



di.unito.it

DIPARTIMENTO DI INFORMATICA

# FEDERATED LEARNING

## *An Overview*

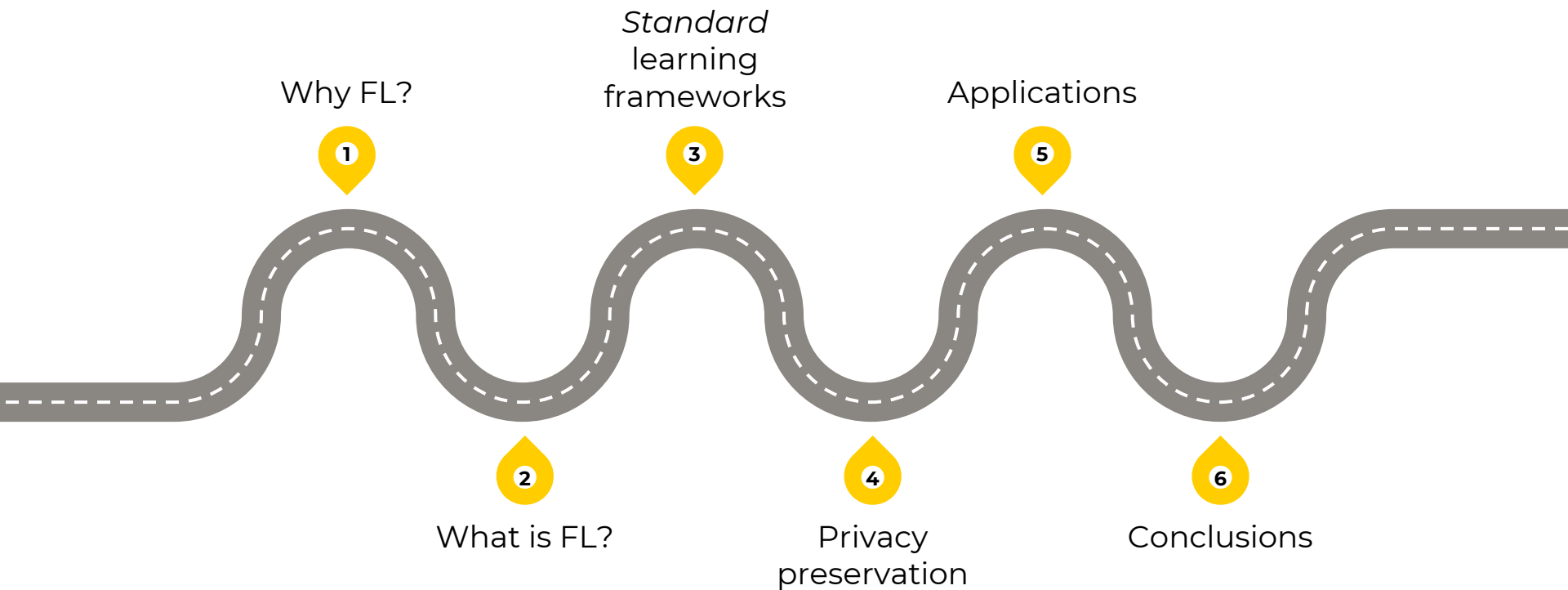


Mirko **Polato**, Ph.D.

Assistant Professor @ Dept. of Computer Science, University of Turin



# Roadmap



---

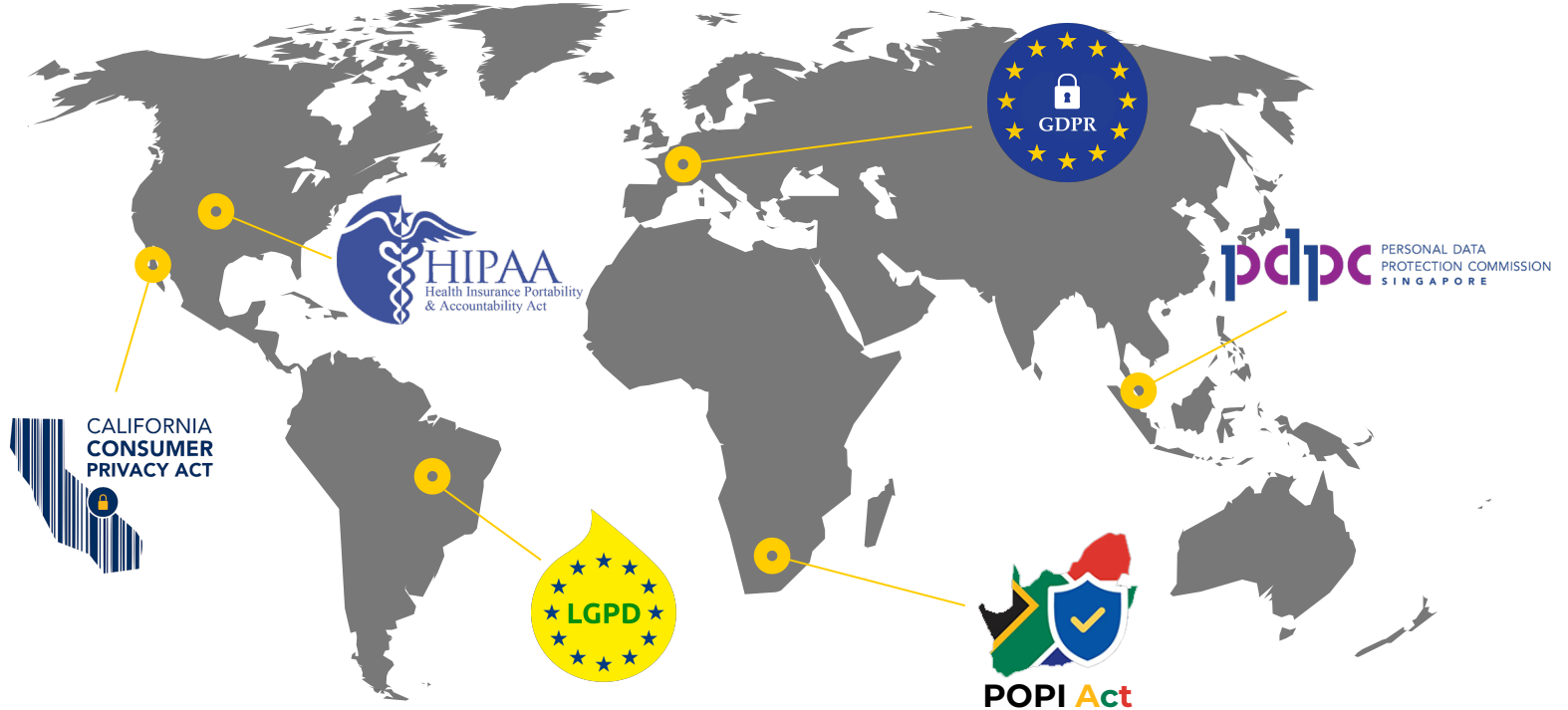
1

# Why FL?

The main reasons behind this new learning paradigm

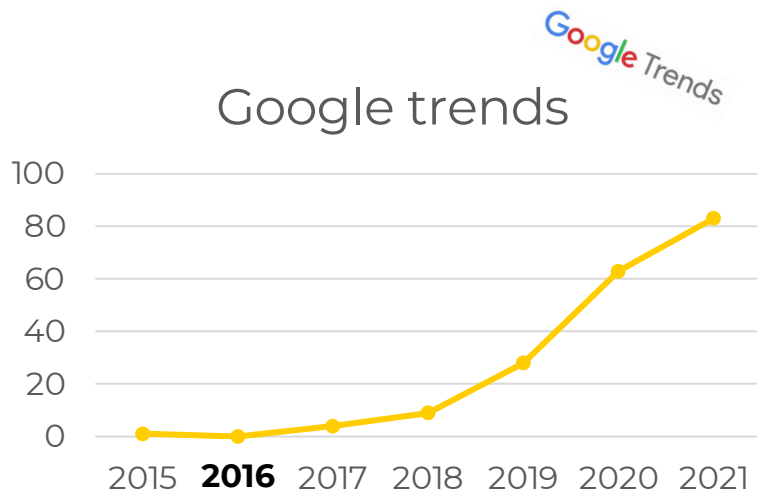
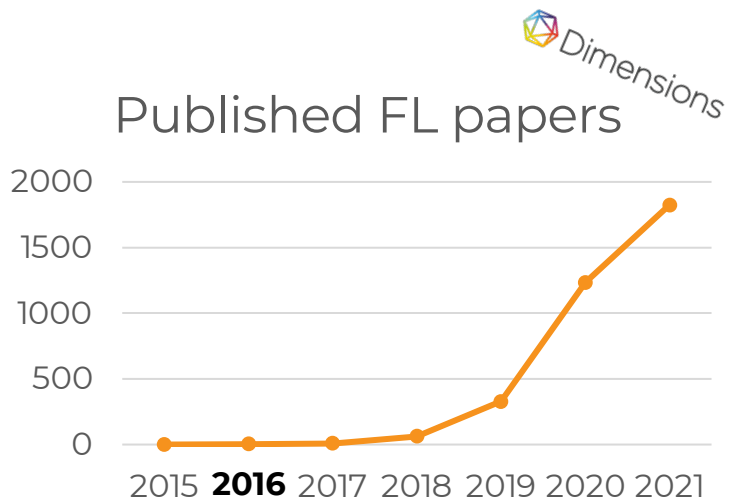


# Data protection regulations





## FL is on fire!



Google publishes the seminal work on FL

---

2

## What is FL?

Definition and overview of the FL taxonomy



## Informal definition

**Federated Learning** is a machine learning setting where *multiple entities (clients) collaborate* in solving a machine learning problem, under the coordination of a central server. *Each client's raw data is stored locally and not exchanged or transferred; instead, focused updates intended for immediate aggregation are used to achieve the learning objective.*



## “Formal” definition

	Parties	Data	Learning	Performance
<b>Centralized</b>	Server	$\mathcal{D}$	On the server	$P$
<b>Non-federated</b>	Clients	$\mathcal{D} \equiv \cup_{i=1}^N \mathcal{D}_i$	On clients	$P_i \leq P$
<b>Federated</b>	Clients Server	$\mathcal{D} \equiv \cup_{i=1}^N \mathcal{D}_i$	Collaborative	$\hat{P}$

### Goals

*Clients do not share their data*

*Clients benefit from the federation*

*The federated model is close to the “ideal” one*

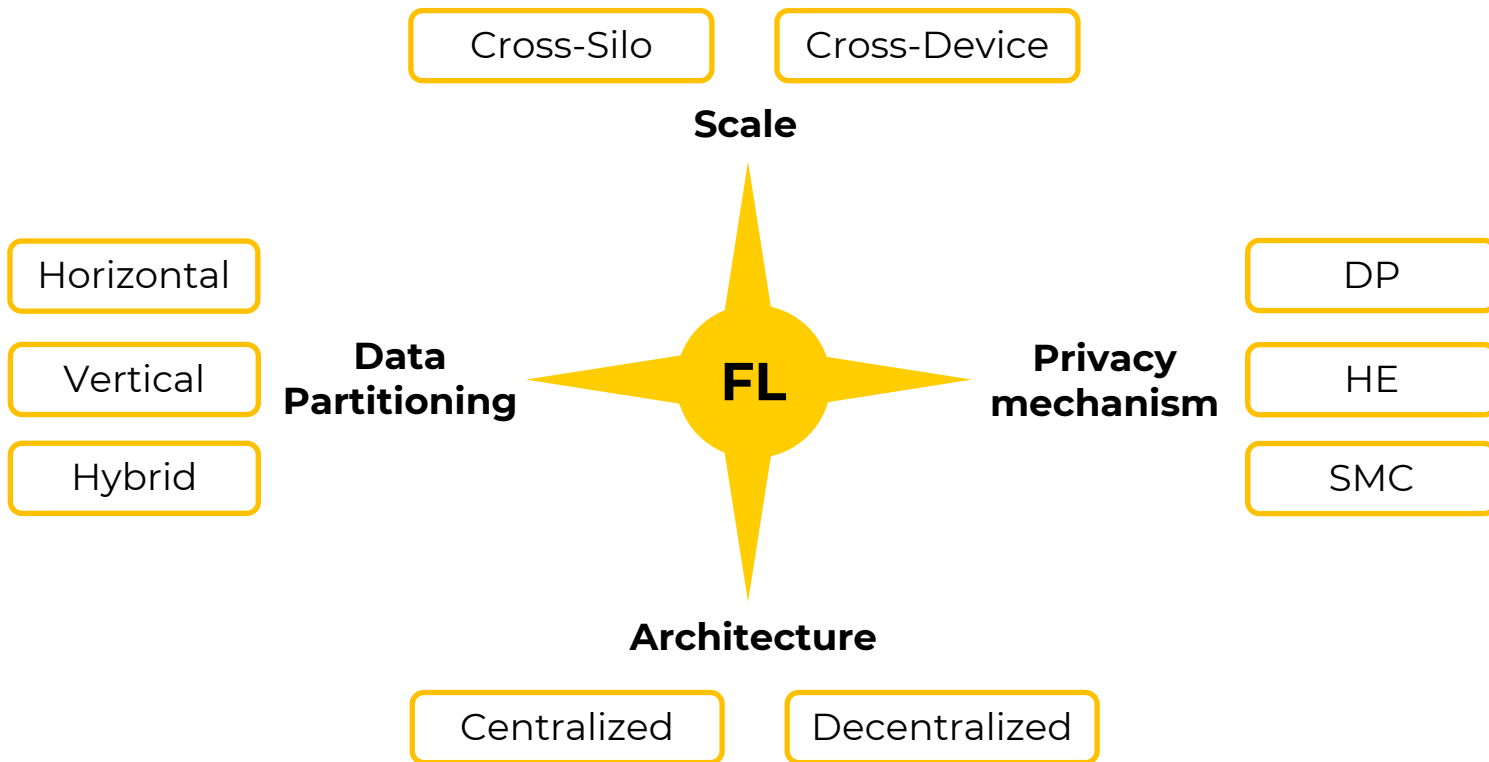
$$P_i \leq \hat{P}$$

$$P - \hat{P} \leq \delta \quad \delta \rightarrow 0$$



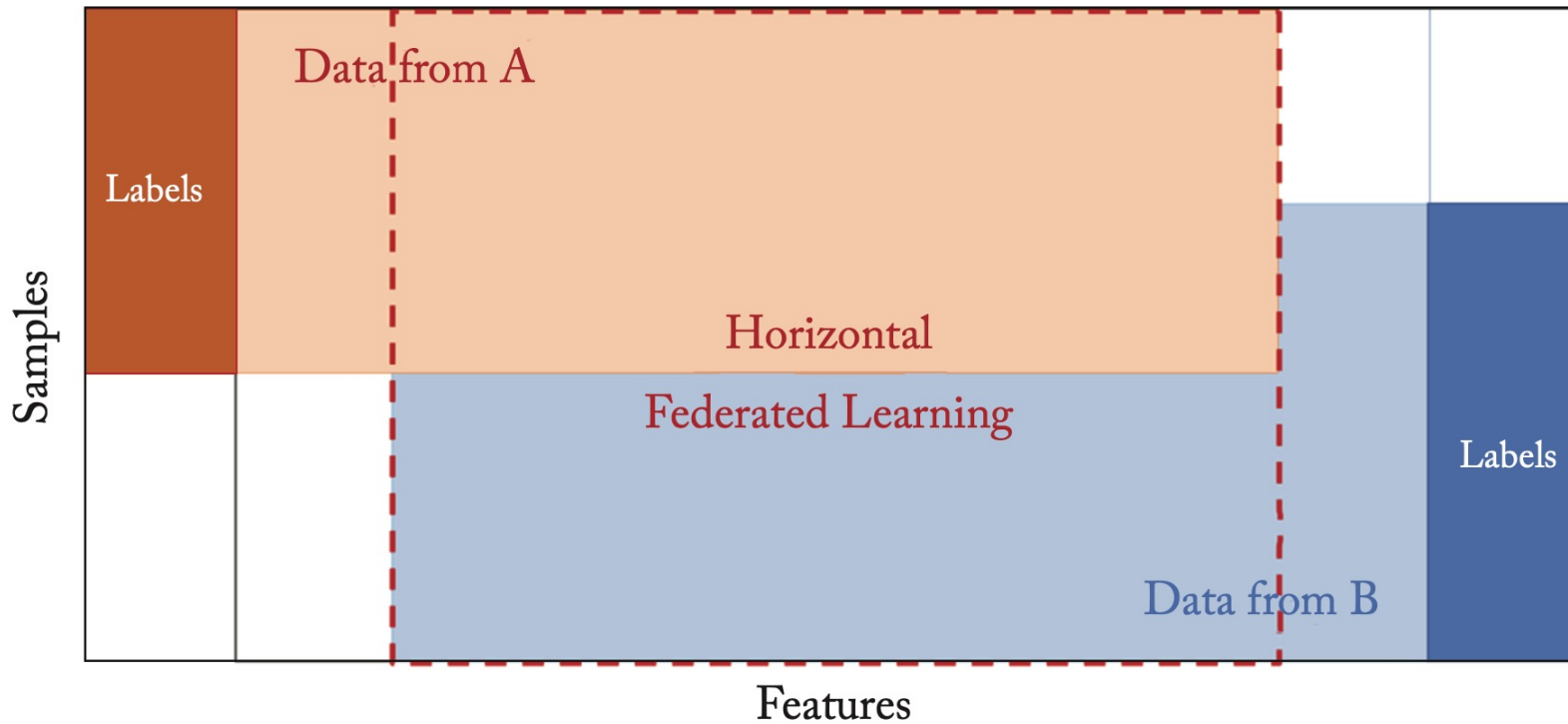


# FL settings landscape



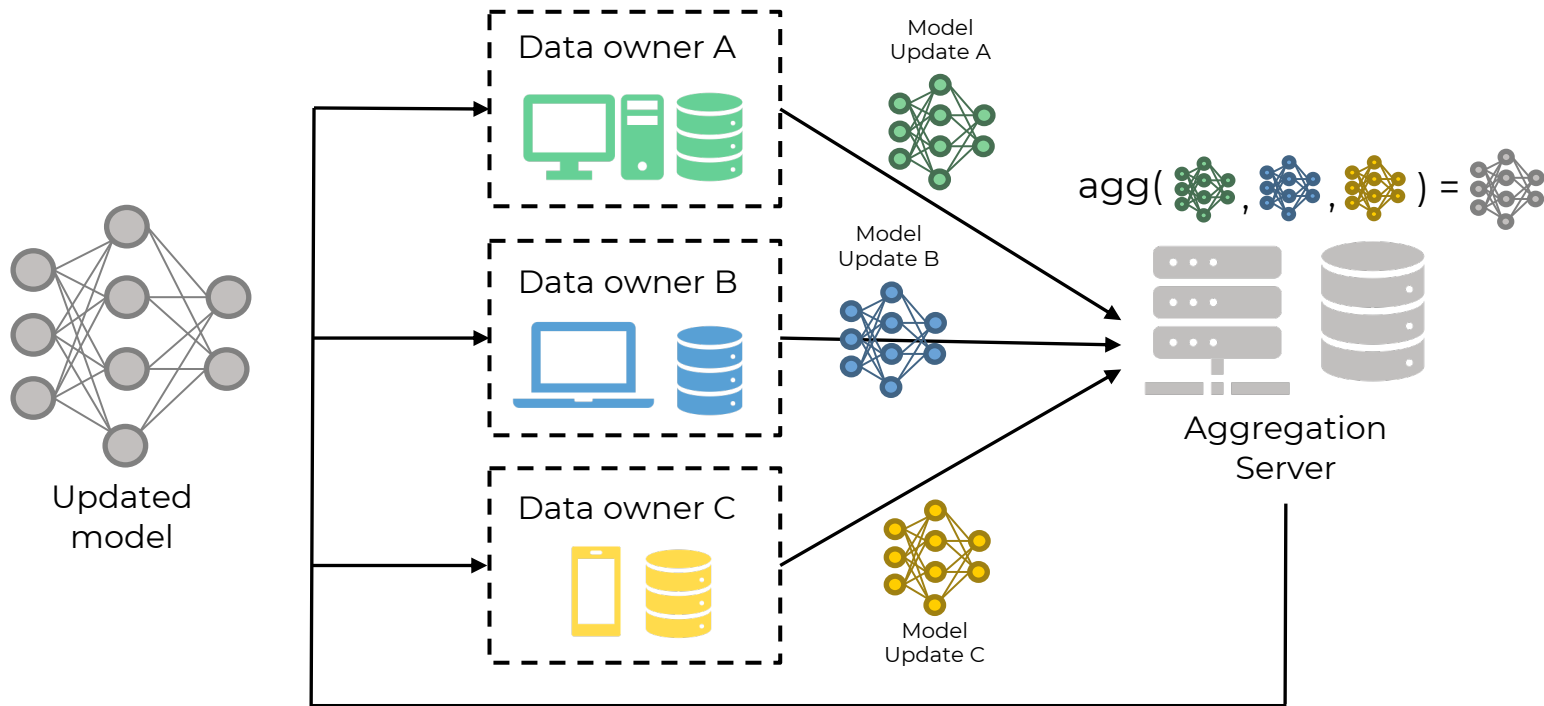


# Horizontal FL





# Standard HFL architecture



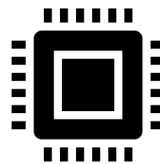


## (Cross-device) FL characteristics

**Large scale:  
thousands of  
devices**



**Limited  
computational  
capacity**



**Unreliable  
connections**



**Power  
consumption**

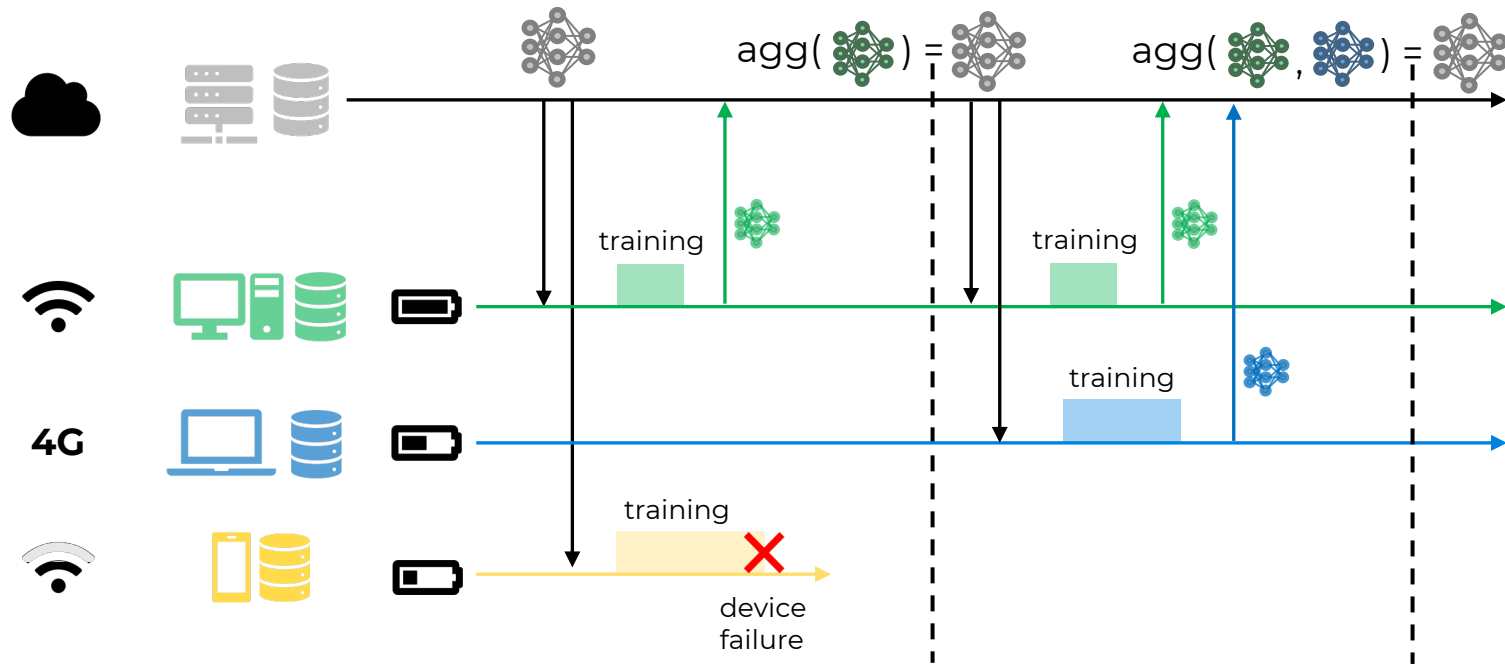


---

*System heterogeneity*

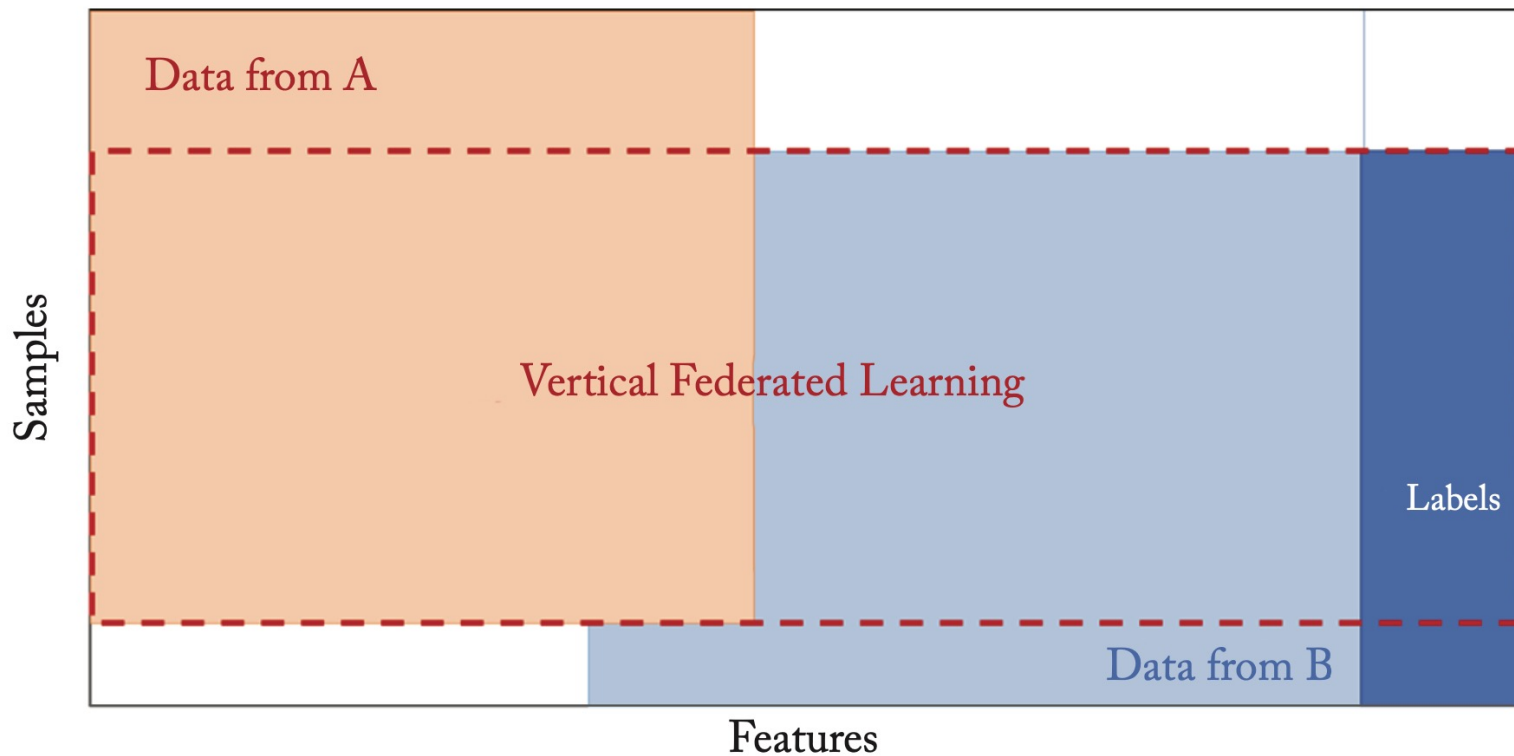


# System heterogeneity



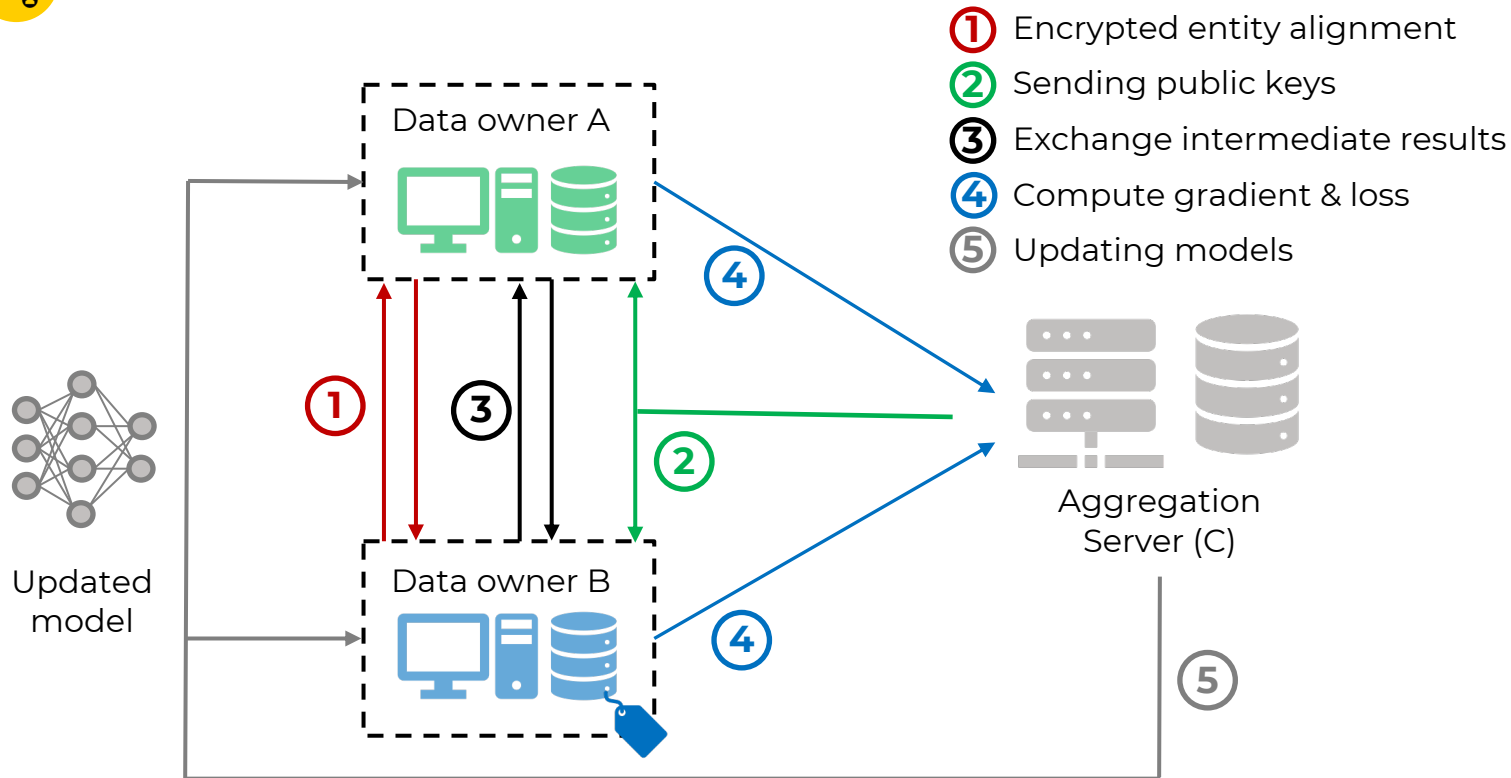


## Vertical FL



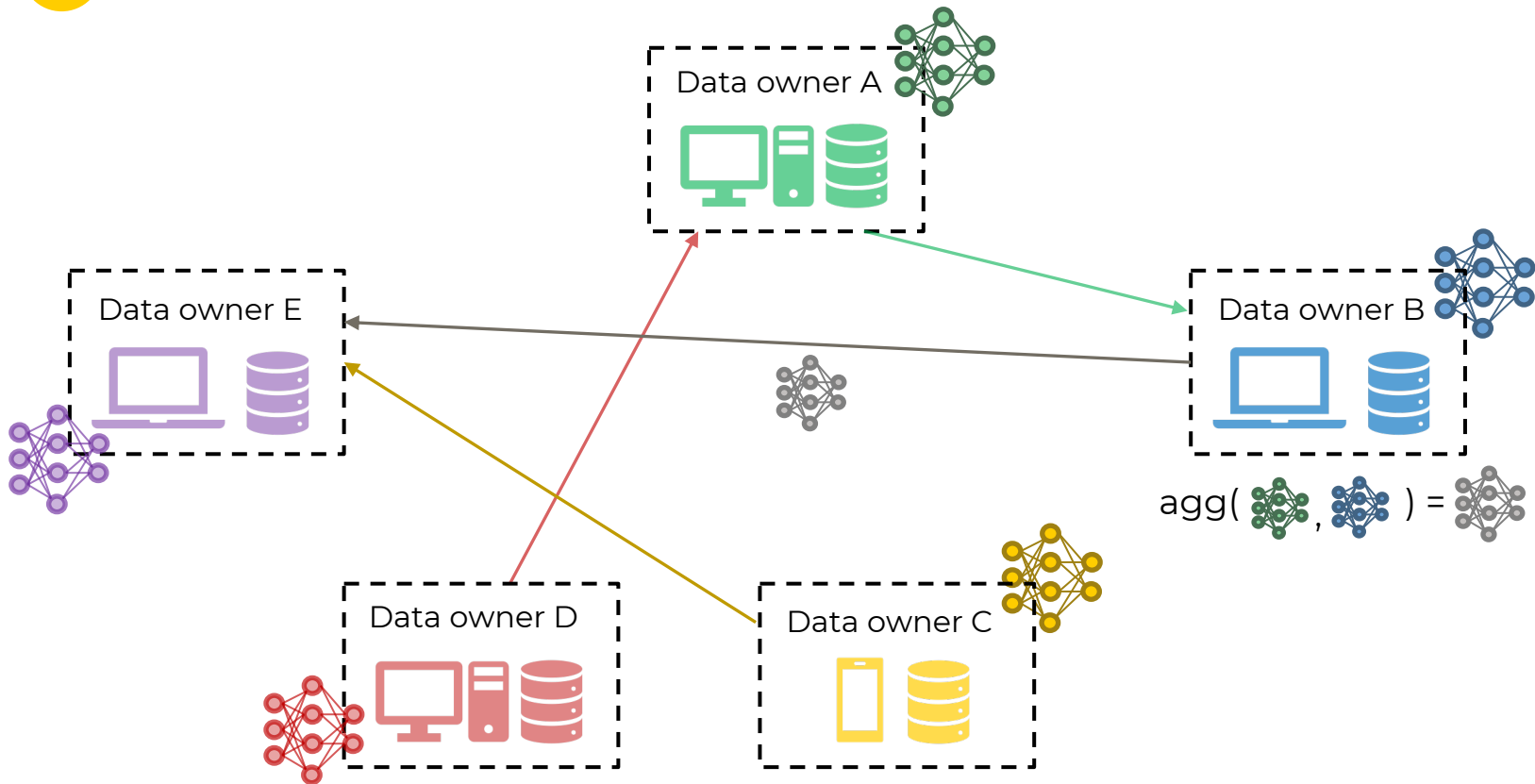


# VFL Architecture





# Decentralized FL







## Centralized vs Decentralized

	<b>Centralized</b>	<b>Decentralized</b>
<b>Orchestration</b>	Server	No centralized orchestration*
<b>Topology</b>	Hub-and-spoke	Peer-to-peer
<b>Global model</b>	Single	Many
<b>Setup</b>	Centralized	Consensus*

\* A central authority might be needed!



## Cross-Silo FL

**Limited  
number of  
collaborators**

**Big local  
datasets**

**Reliable  
connection/  
participation**

- **Few organizations** share incentives to train a model without sharing their data
- Or, same organization cannot centralized its data (e.g., legal constraints)



## Cross-Silo: incentives

### *Free-rider problem*

Organizations/compatitors may benefit from the federation without contributing as much



Monetary payout



Assign FL model with performance commensurate with the contributions (game theory)



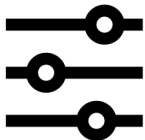
## Cross-Silo vs Cross-Device

	<b>Cross-Silo</b>	<b>Cross-Device</b>
<b>Availability</b>	~Always	Small fraction available
<b>Scale</b>	2-100 clients	Up to $10^{10}$ clients
<b>Addressability</b>	Direct	No client identifier
<b>Reliability</b>	Few failures	Highly unreliable
<b>Dataset size</b>	Big	Relatively small



# Challenges

**Hyper  
parameters  
tuning**



**Avoid  
overfitting**



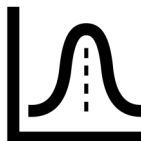
**Model  
debugging**



**Malicious  
participants**



**Non-iidness**



**System  
heterogeneity**





## Non-identical client distributions

$$P_i(x, y) = P_i(y|x)P_i(x) = P_i(x|y)P_i(y)$$

### Covariate shift

$$P_i(y|x) = P_j(y|x)$$

$$P_i(x) \neq P_j(x)$$

### Prior shift

$$P_i(x|y) = P_j(x|y)$$

$$P_i(y) \neq P_j(y)$$

### Concept drift

$$P_i(y) = P_j(y)$$

$$P_i(x|y) \neq P_j(x|y)$$

### Concept shift

$$P_i(y|x) \neq P_j(y|x)$$

$$P_i(x) = P_j(x)$$

### Quantity skew

Clients hold  
hugely different  
amounts of data

3

# *Standard* Learning Frameworks

Learning in a federation



## General Learning Scheme

---

1. Setup task – Model initialization
2. For each federated round:
  - a. Broadcast the current global model
  - b. In parallel:  
clients update and send back the local models
  - c. Update global model with the users' ones





# FedAvg (model averaging)



Aggregation Server

Algorithm: **FedAvg**

1. initialize model  $\bar{w}_0$
2. for each round  $t=1, \dots$ :
3. Broadcast  $\bar{w}_{t-1}$
4. select  $C$  eligible participants
5. foreach|| participant  $p$ :
6.  $w_t^p \leftarrow \text{LocalUpdate}(p)$
7.  $\bar{w}_t \leftarrow \text{aggregate}(\forall p w_t^p)$



Does not guarantee converge



Collaborator/Client

Algorithm: **LocalUpdate**

1.  $w \leftarrow$  global model from Server
2. for each epoch  $s \in 1, \dots, S$ :
3. for each batch  $b$ :
4.  $g_b \leftarrow$  compute gradient for  $b$
5.  $w \leftarrow w - \eta g_b$
6. send  $w$  to the Server



Practically works most of the time



Efficient (communication-wise)



Not bound to SGD



# FedSgd (gradient averaging)



Aggregation Server



Collaborator/Client

Algorithm: **FedSgd**

1. initialize model  $\bar{w}_0$
2. for each round  $t=1, \dots$ :
3. Broadcast  $\bar{w}_{t-1}$
4. select  $C$  eligible participants
5. foreach|| participant  $p$ :
6.  $g_t^p \leftarrow \text{LocalUpdate}(p)$
7.  $g_t \leftarrow \text{aggregate}(\forall p g_t^p)$
8.  $\bar{w}_t \leftarrow \bar{w}_{t-1} - \eta g_t$

Algorithm: **LocalUpdate**

1.  $w \leftarrow$  global model from Server
2. select batch  $b$
3.  $g_b \leftarrow$  compute gradient for  $b$
4. send  $g_b$  to the Server



Inefficient (communication-wise)



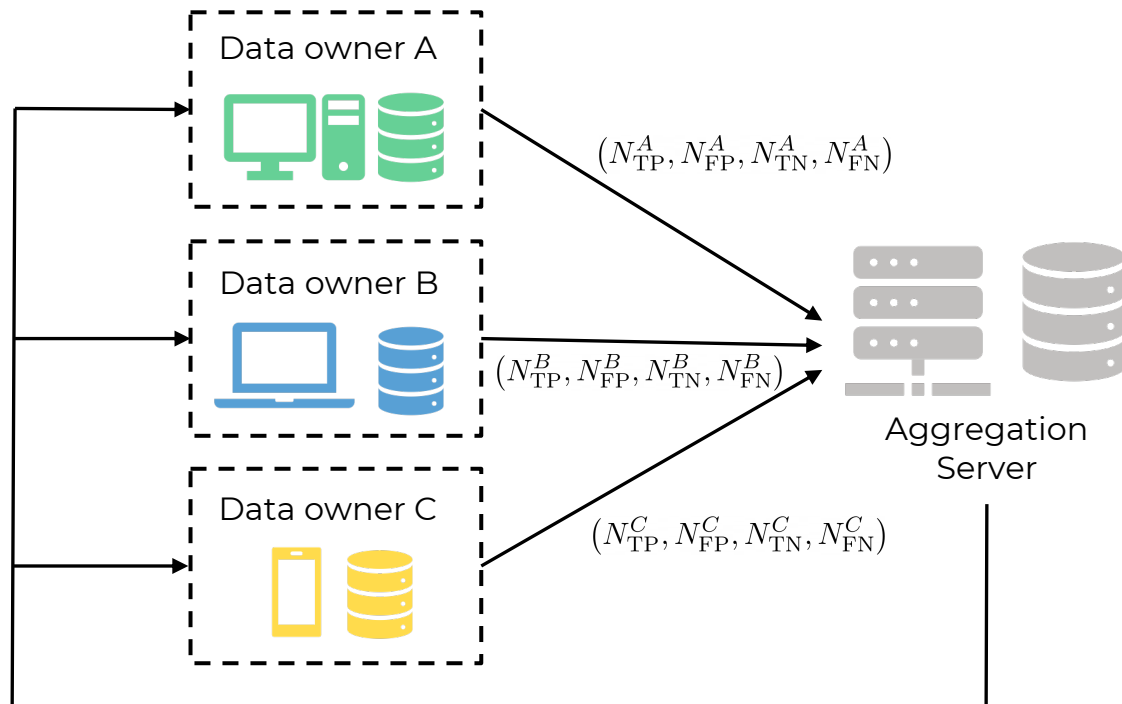
Guaranteed convergence



# Model evaluation (classification)

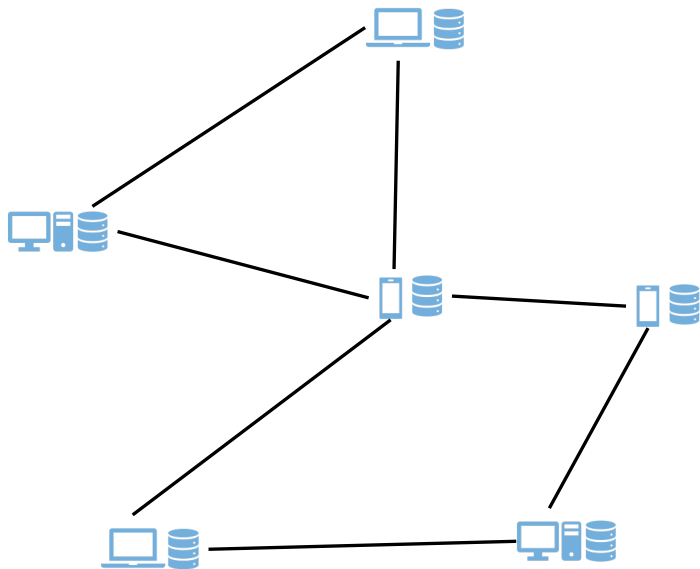
$$\frac{\sum_{k \in \{A, B, C\}} N_{TP}^k}{\sum_{k \in \{A, B, C\}} (N_{TP}^k + N_{FN}^k)}$$

e.g., global recall





# Gossip Learning



## Algorithm: **Main gossip loop (Push)**

1. initialize local model  $w$
2. loop (forever)
3. wait for a fixed time  $\Delta$
4. select neighbor peer  $p$
5. send  $w$  to  $p$

## Algorithm: **On receive model**

1. receive  $w_p$  from  $p$
2. merge( $w$ ,  $w_p$ )
3. update( $w$ )

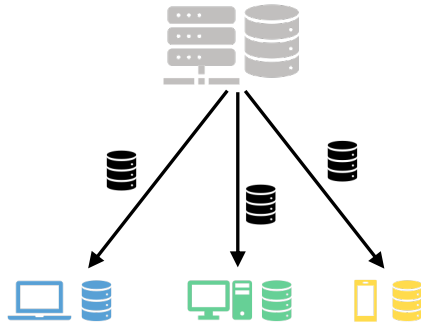
↑  
e.g., gradient  
descent step  
on local data

↖  
e.g., model  
averaging

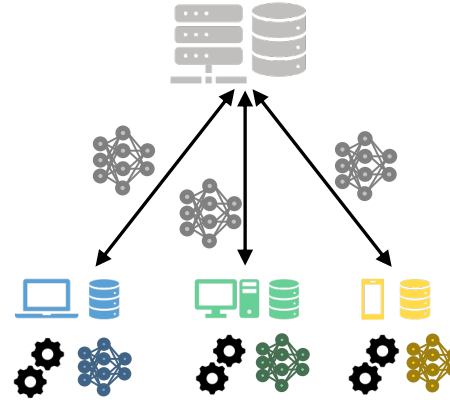


# Dealing with non-iidness

## Data augmentation



## Personalization: It's a feature not a bug!



Ad hoc learning methods



Hyper-parameter tuning

---

4

# Privacy preservation

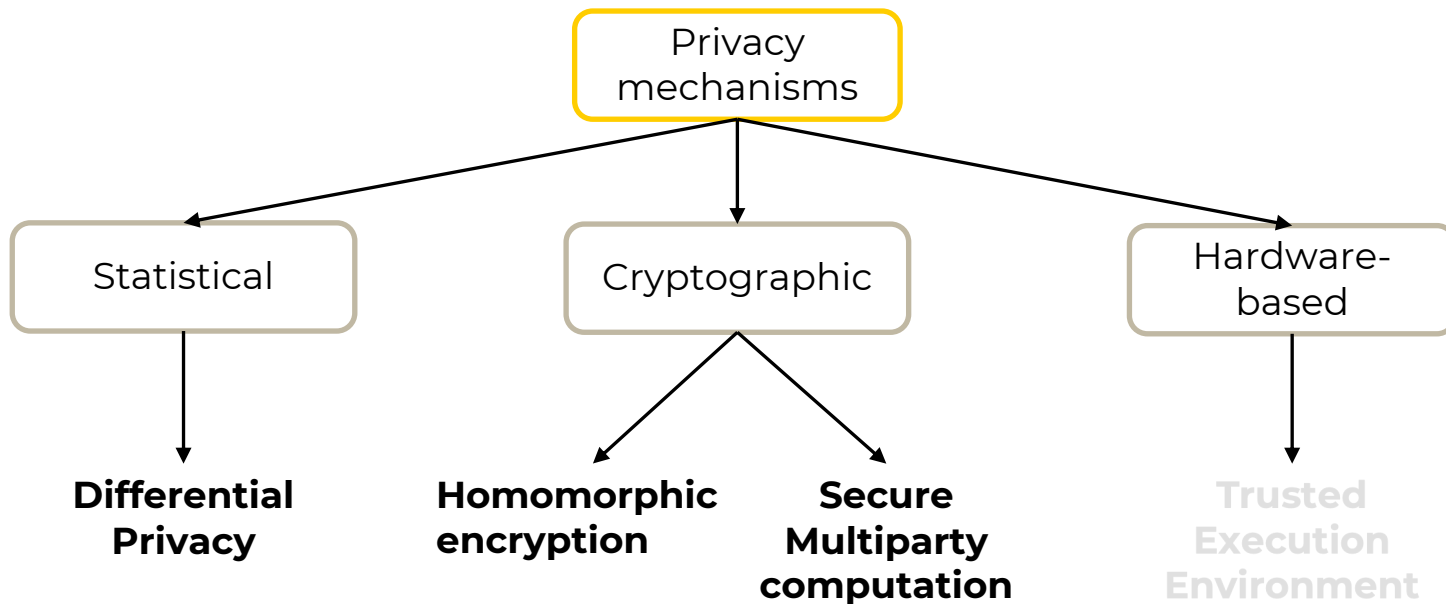
How to improve privacy in FL

---



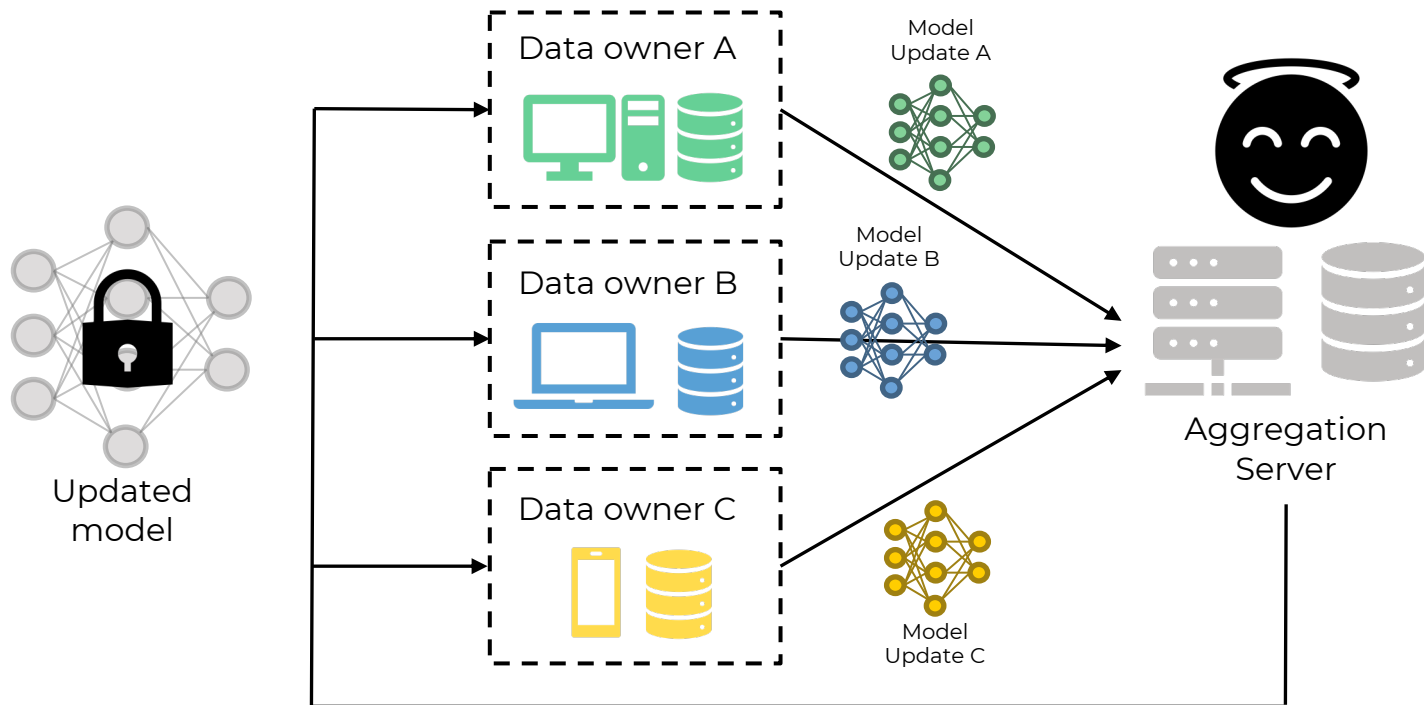
## FL may be not enough

Gradient/model updates may leak information about the user data!





# Global Differential Privacy







# FedAvg + Global DP



Aggregation Server

Algorithm: **FedAvg**

1. Initialize model  $\bar{w}_0$
2. for each round  $t=1, \dots$ :
3. Broadcast  $\bar{w}_{t-1} + \text{noise}$
4. select  $C$  eligible participants
5. foreach|| participant  $p$ :
6.  $w_t^p \leftarrow \text{LocalUpdate}(p)$
7.  $\bar{w}_t \leftarrow \text{aggregate}(\forall p w_t^p)$



Collaborator/Client

Algorithm: **LocalUpdate**

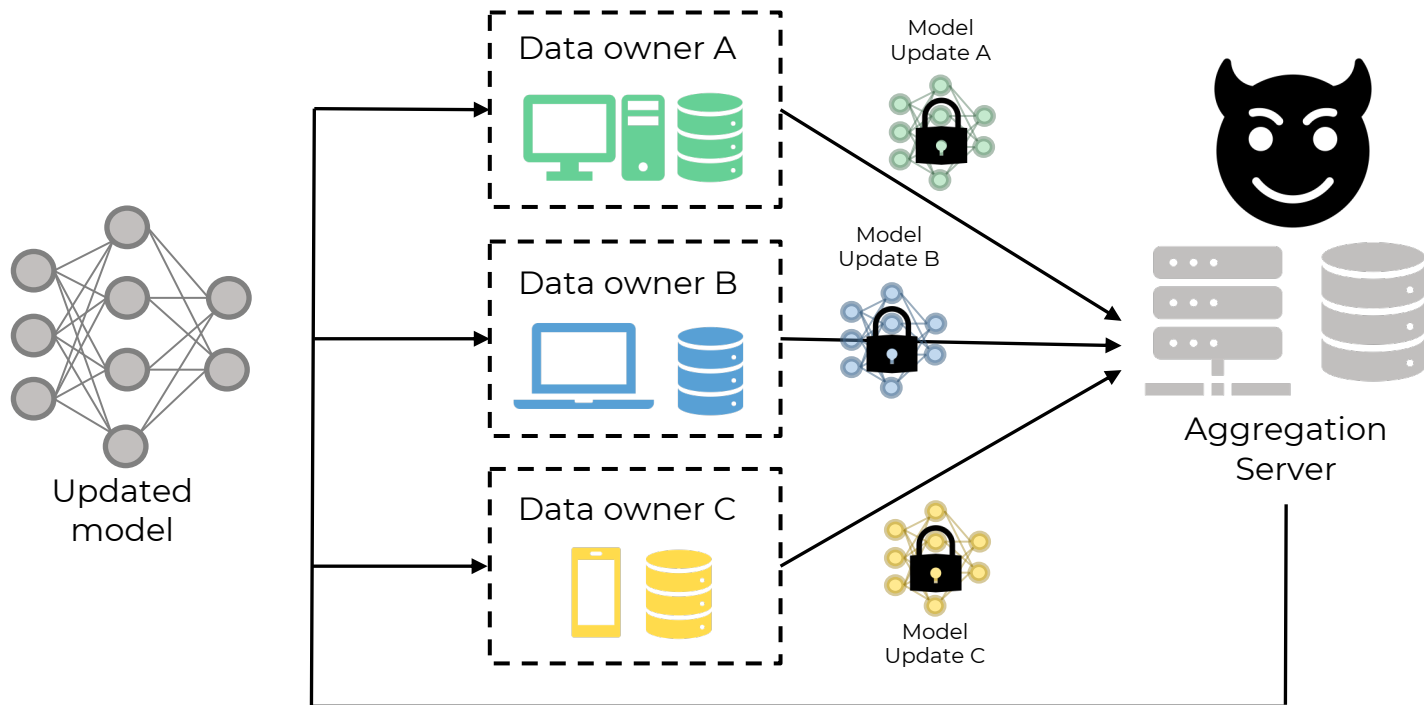
1.  $w \leftarrow$  global model from Server
2. for each epoch  $s \in 1, \dots, S$ :
3. for each batch  $b$ :
4.  $g_b \leftarrow$  compute gradient for  $b$
5.  $w \leftarrow w - \eta g_b$
6. send  $w$  to the Server



Not safe with honest-but-curious servers



# Local Differential Privacy





# FedAvg + Local DP



Aggregation Server

Algorithm: **FedAvg**

1. Initialize model  $\bar{w}_0$
2. for each round  $t=1, \dots$ :
3. Broadcast  $\bar{w}_{t-1}$
4. select  $C$  eligible participants
5. foreach || participant  $p$ :
6.  $w_t^p \leftarrow \text{LocalUpdate}(p)$
7.  $\bar{w}_t \leftarrow \text{aggregate}(\forall p w_t^p)$



Collaborator/Client

Algorithm: **LocalUpdate**

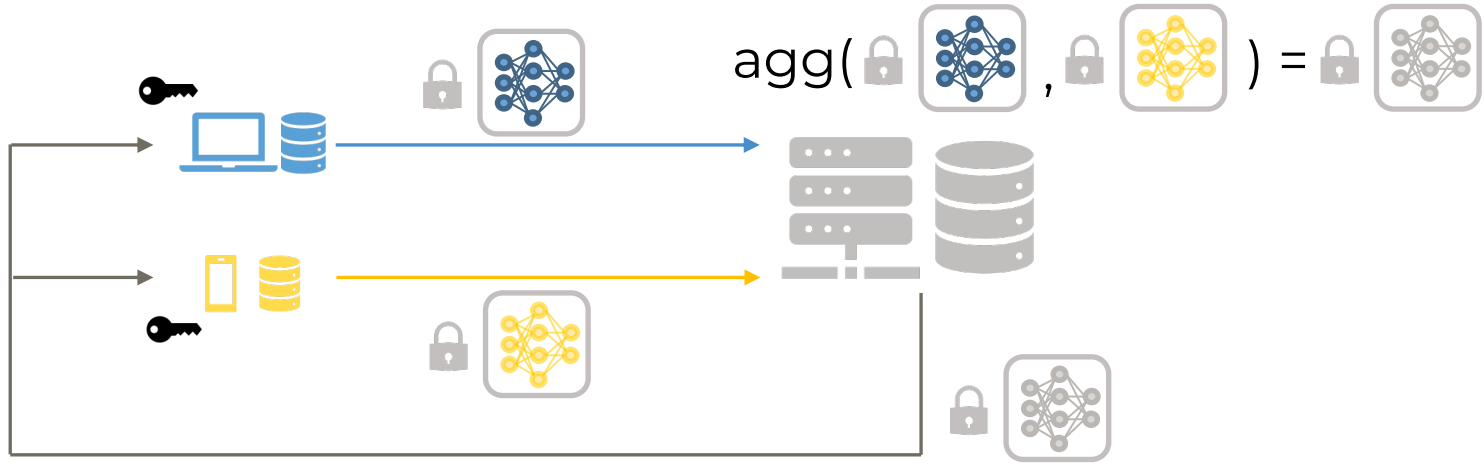
1.  $w \leftarrow$  global model from Server
2. for each epoch  $s \in 1, \dots, S$ :
3. for each batch  $b$ :
4.  $g_b \leftarrow$  compute gradient for  $b$
5.  $w \leftarrow w - \eta g_b$
6. **send ( $w + \text{noise}$ ) to the Server**



It may affect the performance



# Homomorphic Encryption





# FedAvg + HE



Aggregation Server

Algorithm: **FedAvg**

1. Initialize model  $\bar{w}_0$
2. for each round  $t=1, \dots$ :
3. Broadcast  $\bar{w}_{t-1}$
4. select  $C$  eligible participants
5. foreach || participant  $p$ :
6.  $enc(w_t^p) \leftarrow \text{LocalUpdate}(p)$
7.  $\bar{w}_t \leftarrow \text{aggregate\_he}(\forall p \text{ enc}(w_t^p))$



Collaborator/Client

Algorithm: **LocalUpdate**

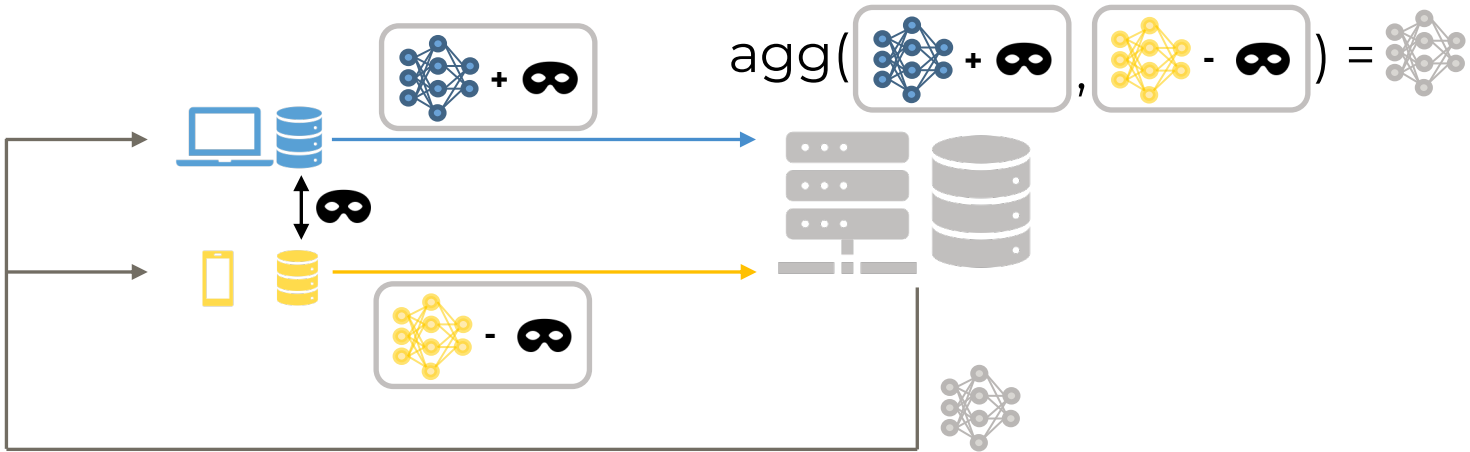
1.  $w_{enc} \leftarrow$  global model from Server
2.  $w \leftarrow \text{decrypt}(w_{enc})$
3. for each epoch  $s \in 1, \dots, S$ :
4. for each batch  $b$ :
5.  $g_b \leftarrow$  compute gradient for  $b$
6.  $w \leftarrow w - \eta g_b$
7. send  $\text{encrypt}(w)$  to the Server



Computationally expensive



# Secure Multiparty Computation





# FedAvg + SMC



Aggregation Server

Algorithm: **FedAvg**

1. initialize model  $\bar{w}_0$
2. for each round  $t=1, \dots$ :
3. Broadcast  $\bar{w}_{t-1}$
4. select  $C$  eligible participants
5. foreach || participant  $p$ :
6.  $\text{mask}(w_t^p) \leftarrow \text{LocalUpdate}(p)$
7.  $\bar{w}_t \leftarrow \text{agg\_smc}(\forall p \text{ mask}(w_t^p))$



**Weigh down the communication protocol**



Collaborator/Client

Algorithm: **OneTimePadAgreement**

1. For each active client  $p$ :
2. agree on perturbation  $s_p$

Algorithm: **LocalUpdate**

1.  $w \leftarrow$  global model from Server
2. for each epoch  $s \in 1, \dots, S$ :
3. for each batch  $b$ :
4.  $g_b \leftarrow$  compute gradient for  $b$
5.  $w \leftarrow w - \eta g_b$
6. send  $\text{mask}(w, \forall p s_p)$  to the Server

---

**5**

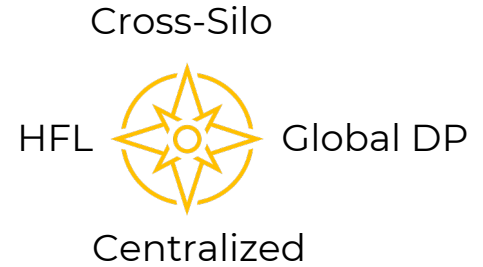
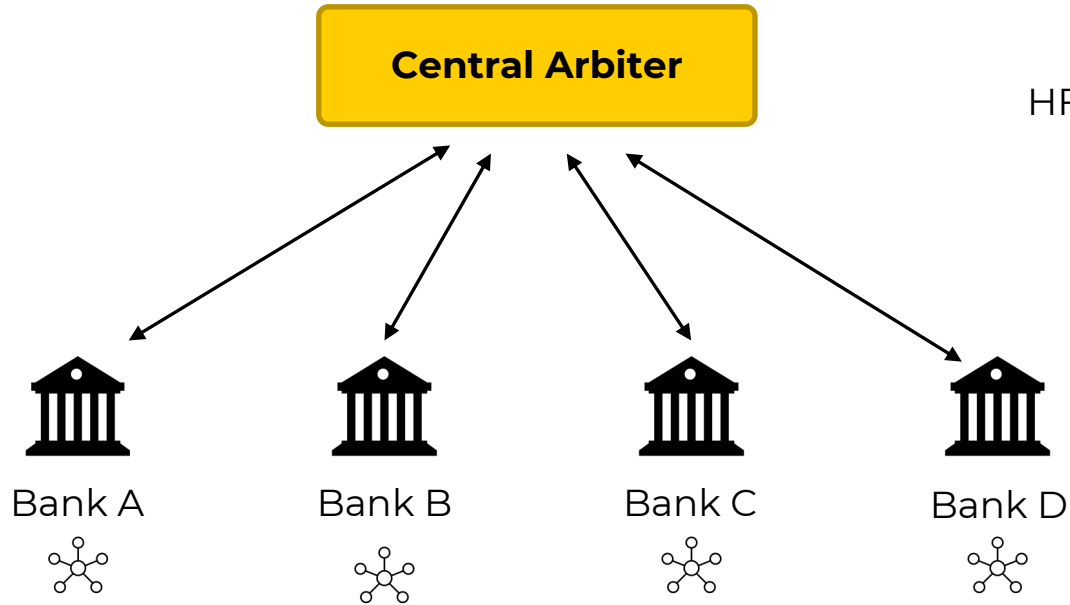
# Applications

Some use cases and real-world examples of FL



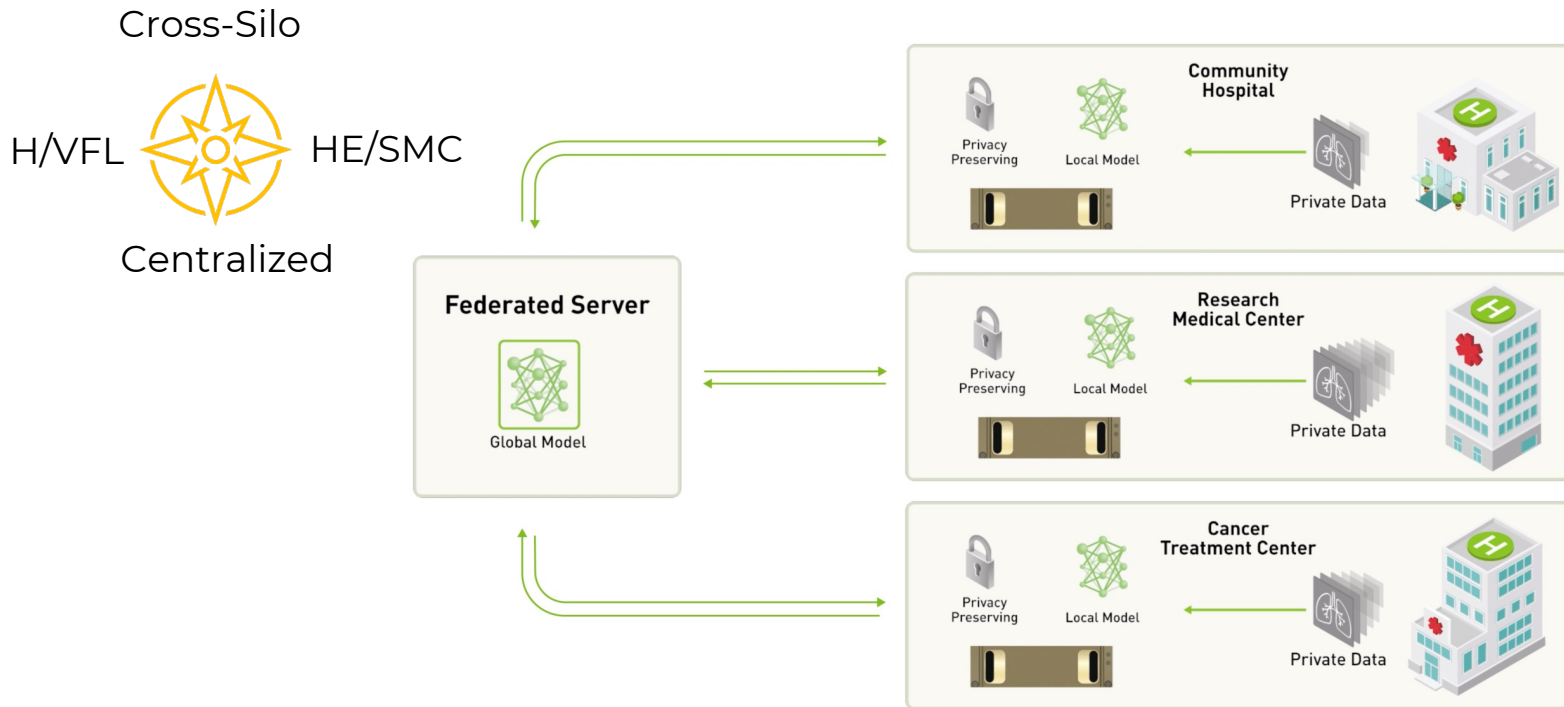


# Use case 1 Anti-money laundering





# Use case 2 Medical diagnosis



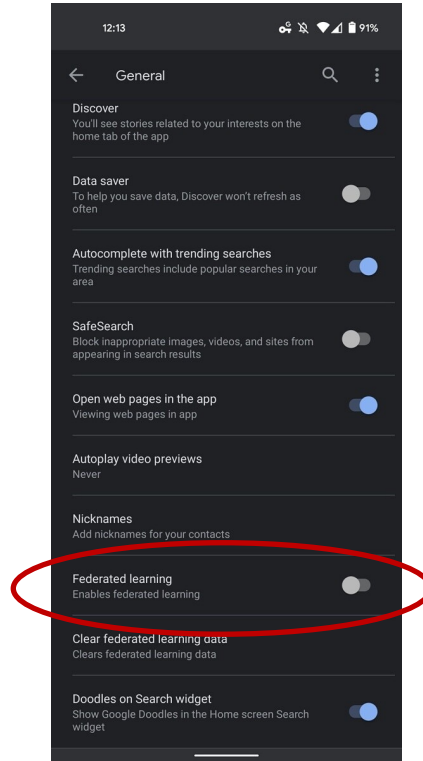
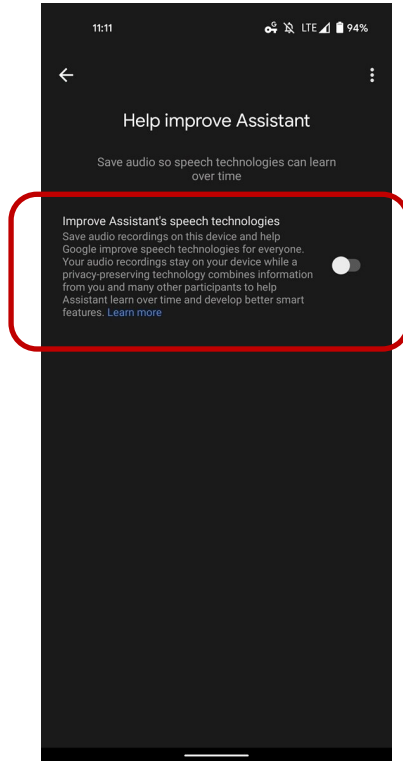


# Google's GBoard



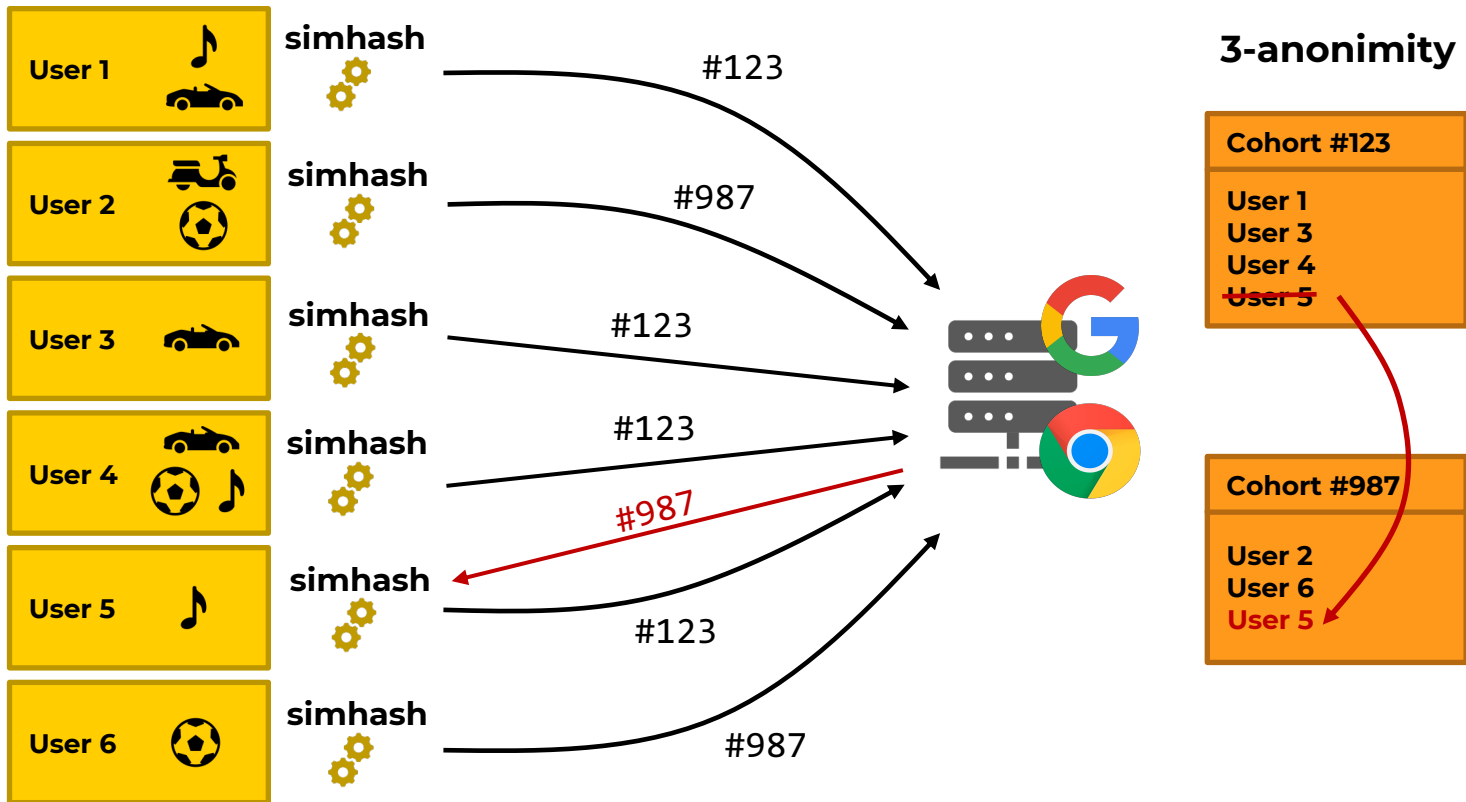


# Google's "Hey Google!" recognition





# Google's Federated Learning Of Cohorts



---

6

# Conclusions

The glorious “take home message”



## What we did not cover

---

- ⦿ Attacks to FL systems
- ⦿ Federated Transfer Learning
- ⦿ Improve communication efficiency, e.g., model quantization
- ⦿ Fairness



## Take home message 😊

- 🕒 FL is a “novel” yet **interesting framework for privacy-preserving ML**
- 🕒 FL methods must be designed considering the **communication-computation-privacy-effectiveness trade-off**
- 🕒 FL is still in its infancy and there are **many open problems**





---

# Thanks!

*Any Questions?*

*The only stupid question is the one you were  
afraid to ask but never did.*

*Richard Sutton*

“

---



## Resources

- Li, et al. 'A Survey on Federated Learning Systems: Vision, Hype and Reality for Data Privacy and Protection', 2021. <http://arxiv.org/abs/1907.09693>.
- Kairouz, et al. 'Advances and Open Problems in Federated Learning'. Foundations and Trends in Machine Learning 14, 2021. <https://doi.org/10.1561/22000000083>.
- Yang et al. 'Federated Learning'. Synthesis Lectures on Artificial Intelligence and Machine Learning 13, 2019. <https://doi.org/10.2200/S00960ED2V01Y201910AIM043>.
- Li, et al. 'Federated Learning: Challenges, Methods, and Future Directions', 2019. <https://arxiv.org/abs/1908.07873>
- Bonawitz, et al. 'Practical Secure Aggregation for Federated Learning on User-Held Data'. 2016. <http://arxiv.org/abs/1611.04482>.
- McMahan, et al. 'Communication-efficient learning of deep networks from decentralized data', 2016. <https://arxiv.org/abs/1602.05629>