

## I- FedAvg (2016)

### Communication-Efficient Learning of Deep Networks from Decentralized Data

H. Brendan McMahan Eider Moore Daniel Ramage Seth Hampson Blaise Agüera y Arcas

Google, Inc., 651 N 34th St., Seattle, WA 98101 USA

#### Abstract

Modern mobile devices have access to a wealth of decentralized learning models, which in turn can greatly improve the user experience on the device. For example, language models can improve speech recognition and text entry, and image models can automatically select good photos. However, rich decentralized learning models are large in quantity or both, which may preclude logging to the central device and training there using conventional approaches. We advocate an alternative that leaves the model distributed on the mobile devices and learns a shared model by aggregating locally-computed updates. We term this decentralized approach *Federated Learning*. We present a practical method for the federated learning of deep networks based on iterative model averaging, and an extensive model architectural evaluation, considering five different model architectures and four datasets. These experiments demonstrate the approach is robust to the unbalanced and non-IID data distributions that are a defining characteristic of this setting. Communication costs are the principal constraint, and we show a reduction in required communication rounds by 10–100x as compared to synchronized stochastic gradient descent.

#### 1 Introduction

Increasingly, phones and tablets are the primary computing devices for many [3], and billions of users on these devices (including cameras, microphones, and GPS), combined with the fact they are frequently carried, means they have access to an unprecedented amount of data, much of it private in nature. Models learned on such data hold the

Proceedings of the 20<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2017, Fort Lauderdale, Florida, USA. JMLR: W&CP volume 54. Copyright 2017 by the author(s).

promise of greatly improving usability by powering more intelligent applications, but the sensitive nature of the data means there are risks and responsibilities to storing it in a centralized location.

We investigate a learning technique that allows users to collectively reap the benefits of shared models trained from this rich data, without the need to centrally store it. We term our approach *Federated Learning*, since the learning task is solved by a local client, and the global model is aggregated (which were the original terms which are coordinated by a central server). Each client has a local training dataset which is never uploaded to the server. Instead, each client computes an update to the current global model maintained by the server, and only uploads this update to the server. This is a direct application of the principle of *focused optimization* or *data minimization* proposed by the 2012 White House report on privacy of consumer data [39]. Since these updates are specific to improving the current model, there is no reason to store them once they have been applied.

A primary challenge of federated learning is the decoupling of model training from the need for direct access to the raw training data. Clearly, some trust of the server coordinating the training is still required. However, for applications where the training objective can be specified on the basis of data available at each client, federated learning can significantly reduce privacy and security risks by limiting the attack surface to only the device, rather than the device and the cloud.

Our primary contributions are 1) the identification of the problem of training on decentralized data from mobile devices as an important research direction; 2) the selection of a simple and effective algorithm, FedAvg, which can be applied to this setting; and 3) an extensive empirical evaluation of the proposed approach. More concretely, we introduce the *FederatedAveraging* algorithm, which combines local stochastic gradient descent (SGD) on each client with a server that performs model averaging. We perform experiments to show that FedAvg is robust to the fact that it is robust to unbalanced and non-IID data distributions, and reduce the rounds of communication needed to train a deep network on decentralized data by orders of magnitude.

<sup>1</sup>Carnegie Mellon University <sup>2</sup>Bosch Center for Artificial Intelligence <sup>3</sup>Google Research <sup>4</sup>Facebook AI <sup>5</sup>Determined AI. Correspondence to: H. Brendan McMahan (hbmcmah@cs.cmu.edu).

Proceedings of the 37<sup>th</sup> ICML Conference, Austin, TX, USA, 2020. Copyright 2020 by the authors.

Federated learning is a key challenge in the federated setting. While not the focus of this work, standard privacy-preserving approaches such as differential privacy and secure multiparty computation can naturally be combined with methods proposed herein, particularly since our framework proposes only lightweight algorithmic modifications to prior work.

## II- FedProx (2020)

### FEDERATED OPTIMIZATION IN HETEROGENEOUS NETWORKS

Tian Li<sup>1</sup> Anit Kumar Sahu<sup>2</sup> Manzil Zaheer<sup>3</sup> Maziar Sanjabi<sup>4</sup> Ameet Talwalkar<sup>1,5</sup> Virginia Smith<sup>1</sup>

#### ABSTRACT

Federated Learning is a distributed learning paradigm with two key challenges that differentiate it from traditional distributed optimization: (1) significant variability in terms of the systems characteristics on each device in the network (systems heterogeneity), and (2) non-identically distributed data across the network (statistical heterogeneity). In this work, we introduce a framework, FedProx, to tackle heterogeneity in federated networks. FedProx is the first view of a generalized framework for federated learning that builds upon the state-of-the-art method for federated learning. While this re-parameterization makes only minor modifications to the method itself, these modifications have important ramifications both in theory and in practice. Theoretically, we provide convergence guarantees for our framework when learning over data from non-identical distributions (statistical heterogeneity), and while adhering to device-level systems constraints by allowing each participating device to perform a variable amount of work (systems heterogeneity). Practically, we demonstrate that FedProx allows for more robust convergence than FedAvg across a suite of realistic federated datasets. In particular, in highly heterogeneous settings, FedProx demonstrates significantly more stable and accurate convergence behavior relative to FedAvg—improving absolute test accuracy by 22%.

#### 1 INTRODUCTION

Federated learning has emerged as an attractive paradigm for distributed training of machine learning models in networks of remote devices. While there is a wealth of work on distributed optimization in the context of machine learning, two key challenges distinguish federated learning from traditional distributed optimization: high degrees of *systems and statistical heterogeneity* [McMahan et al., 2017; Li et al., 2019].

In an attempt to handle heterogeneity and tackle high communication costs, optimization methods that allow for local updating and low participation are a popular approach for federated learning [McMahan et al., 2017; Smith et al., 2017]. In particular, FedAvg [McMahan et al., 2017] is an iterative method that has emerged as the de facto optimization method in the federated setting. At each iteration, FedAvg first locally performs  $E$  epochs of stochastic gra-

dient descent (SGD) on  $K$  devices—where  $E$  is a small constant and  $K$  is a small fraction of the total devices in the network. The devices then communicate their model updates to a central server, where they are averaged.

While FedAvg has demonstrated empirical success in heterogeneous settings, it does not fully address the underlying challenges associated with heterogeneity. In the context of statistical heterogeneity, FedAvg is prone to having participating devices to perform variable amounts of local work based on their underlying systems constraints; instead it is common to simply drop devices that fail to compute  $E$  epochs within a specified time window [Bonawitz et al., 2019]. From a statistical perspective, FedAvg has been shown to be statistically suboptimal in settings where data is non-identically distributed [McMahan et al., 2017; Li et al., 2019].

In this work, we propose FedProx, a federated optimization algorithm that addresses the challenges of heterogeneity both theoretically and empirically. A key insight we have in developing FedProx is that a trade-off exists between systems and statistical heterogeneity in federated learning. Indeed, both dropping stragglers (as in FedAvg) or naively incorporating partial information from stragglers (as in FedProx with the proximal term set to 0) implicitly increases statistical heterogeneity and can adversely impact

## III- FedALA (2023)

The Thirty-Seventh AAAI Conference on Artificial Intelligence (AAAI-23)

#### 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>4</sup>, Haibing Guan<sup>1</sup>

<sup>1</sup>Shanghai Jiao Tong University  
<sup>2</sup>Queen's University Belfast  
<sup>3</sup>Louisiana State University  
<sup>4</sup>{zqz, songt333, zhengguix, ruhuihua, hbguan}@sjtu.edu.cn, YHua@qub.ac.uk, hwang@sjtu.edu.cn

#### Abstract

A key challenge in federated learning (FL) is the statistic heterogeneity that impairs the generalization of the global model on each client. To address this, we propose a method *FedALA* learning with *Adaptive Local Aggregation (ALA)* for FL. FedALA is an *Adaptive Local Aggregation (ALA)* model, which consists of a local model initialized by a global model and local model toward the local objective on each client to initialize the local model before training in each iteration. FedALA is a personalized FL method that can conduct extensive experiments with five benchmark datasets in computer vision and natural language processing domains.

FedALA outperforms the baseline methods and achieve up to 24.19% improvement in test accuracy. Furthermore, we also apply ALA model to the personalized FL methods and FedAvg to capture the desired information in the global model through personalized aggregation.

However, the methods in Category (3) still have shortcomings. FedAMP/FedPHP performs personalized aggregation on the server/clients without considering the local objective. FedFomo/APPLE download other client models and locally aggregate the aggregated weights on each client to obtain the final global model. FedFomo only uses the desired information in each iteration to improve the quality of the local model is beneficial for the client. The global model has poor generalization ability since it has desired and noisy information for all individual client simultaneously. Their methods in Category (3) intend to capture the desired information in the global model through personalized aggregation.

However, the methods in Category (3) still have shortcomings. FedAMP/FedPHP performs personalized aggregation on the server/clients without considering the local objective. FedFomo/APPLE download other client models and locally aggregate the aggregated weights on each client to obtain the final global model. FedFomo only uses the desired information in each iteration to improve the quality of the local model is beneficial for the client. The global model has poor generalization ability since it has desired and noisy information for all individual client simultaneously. Their methods in Category (3) intend to capture the desired information in the global model through personalized aggregation.

However, the methods in Category (3) still have shortcomings. FedAMP/FedPHP performs personalized aggregation on the server/clients without considering the local objective. FedFomo/APPLE download other client models and locally aggregate the aggregated weights on each client to obtain the final global model. FedFomo only uses the desired information in each iteration to improve the quality of the local model is beneficial for the client. The global model has poor generalization ability since it has desired and noisy information for all individual client simultaneously. Their methods in Category (3) intend to capture the desired information in the global model through personalized aggregation.

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>4</sup>, Haibing Guan<sup>1</sup>

<sup>1</sup>Shanghai Jiao Tong University  
<sup>2</sup>Queen's University Belfast  
<sup>3</sup>Louisiana State University  
<sup>4</sup>{zqz, songt333, zhengguix, ruhuihua, hbguan}@sjtu.edu.cn, YHua@qub.ac.uk, hwang@sjtu.edu.cn

11237



Fair Federated AI  
Summer School



<https://basira-lab.com/>

RISE  
MICCAI

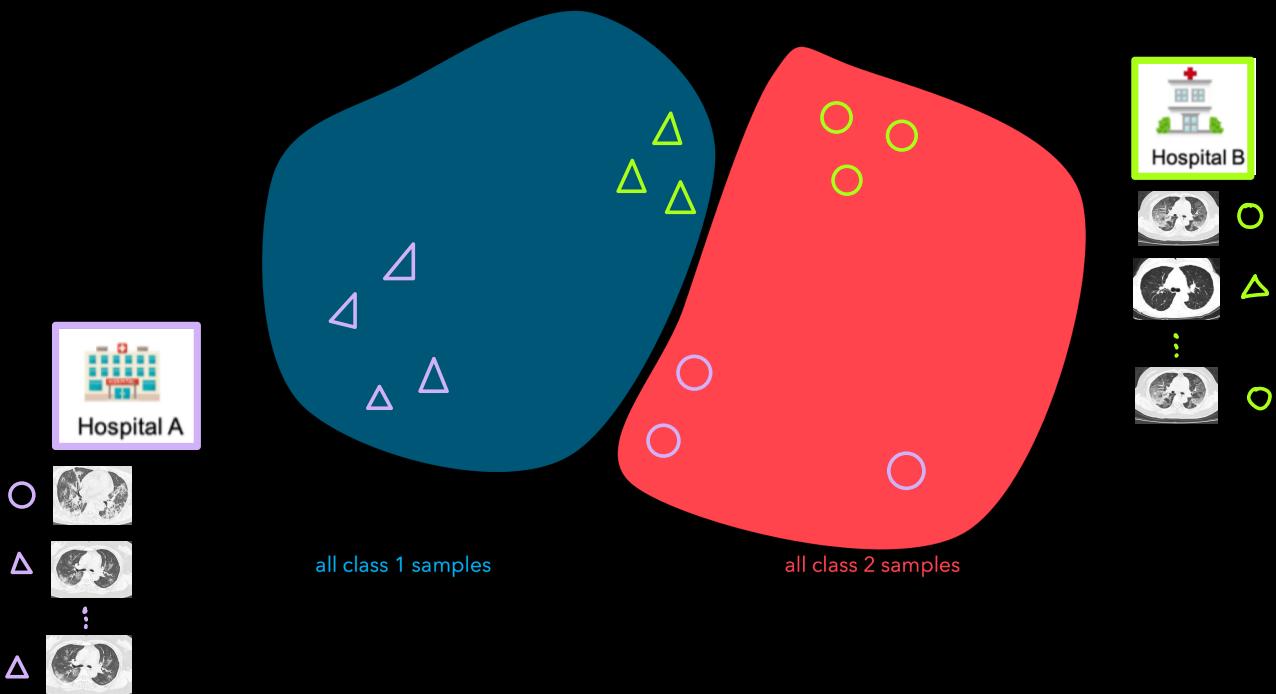
FAIMI



<https://github.com/RISE-MICCAI/FIFAI-Summer-School-2024>

2) What is federated learning?

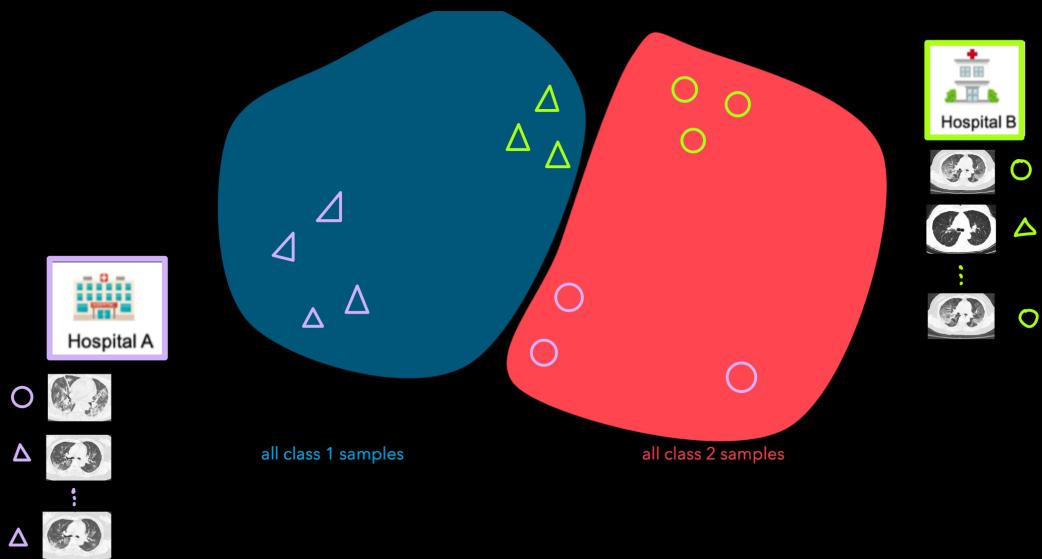
why federate?  
How to federate?



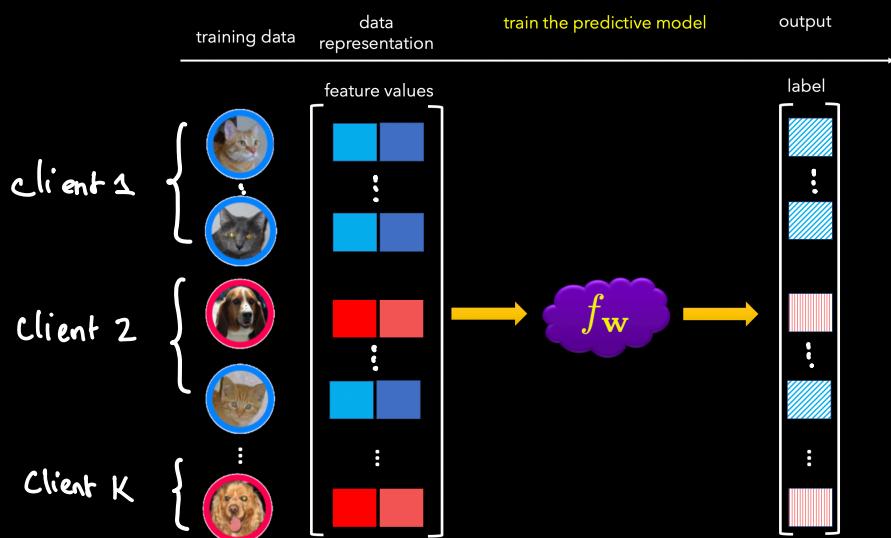
现代社会的机器学习算法（分类，回归）  
是数据和计算饥渴的。

现代机器学习算法（分类，回归）  
的数据和计算需求很高。  
训练数据的大小、分布等决定了模型的性能。

→ 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.

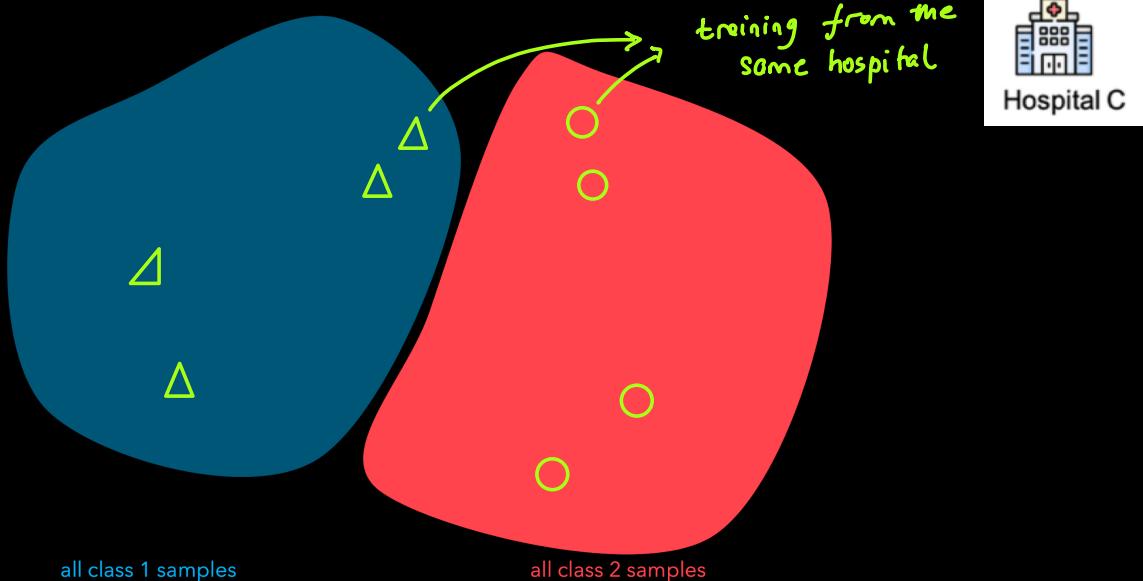


- || Concerns about data security and privacy.
- || Such a data "sharing" process is not straightforward.
- || Large-scale datasets are computationally expensive to upload.  
(e.g., on mobile phones & apps)
- || Intrusive and difficult if data is spread across multiple devices or clients.

→ How can we design powerful & generalizable ML models ?

## Solution 2

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



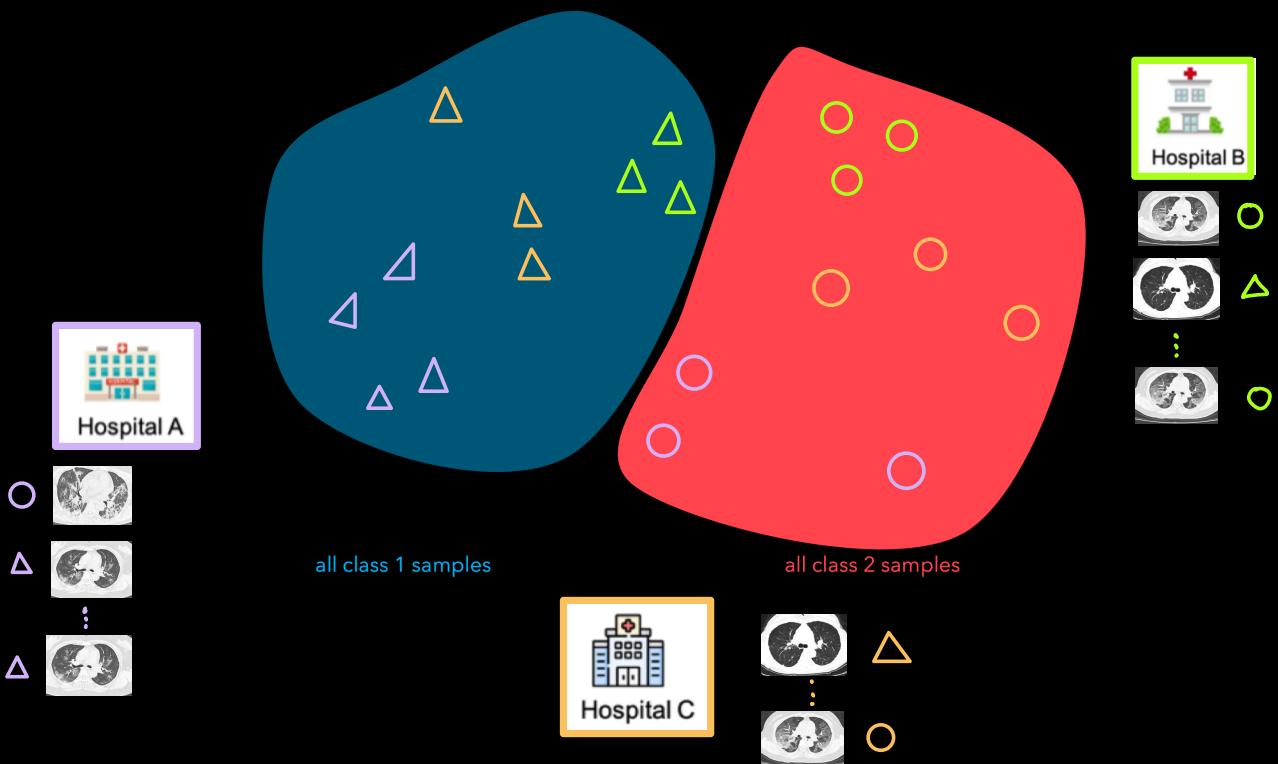
single model (classifier)



- boost the model's performance
  - high computations (many models to train on a single device)
  - training only benefits from a single dataset (here Hospital C)

→ 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)

 Communication-Efficient Learning of Deep Networks  
from Decentralized Data

---

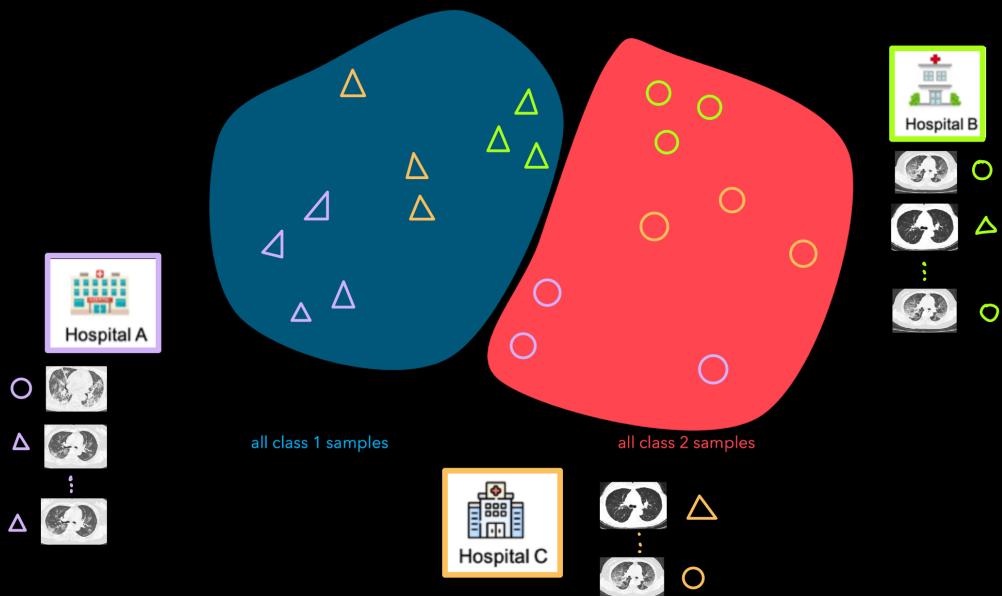
H. Brendan McMahan   Eider Moore   Daniel Ramage   Seth Hampson   Blaise Agüera y Arcas  
Google, Inc., 651 N 34th St., Seattle, WA 98103 USA

<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 ( ? )

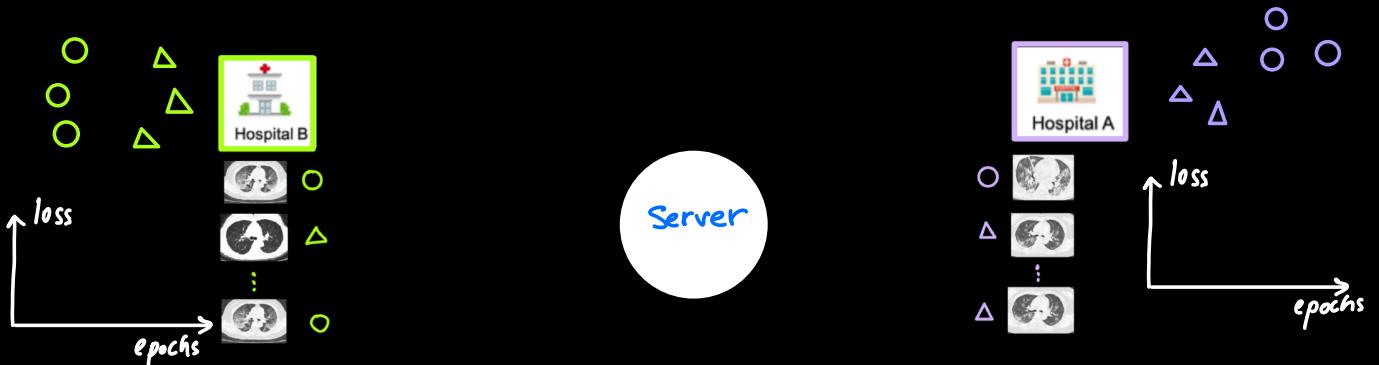
**THINK!**



What's the trick here ?

## 2) What is federated learning?

why federate?      How to federate?



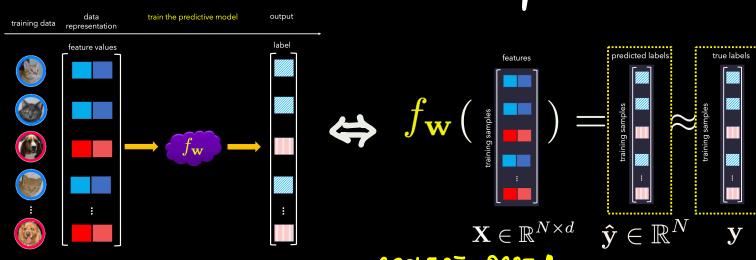
In each round

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

**FED AVG**

\* Hyperparameters :

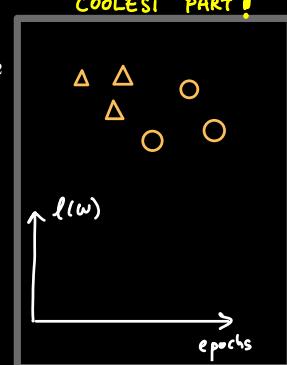
- K : total number of clients
- T : total communication rounds
- N : local iteration steps per round
- M : clients per round (randomly selected)



$$\min_w l(w) = \sum_{i=1}^n (\hat{y}^i - y^i)^2$$

↓  
solving the problem using cool math

$$\tilde{w}^* = [\tilde{X} \tilde{X}^T]^{-1} \tilde{X} y$$

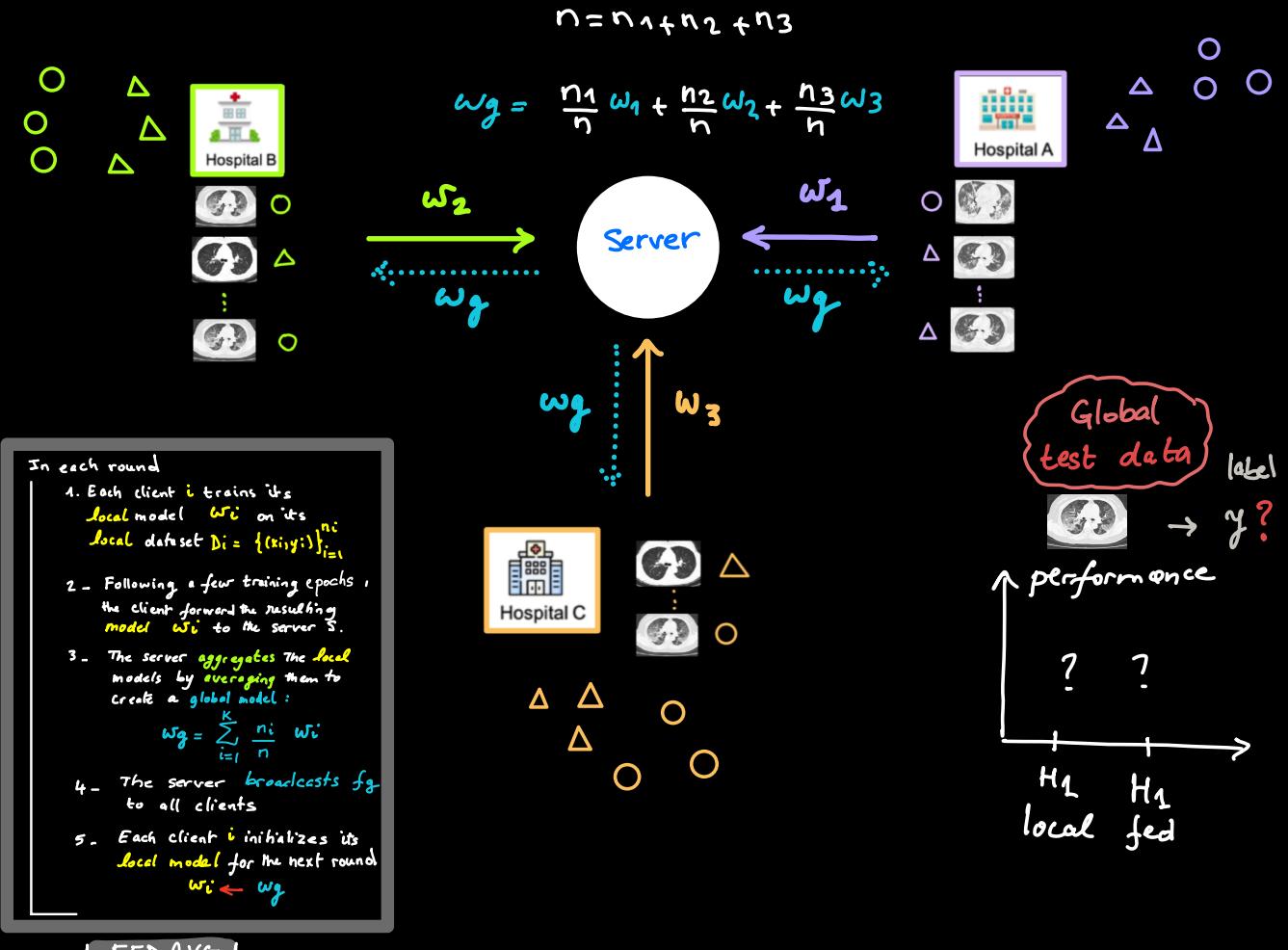


$$\tilde{w}^* = [\tilde{X} \tilde{X}^T]^{-1} \tilde{X} y$$

training samples      training samples  
 $\tilde{X}$        $y$

## 2) What is federated learning?

why federate?      How to federate?



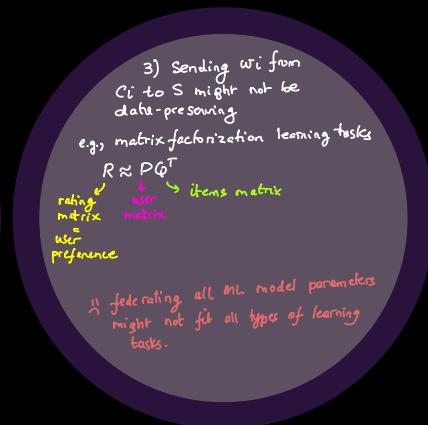
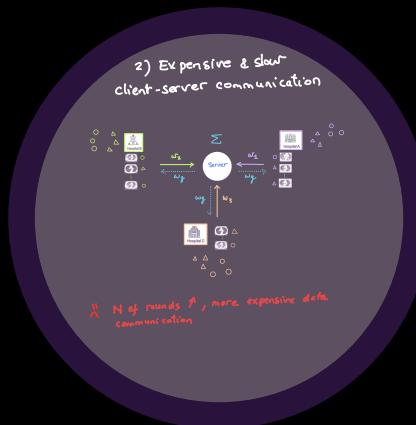
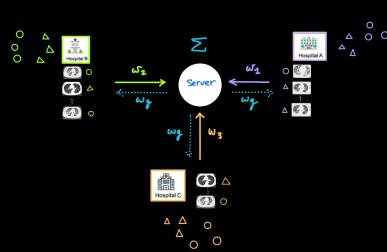
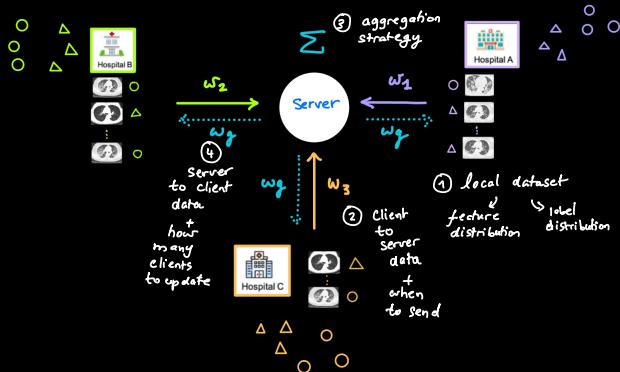
### \* Hyperparameters :

- K : total number of clients
- T : total communication rounds
- N : Local epochs per round
- M : clients per round (randomly selected)

### 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



design new solutions

design new solutions

design new solutions

I- FedAvg (2016) [1]

a) IID label distribution across clients in classification

b) Non-IID scenario

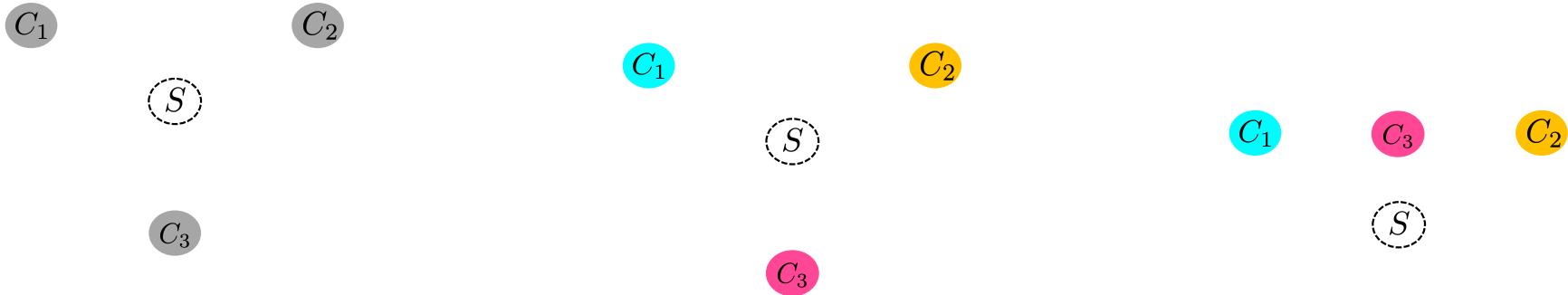
FL experiment setup for evaluation

[Scenario 1] Non-IID label distribution in classification

Random data split

{ global test / eval  
local test / eval

•  $d_i = \{x_i, y_i\}$   
data point

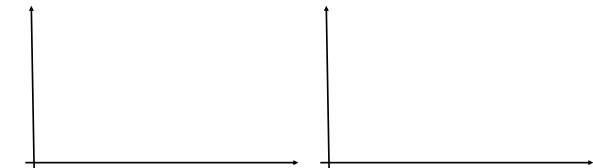
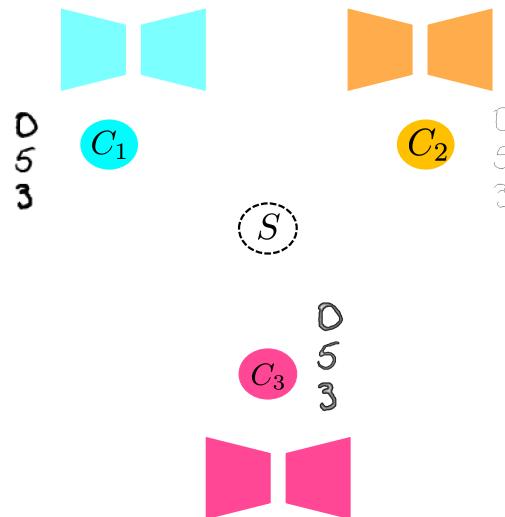


Perturb initial model and data distributions by changing the random seed



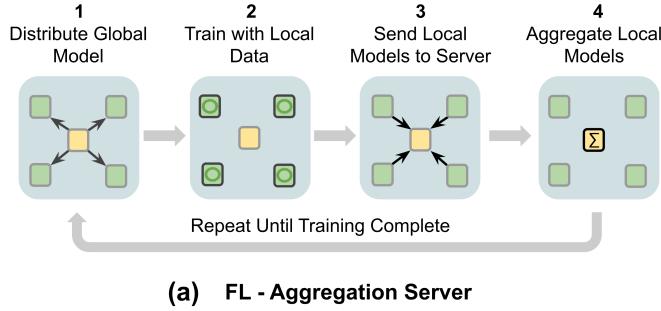
[Scenario 2] Non-IID feature distribution in generative learning

Evaluation measures and benchmarks



Diverse learning paradigms: centralized versus federated [2]

Personal reflections and notes



**(b) FL - Peer to Peer**

**Key**

- |                            |                         |
|----------------------------|-------------------------|
| ■ Aggregation Server       | ■ Central Data Lake     |
| ■ Training Node            | ■ Data Donor            |
| Σ Model Aggregation        | → Time                  |
| ← Weight/Gradient Exchange | → Medical Data Exchange |
|                            | ○ 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

Existing solutions:

- 
- 

System heterogeneity:

Existing solutions:

- 
- 

### 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  $\gamma_k^t$ -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)

**Challenge 1 to solve:** With dissimilar (heterogeneous) local objectives  $F_k(\cdot)$ , 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:**

Proposed solution:

local updates. In particular, instead of just minimizing the local function  $F_k(\cdot)$ , device  $k$  uses its local solver of choice to approximately minimize the following objective  $h_k$ :

$$\min_w h_k(w; w^t) = F_k(w) + \frac{\mu}{2} \|w - w^t\|^2. \quad (2)$$

The proximal term is beneficial in two aspects: (1) It addresses the issue of statistical heterogeneity by restricting the local updates to be closer to the initial (global) model without any need to manually set the number of local epochs. (2) It allows for safely incorporating variable amounts of local work resulting from system heterogeneity. We summarize the steps of FedProx in Algorithm 2.

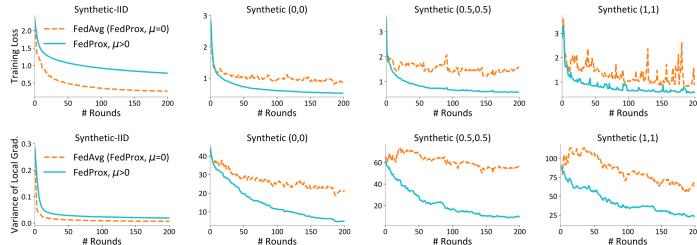
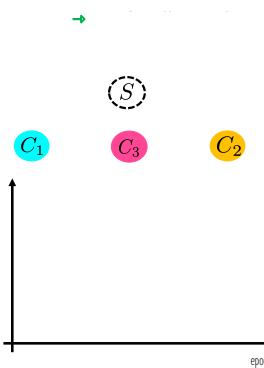


Figure 2. Effect of data heterogeneity on convergence. We remove the effects of systems heterogeneity by forcing each device to run the same amount of epochs. In this setting, FedProx with  $\mu = 0$  reduces to FedAvg. (1) Top row: We show training loss (see results on testing accuracy in Appendix C.3, 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.

## 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>

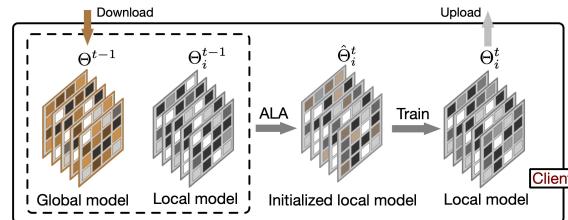


Figure 1: Local learning process on client  $i$  in the  $t$ -th iteration. Specifically, client  $i$  downloads the global model from the server, locally aggregates it with the old local model by ALA module for local initialization, trains the local model, and finally uploads the trained local model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

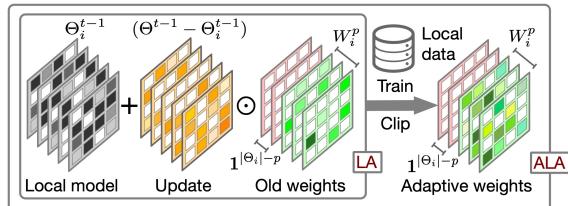


Figure 2: The learning process in ALA. LA denotes “local aggregation”. Here, we consider a five-layer model and set  $p = 3$ . The lighter the color, the larger the value.

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.

$$\min_w h_k(w; w^t) = F_k(w) + \frac{\mu}{2} \|w - w^t\|^2$$

$$\hat{\Theta}_i^t := \Theta_i^{t-1} + (\Theta^{t-1} - \Theta_i^{t-1}) \odot W_i$$

#### 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\\_files/paper/2020/file/1f5fe83998a09396ebe477d9475ba0c-Paper.pdf](https://proceedings.mlsys.org/paper_files/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/aaai.v37i9.76330>; <https://github.com/TsingU/FedALA>