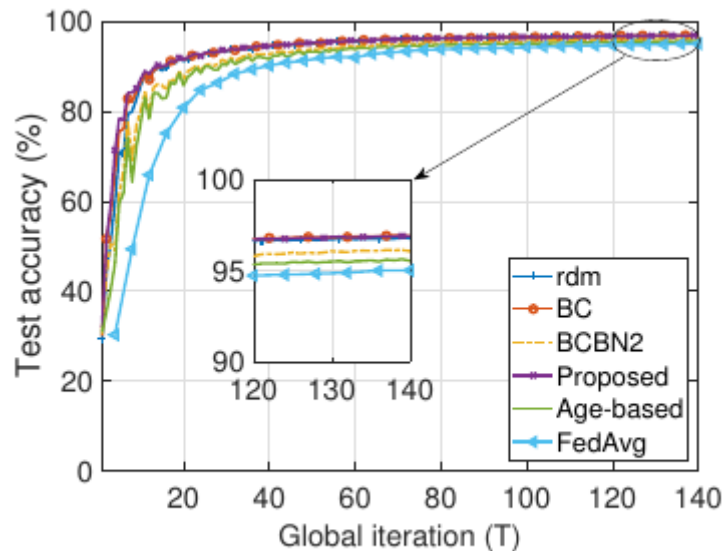
- Problem:
  - Different “age” of previously received global model
- Solution:
  - Age-aware aggregation policy to balance update freshness

# Simulation Results

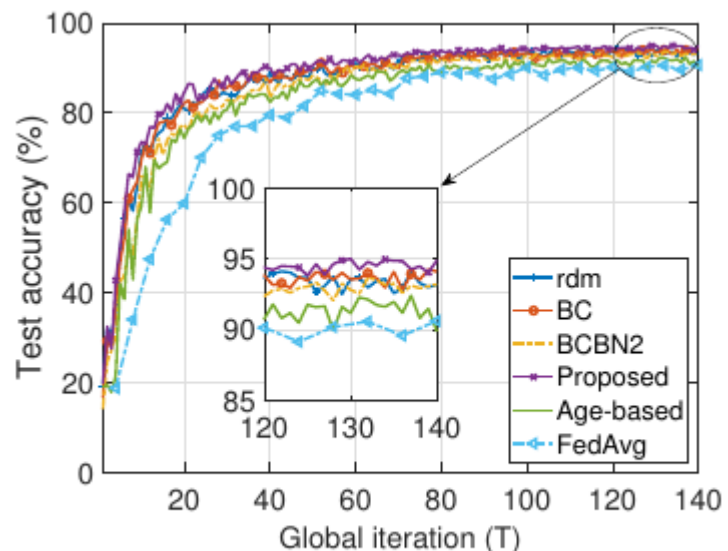
## Simulation setup:

MNIST dataset, CNN, model dimension  $d = 21840$ , Rayleigh fading, 40 users in total, scheduling 20% in every round

Comparison with alternative methods in [5] and [6]

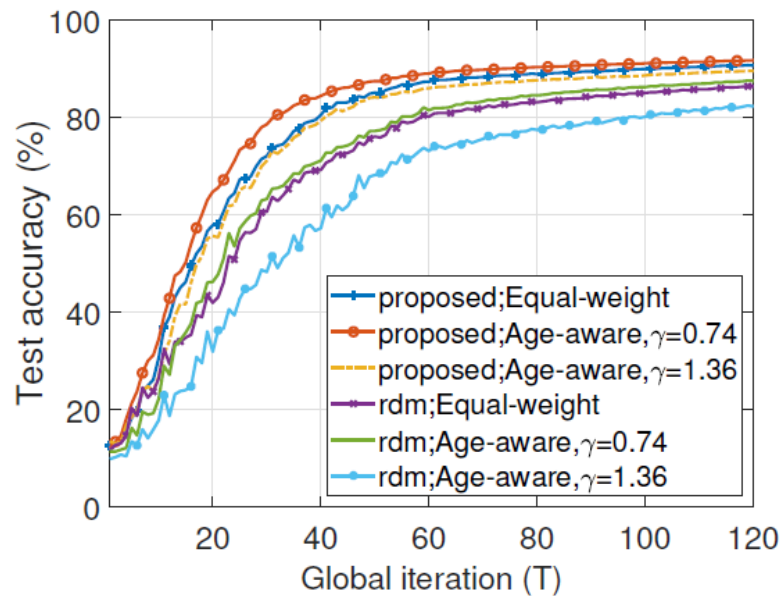


(a) IID data distribution.

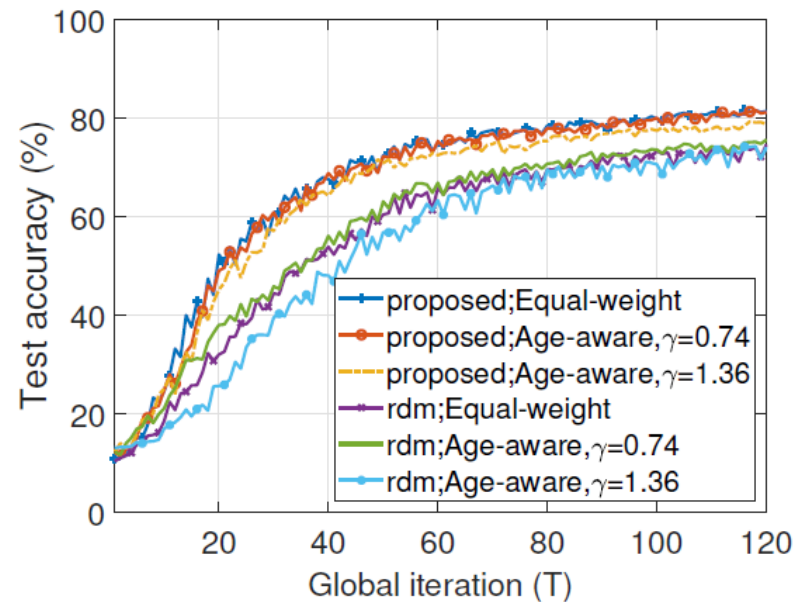


(b) Non-IID data distribution.

# Simulation Results



(a) IID data distribution.



(b) Non-IID data distribution.

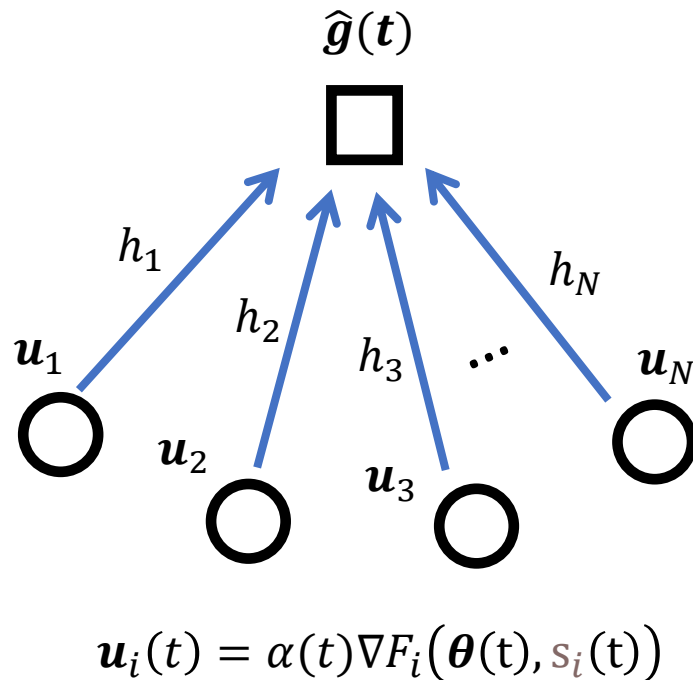
# From Higher Rate to Better Task Performance

Why do we need new wireless design for distributed intelligent systems?

Properties of AI data traffic:

- Higher data rate  $\neq$  higher impact on the task
- Temporal correlation/evolution of data
- Error-free transmission may not be necessary
- Importance matters more than fairness

# Model Aggregation over Distributed Nodes



## Communication Goal:

Compute  $\mathbf{g}(t) = \sum_{i=1}^N w_i \mathbf{u}_i(t)$

- Do we need to decode each stream correctly?
- What happens when N is very large?



# Goal-Oriented Communication Design

## Data property:

- Stochastic gradient vector is a noisy measure itself

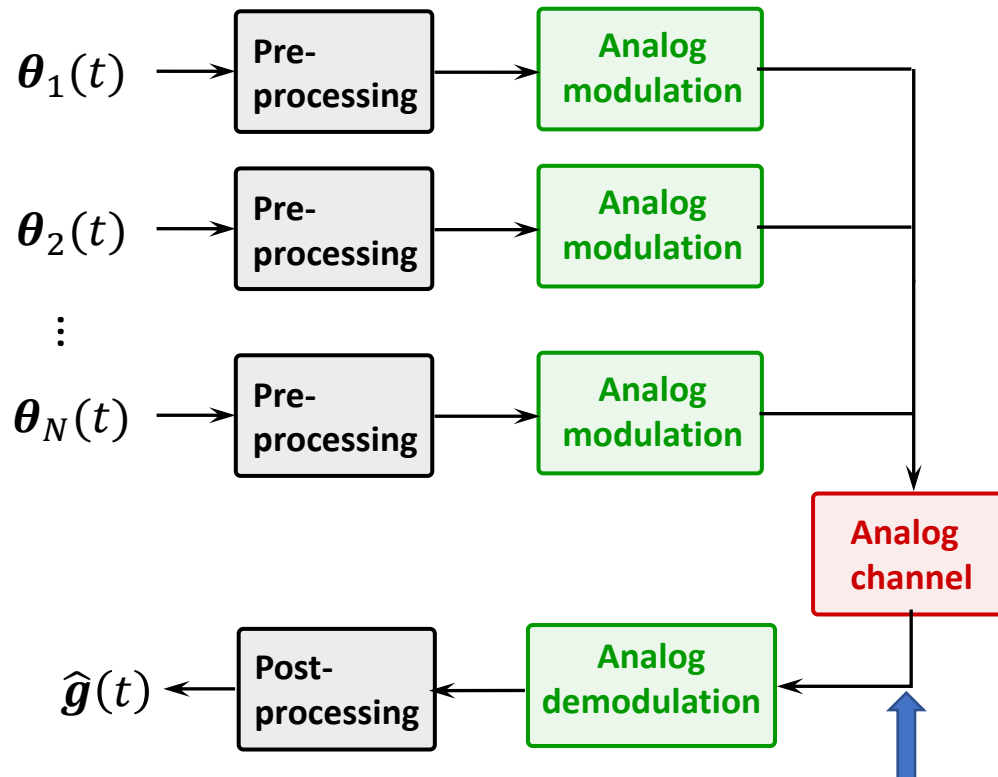
## Objective:

- Received aggregated parameter vector  $\hat{\mathbf{g}}(t)$  as close to the ground-truth  $\mathbf{g}(t)$

## Communication design:

- Exploit signal superposition in wireless channels
- Create a “linear” mapping between data value and received signal amplitude

# Over-the-Air Computation for Federated Learning



- All users share the same frequency-time resources
- Linear processing with analog modulation
- Noise and power limitation affect the reconstruction error

$$\|\hat{g}(t) - g(t)\|^2$$

Noise is directly added to the received superimposed signal

# Computation of Nomographic Functions over Multiple Access Channels

A function  $f$  of  $N$  variables is nomographic if it can be represented in the form

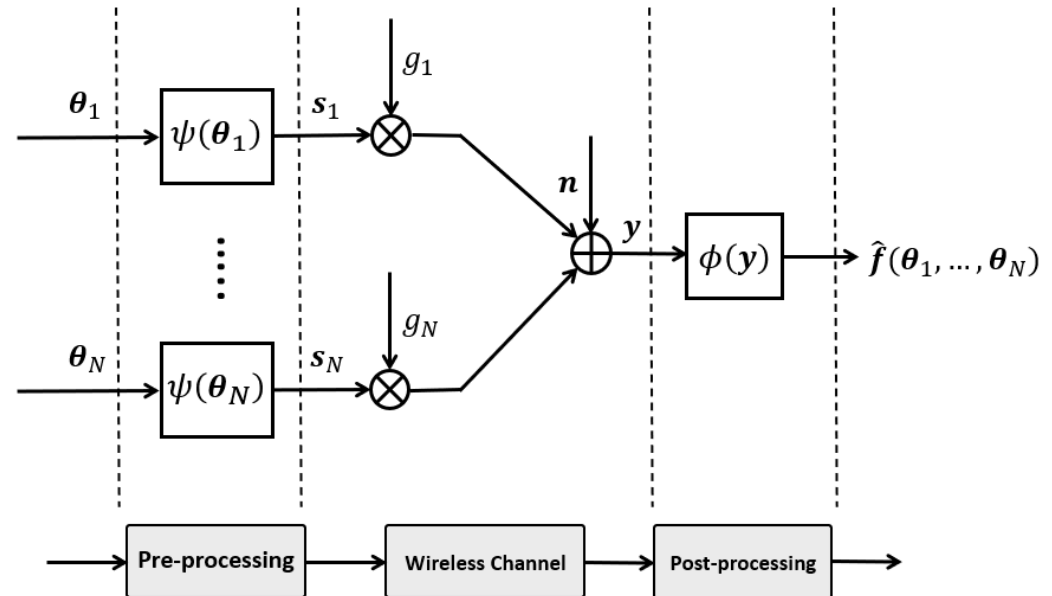
$$f(s_1, \dots, s_N) = \phi \left( \sum_{i=1}^N \psi_i(s_i) \right)$$

- Joint source-channel design
- Exploit signal superposition in MAC channels
- Communication is part of the computation process
- Harness interference instead of fighting it

# OtA Computation over Fading Channels

- N senders, one receiver
- Channel gain of link  $i$ :  $g_i$
- Power constraint:  $P_{\max}$
- Computation objective

$$\hat{f}(\theta_1, \dots, \theta_N) = \sum_{i=1}^N w_i \theta_i$$



- Pre-processing:  $\psi(\theta_i) = \frac{w_i \eta}{g_i} \theta_i$
- Wireless channel:  $y = \sum_{i=1}^N \psi(\theta_i) \cdot g_i + n$
- Post-processing:  $\phi(y) = y/\eta$

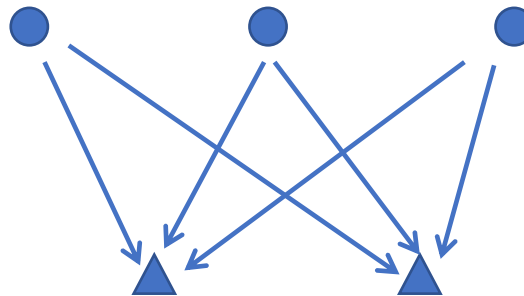
$$\eta = \sqrt{P_{\max}} \min_{i=1 \dots N} \left\{ \frac{|g_i|}{w_i \sqrt{E_i}} \right\}$$

**Bottleneck effect!**

# OtA Computation with Multiple Receivers

- Set of senders:  $\mathcal{S} = \{1, 2, \dots, N_s\}$ , set of receivers:  $\mathcal{R} = \{1, 2, \dots, N_r\}$
- Set of directed links:  $\mathcal{E}$
- Set of senders connected to receiver  $j$ :  $\mathcal{N}_j = \{i \in \mathcal{S} | (i, j) \in \mathcal{E}\}$
- Channel gain of link  $(i, j)$ :  $h_{ij}$
- Data sample of sender  $i$ :  $s_i$
- Computation objective at each receiver  $j$ :

$$f_j = \frac{1}{|\mathcal{N}_j|} \sum_{i \in \mathcal{N}_j} s_i$$

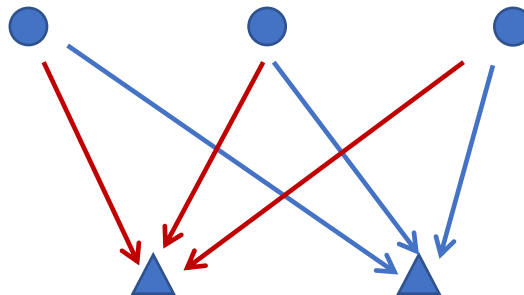


# OtA Computation with Multiple Receivers

Baseline approach:

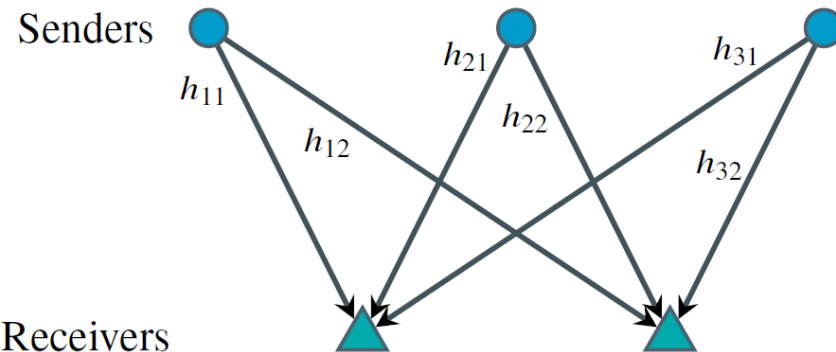
- Slot 1: first receiver active
- Slot 2: second receiver active

Number of required slots=number of receivers



# Multi-Slot Joint Precoding and Decoding Design

Precoder  $\mathbf{p}_i = [p_{i,1}, \dots, p_{i,T}]^T$



Decoder  $\mathbf{q}_j = [q_{j,1}, \dots, q_{j,T}]^T$

At slot  $t$ :

- Sender  $i$  transmits  $s_i p_{i,t}$
- Receiver  $j$  receives

$$\sum_{i \in \mathcal{N}_j} s_i p_{i,t} h_{ij} + n_j$$

Multi-slot decoding:

$$\hat{\theta}_j = \sum_{t=1}^T \left( \sum_{i \in \mathcal{N}_j} s_i p_{i,t} h_{ij} + n_j \right) q_{j,t}$$

# Multi-Slot Joint Precoding and Decoding Design

Unbiased estimation if:

$$\sum_{t=1}^T p_{i,t} q_{j,t} = \frac{1}{|\mathcal{N}_j| h_{ij}} = w_{ij}$$

Conditioning on unbiased estimation, the MSE averaged over all receivers is

$$\begin{aligned} \text{MSE} &= \frac{1}{N_r} \sum_{j=1}^{N_r} \mathbb{E} [|\hat{\theta}_j - \theta_j|^2] \\ &= \frac{\sigma^2}{N_r} \sum_{j=1}^{N_r} \sum_{t=1}^T q_{j,t}^2. \end{aligned}$$



# Problem Formulation

Define  $\mathbf{P} = [\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_{N_s}]$  and  $\mathbf{Q} = [\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{N_r}]$

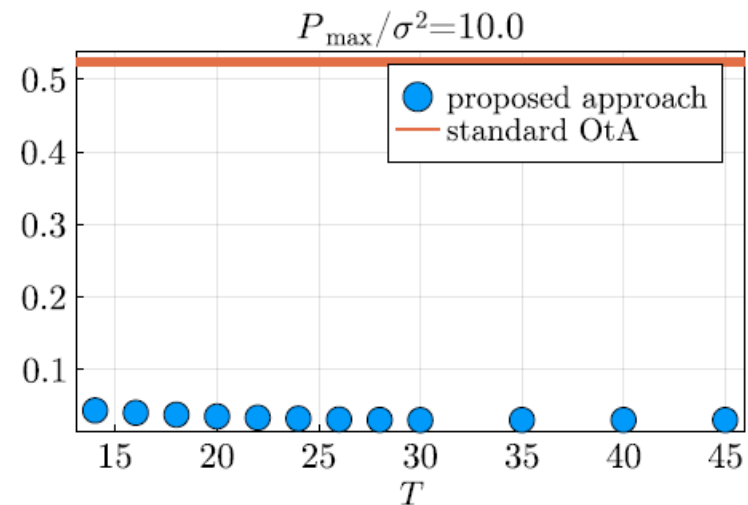
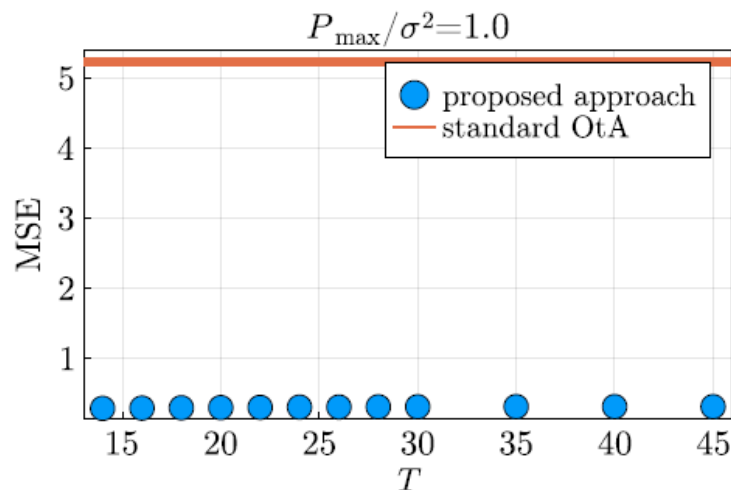
$$\begin{array}{ll}
 \text{minimize} & \sum_{j=1}^{N_r} \|\mathbf{q}_j\|^2 \\
 \text{subject to} & \mathbf{p}_i \cdot \bar{\mathbf{q}}_j = w_{ij}, \forall (i, j) \in \mathcal{E} \\
 & \|\mathbf{p}_i\|^2 \leq C_i, \forall i \in \mathcal{S}.
 \end{array}
 \qquad
 \begin{array}{ll}
 \text{minimize} & \|\mathbf{Q}\|_F^2 \\
 \text{subject to} & \mathbf{P}^\top \mathbf{Q} = \mathbf{W} \\
 & \|\mathbf{p}_i\|^2 \leq C_i, i = 1, \dots, N_s.
 \end{array}$$

$$\begin{array}{ll}
 \text{minimize} & \|\mathbf{P}^\top \mathbf{Q} - \mathbf{W}\|_F^2 + \lambda \|\mathbf{Q}\|_F^2 \\
 \text{subject to} & \|\mathbf{p}_i\|^2 \leq C_i, i = 1, \dots, N_s,
 \end{array}$$

# Simulation Results

## Simulation setup:

- 50 senders, 30 receivers
- Each sender randomly connected to 20 receivers
- Data and channel gain randomly generated from  $\mathcal{CN}(0,1)$



# Simulation Results

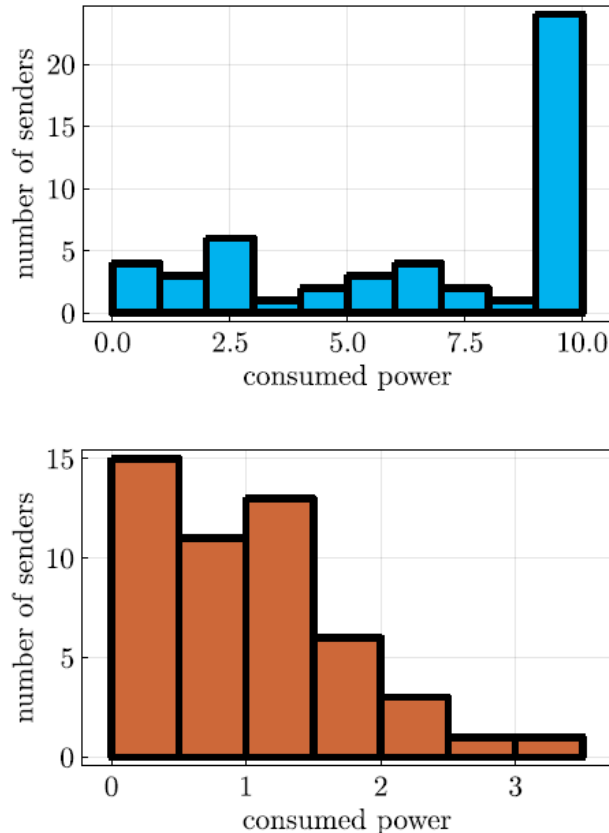


Fig. 3: Power consumption: proposed (above) and standard OtA (bottom) approaches.

## Remark:

- With proposed approach, many senders reach the power limit
- With standard (parallel) OtA, much available power is unconsumed
- Standard OtA can be potentially optimized

## Conclusion:

- With multi-slot joint precoding and decoding, we can save communication resources (e.g., time slots) as compared to standard parallel OtA approach

# OtA FL with multiple antennas

If no channel state information, with multiple-antenna receiver, what can we do?

- Users send orthogonal pilot signals and data
- Server estimates individual channels
- Server decodes the signal using estimated individual channels

Or

- Users send one common pilot and data
- Server estimates the sum channel
- Server decodes the signal using estimated sum channel

# Estimating Individual Channels vs. Estimating Sum Channel

Estimating  $\hat{G}$

$$\mathbf{Y}_{p,\text{orth}} = \sum_{k=1}^K \sqrt{\rho\tau_p} \mathbf{g}_k \phi_k^H + \mathbf{N}_p$$

$$\hat{\mathbf{g}}_k = \frac{\sqrt{\rho\tau_p} \beta_k}{1 + \rho\tau_p \beta_k} \mathbf{Y}_{p,\text{orth}} \phi_k$$

$$\gamma_k = \frac{\rho\tau_p \beta_k^2}{1 + \rho\tau_p \beta_k}$$

Estimating  $\hat{h}_{\text{sum}}$

$$\begin{aligned} \mathbf{Y}_{p,\text{sum}} &= \sum_{k=1}^K \sqrt{\rho\tau_p \frac{\beta_{\min}}{\beta_k}} \mathbf{g}_k \phi^H + \mathbf{N}_p \\ &= \sum_{k=1}^K \sqrt{\rho\tau_p \beta_{\min}} \mathbf{h}_k \phi^H + \mathbf{N}_p \\ &= \sqrt{\rho\tau_p \beta_{\min}} \mathbf{h}_{\text{sum}} \phi^H + \mathbf{N}_p \end{aligned}$$

$$\hat{h}_{\text{sum}} = \frac{\sqrt{\rho\tau_p \beta_{\min} K}}{1 + \rho\tau_p \beta_{\min} K} \mathbf{Y}_{p,\text{sum}} \phi$$

$$\bar{\gamma} = \frac{\rho\tau_p \beta_{\min} K^2}{1 + \rho\tau_p \beta_{\min} K}$$

# Best Linear Unbiased Estimator for Signal Decoding

$$\mathbf{Y} = \sum_{k=1}^K \sqrt{\rho\eta_k} \mathbf{g}_k \mathbf{x}_k^T + \mathbf{N}$$

BLUE

$$[\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_K] \\ = \left( \frac{1}{\sqrt{\rho}} \mathbf{D}_\eta^{-1/2} (\hat{\mathbf{G}}^H \hat{\mathbf{G}})^{-1} \hat{\mathbf{G}}^H \mathbf{Y} \right)^T$$

$$\mathbf{D}_\eta = \text{diag}(\eta_1, \dots, \eta_K)$$

Full power!

Benchmark

$$\hat{\mathbf{x}} = \sum_{k=1}^K w_k \mathbf{x}_k = c (\hat{\mathbf{h}}_{\text{sum}}^H \mathbf{Y})^T$$

$$\max_k \|\sqrt{\eta_k} \mathbf{x}_k\|_2^2 = \max_k \eta \frac{w_k^2}{\beta_k} \|\mathbf{x}_k\|_2^2 = T$$

$$c = \frac{K}{M \sqrt{\eta \rho \gamma}}$$

# What Happens at Client Side

## Sparse BLUE

- 1: **for** client  $k \in \{1, \dots, K\}$  **in parallel do**
- 2:   receive (error free)  $\theta(t)$
- 3:   obtain  $\Delta\theta_k(t)$  from stochastic gradient descent
- 4:    $\mathbf{x}_k^{\text{full}} \leftarrow \text{SPLIT}(\Delta\theta_k(t))$
- 5:    $\mathbf{x}_k^{\text{full}} \leftarrow \mathbf{x}_k^{\text{full}} + \mathbf{r}_k$
- 6:    $\mathbf{x}_k^{\text{sparse}} \leftarrow \text{SPARSE}(\mathbf{x}_k^{\text{full}}, S)$
- 7:    $\mathbf{r}_k \leftarrow \mathbf{x}_k^{\text{full}} - \mathbf{x}_k^{\text{sparse}}$
- 8:    $\mathbf{x}_k \leftarrow \mathbf{A}_k^t \mathbf{x}_k^{\text{sparse}}$
- 9:   transmit  $\phi_k, \sqrt{\eta_k} \mathbf{x}_k$
- 10: **end for**

## Sparse SUM

- 1: **for** client  $k \in \{1, \dots, K\}$  **in parallel do**
- 2:   receive (error free)  $\theta(t)$
- 3:   obtain  $\Delta\theta_k(t)$  from stochastic gradient descent
- 4:    $\mathbf{x}_k^{\text{full}} \leftarrow \text{SPLIT}(\Delta\theta_k(t))$
- 5:    $\mathbf{x}_k^{\text{full}} \leftarrow \mathbf{x}_k^{\text{full}} + \mathbf{r}_k$
- 6:    $\mathbf{x}_k^{\text{sparse}} \leftarrow \text{SPARSE}(\mathbf{x}_k^{\text{full}}, S)$
- 7:    $\mathbf{r}_k \leftarrow \mathbf{x}_k^{\text{full}} - \mathbf{x}_k^{\text{sparse}}$
- 8:    $\mathbf{x}_k \leftarrow \mathbf{A}_k^t \mathbf{x}_k^{\text{sparse}}$
- 9:   transmit  $\phi, \sqrt{\eta_k} \mathbf{x}_k$
- 10: **end for**

# What Happens at Server Side

## Sparse BLUE

- 11: receive  $\mathbf{Y}_{p,orth}$  and  $\mathbf{Y}$
- 12: estimate  $\hat{\mathbf{G}}$
- 13: **for**  $k \in \{1, \dots, K\}$  **do**
- 14:  $\hat{\mathbf{x}}_k \leftarrow \frac{1}{\sqrt{\eta_k \rho}} [\hat{\mathbf{G}}(\hat{\mathbf{G}}^H \hat{\mathbf{G}})^{-1}]_k^H \mathbf{Y}$
- 15: solve sparsity problem to get  $\hat{\mathbf{x}}_k^{sparse}$
- 16:  $\widehat{\Delta \boldsymbol{\theta}}_k(t, \tau) \leftarrow \text{UNSPLIT}(\hat{\mathbf{x}}_k^{sparse})$
- 17: **end for**
- 18:  $\widehat{\Delta \boldsymbol{\theta}}(t) \leftarrow \sum_{k=1}^K w_k \widehat{\Delta \boldsymbol{\theta}}_k(t)$
- 19:  $\boldsymbol{\theta}(t+1) \leftarrow \boldsymbol{\theta}(t) + \alpha_t^{\text{global}} \widehat{\Delta \boldsymbol{\theta}}(t)$
- 20: broadcast (error free)  $\boldsymbol{\theta}(t+1)$

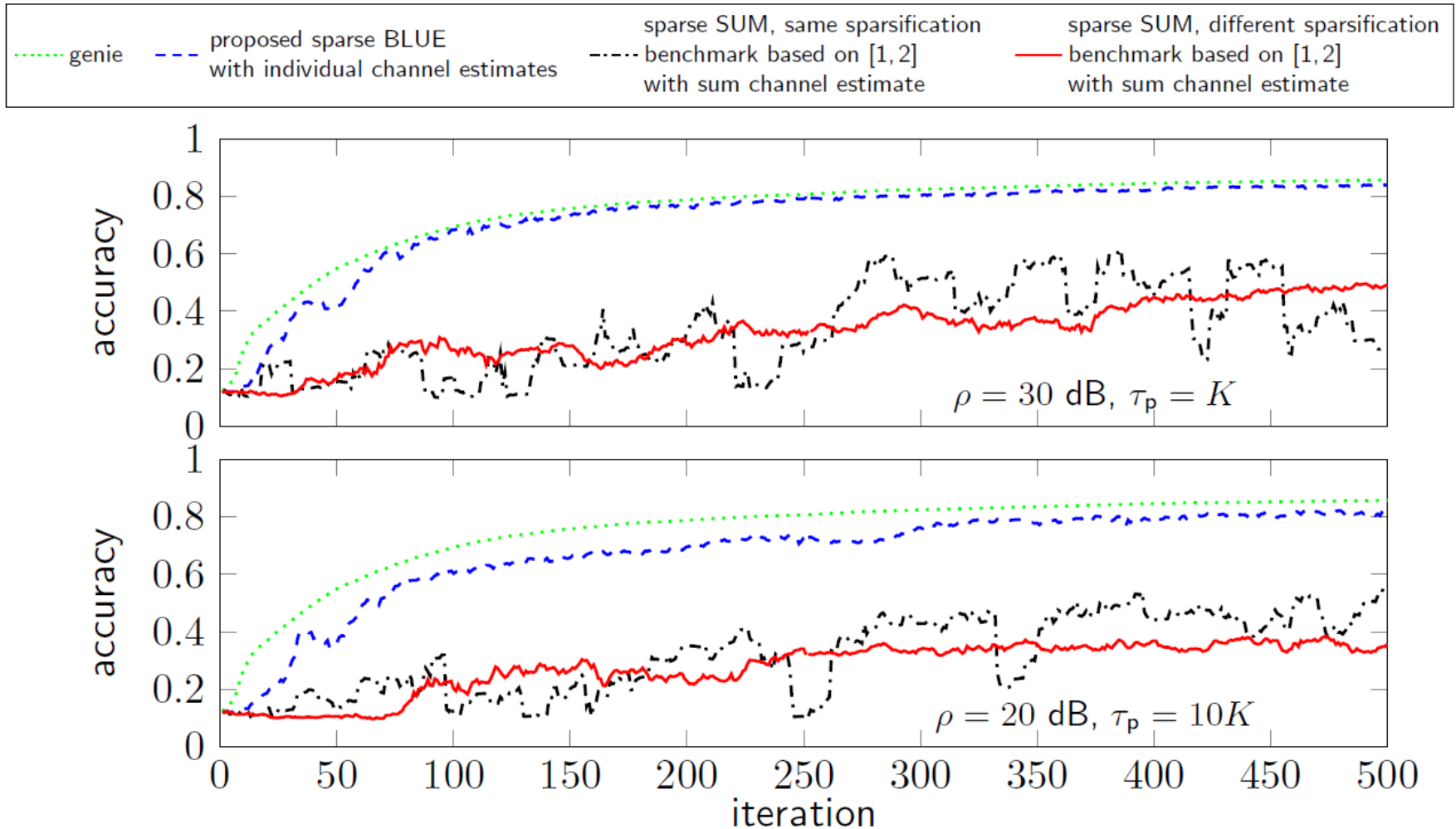
Individual processing of clients!

## Sparse SUM

- 11: receive  $\mathbf{Y}_{p,sum}$  and  $\mathbf{Y}$
- 12: estimate  $\hat{\mathbf{h}}_{sum}$
- 13:  $\hat{\mathbf{x}} \leftarrow \frac{K}{M\sqrt{\rho\eta\gamma}} \hat{\mathbf{h}}_{sum}^H \mathbf{Y}$
- 14: solve sparsity problem to get  $\widehat{\left( \sum_{k=1}^K w_k \mathbf{x}_k^{sparse} \right)}$
- 15:  $\widehat{\Delta \boldsymbol{\theta}}(t) \leftarrow \text{UNSPLIT}\left(\widehat{\left( \sum_{k=1}^K w_k \mathbf{x}_k^{sparse} \right)}\right)$
- 16:  $\boldsymbol{\theta}(t+1) \leftarrow \boldsymbol{\theta}(t) + \alpha_t^{\text{global}} \widehat{\Delta \boldsymbol{\theta}}(t)$
- 17: broadcast (error free)  $\boldsymbol{\theta}(t+1)$



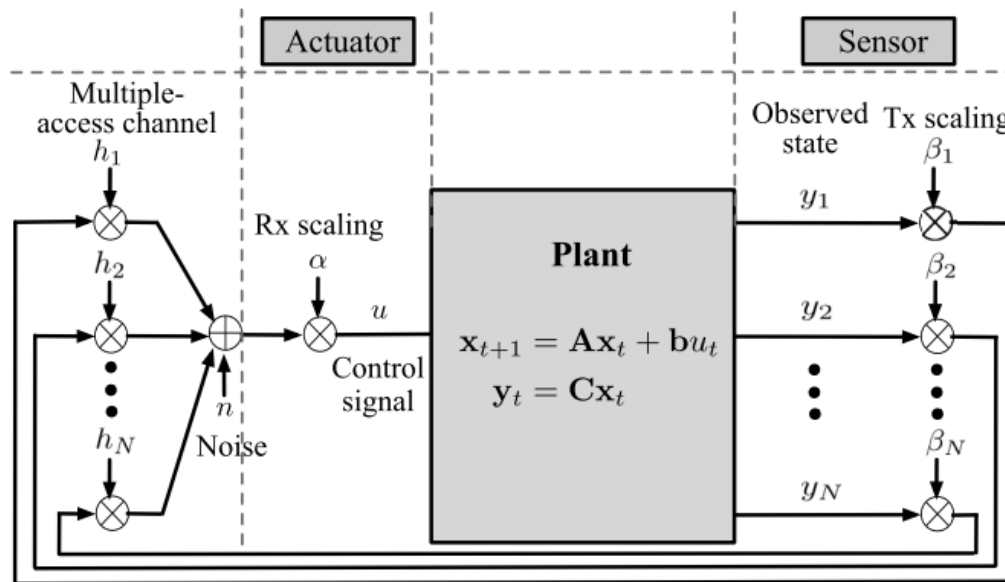
# Simulation Results



# Applications of OtA Computation

- Distributed learning
  - Distributed control
  - Distributed inference and estimation
- 
- Key feature: the goal of communication is to compute some **functions of data** from a set of distributed nodes **approximately correct**, instead of receiving each independent data stream without errors

# Example: Wireless Control Systems



- $\mathbf{x}_t$ : system state
- $\mathbf{y}_t$ : observed output (measured at different sensors)
- $u_t$ : control input signal

# Challenges of OtA Aggregation in FL

MSE of aggregated data may not be the most suitable performance metric, need task-specific metric design

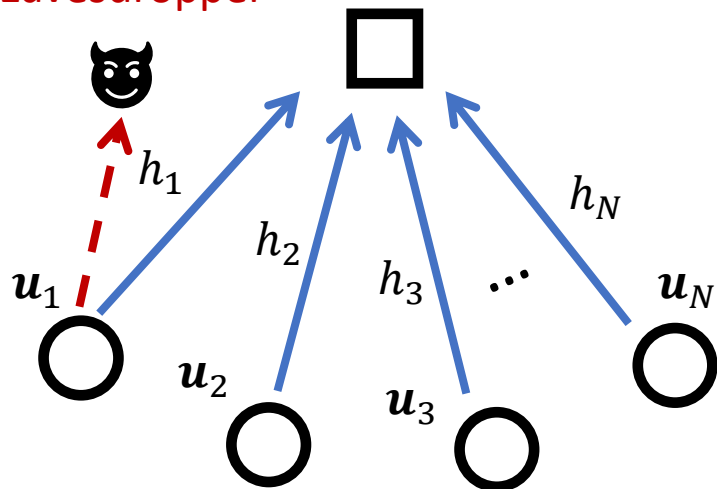
Scheduled users	50	40	30	20	10
Test accuracy	92%	71.65%	56.86%	42.45%	27.37%
MSE	$7.75e-14$	$3.70e-15$	$7.54e-17$	$7.60e-18$	$1.81e-18$

- Opting out bottleneck users can reduce the MSE of aggregated data
- Aggregated model might be biased (especially with non-IID training data) when opting out users

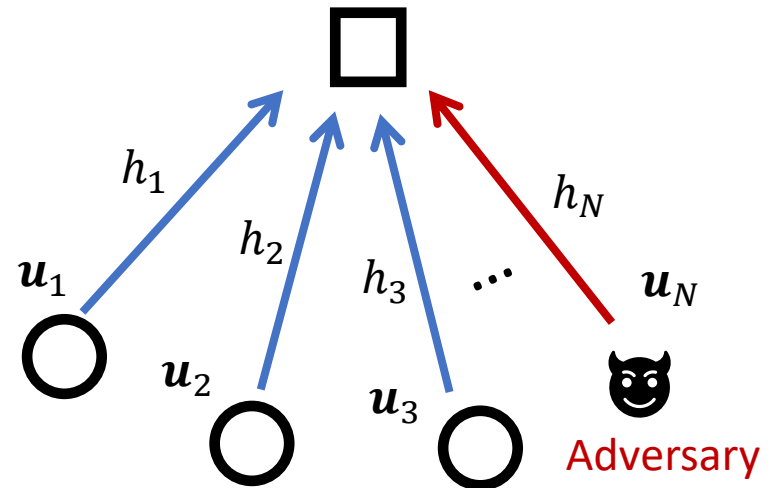
# Challenges of OtA Aggregation in FL

## Privacy, security and robustness

Eavesdropper



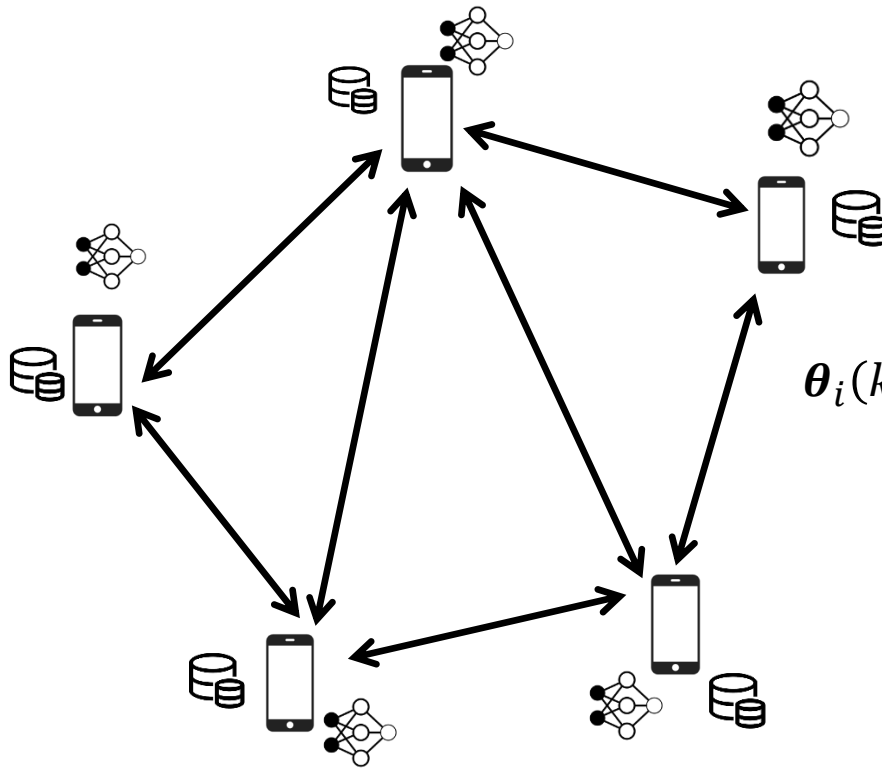
Membership inference attack



Model poisoning attack

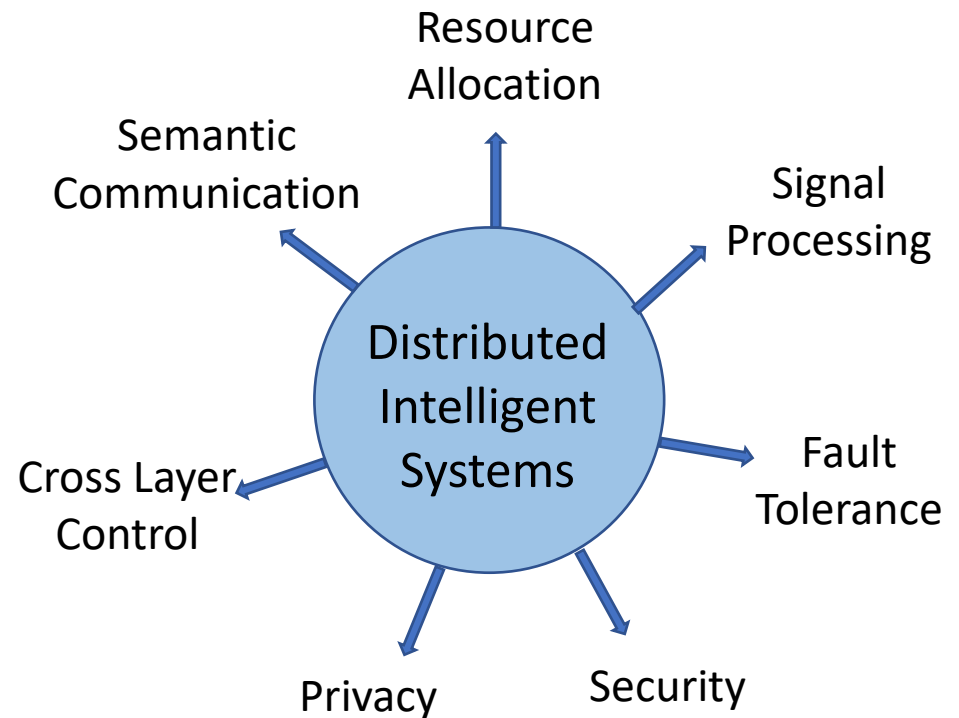
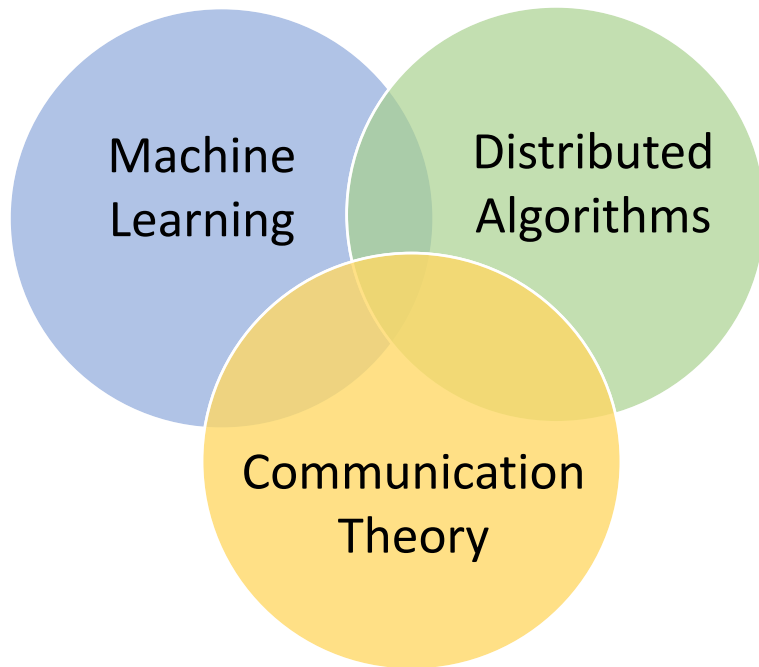
# Challenges of OtA Computation

## Communication design in fully decentralized ML



$$\boldsymbol{\theta}_i(k+1) = \underbrace{\sum_{j=1}^N w_{ij} \boldsymbol{\theta}_j(k)}_{\text{Consensus}} - \underbrace{\alpha \nabla f_i(\boldsymbol{\theta}_i(k))}_{\text{Innovation}}$$

# Distributed Intelligence over Wireless Networks: An Interdisciplinary View



Thank you!