# 15-884: Machine Learning Systems
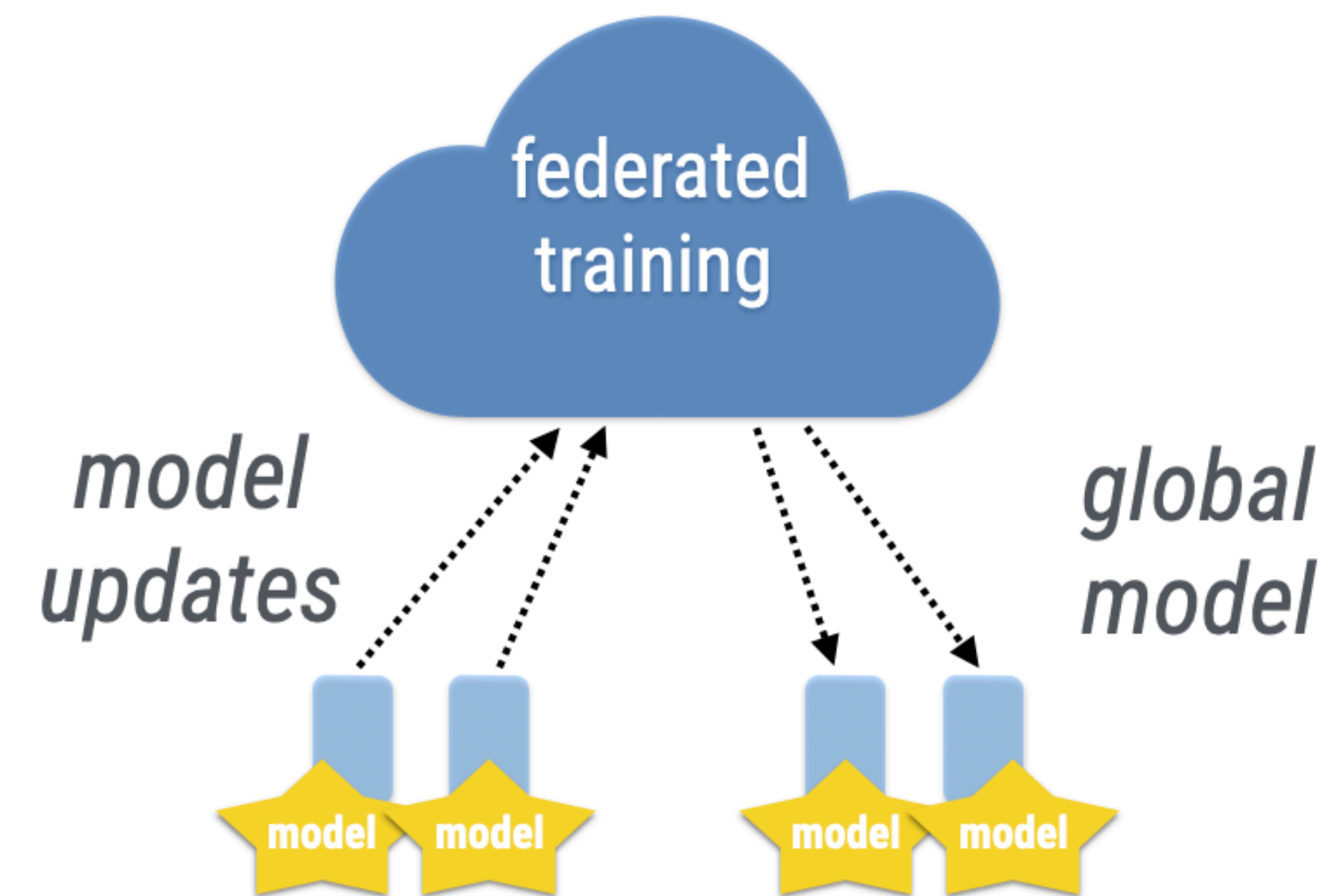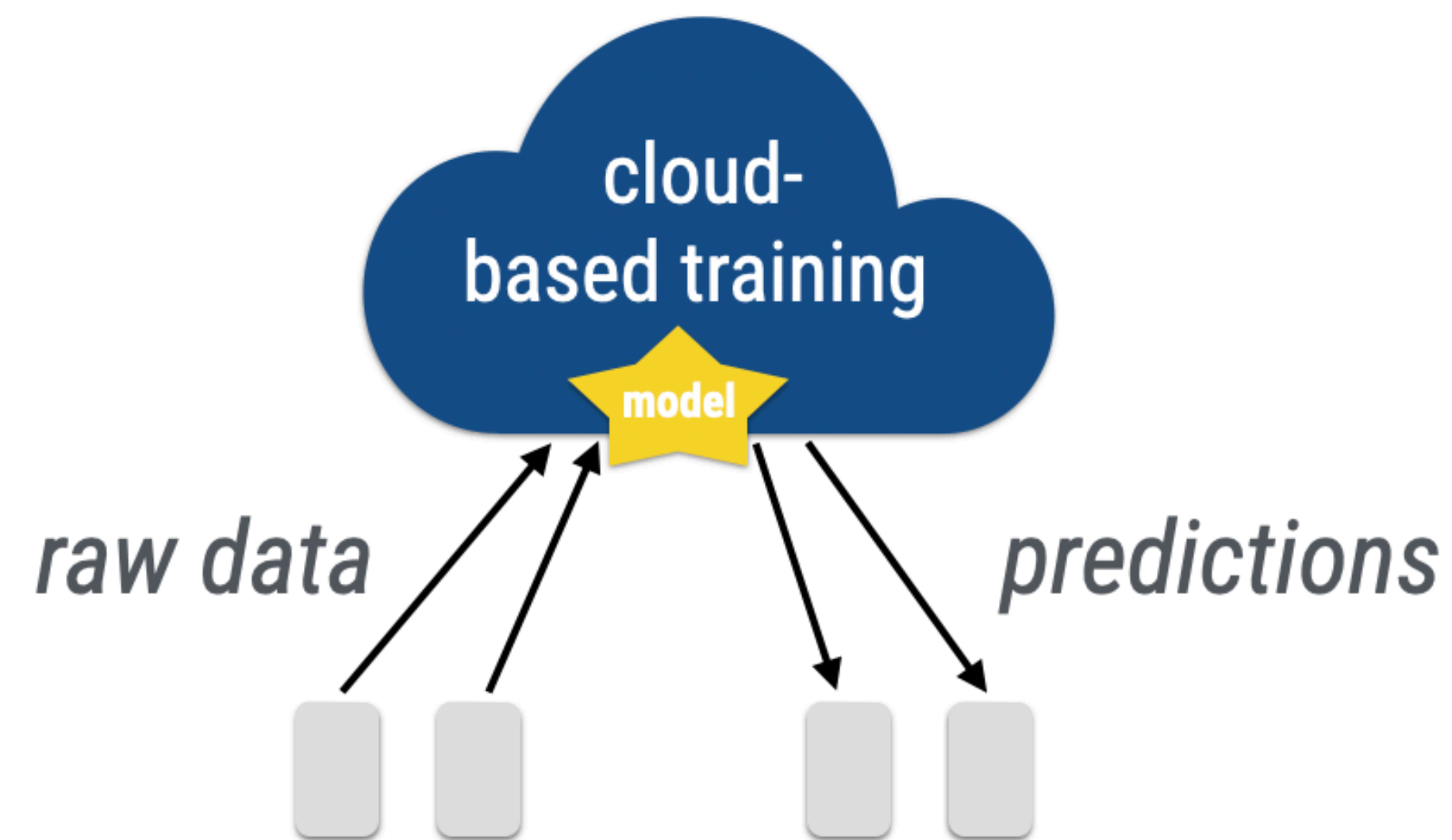
## *Federated Learning*

Tian Li

tianli@cmu.edu

**Carnegie Mellon University**
School of Computer Science

# Federated Learning

**Privacy-preserving *training* in heterogeneous, (potentially) massive networks**
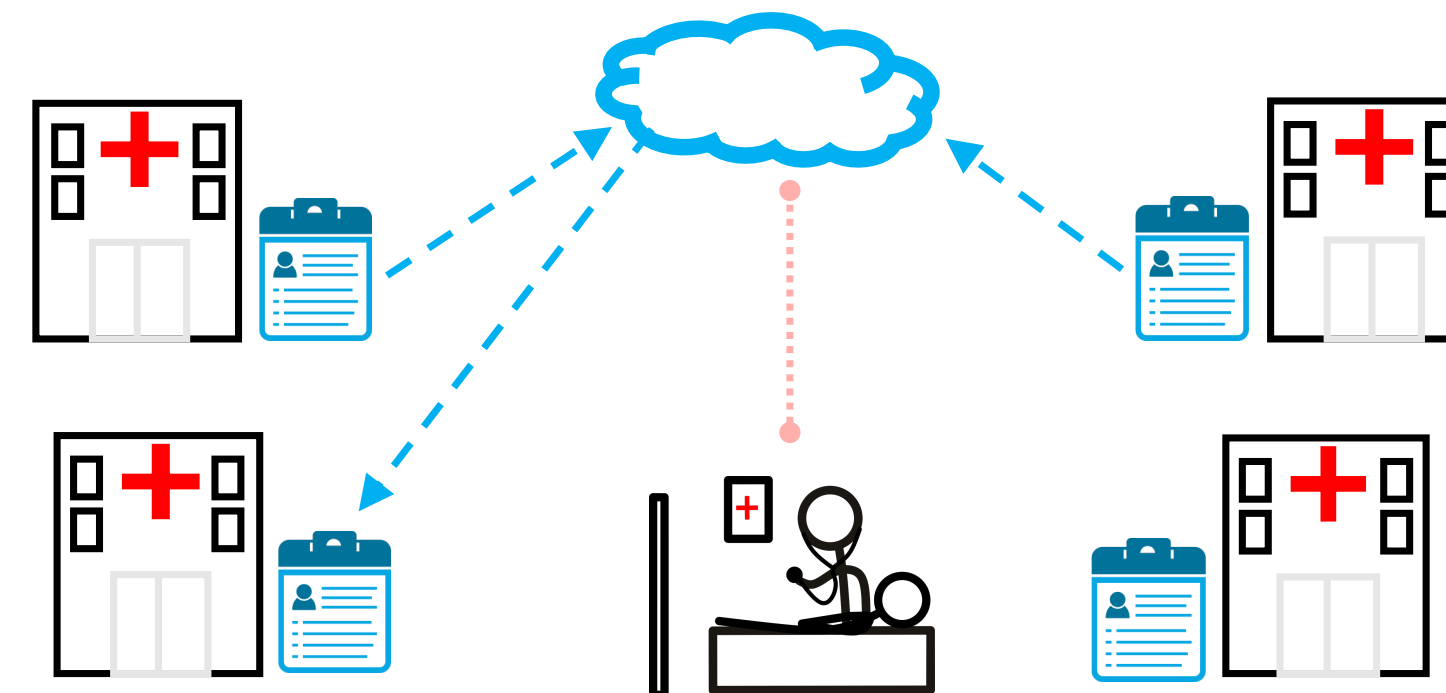
# Federated Learning

**Privacy-preserving _training_ in heterogeneous, (potentially) massive networks**

Networks of remote devices

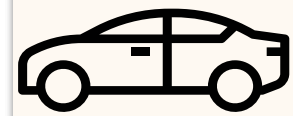Networks of isolated organizations



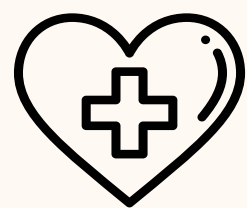cross-device setting

cross-silo setting
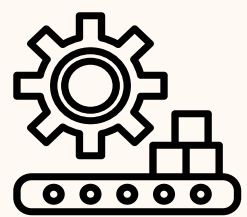
# Example Applications

Anomaly detection in IoT devices

Adapting to pedestrian behavior on autonomous vehicles

Personalized healthcare on wearable devices

Predictive maintenance for industrial machines

Assumptions: (1) local data is important (2) labels are available (3) privacy is a concern

# Workflow & Challenges

**Objective:**

$$\min_w f(w) = \sum_{k=1}^{N} p_k F_k(w)$$

*loss on device k*

*Training setup:*

server

devices

$\mathbf{w}'$      $\mathbf{w}''$

$\mathbf{w_t}$      $\mathbf{w_t}$

**Systems heterogeneity**
variable hardware, network connectivity, power, etc

**Statistical heterogeneity**
highly non-identically distributed data

**Expensive communication**
massive, slow networks

**Privacy & security**
user privacy constraints

# Federated Optimization: Challenges

Systems and statistical heterogeneity (non-identical data) can bias the optimization procedure;
can affect the modeling approach

**Systems heterogeneity**
variable hardware, network connectivity, power, etc

**Statistical heterogeneity**
highly non-identically distributed data

**Expensive communication**
massive, slow networks

**Privacy & security**
user privacy constraints

# Federated Optimization: Challenges

1) reduce the size of messages per round
2) reduce the communication rounds
3) reduce the number of selected devices per round

**Systems heterogeneity**
variable hardware, network connectivity, power, etc

**Statistical heterogeneity**
highly non-identically distributed data

**Expensive communication**
massive, slow networks

**Privacy & security**
user privacy constraints

# Federated Optimization: Challenges

1) keep data on local devices
2) differentially private mechanisms
3) crypto-based methods

*(not the focus today)*

**Systems heterogeneity**
variable hardware, network connectivity, power, etc

**Statistical heterogeneity**
highly non-identically distributed data

**Expensive communication**
massive, slow networks

**Privacy & security**
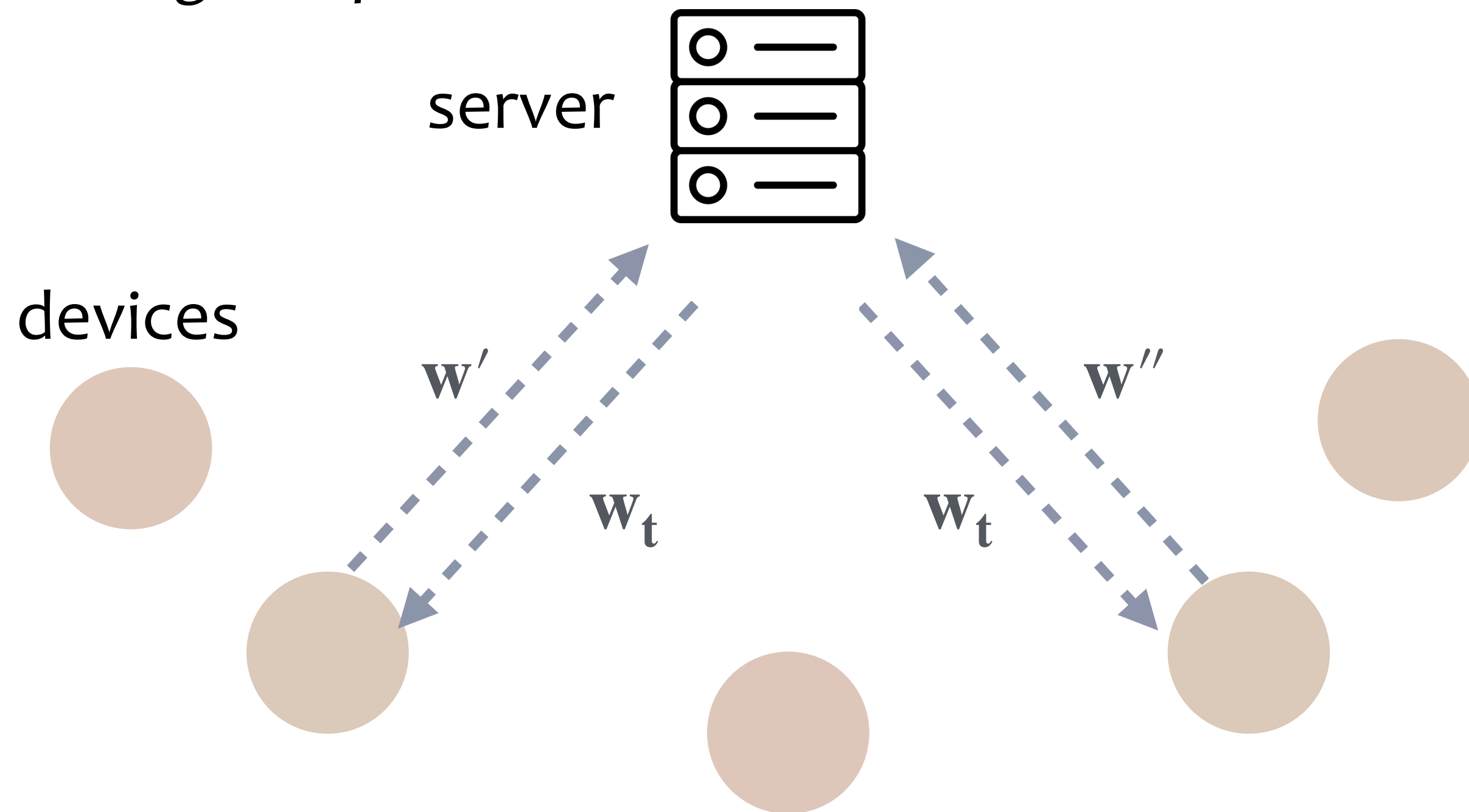user privacy constraints

# How does heterogeneity affect federated optimization methods?

# Federated Optimization: Formulation

**Objective:**

$$\min_w f(w) = \sum_{k=1}^{N} p_k F_k(w)$$

*loss on device k*

*Training setup:*



server

devices

$\mathbf{w}'$  $\mathbf{w}''$

$\mathbf{w_t}$  $\mathbf{w_t}$

Typically solving an empirical risk minimization (ERM) objective:

$$\min_w \sum_{k=1}^{N} p_k \sum_{i=1}^{n_k} \ell(h(x_k^{(i)}; w), y_k^{(i)})$$

# Federated Optimization: Formulation

**Risk:**

$$R(h) = \mathbb{E}_{k \sim Q} \mathbb{E}_{(x,y) \sim P_k} [\ell(h(x; w), y)]$$

**Empirical Risk:**

$$R_{\text{emp}}(h) = \sum_{k=1}^{N} p_k \sum_{i=1}^{n_k} \ell(h(x_k^{(i)}; w), y_k^{(i)})$$

Typically solving an empirical risk minimization (ERM) objective:

$$\min_w \sum_{k=1}^{N} p_k \sum_{i=1}^{n_k} \ell(h(x_k^{(i)}; w), y_k^{(i)})$$

# Optimization for FL: Federated Averaging (FedAvg*)

At each communication round:

- Server randomly selects a subset of devices & sends the current global model $w^t$
- Each selected device $k$ updates $w^t$ for $E$ epochs of SGD to optimize $F_k$ & sends the new local model back
- Server aggregates local models to form a new global model $w^{t+1}$

- Simple method
- Using local updates can lead to much faster convergence empirically
- Works well in many settings (especially non-convex)

* McMahan, H. Brendan, et al. "Communication-efficient learning of deep networks from decentralized data." AISTATS, 2017.

# [Aside] How does FedAvg Differ from Distributed SGD?

**Local updating is not new\***

- one-shot averaging
- ADMM
- COCOA
- Local SGD

Federated settings defer in terms of:

- heterogeneous data
- partial device participation
- often for non-convex objectives

\*   [Zhang, Duchi, Wainwright, Communication-Efficient Algorithms for Statistical Optimization, JMLR 2013]
\*   [Boyd et al, Distributed Optimization and Statistical Learning via ADMM, FnT in ML, 2010]
\*   [Jaggi & Smith et al, Communication-Efficient Distributed Dual Coordinate Ascent, NeurIPS 2014]
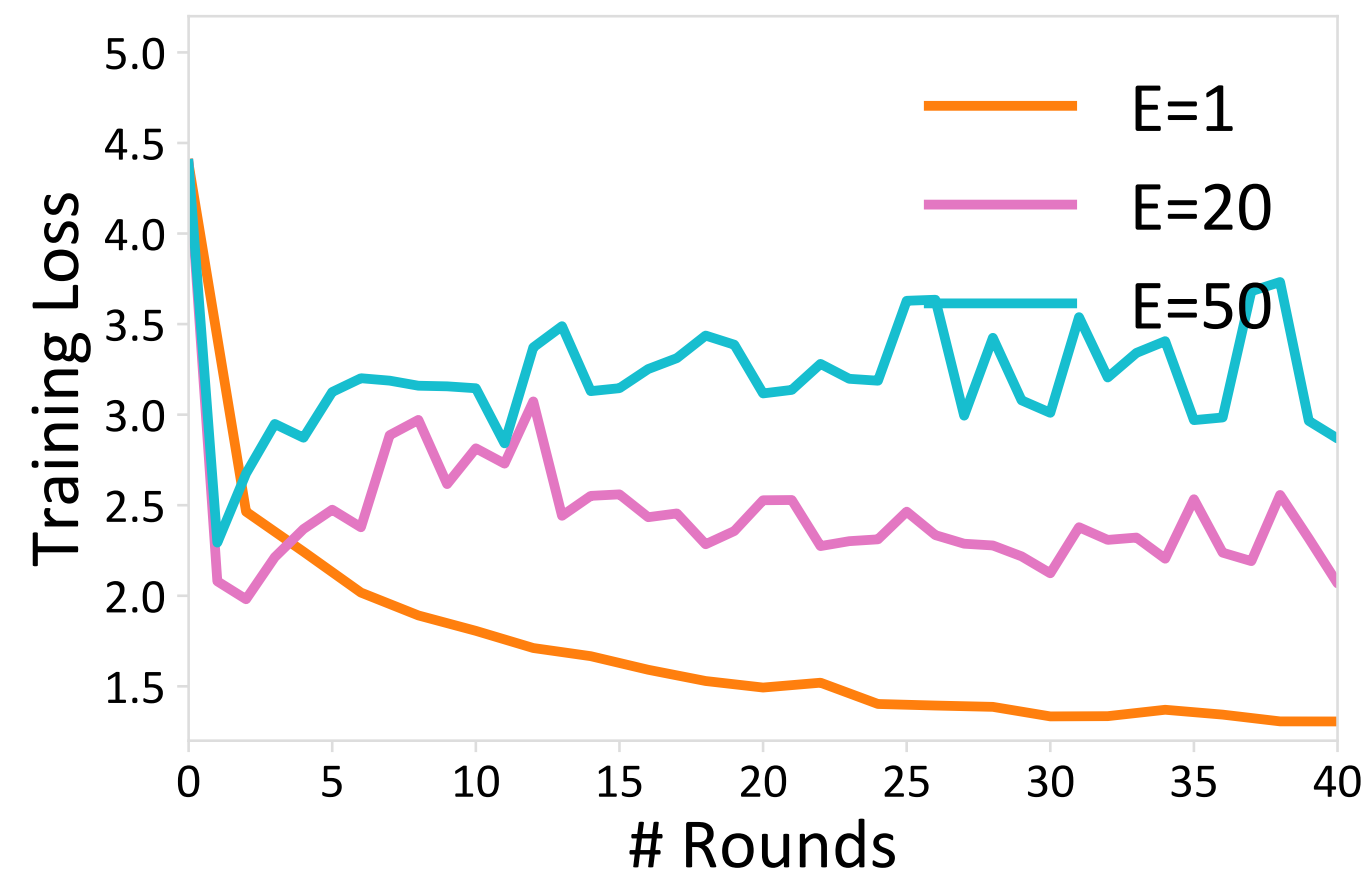\*   [MacDonald et al, Efficient large-scale distributed training of conditional maxent models, NeurIPS 2009]

# Challenge: Heterogeneity

## statistical heterogeneity
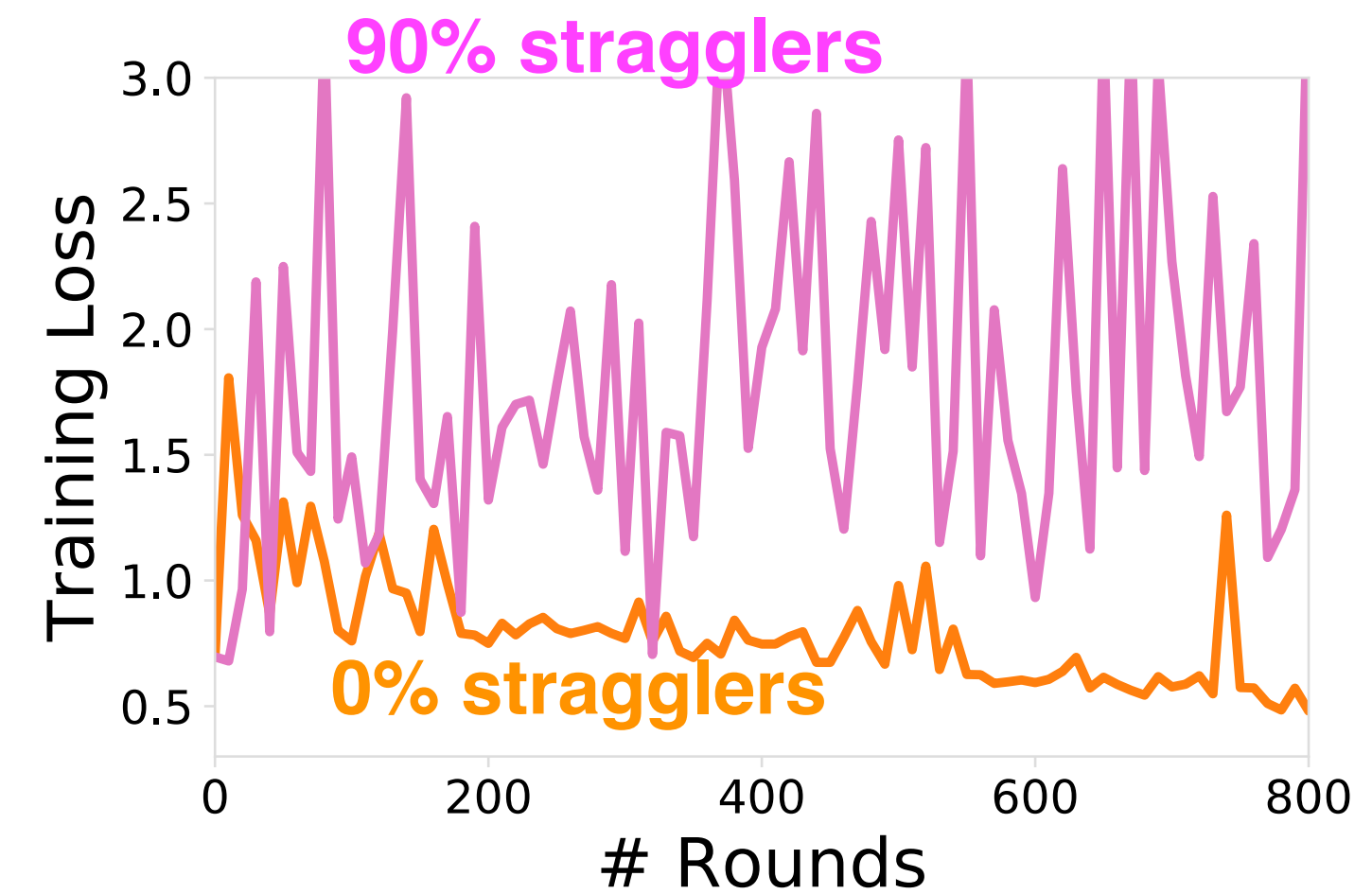
highly non-identically distributed data

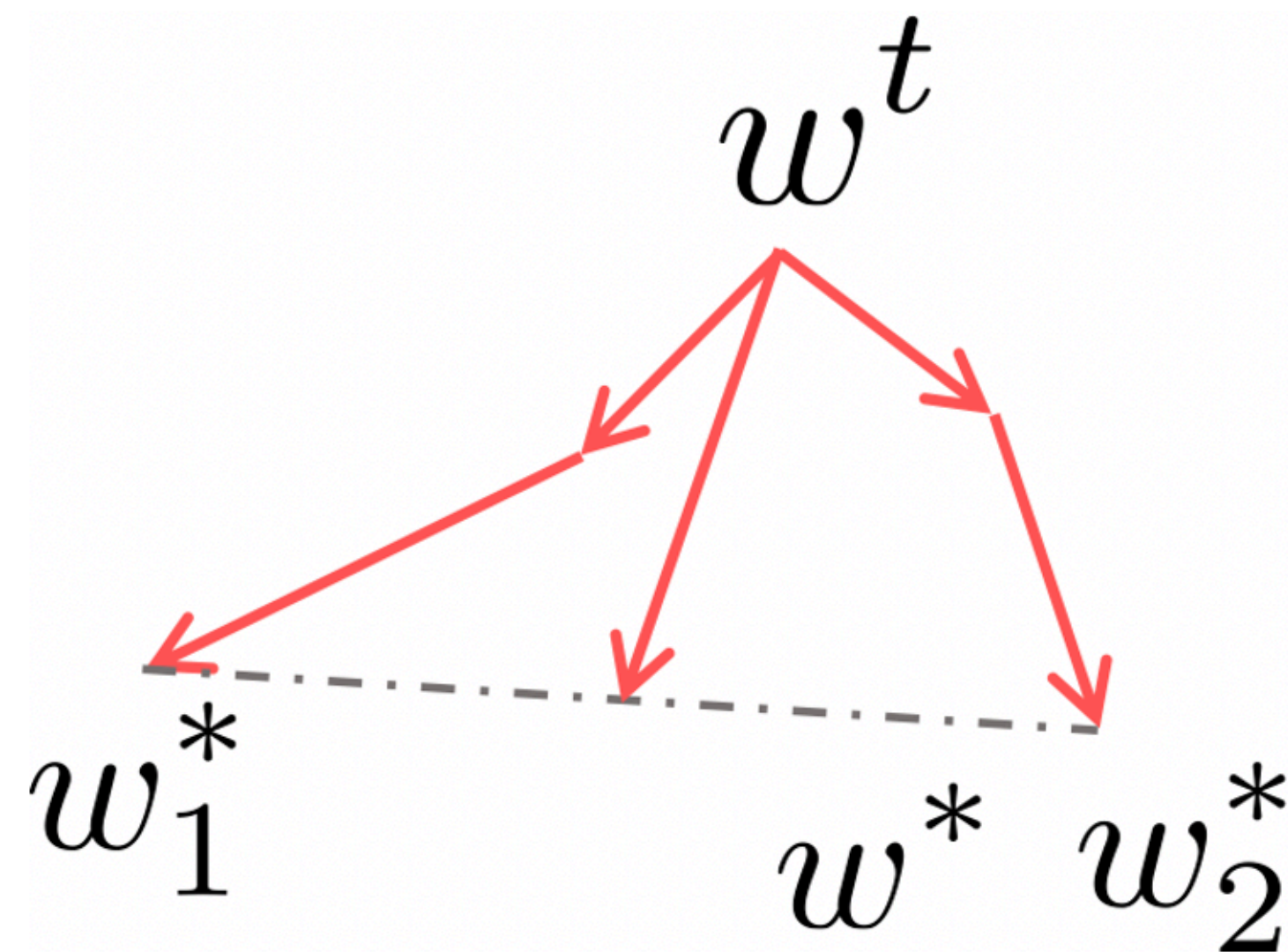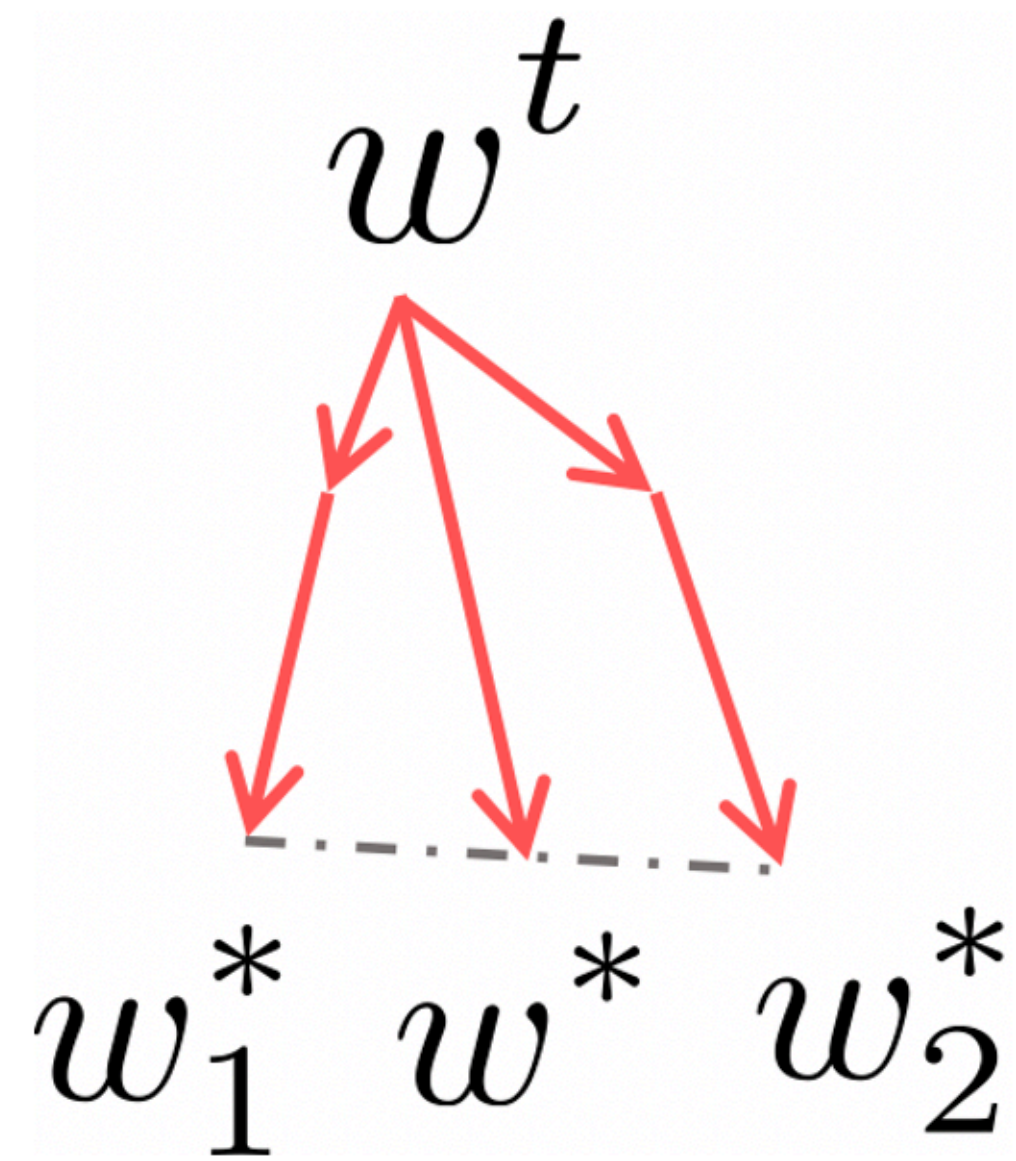too much local work can hurt convergence

## systems heterogeneity

stragglers

dropping slow devices can exacerbate convergence issues



Bonawitz, Keith, et al. "Towards Federated Learning at Scale: System Design." MLSys, 2019.
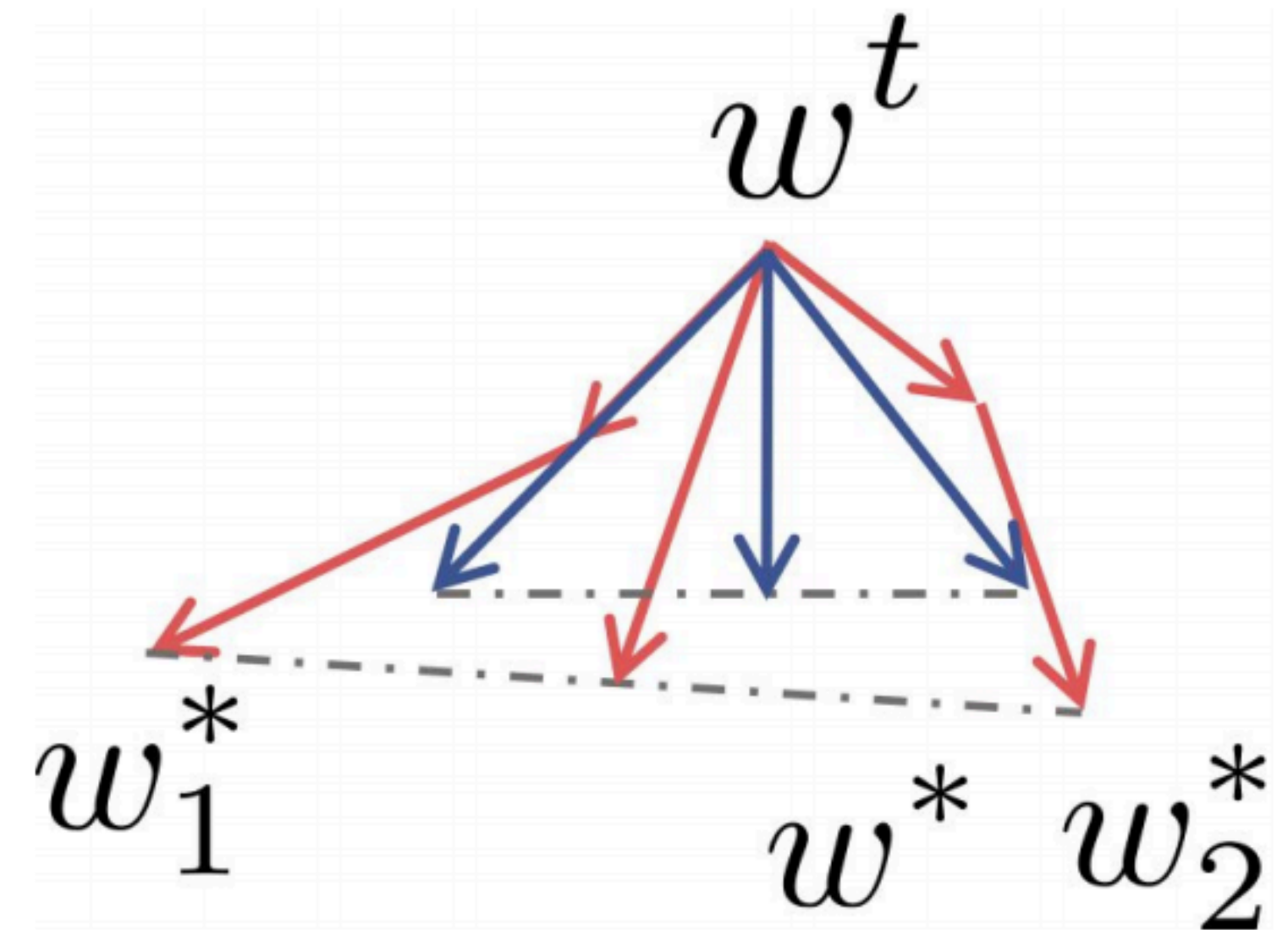
# Challenge: Heterogeneity

# FedProx: A Framework For Federated Optimization

**_Modified_ Local Subproblem:** $\min\limits_{w_k} F_k(w_k) + \dfrac{\mu}{2} \left\| w_k - w^t \right\|^2$

_a proximal term_

- The proximal term explicitly limits the impact of heterogeneous local updates
- Don't drop devices: instead [safely] incorporate partial work
- Generalization of FedAvg; Allows for any local solver
- Theoretical guarantees (with a dissimilarity assumption)

# FedProx: Convergence Analysis

- High-level: **converges** despite non-IID data, local updating, and partial device participation

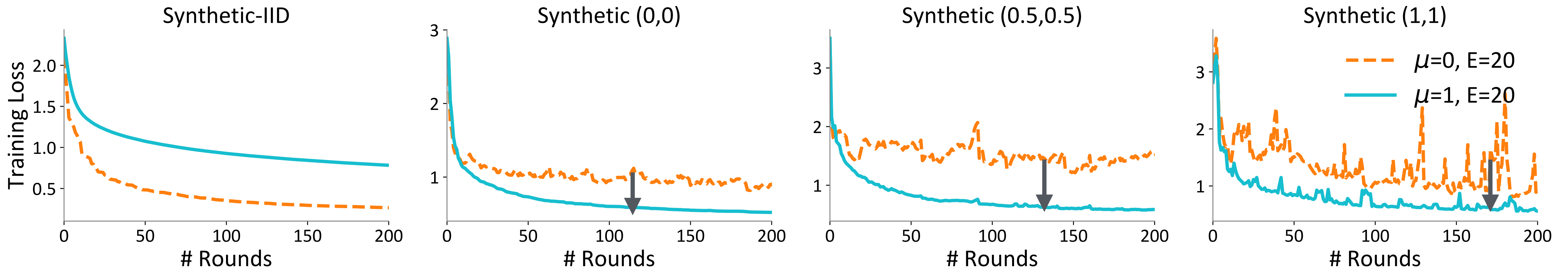- Introduces notion of **B-dissimilarity** in to characterize statistical heterogeneity:

$$\mathbb{E}\left[\|\nabla F_k(w)\|^2\right] \leq \|\nabla f(w)\|^2 B^2$$

IID data: $B = 1$

non-IID data: $B > 1$

*\* used in other contexts, e.g., gradient diversity to quantify the benefits of scaling distributed SGD*

Yin, Dong, et al. "Gradient Diversity: a Key Ingredient for Scalable Distributed Learning." AISTATS, 2018.

# Impact of Statistical Heterogeneity



Increasing heterogeneity leads to worse convergence
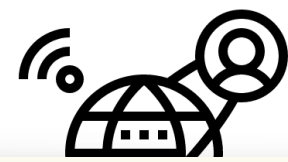
Setting μ > 0 can help to combat this

# How does heterogeneity affect federated optimization methods?

- Heterogeneity can lead to:
  - Slower convergence, reduced stability, divergence
- Critical to analyze and evaluate federated methods with:
  - Non-IID data, partial / variable participation

# Can we equalize performance across heterogeneous networks?

# FL: Traditional Empirical Risk Minimization

ERM: $\min\limits_{w} \left( p_1 \; F_1 \quad +p_2 \; F_2 + \quad \cdots \quad +p_N \; F_N \right)$

no accuracy guarantees for individual devices

Can we encourage a more **fair** (i.e., more **uniform**) distribution
of the model performance across devices?

\#

test accuracy

0.2      0.4      0.6      0.8

# Fair Resource Allocation Objective

$$\textbf{q-FFL:} \quad \min_{w} \frac{1}{q+1}\left( p_1 F_1^{q+1} + p_2 F_2^{q+1} + \cdots + p_N F_N^{q+1} \right)$$

- Inspired by $\alpha$-fairness for fair resource allocation in wireless networks

- A tunable **framework** ($q = 0$: previous objective; $q = \infty$: minimax fairness*)

*Fairness without Demographics in Repeated Loss Minimization, Hashimoto et al, ICML 2018
*Agnostic Federated Learning, Mohri, Sivek, Suresh, ICML 2019

# Fair Resource Allocation Objective

**q-FFL:** $\min_{w} \dfrac{1}{q+1} \left( p_1 \, F_1^{\,\textcolor{red}{q+1}} + p_2 \, F_2^{\,\textcolor{red}{q+1}} + \cdots + p_N \, F_N^{\,\textcolor{red}{q+1}} \right)$

- Theory

✓ Generalization guarantees (recover the known case of $q \rightarrow \infty$)

✓ Increasing $q$ results in more 'uniform' accuracy distributions (in terms of various uniformity measures such as variance)

# Fair Resource Allocation Objective

**q-FFL:** $\min\limits_{w} \dfrac{1}{q+1}\left(p_1\, F_1^{\,q+1} + p_2\, F_2^{\,q+1} + \cdots + p_N\, F_N^{\,q+1}\right)$



test accuracy

0.2    0.4    0.6    0.8

Baseline

q-FFL

#

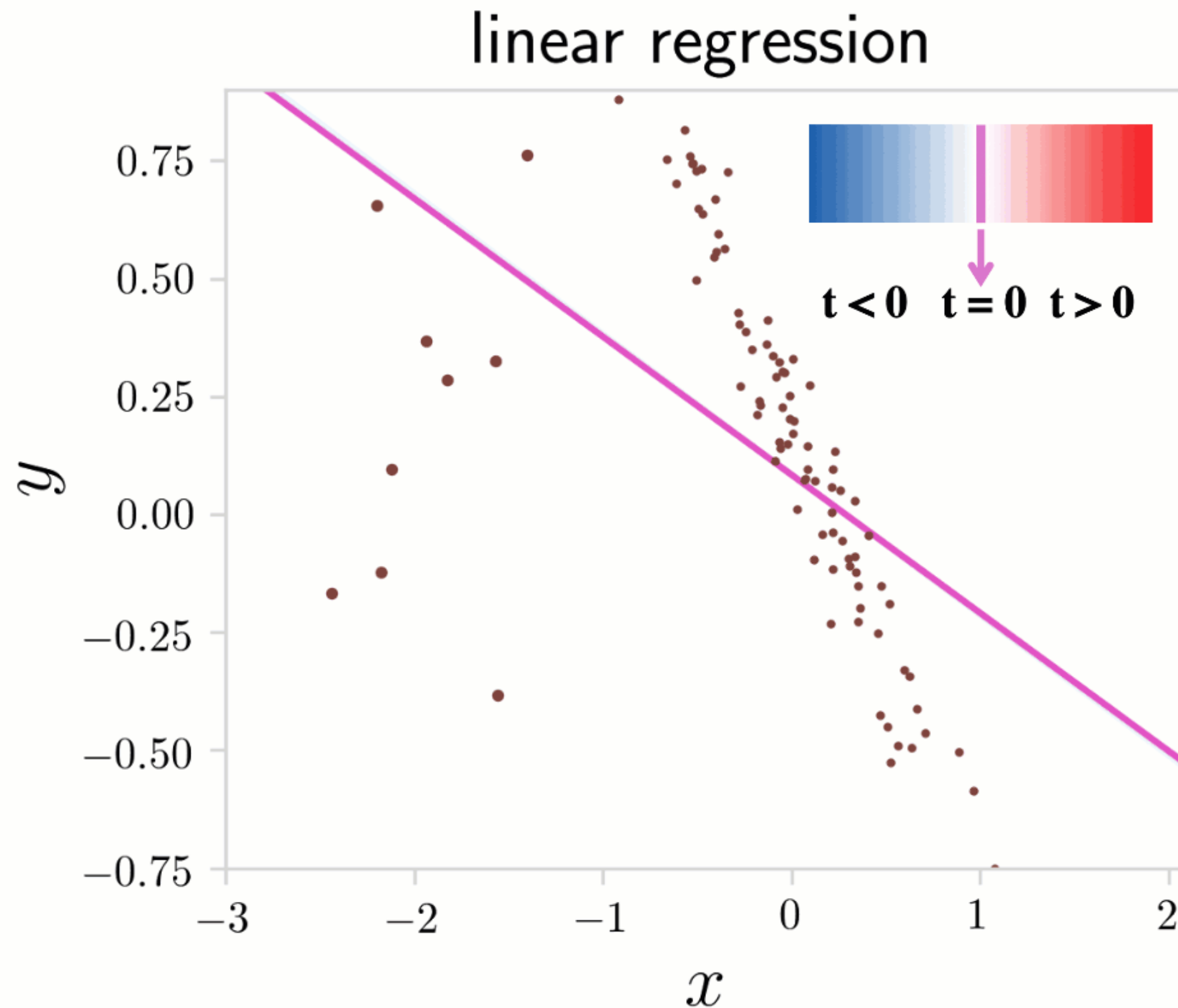# Empirical Results



on average, cut variance of accuracy by 45% while maintaining mean accuracy

# Tilted ERM (TERM) Objective

## linear regression



t < 0   t = 0   t > 0

### Empirical Risk Minimization

$$\min_{w} \frac{1}{n} \sum_{i=1}^{n} f(x_i; w)$$

### Tilted ERM

$$\min_{w} \frac{1}{t} \log \left( \frac{1}{n} \sum_{i=1}^{n} e^{t f(x_i; w)} \right)$$

TERM can increase or decrease the influence of outliers to enable fairness or robustness
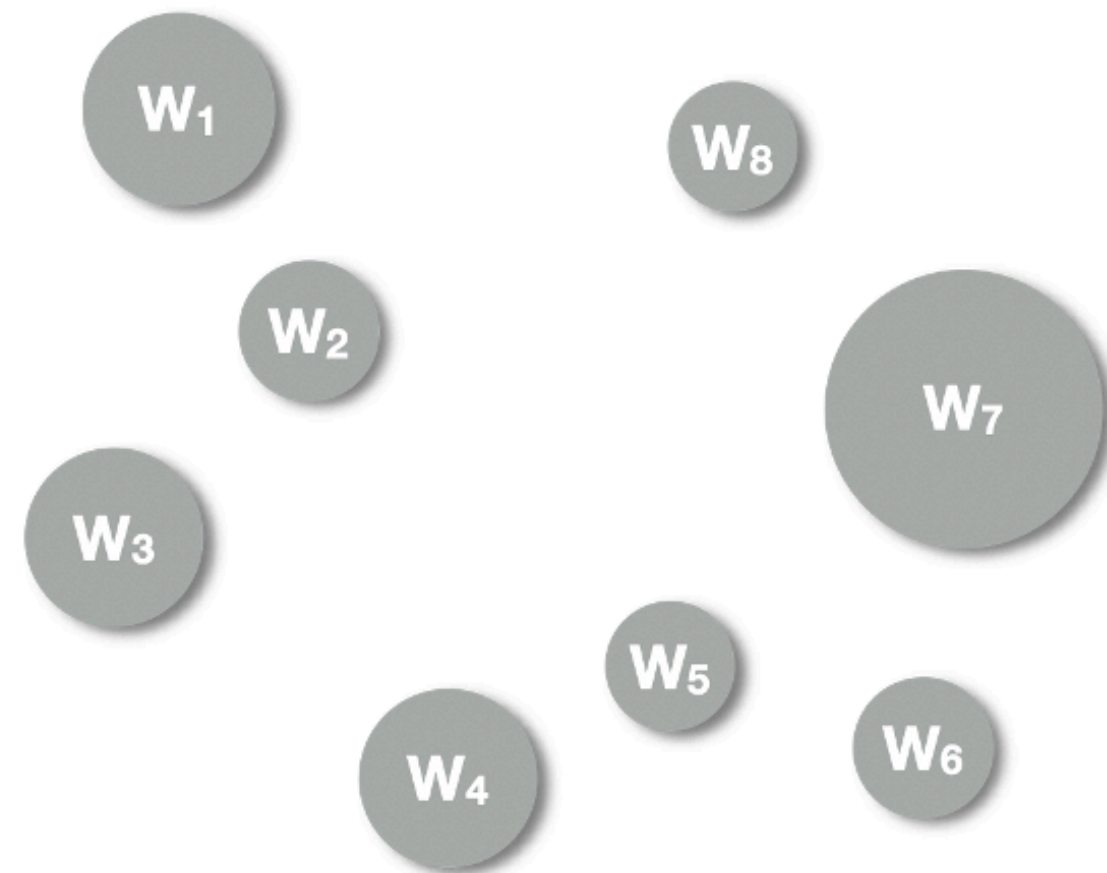
26

# Can we <span style="color:red">equalize</span> performance across heterogeneous networks?

- Vanilla ERM may deliver poor quality of service in heterogeneous networks
- q-FFL/TERM allows for flexible trade-off between fairness and accuracy
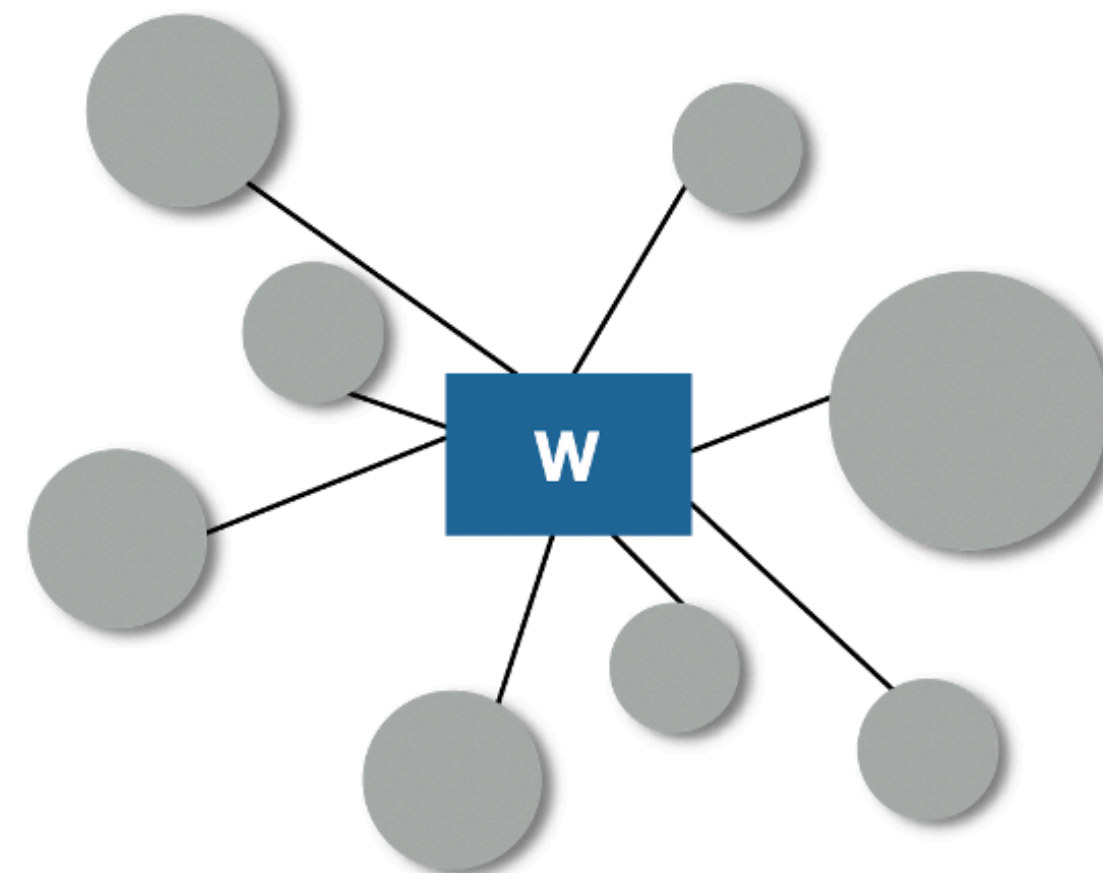
# How to model federated data?

# Personalization for Federated Learning



local

global

??

personalized models
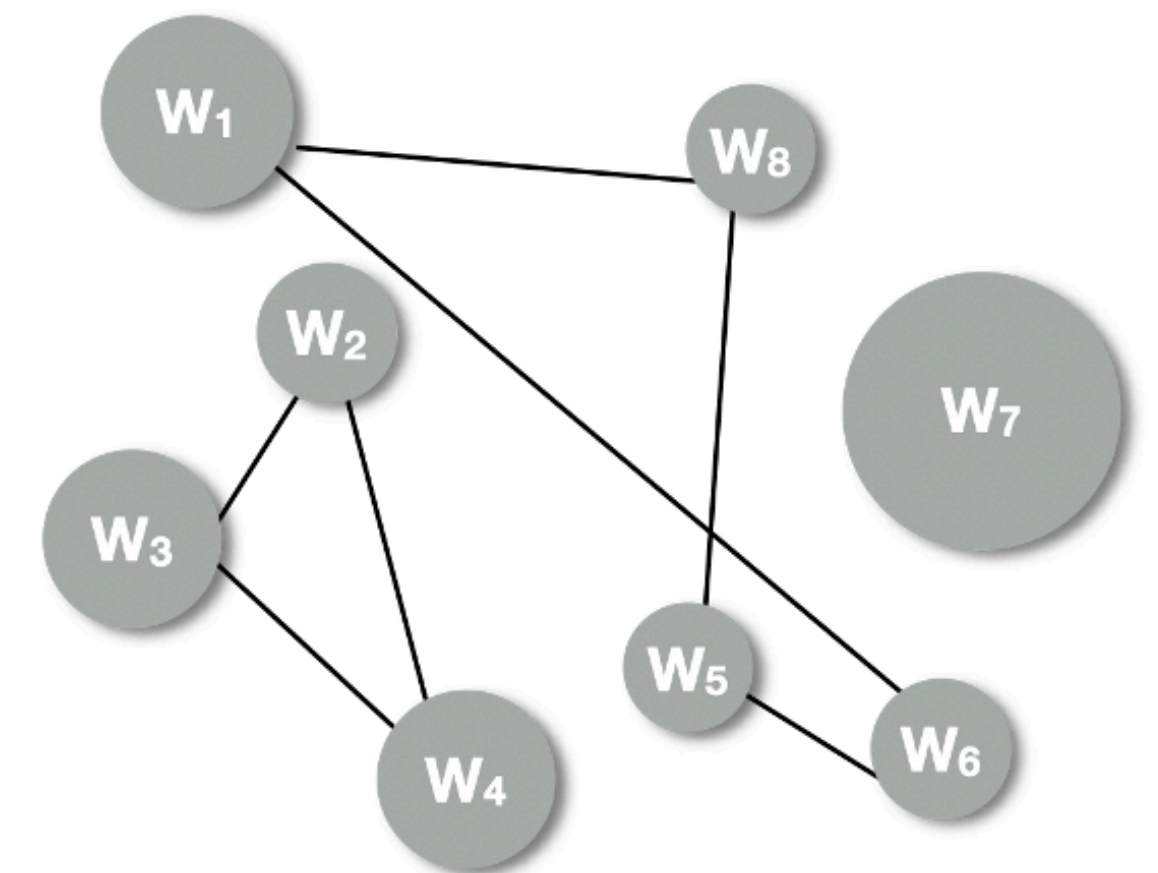not learn from peers
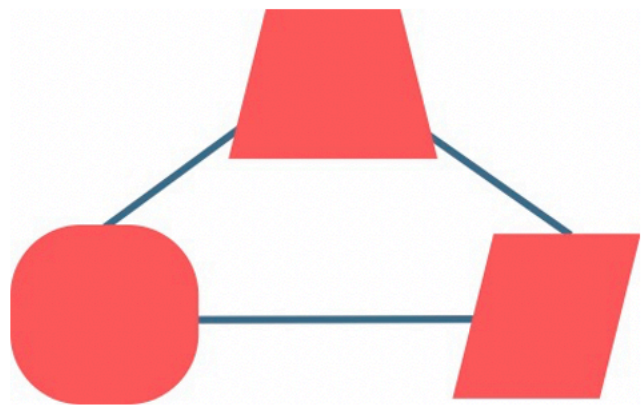
non-personalized models
learn from peers
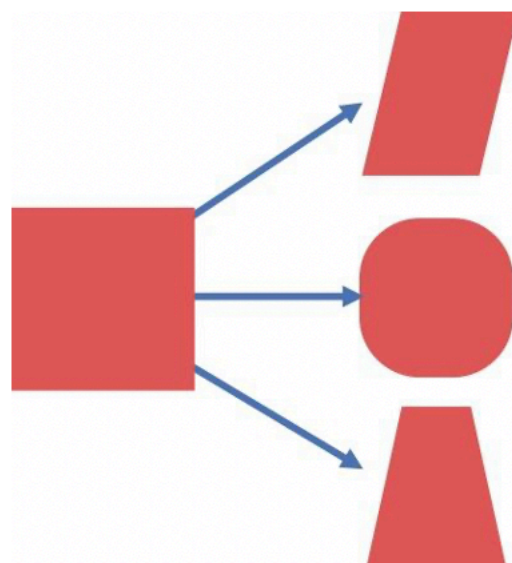
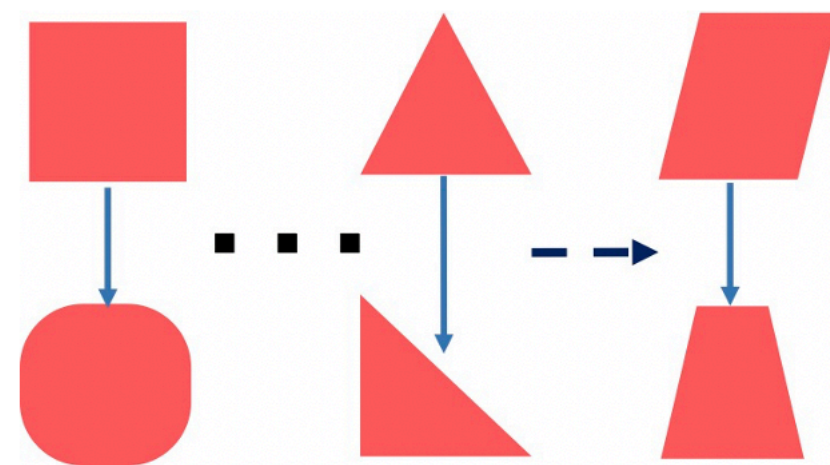personalized models
learn from peers

# Approaches for Personalization

Multi-task Learning

Jointly learn shared, yet personalized models

Fine-tuning
- Learn a global model, then "fine-tune"/adapt it on local data
- See also: transfer learning, domain adaptation

Meta-learning
- Learn initialization over multiple tasks, then train locally

# Meta-learning & Federated learning

**Algorithm 1** Connects FL and MAML (left), Reptile Batch Version(middle), and FedAvg (right).

OuterLoop/Server learning rate $\alpha$
InnerLoop/Client learning rate $\beta$
Initial model parameters $\theta$
**while** not done **do**
    Sample batch of tasks/clients $\{T_i\}$
    **for** Sampled task/client $T_i$ **do**
        **if** FL **then**
            $g_i, w_i = ClientUpdate(\theta, T_i, \beta)$
        **else if** MAML **then**
            $g_i = InnerLoop(\theta, T_i, \beta)$
        **end if**
    **end for**
    **if** FL **then**
        $\theta = ServerUpdate(\theta, \{g_i, w_i\}, \alpha)$
    **else if** MAML **then**
        $\theta = OuterLoop(\theta, \{g_i\}, \alpha)$
    **end if**
**end while**

**Require:** : Reptile Step $K$.
**function** $InnerLoop(\theta, T_i, \beta)$
    Sample $K$-shot data $D_{i,k}$ from $T_i$.
    $\theta_i = \theta$
    **for** each local step i from 1 to K **do**

$$\theta_i = \theta_i - \beta \nabla_\theta L(\theta_i, D_{i,k})$$

    **end for**
    Return $g_i = \theta_i - \theta$
**end function**
**Require:** : Meta Batch Size $M$.
**function** $OuterLoop(\theta, \{g_i\}, \alpha)$

$$\theta = \theta + \alpha \frac{1}{M} \sum_{i=1}^{M} g_i$$

    Return $\theta$
**end function**

**Require:** FedAvg Local Epoch $E$.
**function** $ClientUpdate(\theta, T_i, \beta)$
    Split local dataset into batches $B$
    $\theta_i = \theta$
    **for** each local epoch i from 1 to E **do**
        **for** batch $b \in B$ **do**
            $\theta_i = \theta_i - \beta \nabla_\theta L(\theta_i, b)$
        **end for**
    **end for**
    Return $g_i = \theta_i - \theta$
**end function**
**Require:** Clients per training round $M$.
**function** $ServerUpdate(\theta, \{g_i, w_i\}, \alpha)$

$$\theta = \theta + \alpha \sum_{i=1}^{M} w_i g_i / \sum_{i=1}^{M} w_i$$

    Return $\theta$
**end function**

[Jiang et al, Improving federated learning personalization via model agnostic meta learning, 2019]
[Khodak, Balcan, Talwalkar, Adaptive gradient-based meta-learning methods, NeurIPS 2019

# Personalization for Practical Constraints

constraints in federated learning

| fairness | *representation disparity* |
| robustness | *against data and model poisoning attacks* |

privacy

security

communication

……

competing with each other

$$w* \in \arg\min_{w} G\left(F_1(w), \dots, F_k(w)\right)$$



Subject
Thank you for the feedback `tab`

Ditto: Fair and Robust Federated Learning Through Personalization
Li, Hu, Beirami, Smith, ArXiv 2021
Best paper at ICLR Secure ML Workshop

# Ditto: Global-regularized Federated MTL

*personalization* to achieve robustness and fairness simultaneously

for each device k,

global-regularized

local loss

Ditto:

$$\min_{v_k} \quad h_k(v_k; w^*) := F_k(v_k) + \frac{\lambda}{2}\|v_k - w^*\|^2$$

$$\text{s.t.} \quad w^* \in \arg\min_w G\big(F_1(w), \ldots, F_k(w)\big)$$

✳ simple form of MTL: ensure personalized models are close to global model
✳ easy to implement in federated settings
✳ accurate, robust, and fair

# Ditto Solver

solver for the global model $w^*$ <span style="background-color:#f9dede">+ personalization add-on</span>
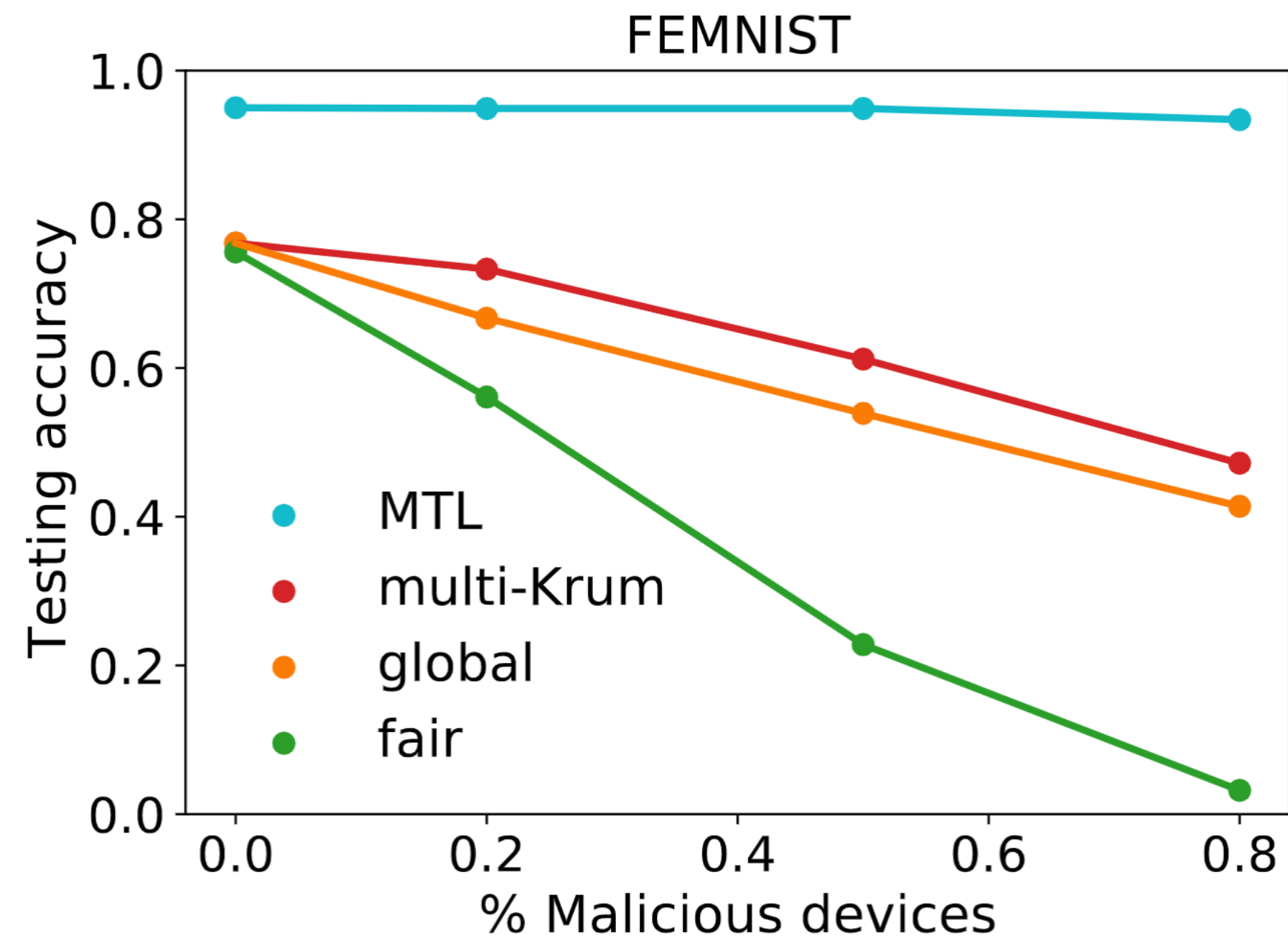
---

**Algorithm 1:** `Ditto` for Personalized FL

---

1   **Input:** $K$, $T$, $s$, $\lambda$, $\eta$, $w^0$, $\{v_k^0\}_{k\in[K]}$

2   **for** $t = 0, \cdots, T-1$ **do**

3     Server randomly selects a subset of devices $S_t$, and sends the current global model $w^t$ to them

4     **for** *device* $k \in S_t$ *in parallel* **do**

5       Solve the local sub-problem of $G(\cdot)$ inexactly starting from $w^t$ to obtain $w_k^t$:

$$w_k^t \leftarrow \text{UPDATE\_GLOBAL}(w^t, \nabla F_k(w^t))$$

/* Solve $h_k(v_k; w^t)$ */

6       Update $v_k$ for $s$ local iterations:

$$v_k = v_k - \eta(\nabla F_k(v_k) + \lambda(v_k - w^t))$$

      Send $\Delta_k^t := w_k^t - w^t$ back

7     Server aggregates $\{\Delta_k^t\}$:

$$w^{t+1} \leftarrow \text{AGGREGATE}\left(w^t, \{\Delta_k^t\}_{k\in\{S_t\}}\right)$$

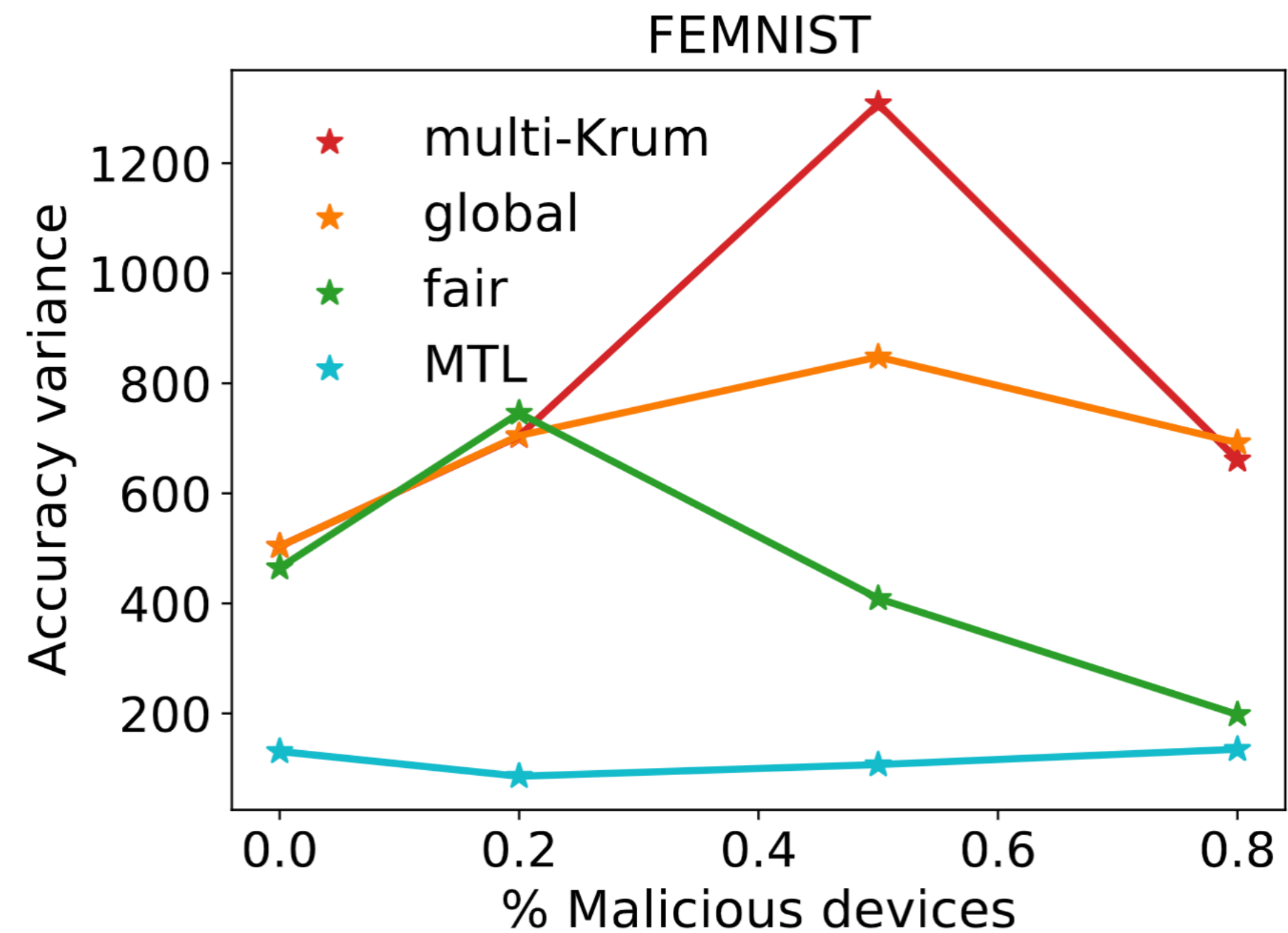8   **return** $\{v_k\}_{k\in[K]}$ *(personalized models)*, $w^T$ *(global model)*

---

✱ a scalable, simple personalization add-on for any federated global solver
✱ preserves the practical properties of the global FL solver (e.g., communication, privacy)
✱ with convergence guarantees
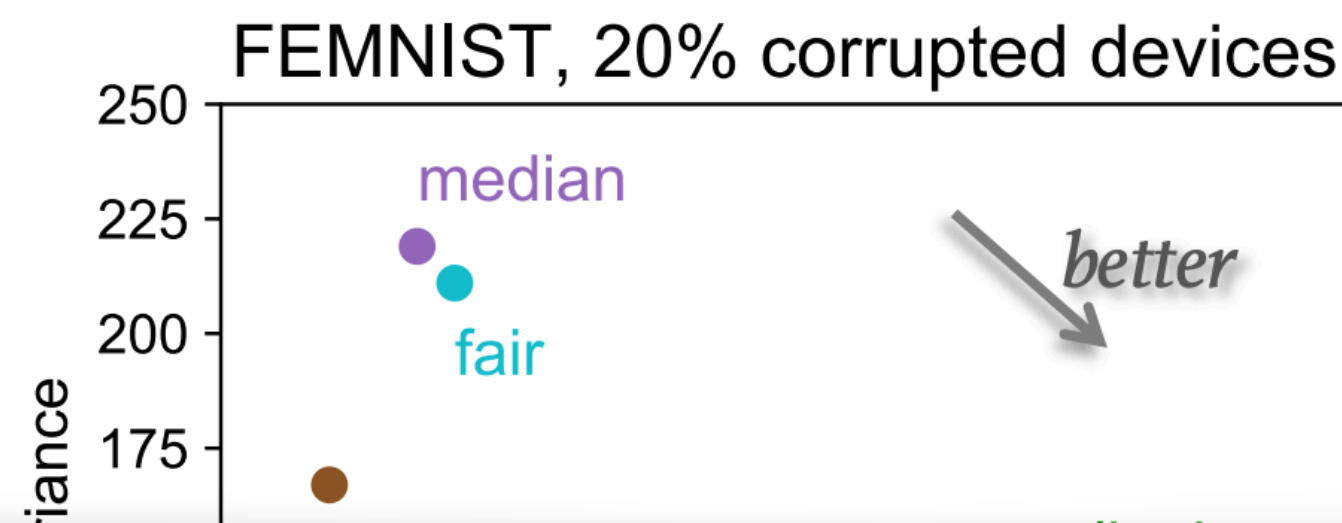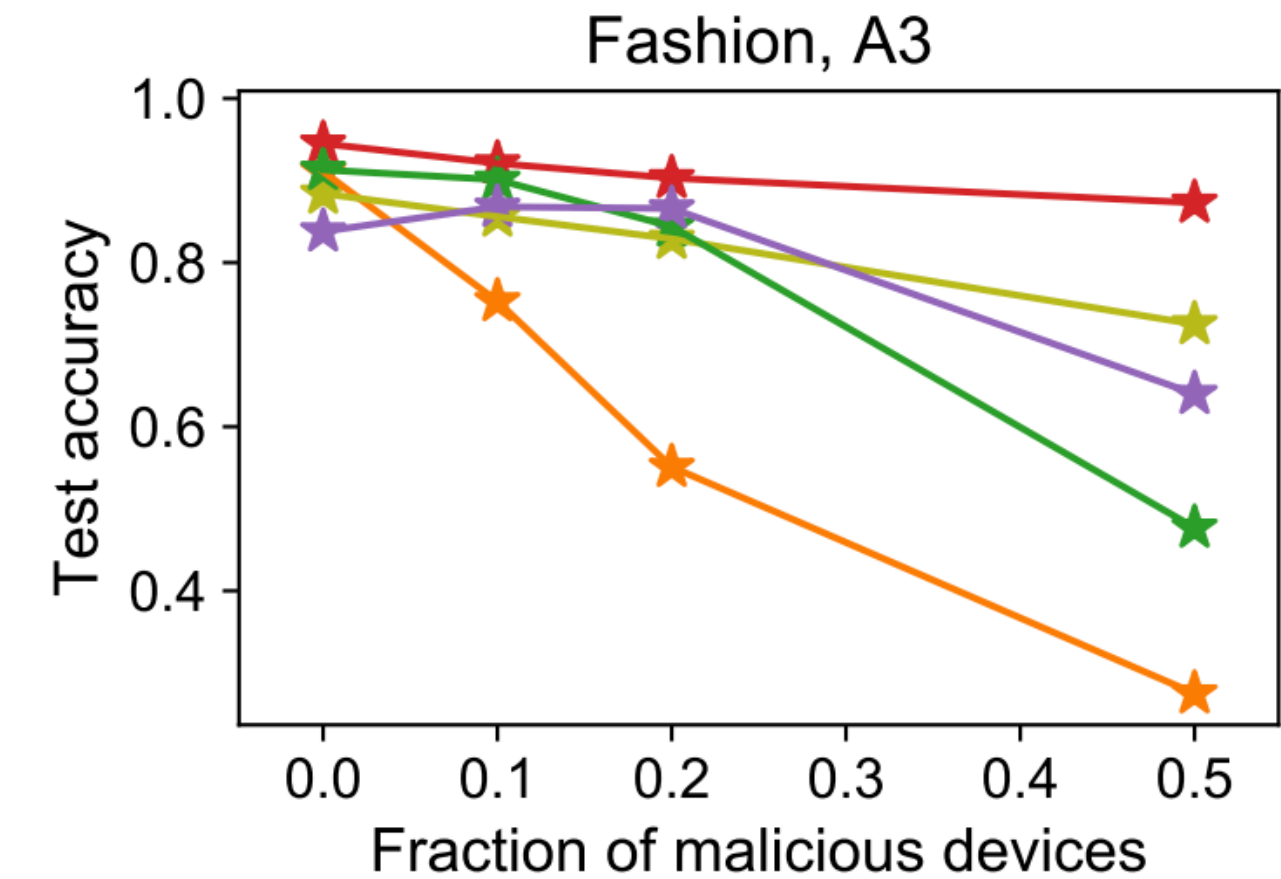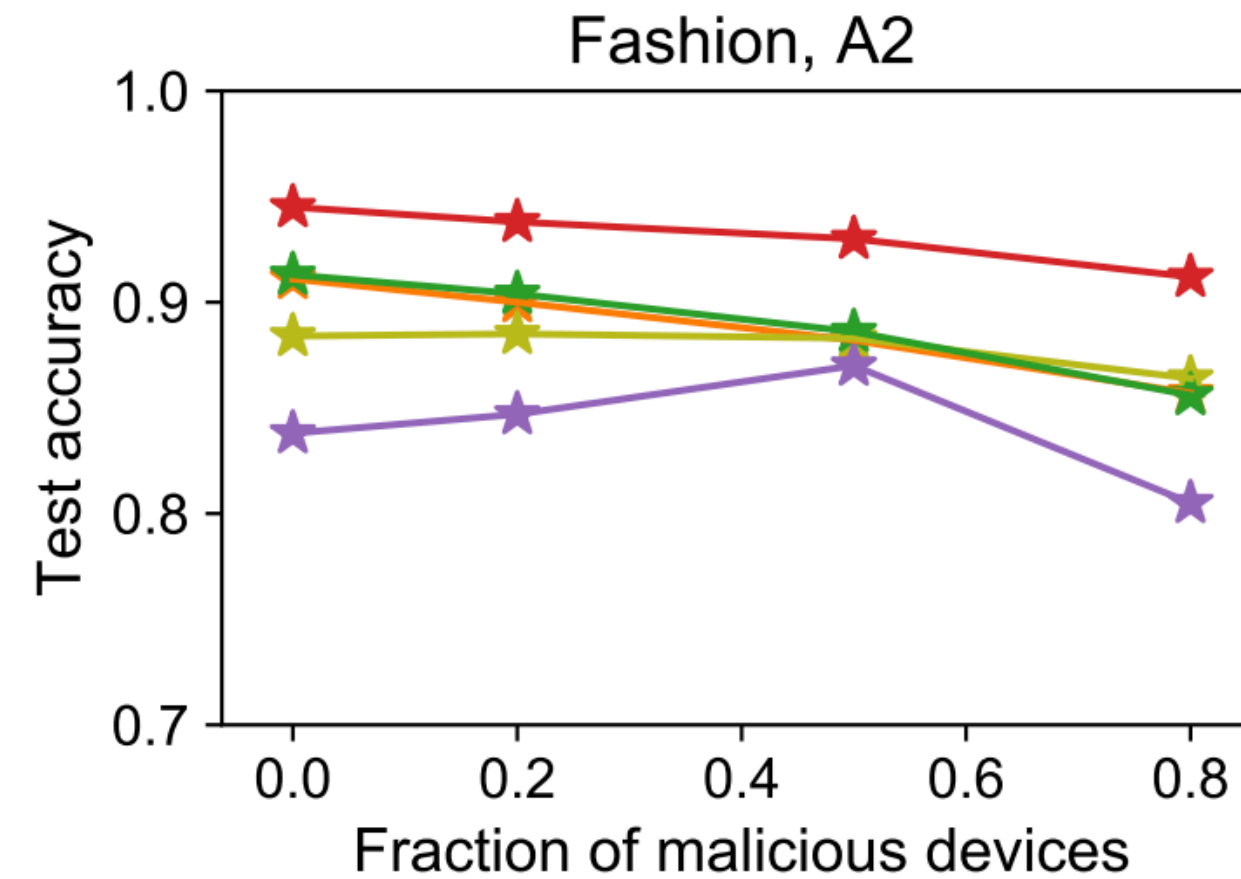
# Experiments: Competing Constraints



FEMNIST

fair methods are not robust

robust methods are not fair (with high variance)

# Experiments: Benefits of Personalization



Fashion, A1

Fashion, A2

Fashion, A3

Legend:
- Ditto
- global
- median
- clipping
- Krum

FEMNIST, 20% corrupted devices

better

Ditto is also more fair

Ditto is more robust than strong baselines under various attacks

on average, improve absolute accuracy by ~6% over the strongest robust baseline

reduce variance by ~10% over STOA fair methods

# How to model federated data?
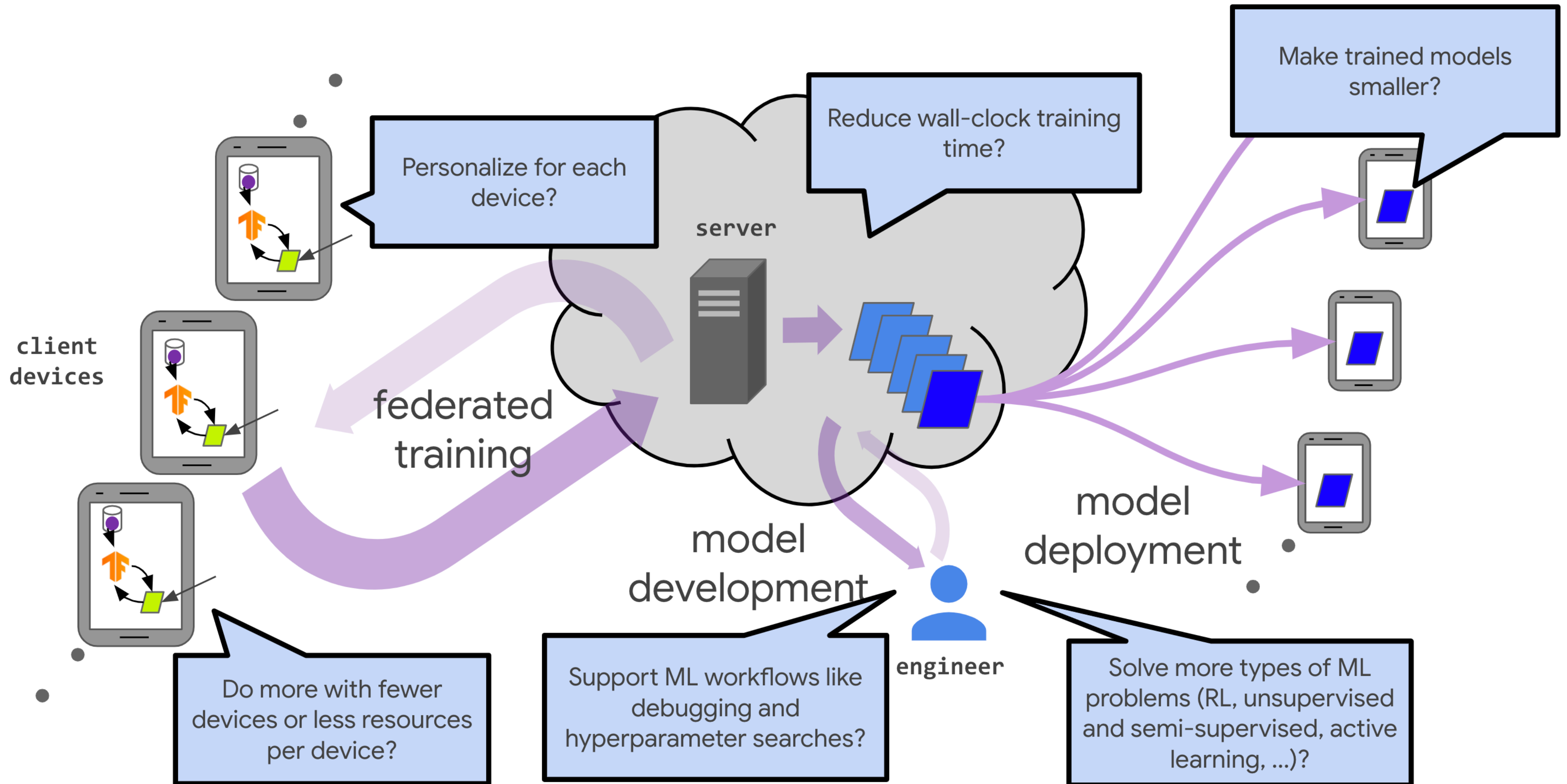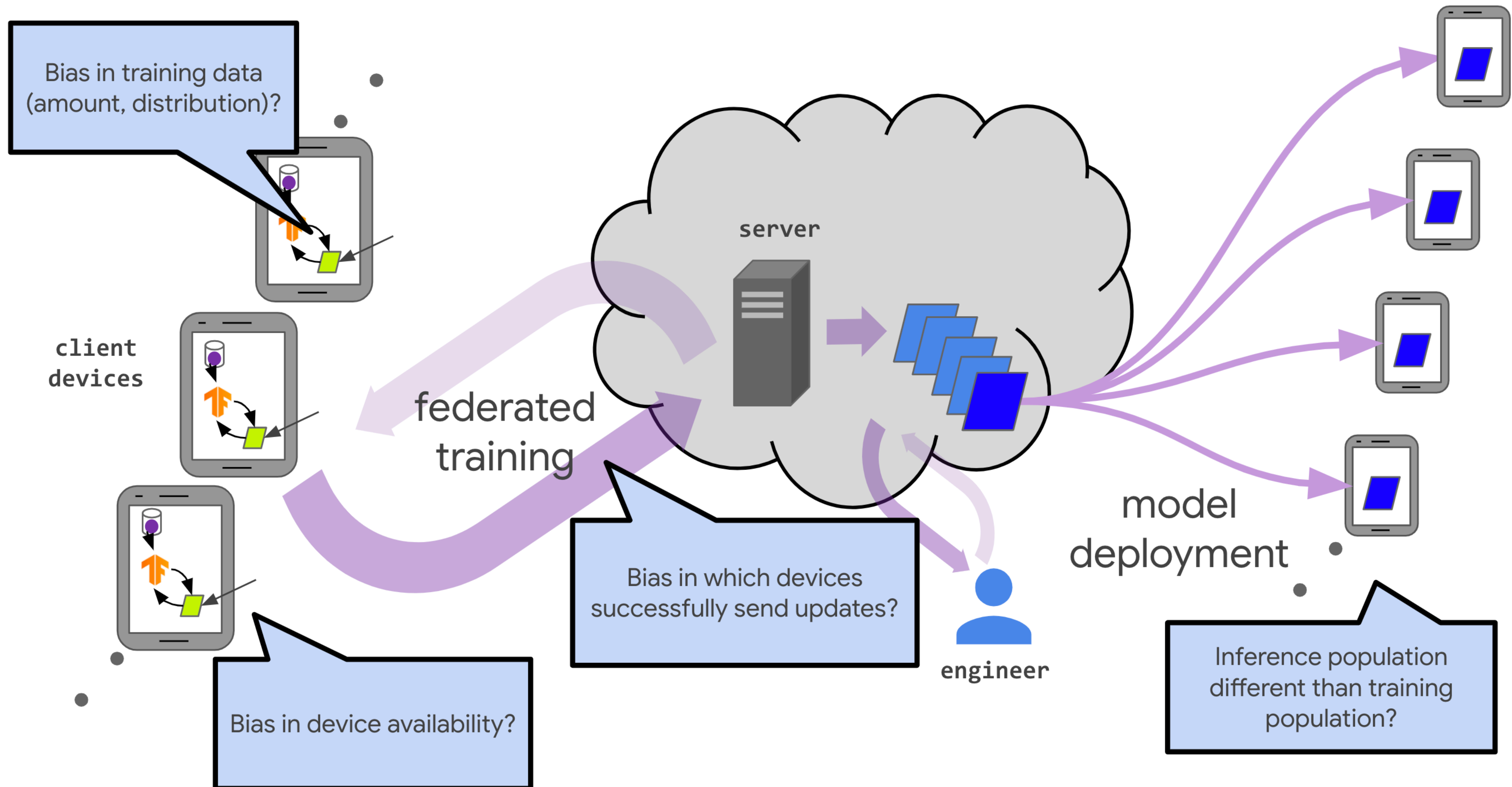
- Personalization is a promising approach (need to be scalable, automated)
- Personalization has additional benefits beyond accuracy, e.g., fairness, robustness, etc.

# What's next??

# Improving efficiency and effectiveness

# Ensuring fairness and addressing sources of bias



Bias in training data (amount, distribution)?

client devices

Bias in device availability?

federated training

server

Bias in which devices successfully send updates?

engineer

model deployment

Inference population different than training population?

[Credit: B. McMahan, FL Tutorial, NeurIPS 2020]

# Robustness to attacks and failures



Compromised device sending malicious updates

Inference-time evasion attacks

server

client devices

federated training

Devices training on compromised data (data poisoning)

Device dropout, data corruption in transmission

model development

engineer

model deployment

[Credit: B. McMahan, FL Tutorial, NeurIPS 2020]

# Additional Reading

- **FedAvg: Communication-Efficient Learning of Deep Networks from Decentralized Data**, McMahan et al, AISTATS 2017

- **MOCHA: Federated Multi-Task Learning**, Smith et al, NeurIPS 2017

- [White Paper] **Federated Learning: Challenges, Methods, and Future Directions**, Li et al, IEEE Signal Processing Magazine, 2020

- **NeurIPS 2020 federated learning tutorial**, https://sites.google.com/view/fl-tutorial

**Carnegie Mellon University**
School of Computer Science

# Questions?

Tian Li

tianli@cmu.edu