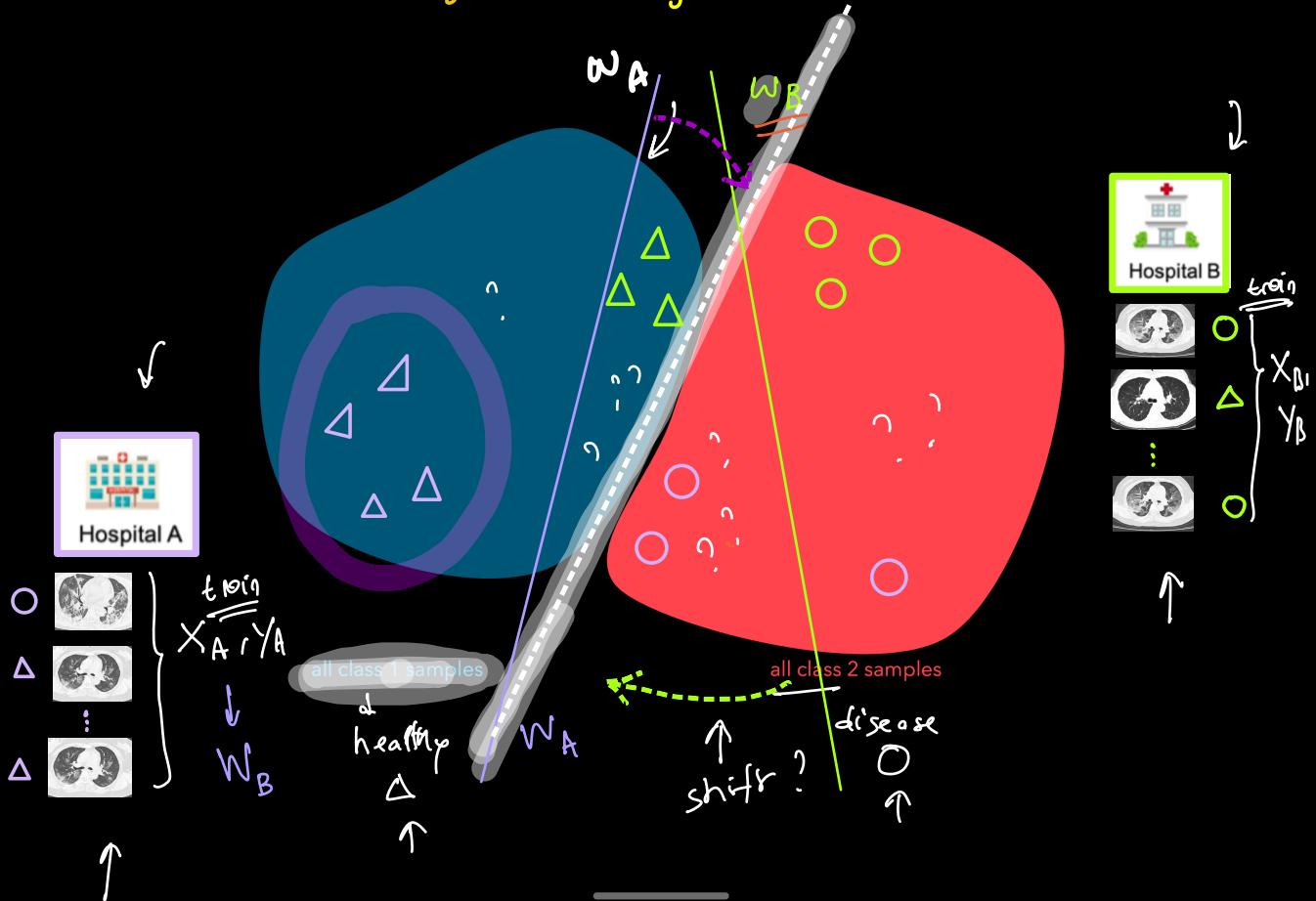




2) What is federated learning?

why  
federate?

How to  
federate?

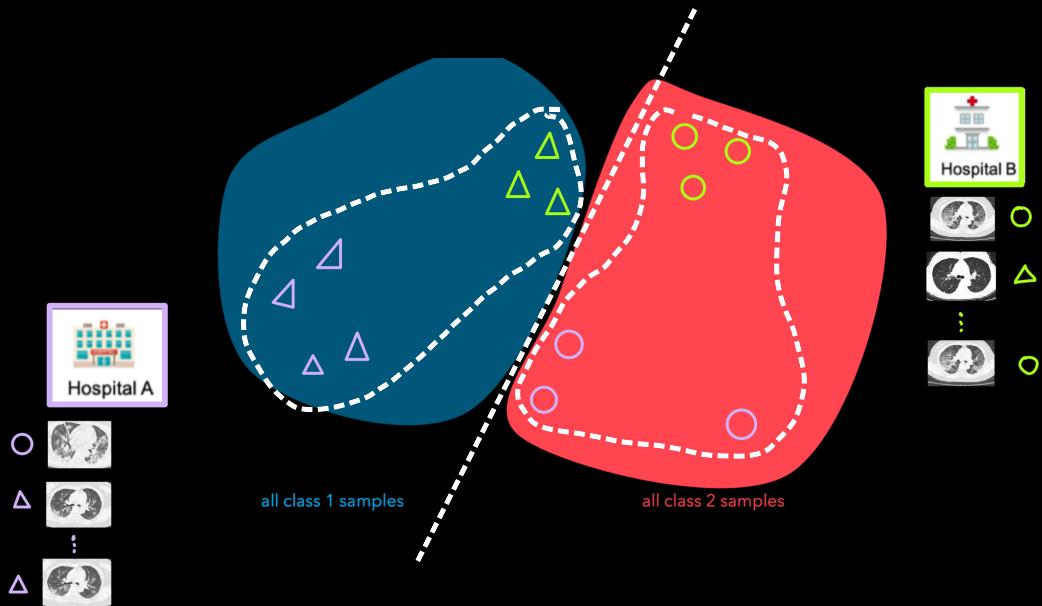


Modern ML algorithms (classification, regression)  
are data & computation hungry.

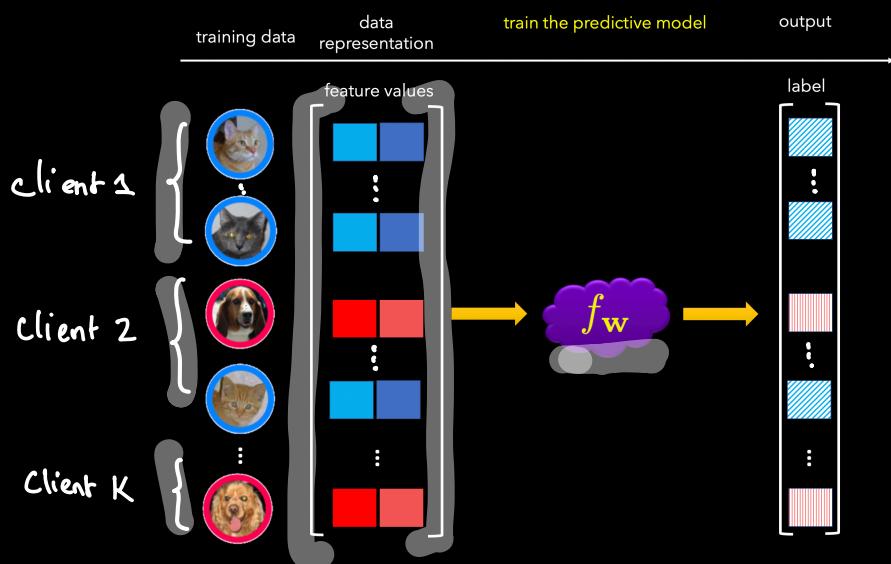
The performance of an ML model is dictated by the data it is trained on (size, distribution, etc)

→ How can we design powerful & generalizable ML models?

Solution 1 Build a single model using centralized data.



- \* The data is uploaded into a single device where a model is trained.

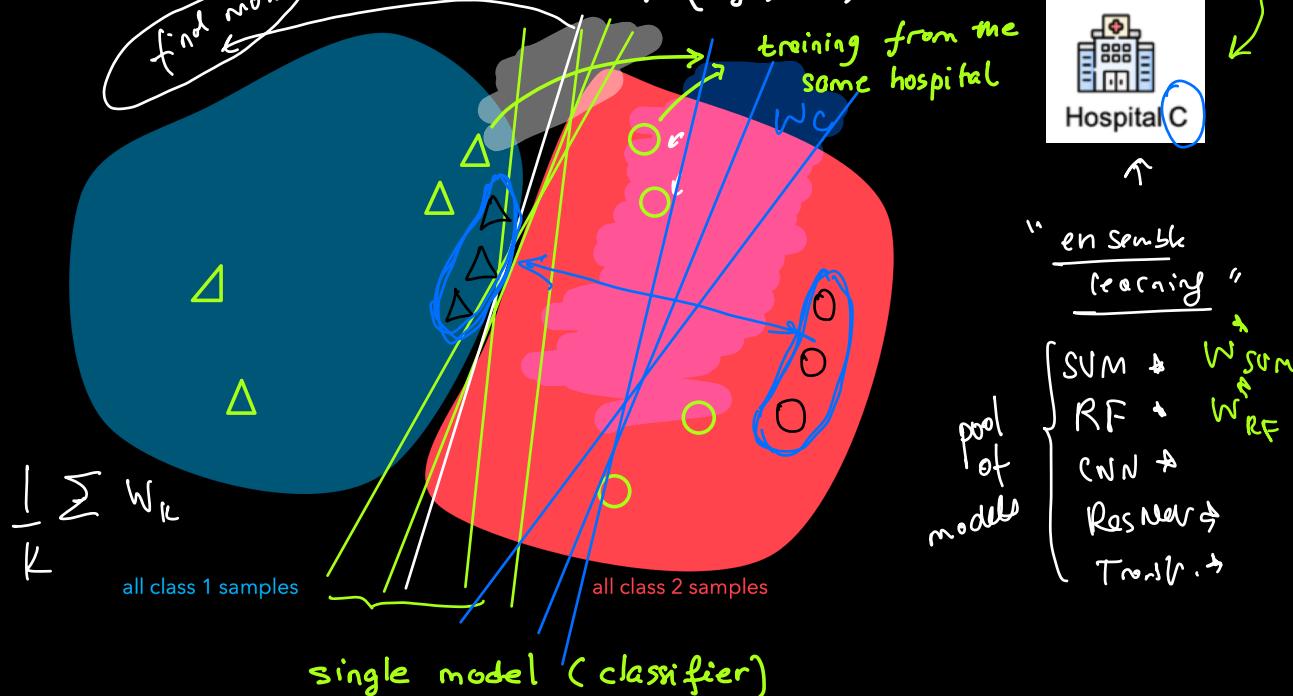


- Concerns about data security and privacy
- Such a data "sharing" process is not straightforward.
- large-scale datasets are computationally expensive to upload.  
(e.g., on mobile phones & apps)
- Intrusive and difficult if data is spread across multiple devices or clients.

→ How can we design powerful & generalizable ML models ?

## Solution 2

Build an "ensemble" model on a single local dataset (e.g., H.C)



**Input data** →  $f_w$  → **Output data**

ensemble

1) average models

2) average predictions

boost the model's performance ↗

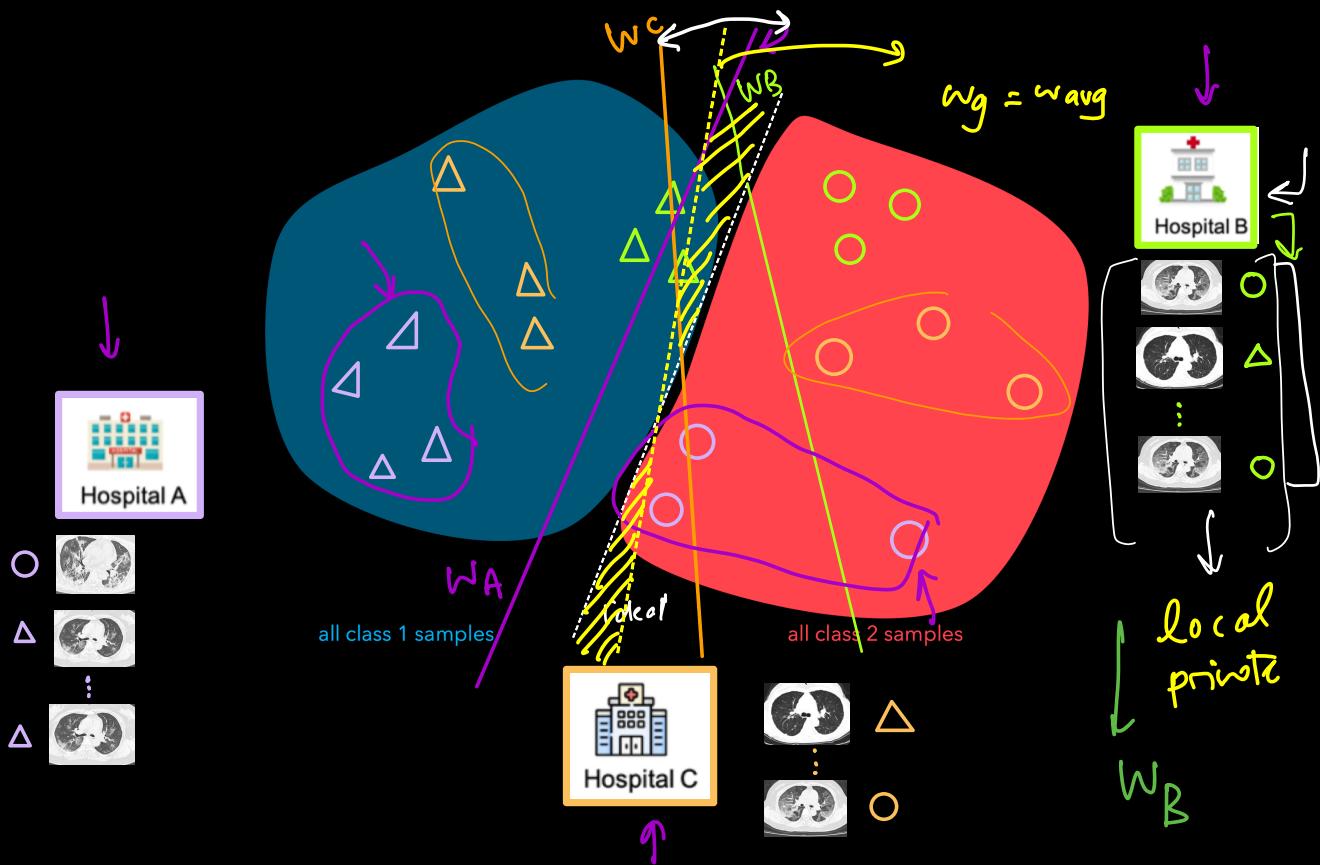
- high computations (many models to train on a single device)
- training only benefits from a single dataset (here 

# Distribution

where your points  
are located ?

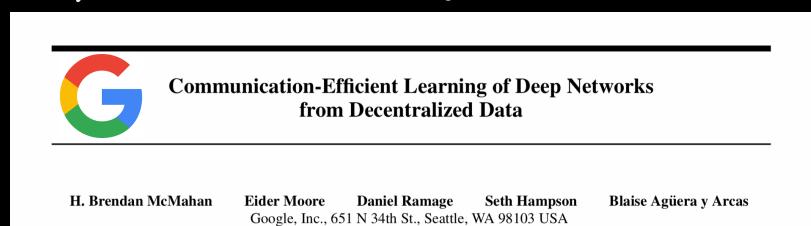
→ How can we design powerful & generalizable ML models?

### Solution 3 Data-preserving federated learning



**Goal:** boost the local performance of a model of particular client/hospital without any data sharing.

Simple but creative Google paper by McMahan (2016)



<https://arxiv.org/> cs ::  
Communication-Efficient Learning of Deep Networks - arXiv  
by HB McMahan · 2016 · Cited by 4039 — [Submitted on 17 Feb 2016 (v1), last revised 28 Feb 2017 (this version, v3)] ... Submission history. From: Hugh Brendan McMahan...

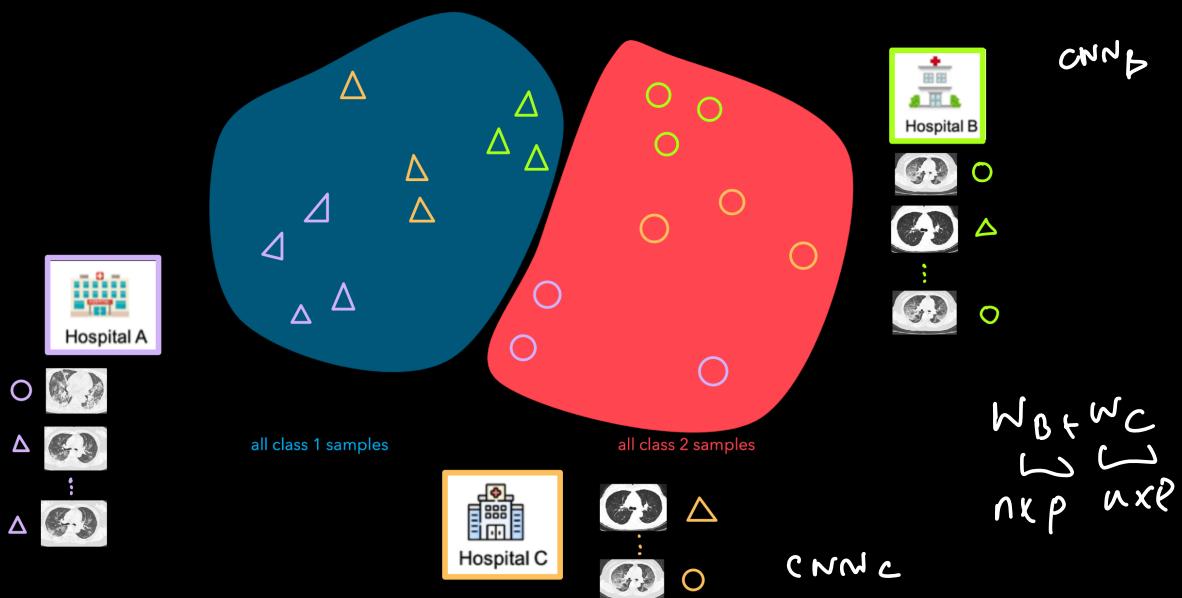
**Goal**: boost the local performance of a model of particular client/hospital without any data sharing.

{ Sol° 1 → centralized (one model on many datasets)

Sol° 2 → ensemble (many models on a single dataset)

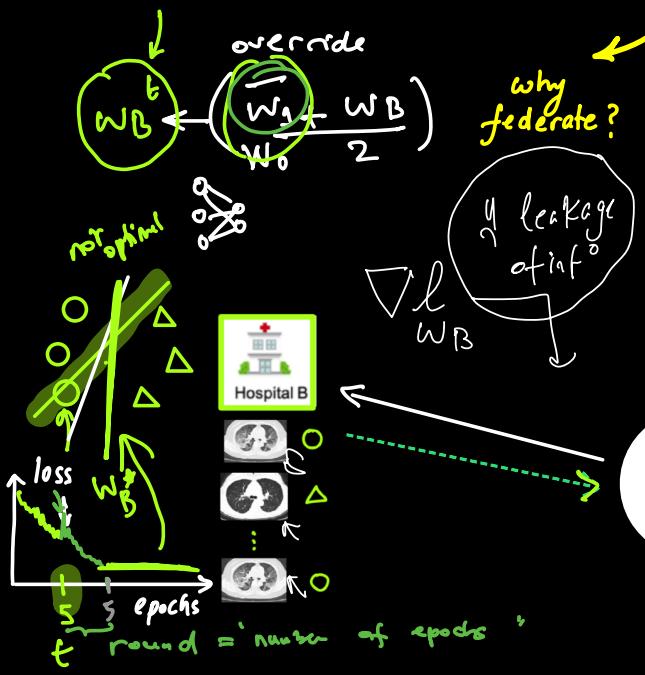
Sol° 3 → federated (many models on many datasets)

**THINK!**



What's the trick here ?

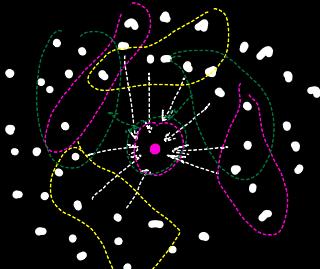
2) What is federated learning?



In each round,

1. Each client  $i$  trains its local model  $w_i^t$  on its local dataset  $D_i = \{(x_{ij}, y_j)\}_{j=1}^{n_i^t}$
2. Following a few training iterations, the client forward the resulting model  $w_i^t$  to the server  $S$ .
3. The server aggregates the local models by averaging them to create a global model:  
 $w_g^t = \frac{1}{n} \sum_{i=1}^n w_i^t$  total across all clients
4. The server broadcasts  $w_g^t$  to all clients
5. Each client  $i$  initializes its local model for the next round  
 $w_i^{t+1} = w_g^t$

- \* Hyperparameters :
  - K: total number of clients
  - T: total communication rounds  $5 \times 10^6$
  - N: local iteration steps per round
  - M: clients per round (randomly selected)



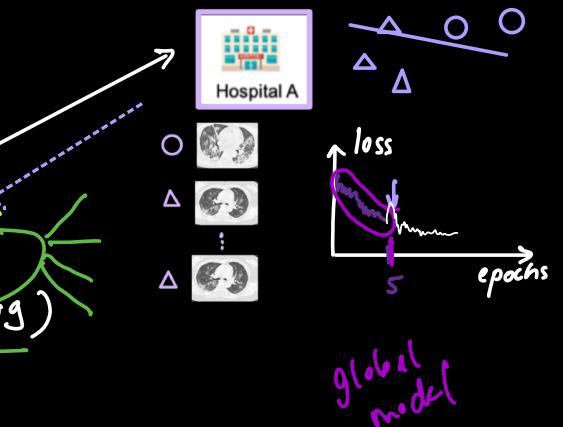
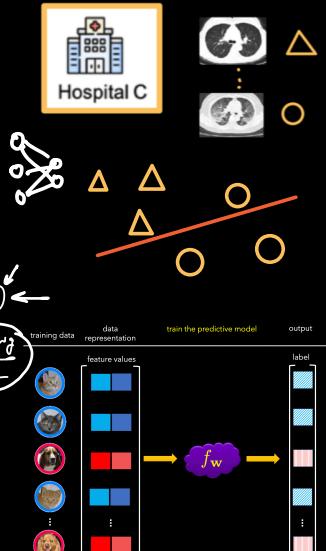
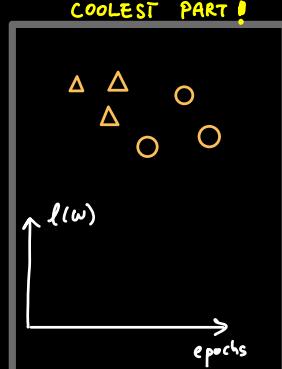
$$\min_{\mathbf{w}} \quad l(\mathbf{w}) = \sum_{i=1}^n (\underbrace{\mathbf{f}_{\mathbf{w}}(x^i) - y^i}_{{\hat{y}}^i})^2$$


solving the problem  
using cool math

## COOLEST PART

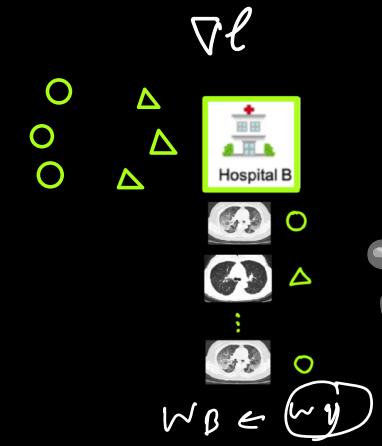
The diagram shows the decomposition of the augmented matrix  $\tilde{\mathbf{X}}^*$  into three components:

- Training samples**: A red rectangular block containing several colored rectangles (blue, red, blue, blue, red, blue, ...).
- $\tilde{\mathbf{X}}$** : A large black rectangular block below the training samples.
- Learned feature matrix**: A blue rectangular block on the right side containing several colored rectangles (blue, pink, blue, blue, pink, blue, ...).



## 2) What is federated learning?

why federate?      How to federate?



$$n = n_1 + n_2 + n_3$$

$$w_g = \frac{n_1}{n} w_1 + \frac{n_2}{n} w_2 + \frac{n_3}{n} w_3$$

global

Server

$w_1$

$w_2$

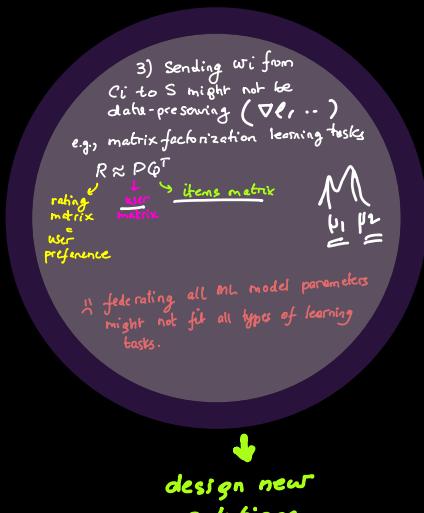
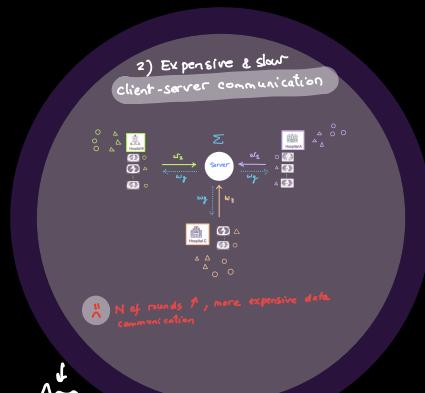
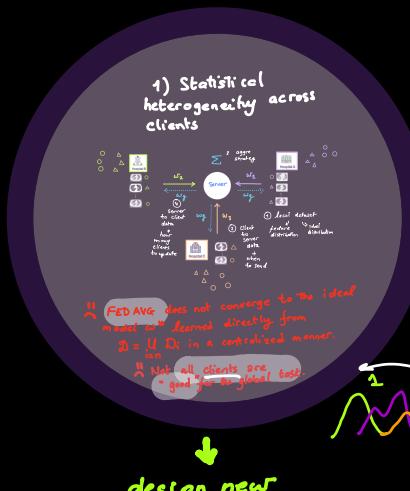
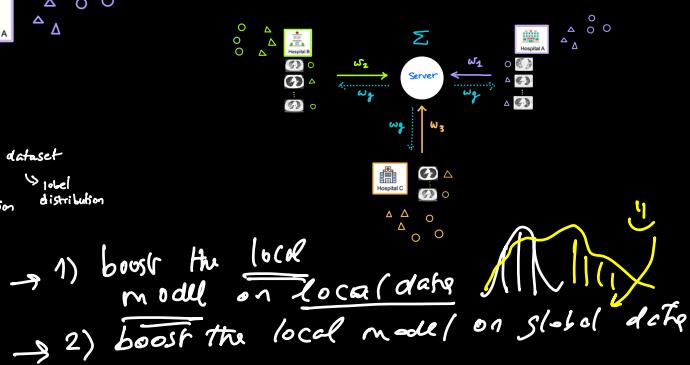
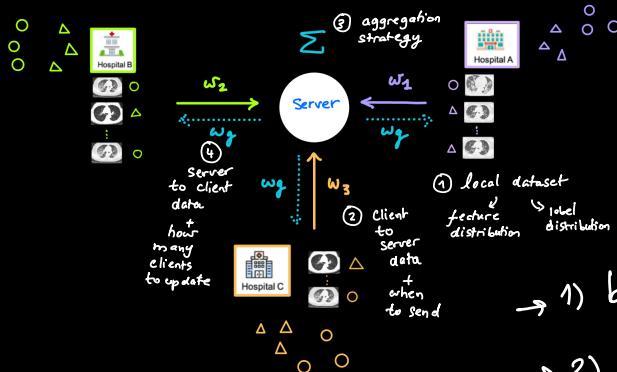
$w_3$

$w_g$

### 3) Exciting challenges & future AI problems TO SOLVE!

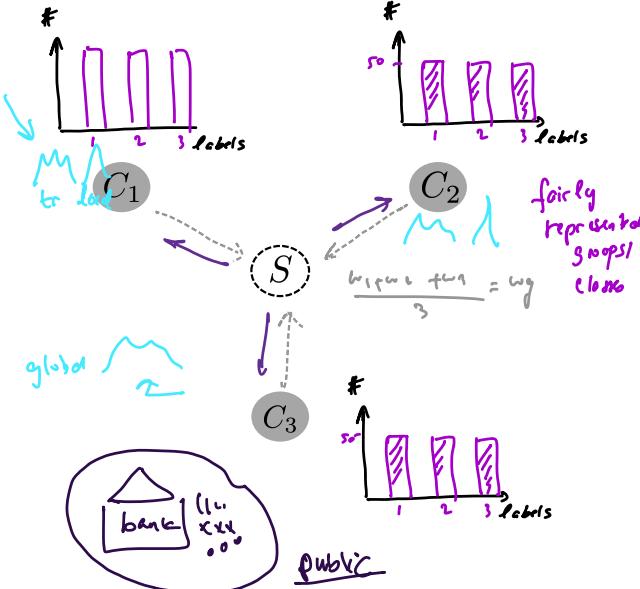
- ① formalize your FL problem
- ② describe a potential solution
- ③ Design your experiments & evaluate

Each FL component = challenge



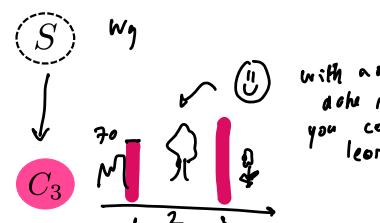
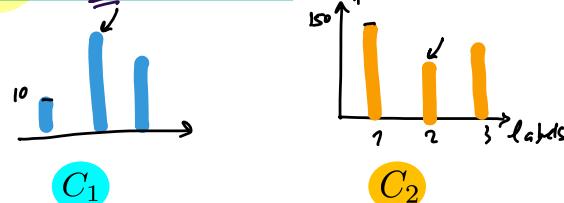
I - FedAvg (2016) [1]

a) IID label distribution across clients in classification

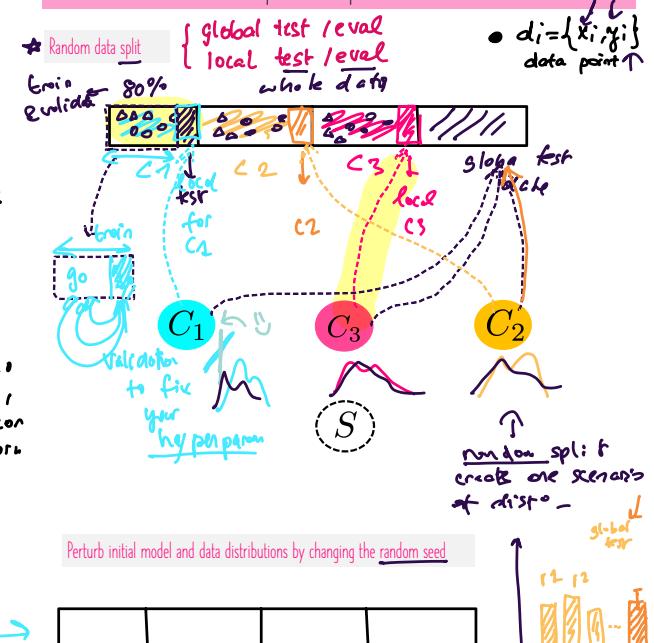


b) Non-IID scenario

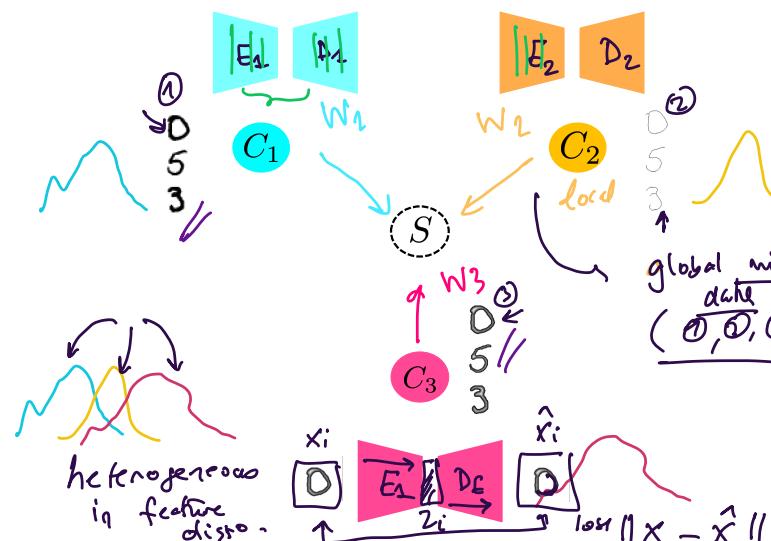
[Scenario 1] Non-IID label distribution in classification



FL experiment setup for evaluation



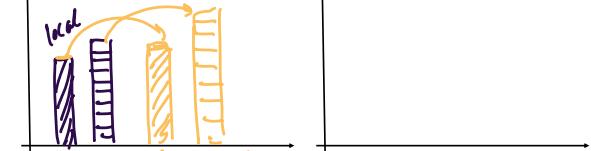
[Scenario 2] Non-IID feature distribution in generative learning



Evaluation measures and benchmarks

- 1) Eval. measures
  - $\rightarrow$  Acc, Sens, Spec
- 2) benchmarks
  - $\rightarrow$  global test set, local test set
  - $\rightarrow$  stand-alone (global trained on local data no fed)
  - $\rightarrow$  Fed. local model

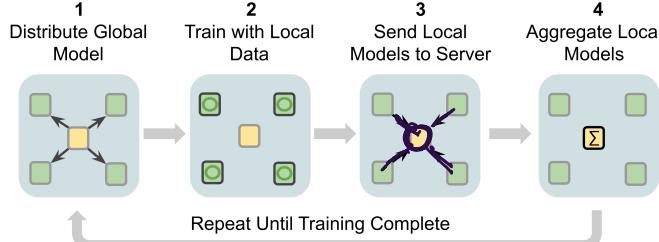
local kit      global



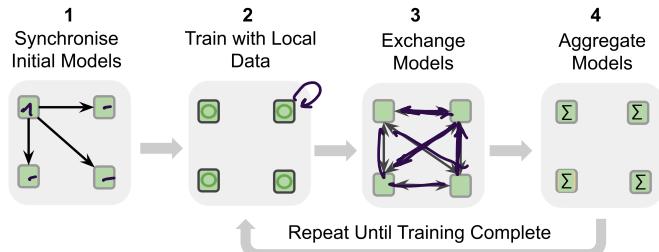
Diverse learning paradigms: centralized versus federated [2]

Personal reflections and notes

Fed Avg



(a) FL - Aggregation Server



(b) FL - Peer to Peer

#### Key

- |  |                          |   |                       |
|--|--------------------------|---|-----------------------|
| <span style="background-color: yellow;">■</span> | Aggregation Server       | <span style="background-color: #c8512e;">■</span> | Central Data Lake     |
| <span style="background-color: green;">■</span>  | Training Node            | <span style="background-color: #6f81bd;">■</span> | Data Donor            |
| $\Sigma$   | Model Aggregation        | $\rightarrow$                                     | Time                  |
| $\rightsquigarrow$                               | Weight/Gradient Exchange | $\textcolor{red}{\leftarrow}$                     | Medical Data Exchange |
|  |                          | <span style="color: green;">○</span>              | Local Training        |

Image source: <https://www.nature.com/articles/s41746-020-00323-1>



### Algorithm 1 Federated Averaging (FedAvg)

```

Input:  $K, T, \eta, E, w^0, N, p_k, k = 1, \dots, N$ 
for  $t = 0, \dots, T - 1$  do
    Server selects a subset  $S_t$  of  $K$  devices at random (each device  $k$  is chosen with probability  $p_k$ )
    Server sends  $w^t$  to all chosen devices
    Each device  $k \in S_t$  updates  $w^t$  for  $E$  epochs of SGD on  $F_k$  with step-size  $\eta$  to obtain  $w_k^{t+1}$ 
    Each device  $k \in S_t$  sends  $w_k^{t+1}$  back to the server
    Server aggregates the  $w$ 's as  $w^{t+1} = \frac{1}{K} \sum_{k \in S_t} w_k^{t+1}$ 
end for

```

② Cons:

Statistical heterogeneity: is common with data being non-identically distributed between devices.



Existing solutions: There are also heuristic approaches that aim to tackle statistical heterogeneity by sharing the local device data or server-side proxy data.

- However, these methods may be unrealistic: in addition to imposing burdens on network bandwidth, sending local data to the server violates the key privacy assumption of federated learning
- Sending globally-shared proxy data to all devices requires effort to carefully generate or collect such auxiliary data.

System heterogeneity: is also a critical concern in federated networks. The storage, computational, and communication capabilities of each device in federated networks may differ due to variability in hardware (CPU, memory), network connectivity (3G, 4G, 5G, WiFi), and power (battery level).

Existing solutions: ignore the more constrained devices (stragglers) failing to complete a certain amount of training (unfair)

- However, this can have negative effects on convergence as it limits the number of effective devices contributing to training,
- May induce bias in the device sampling procedure if the dropped devices have specific data characteristics.

### Algorithm 2 FedProx (Proposed Framework)

```

Input:  $K, T, \mu, \gamma, w^0, N, p_k, k = 1, \dots, N$ 
for  $t = 0, \dots, T - 1$  do
    Server selects a subset  $S_t$  of  $K$  devices at random (each device  $k$  is chosen with probability  $p_k$ )
    Server sends  $w^t$  to all chosen devices
    Each chosen device  $k \in S_t$  finds a  $w_k^{t+1}$  which is a  $\gamma_k$ -inexact minimizer of:  $w_k^{t+1} \approx \arg \min_w h_k(w; w^t) = F_k(w) + \frac{\mu}{2} \|w - w^t\|^2$ 
    Each device  $k \in S_t$  sends  $w_k^{t+1}$  back to the server
    Server aggregates the  $w$ 's as  $w^{t+1} = \frac{1}{K} \sum_{k \in S_t} w_k^{t+1}$ 
end for

```

Proposed solution (FedProx)

$$l_k = \text{loss of client } k$$

**Challenge 1 to solve:** With dissimilar (heterogeneous) local objectives  $F_k$ , a larger number of local epochs may lead each device towards the optima of its local objective as opposed to the global objective—potentially hurting convergence or even causing the method to diverge.

**Challenge 2 to solve:** In federated networks with heterogeneous systems resources, setting the number of local epochs to be high may increase the risk that devices do not complete training within a given communication round and must therefore drop out of the procedure.

**Proposed solution:** merge local models by accounting for their heterogeneity and tolerating partial work → allowing for variable amounts of work to be performed locally across devices based on their available systems resources, and then aggregate the partial solutions sent from the stragglers (as compared to dropping these devices) (fairer)

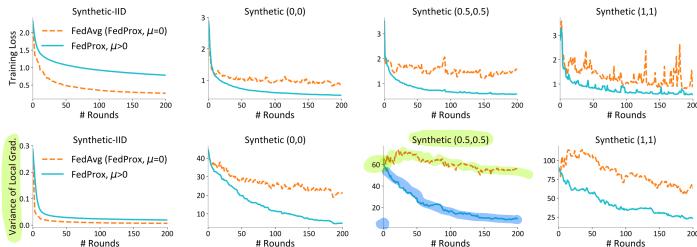
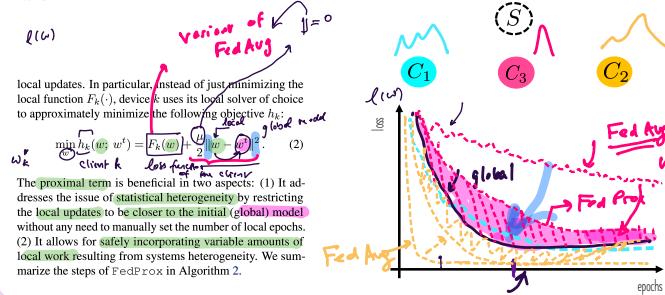
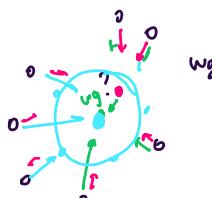


Figure 2. Effect of data heterogeneity on convergence. We remove the effects of systems heterogeneity by forcing each device to run the same amount of epochs. In this setting, FedProx with  $\mu = 0$  reduces to FedAvg. (1) Top row: We show training loss (see results on testing accuracy in Appendix C.3, Figure 6) on four synthetic datasets whose statistical heterogeneity increases from left to right. Note that the method with  $\mu = 0$  corresponds to FedAvg. Increasing heterogeneity leads to worse convergence, but setting  $\mu > 0$  can help to combat this. (2) Bottom row: We show the corresponding dissimilarity measurement (variance of gradients) of the four synthetic datasets. This metric captures statistical heterogeneity and is consistent with training loss — smaller dissimilarity indicates better convergence.

$$\|w - w^t\|_{\text{global}}$$



## FedALA: Adaptive Local Aggregation for Personalized Federated Learning

Jianqing Zhang<sup>1</sup> Yang Hua<sup>2</sup> Hao Wang<sup>3</sup>

Tao Song<sup>1</sup> Zhengui Xue<sup>1</sup> Ruhui Ma<sup>1</sup> Haibing Guan<sup>1</sup>



Source link: <https://github.com/TsingZ0/FedALA/blob/main/FedALA0ral.pdf>

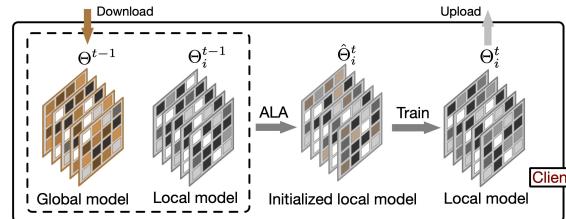


Figure 1: Local learning process on client  $i$  in the  $t$ -th iteration. Specifically, client  $i$  downloads the global model from the server, locally aggregates it with the old local model by ALA module for local initialization, trains the local model, and finally uploads the trained local model to the server.

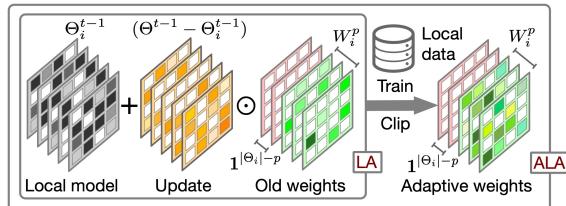
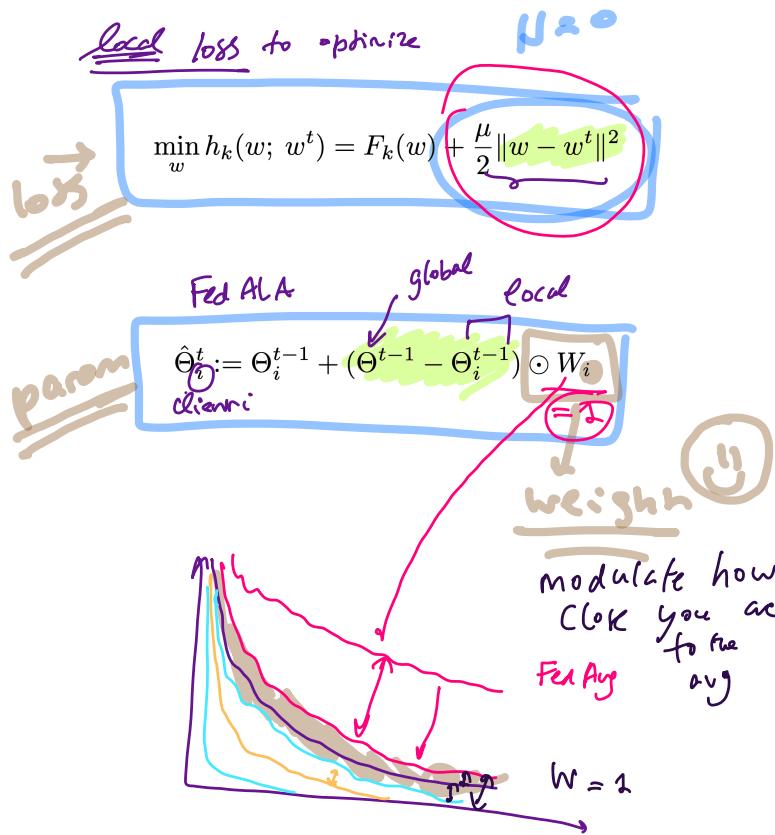


Figure 2: The learning process in ALA. LA denotes “local aggregation”. Here, we consider a five-layer model and set  $p = 3$ . The lighter the color, the larger the value.

FedProx vs FedALA

2020 2023

Key idea: Regularize/control the local model using the global model (how far am I from the average)?  
Stay personalized (retain local information) without diverging too much from the global/universal model.



Personal reflections and notes

References:

1. FedAvg: McMahan et al., 2017. Communication-efficient learning of deep networks from decentralized data. In Artificial intelligence and statistics (pp. 1273–1282). PMLR. <https://proceedings.mlr.press/v54/mcmahan17a/mcmahan17a.pdf>
2. Other federation topologies: Rieke et al., 2020. The future of digital health with federated learning. NPJ digital medicine, 3(1), pp1–7. <https://www.nature.com/articles/s41746-020-00323-1>
3. FedProx: Li et al., 2020. Federated optimization in heterogeneous networks. Proceedings of Machine learning and systems, 2, pp 429–450. <https://proceedings.mlsys.org/paper/2020/file/1f5fe83998a09396ebe477d9475ba0c-Paper.pdf>; <https://github.com/litfan96/FedProx>
4. FedALA: Zhang et al., 2023. FedALA: Adaptive local aggregation for personalized federated learning. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 37 No. 9, pp. 11237–11244). <https://dl.acm.org/doi/10.11609/aaa.v37i9.76330>; <https://github.com/TsingU/FedALA>