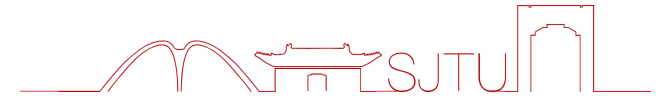




上海交通大学
SHANGHAI JIAO TONG UNIVERSITY



Federated Learning (FL) Lab



饮水思源 · 爱国荣校



0

Backgrounds

1

FedAvg algorithm

2

Benchmark algorithms

3

Advanced algorithms

A photograph of a modern building with a white, angular facade and large glass windows, set against a blue sky with light clouds. The building is the central focus of the background.

00

Backgrounds

- **Governmental support**
- **Develop trend**

- ① 数据最高管理部门——国家数据局成立
- ② 负责协调推进数据基础制度建设，统筹数据资源整合共享和开发利用，统筹推进数字中国、数字经济、数字社会规划和建设等，由国家发展和改革委员会管理。
- ③ 国家数据局的成立或将促进以下行业的发展：
 - 数据收集（各行业数字化转型）
 - 数据储存(国资云，底层算力基础设施)
 - 数据处理（数据安全与隐私计算）
 - 数据定价（数据资产化）
 - 数据流转(数据交易机构、公共数据交易)



新华社

23-3-7 16:46 来自 微博网页版 已编辑

【#组建国家数据局#】根据国务院关于提请审议国务院机构改革方案的议案，组建国家数据局。负责协调推进数据基础制度建设，统筹数据资源整合共享和开发利用，统筹推进数字中国、数字经济、数字社会规划和建设等，由国家发展和改革委员会管理。

将中央网络安全和信息化委员会办公室承担的研究拟订数字中国建设方案、协调推动公共服务和社会治理信息化、协调促进智慧城市建设、协调国家重要信息资源开发利用与共享、推动信息资源跨行业跨部门互联互通等职责，国家发展和改革委员会承担的统筹推进数字经济发展、组织实施国家大数据战略、推进数据要素基础制度建设、推进数字基础设施布局建设等职责划入国家数据局。#全国两会时光#





隐私计算及联邦学习



	安全多方计算	联邦学习	同态加密	可信执行环境	合成数据	差分隐私	秘密共享
波士顿女性劳动力							
欧洲统计系统							
欧盟统计局							
印度尼西亚旅游部							
意大利国家统计研究所							
英国国家统计局							
三星SDS (韩国)							
加拿大统计局							
韩国统计局							
荷兰统计局							
联合国欧洲经济委员会							
美国人口普查局							
美国教育部							



- ① 联邦学习相关国家自然科学基金项目获批数呈现出逐年增加的趋势，这表明联邦学习在学术界和工业界的应用越来越受到重视。联邦学习是多个领域的融合，触及到的研究方向众多，应用场景非常丰富。
- ② 金融领域有复旦大学柴洪峰团队获批的重大研究计划《大数据背景下基于联邦学习的小微企业信用风险评估研究》
- ③ 通信领域有深圳市大数据研究院朱光旭团队获批的青年科学基金项目《面向联邦式边缘学习的高效通信技术研究》
- ④ 医疗领域有上海大学武星团队的面上项目《多中心胶囊内窥镜影像联邦主动学习》



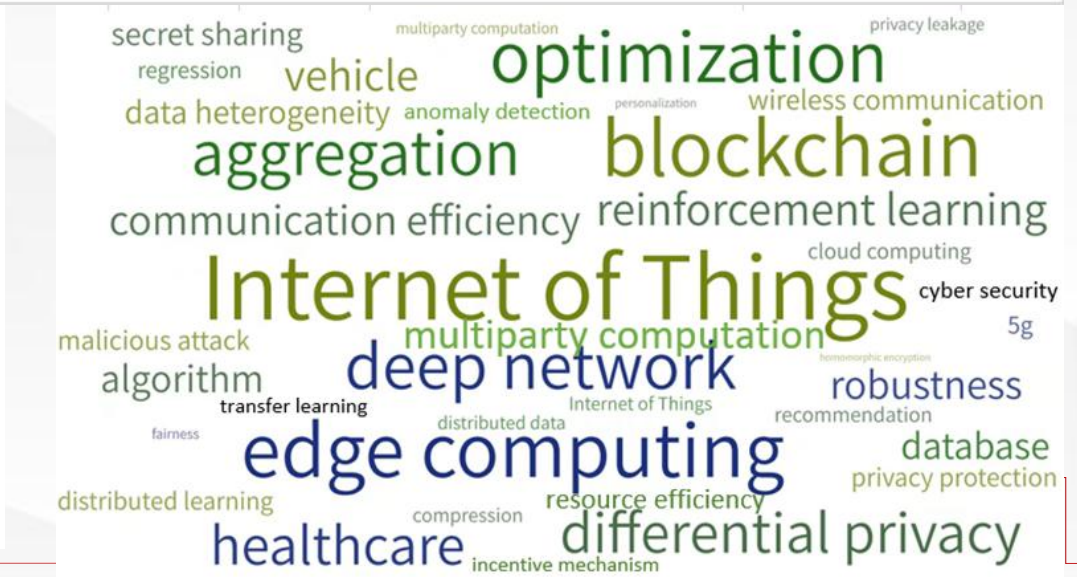
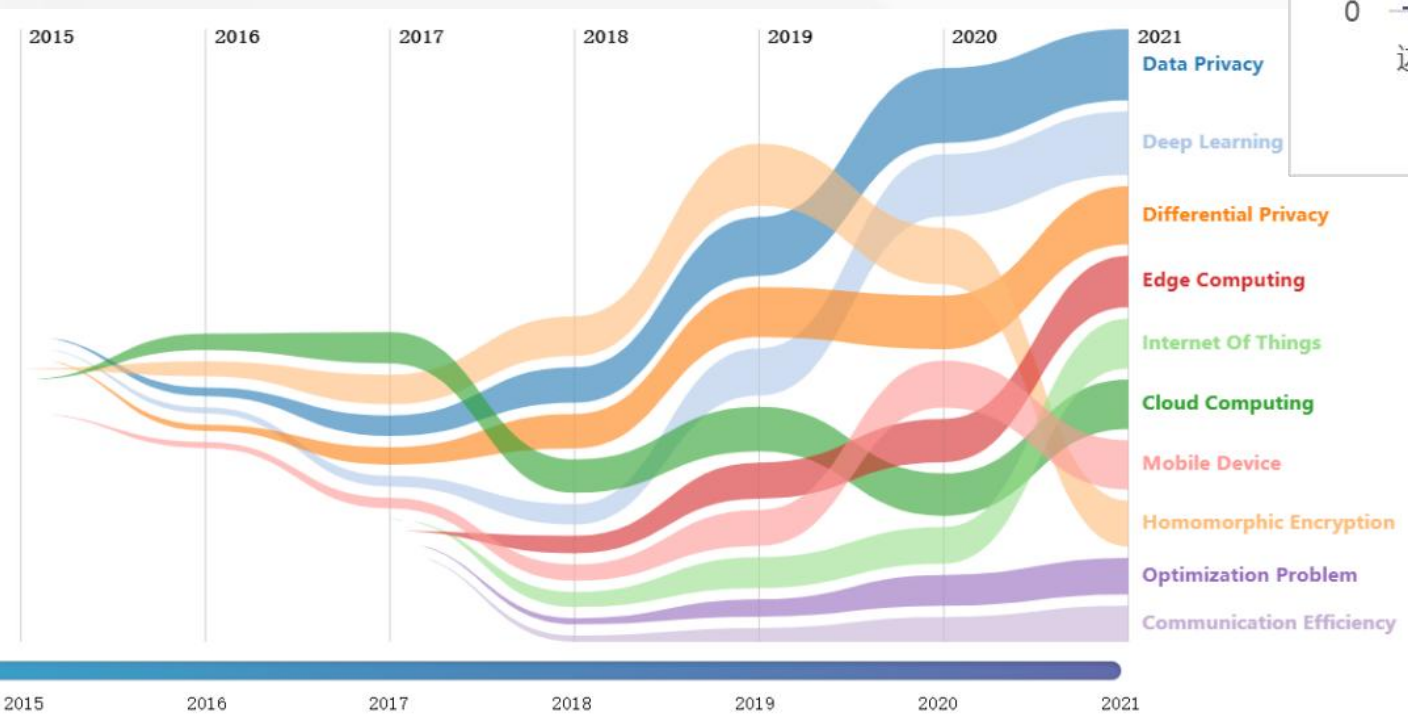
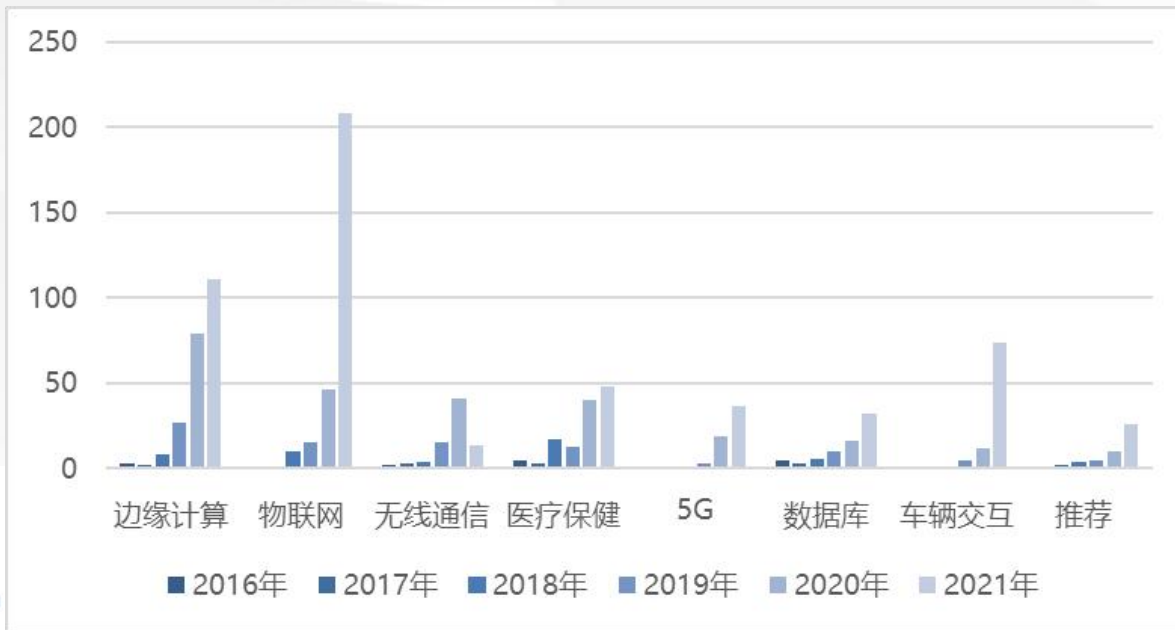


隐私计算及联邦学习



① 涵盖多个领域，主要以边缘计算、物联网和车辆交互为主

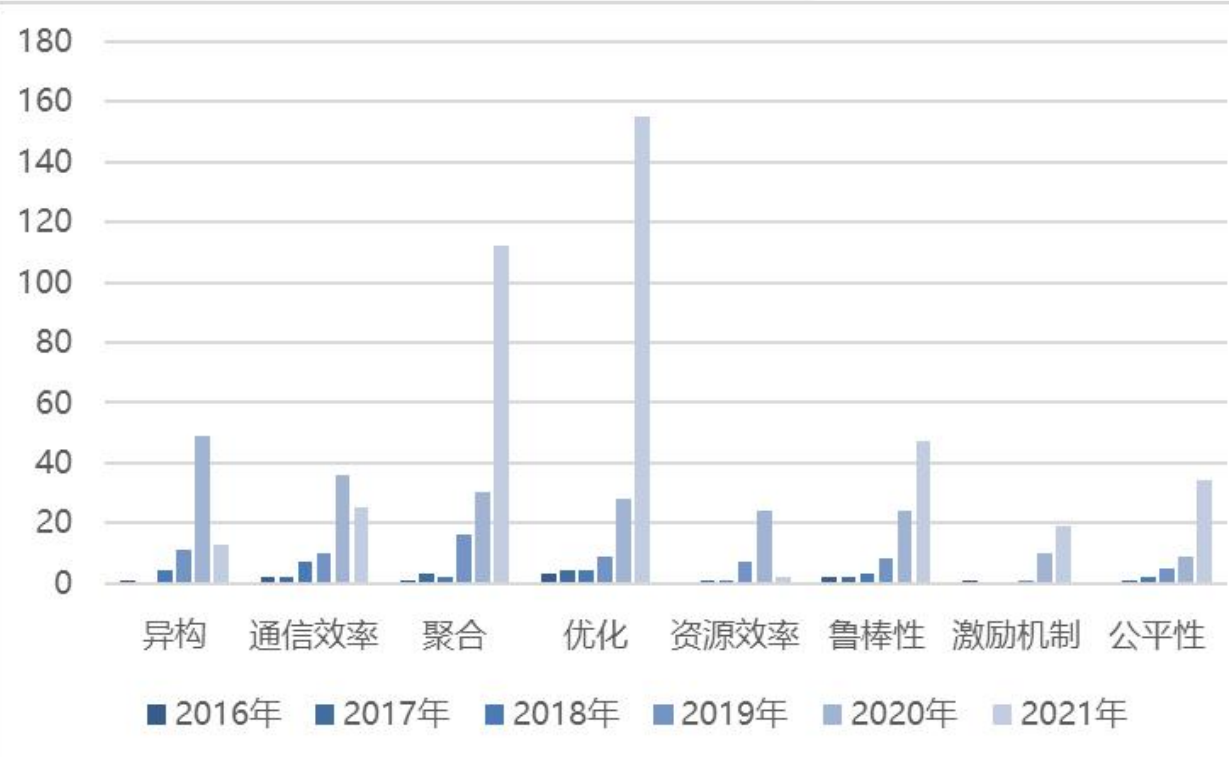
② 各领域联邦学习研究逐年增长，趋势良好





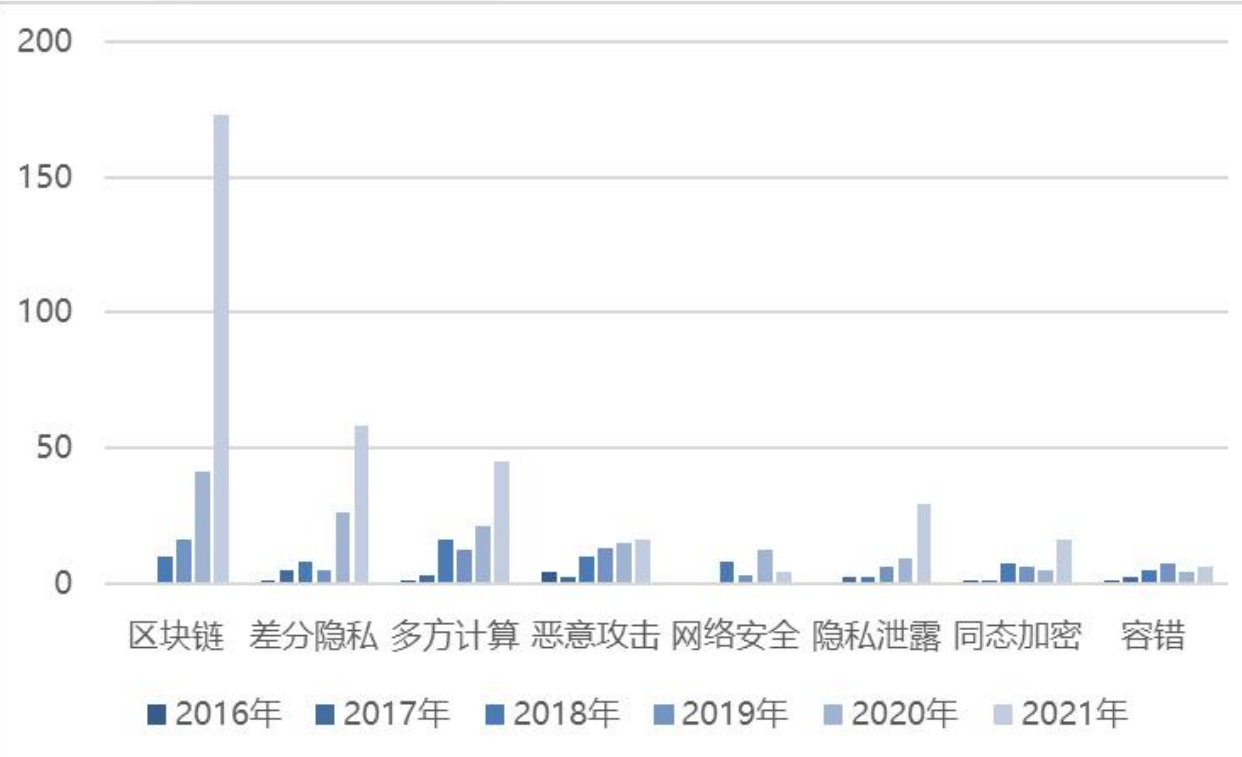
联邦学习性能上以聚合、优化方法和异构为主

各方法研究量逐年增长



联邦学习安全上以区块链、差分隐私和多方计算为主

各方法研究量逐年增长





开源框架

Github热度	发布方	系统名称	开源时间	系统特点
8000	OpenMind	PySyft	2017.7	<ul style="list-style-type: none"> • 一个用于安全和私有深度学习的 Python 库 • 基于 PyTorch, 使用 Unity Game Engine • 安全多方计算 • 联合学习、差异隐私
4100	微众银行	FATE	2019.2	<ul style="list-style-type: none"> • 工业级框架, 采用 Python开发, 底层计算存储基于 EGGROLL、Spark 等高性能计算引擎 • 提供一站式的联邦模型企业级服务解决方案。提供多插件支持联邦学习企业和科研应用 • 支持主流的分类、回归、聚类和迁移学习的联邦化算法 • 提供多种安全计算协议支撑上层应用, 支持同态加密协议、秘密共享协议、不经意传输协议和 DH密钥交换算法等 • 提供20多个联邦算法组件
1800	谷歌	TensorFlow Federated	2019.3	<ul style="list-style-type: none"> • 可以选择 ML 模型架构 • 模型设计理念以数据为主



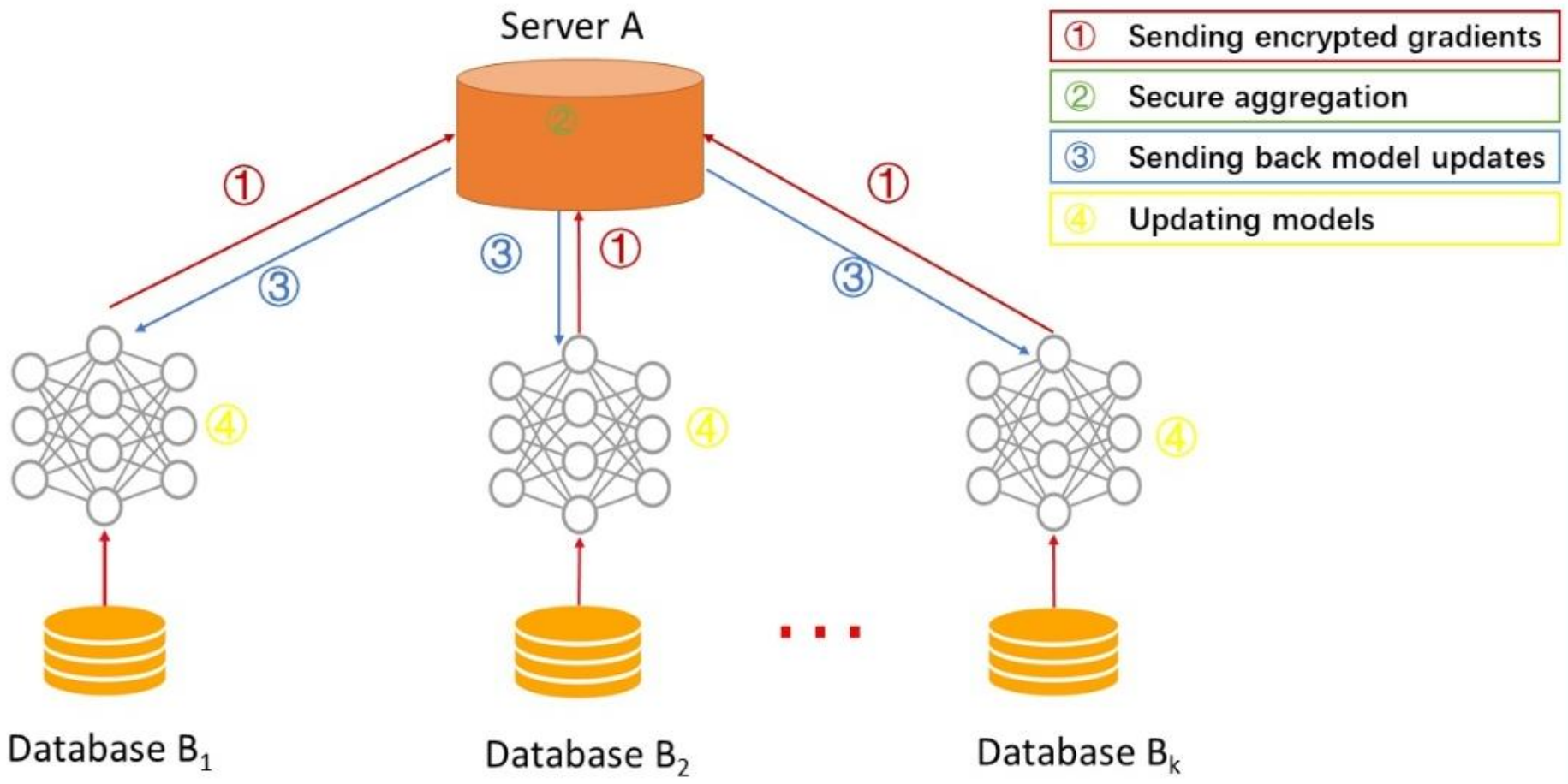
联邦学习介绍



- ① 联邦学习是一种分布式机器学习技术，或机器学习框架。
- ② 在保证**数据隐私安全及合法合规**的基础上，实现**共同建模**，提升AI模型的效果。

③ 本质上是通过多个用户设备共同训练一个代表所有用户设备的**全局模型**

④ 训练过程不需要用户数据的交换，更强调**隐私性**。



01

FedAvg algorithm

- Background and contributions
- Federated learning review
- Federated settings
- Federated Averaging
- Experiments

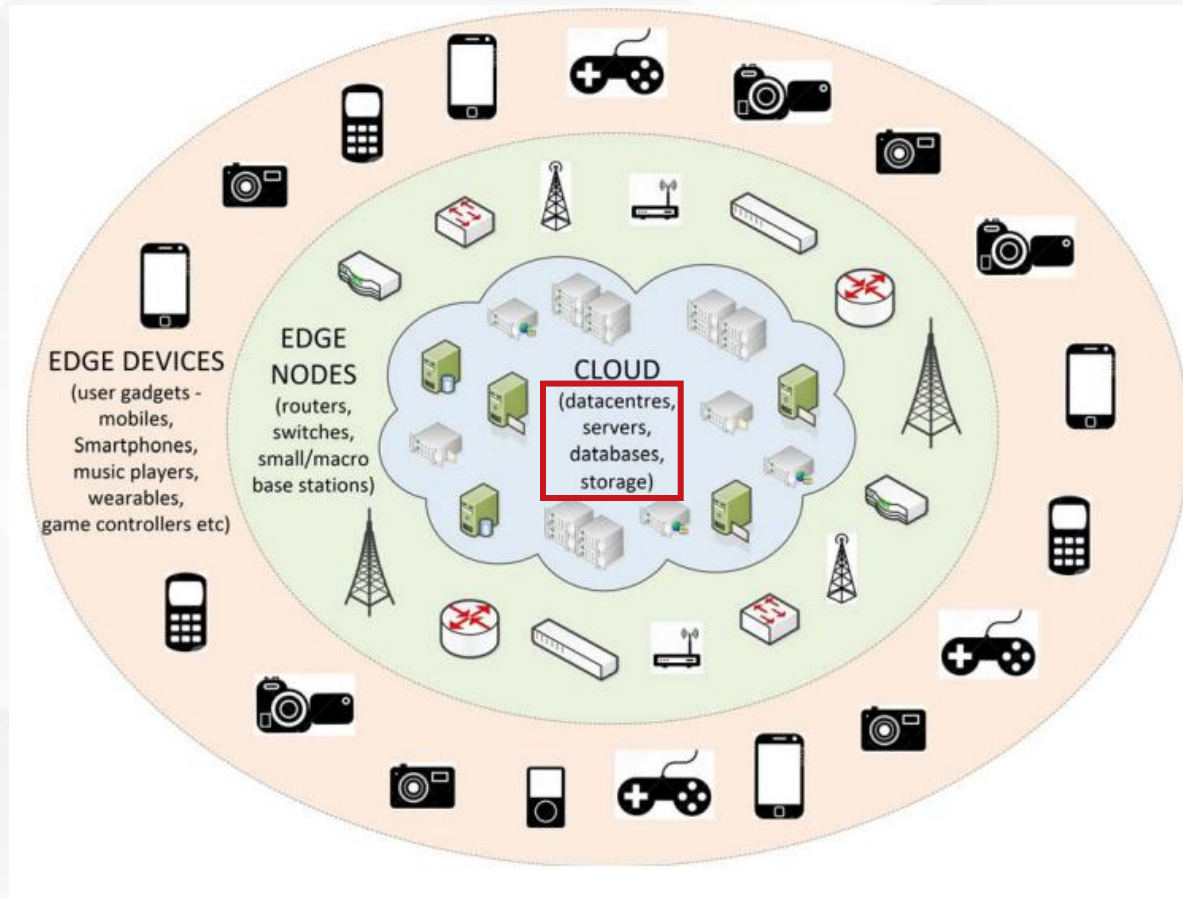


Background

- Recently, attributed to the collection of **massive data** from users or organizations, AI has been thriving for years.
- However, some **private data** is also collected during data collection, such as the shopping behaviors, facial images, house locations, etc.
- Data breach is becoming more and more severe, and many governments have issued **data privacy protecting laws**.

Background

- AI model training based on distributed computing and edge computing is required.



Background

- Different from traditional settings, data cannot be collected and naturally exists locally.

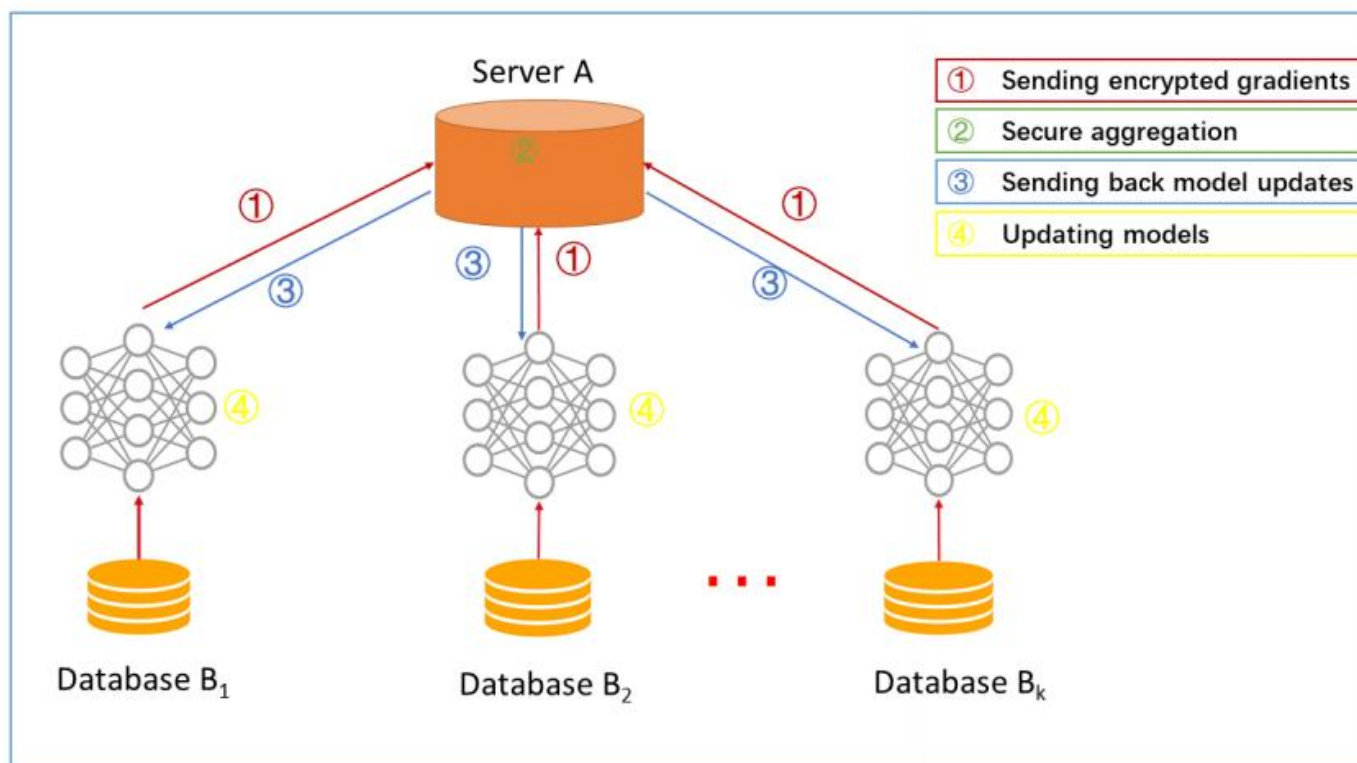


Fig. 3. Architecture for a horizontal federated learning system



Contributions

- Consider the problem of training dispersed data from mobile devices as an important research direction.
- Propose a simple and practical algorithm for Federated Averaging.
- An extensive empirical evaluation of the proposed algorithms shows that they are robust to non-independently identically distributed (Non-IID) and unbalanced data.

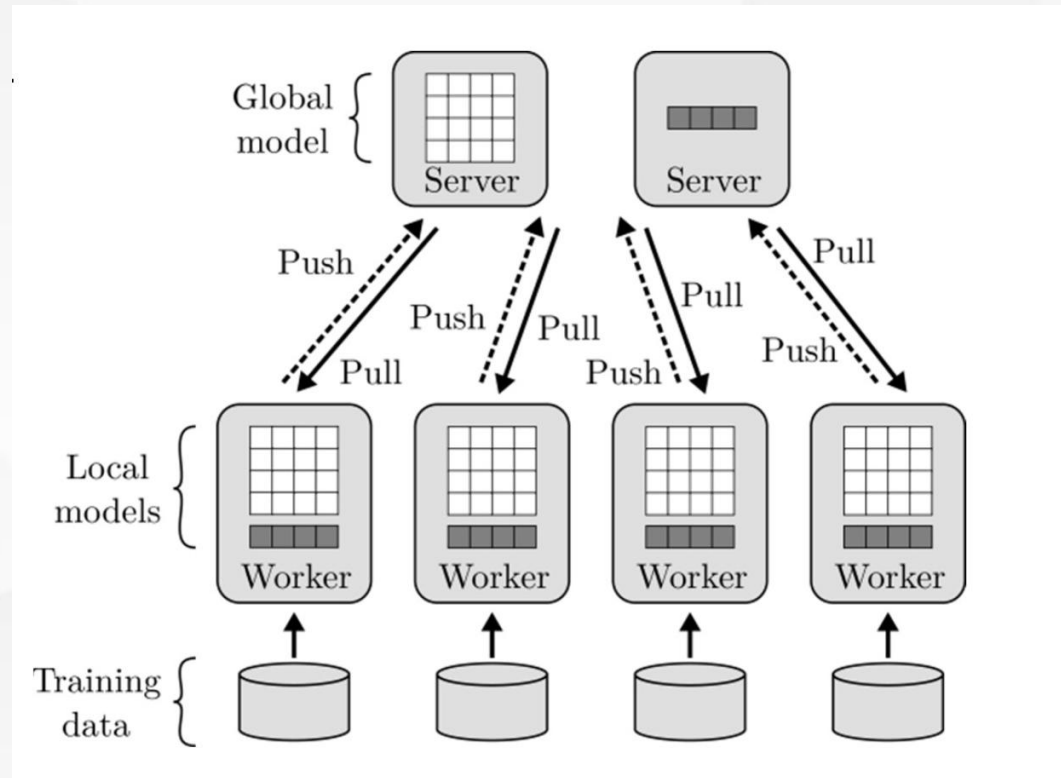


FedAvg algorithm



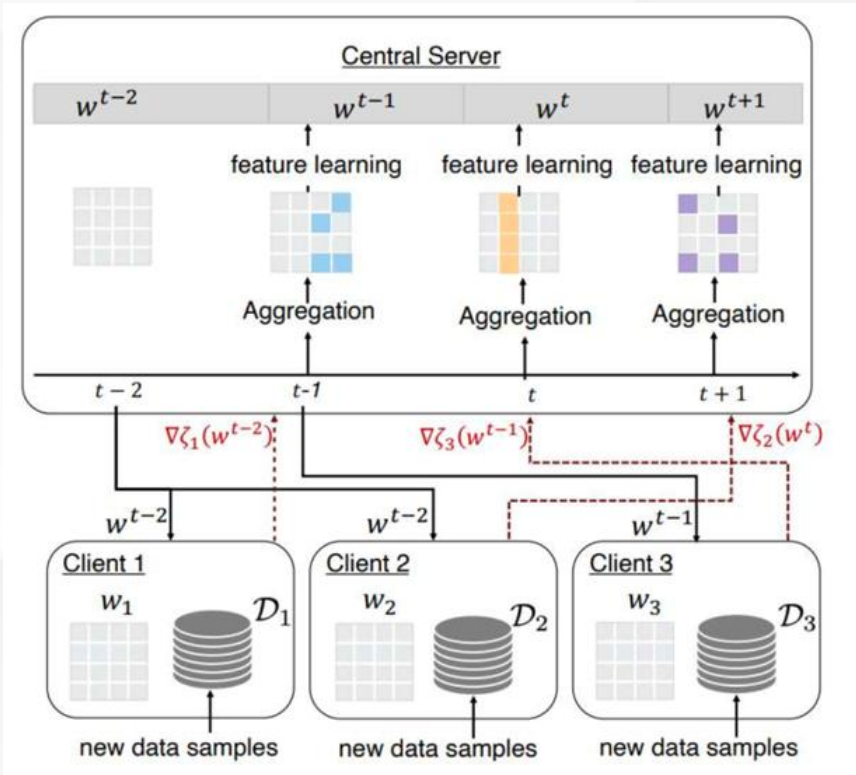
Federated learning review

- Learning tasks are handled by a loose federation of participating devices (clients) coordinated by a central server.



Federated learning review

- Each client has a local training data set that does not need to be uploaded to the server and only sends local model parameters for each update.





- ① Federated learning (FL) review: ideal FL
 - Training on **real-world data** from mobile devices has distinct advantages over the proxy data that is ubiquitous in data centers.
 - This data is **privacy-sensitive** or large (compared to the size of the model) to avoid recording it to the data center for model training.
 - For supervised tasks, labels on data can be **naturally inferred** from user interactions.



⊗ Federated learning (FL) review: privacy

- In the traditional distributed training setting, even if an "anonymous" data set is held, users' privacy will be threatened through the **connection with other data**.
- In contrast, the information transmitted in FL is the **minimum update** (all/part of the model parameters) needed to improve a particular model, and less information means a lower risk of privacy disclosure.
- Combining FL with **secure multi-party computing and differential privacy**.



⊗ Federated settings

- Non-IID data: Training data on a given client is usually based on mobile device usage by a particular user, so any local data set for a particular user **does not represent a group distribution.**
- Imbalance: Some users will use the service or application more than others, resulting in **different amounts of local training data.**
- Massive clients: We expect the number of clients participating in the FL to be **much larger than** the average number of instances per client.
- Limited communication: Mobile devices are often **offline or in a slow or expensive connection.**



Objective

$$\min_{w \in \mathbb{R}^d} f(w) \quad \text{where} \quad f(w) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(w). \quad (1)$$



$$f(w) = \sum_{k=1}^K \frac{n_k}{n} F_k(w) \quad \text{where} \quad F_k(w) = \frac{1}{n_k} \sum_{i \in \mathcal{P}_k} f_i(w).$$

$$f_i(w) = l(x_i, y_i; w)$$

K number of clients

\mathcal{P}_k the distribution of local data

n_k the number of local data samples



⊗ Limited communication

- In data center optimization, communication costs are relatively small while computing costs dominate, and recent emphasis has been on using GPUs to reduce these costs. In contrast, in joint optimization, **communication costs dominate** -- we're typically limited by upload bandwidth of 1MB/s or less.
- Clients typically volunteer for optimization only when charging, plugging in, and using a non-billable Wi-Fi connection, and we expect each client to **participate in a small number of update sessions** per day.
- Modern smartphones have **plenty of local computing power and small datasets** on a single device.



⊗ Reduce communication

- Increasing **parallelism**, that is, using more clients to work independently between rounds of communication.
- Add computations per client, that is, **perform more complex computations** (such as cumulative training) between rounds of communication.



Federated Averaging

Algorithm 1 FederatedAveraging. The K clients are indexed by k ; B is the local minibatch size, E is the number of local epochs, and η is the learning rate.

Server executes:

```
initialize  $w_0$ 
for each round  $t = 1, 2, \dots$  do
   $m \leftarrow \max(C \cdot K, 1)$ 
   $S_t \leftarrow$  (random set of  $m$  clients)
  for each client  $k \in S_t$  in parallel do
     $w_{t+1}^k \leftarrow$  ClientUpdate( $k, w_t$ )
   $w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$ 
```

ClientUpdate(k, w): // Run on client k

```
 $\mathcal{B} \leftarrow$  (split  $\mathcal{P}_k$  into batches of size  $B$ )
for each local epoch  $i$  from 1 to  $E$  do
  for batch  $b \in \mathcal{B}$  do
     $w \leftarrow w - \eta \nabla \ell(w; b)$ 
return  $w$  to server
```

K : Client amount

B : Batch size

E : Local training rounds

η : Local learning rate

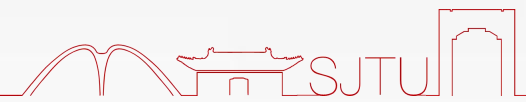
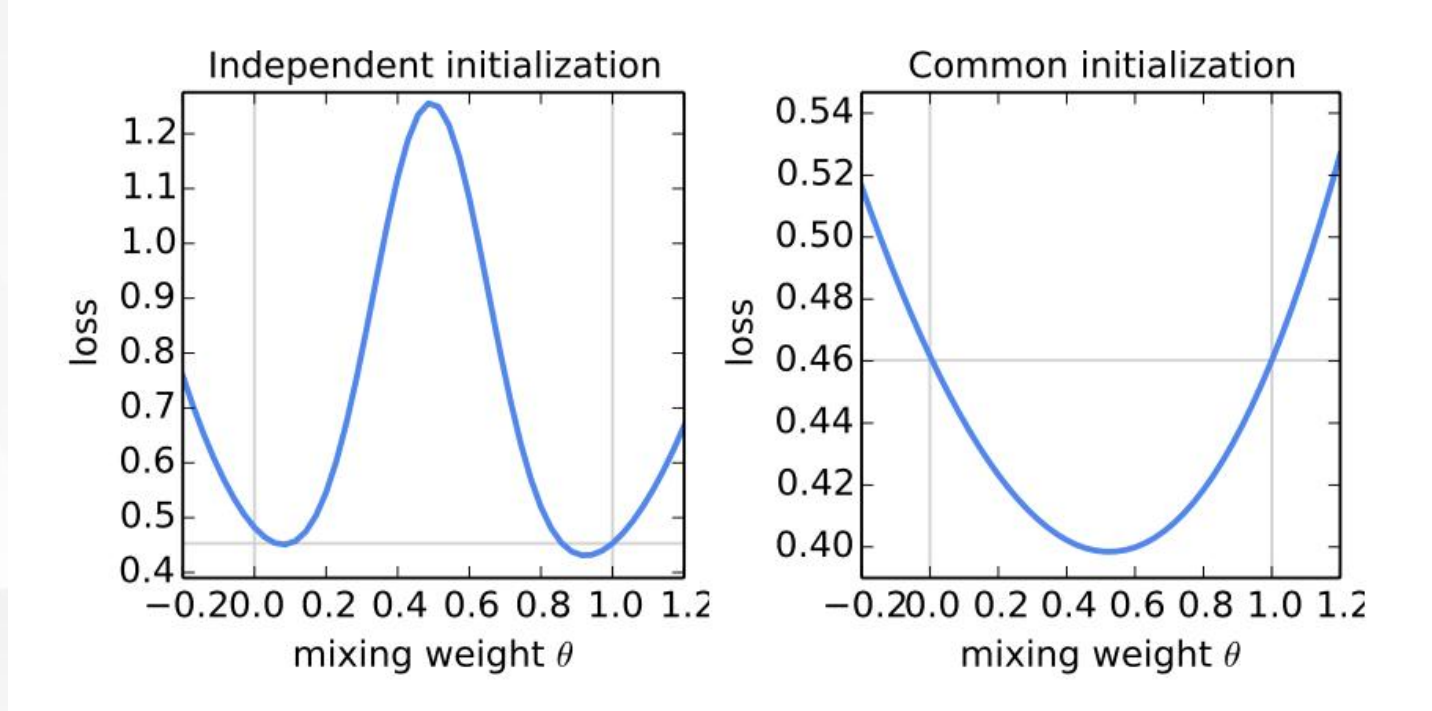


FedAvg algorithm



Federated Averaging

- It is beneficial to initialize local models with the common start point
- Mixed weight = $\theta w + (1 - \theta)w'$





Experiments

- Image classification: MNIST handwritten digital recognition
- Language modeling: Dataset based on the Complete Works of William Shakespeare



Experiments (MNIST)

- IID: shuffling the data and dividing it into 100 clients, each receiving 600 examples.
- Non-IID: (1) sort the data by number label; (2) divide it into 200 shards of size 300; (3) assign 2 shards to each of the 100 clients, most clients will have only two number examples.

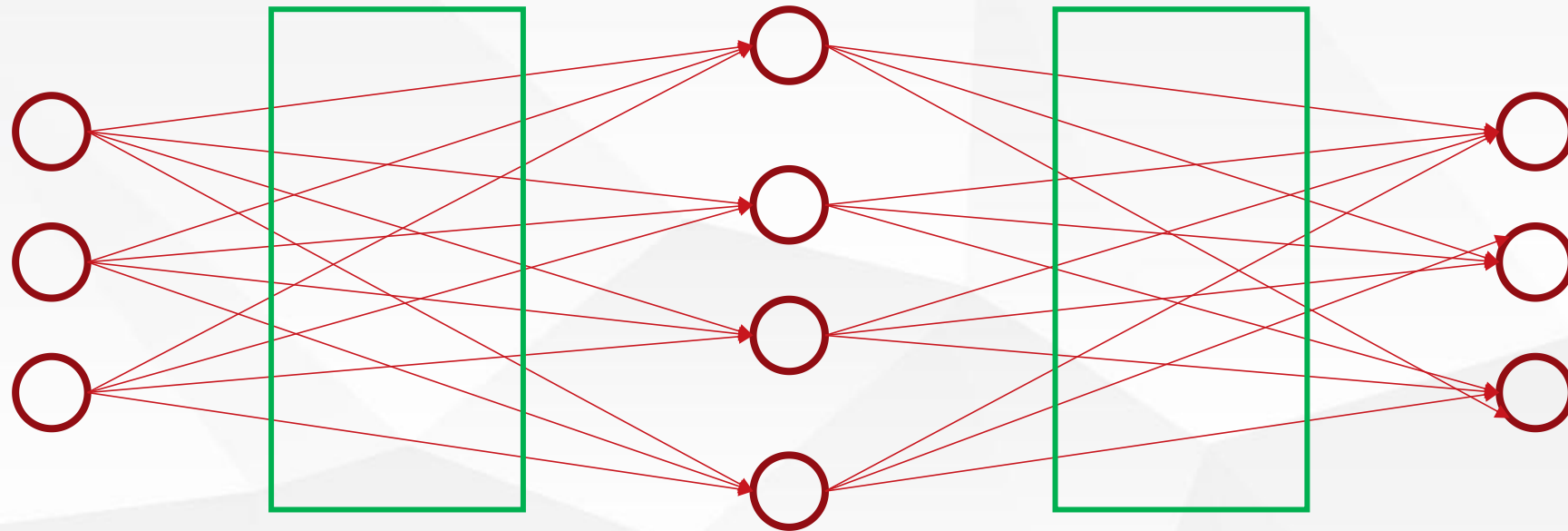


FedAvg algorithm



Experiments (MNIST)

- 1) A simple multilayer-perceptron with 2-hidden layers with 200 units each using ReLU activations (199,210 total parameters), which we refer to as the MNIST 2NN.



Input neurons

Hidden layer

Hidden neurons

Hidden layer

Output neurons





Experiments (MNIST)

- 1) A simple multilayer-perceptron with 2-hidden layers with 200 units each using ReLU activations (199,210 total parameters), which we refer to as the MNIST 2NN.

```
class FedAvgMLP(nn.Module):
    def __init__(self, in_features=784, num_classes=10, hidden_dim=200):
        super().__init__()
        self.fc1 = nn.Linear(in_features, hidden_dim)
        self.fc2 = nn.Linear(hidden_dim, num_classes)
        self.act = nn.ReLU(inplace=True)

    def forward(self, x):
        if x.ndim == 4:
            x = x.view(x.size(0), -1)
        x = self.act(self.fc1(x))
        x = self.fc2(x)
        return x
```



Experiments (MNIST)

- 2) A CNN with two 5x5 convolution layers (the first with 32 channels, the second with 64, each followed with 2x2 max pooling), a fully connected layer with 512 units and ReLU activation, and a final softmax output layer (1,663,370 total parameters).

```
class FedAvgCNN(nn.Module):
    def __init__(self, in_features=1, num_classes=10, dim=1024):
        super().__init__()
        self.conv1 = nn.Sequential(
            nn.Conv2d(in_features,
                      32,
                      kernel_size=5,
                      padding=0,
                      stride=1,
                      bias=True),
            nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=(2, 2))
        )
```





Experiments (Key hyperparameters)

- C: Proportion of clients participating in calculation: 1 indicates that all clients participate in training
- E: Number of training cycles per client between two communications
- B: Mini-batch size of each client. ∞ indicates full-batch
- $u_k = E \frac{n_k}{B}$: The total number of updates in each iteration on client k



FedAvg algorithm



Experiments (communication rounds)

2NN <i>C</i>	IID		NON-IID	
	$B = \infty$	$B = 10$	$B = \infty$	$B = 10$
0.0	1455	316	4278	3275
0.1	1474 (1.0×)	87 (3.6×)	1796 (2.4×)	664 (4.9×)
0.2	1658 (0.9×)	77 (4.1×)	1528 (2.8×)	619 (5.3×)
0.5	— (—)	75 (4.2×)	— (—)	443 (7.4×)
1.0	— (—)	70 (4.5×)	— (—)	380 (8.6×)

CNN, $E = 5$				
0.0	387	50	1181	956
0.1	339 (1.1×)	18 (2.8×)	1100 (1.1×)	206 (4.6×)
0.2	337 (1.1×)	18 (2.8×)	978 (1.2×)	200 (4.8×)
0.5	164 (2.4×)	18 (2.8×)	1067 (1.1×)	261 (3.7×)
1.0	246 (1.6×)	16 (3.1×)	— (—)	97 (9.9×)

C = 0.1 is the best

C: Proportion of clients participating in calculation

E: Number of training cycles per client between two communications

B: Mini-batch size of each client. ∞ indicates full-batch

$u_k = E \frac{n_k}{B}$: The total number of updates in each iteration on client k





FedAvg algorithm



Experiments (communication rounds)

MNIST CNN, 99% ACCURACY						
CNN	E	B	u	IID	NON-IID	
FEDSGD	1	∞	1	626	483	
FEDAVG	5	∞	5	179 (3.5 \times)	1000 (0.5 \times)	
FEDAVG	1	50	12	65 (9.6 \times)	600 (0.8 \times)	
FEDAVG	20	∞	20	234 (2.7 \times)	672 (0.7 \times)	
FEDAVG	1	10	60	34 (18.4 \times)	350 (1.4 \times)	
FEDAVG	5	50	60	29 (21.6 \times)	334 (1.4 \times)	
FEDAVG	20	50	240	32 (19.6 \times)	426 (1.1 \times)	
FEDAVG	5	10	300	20 (31.3 \times)	229 (2.1 \times)	
FEDAVG	20	10	1200	18 (34.8 \times)	173 (2.8 \times)	

SHAKESPEARE LSTM, 54% ACCURACY						
LSTM	E	B	u	IID	NON-IID	
FEDSGD	1	∞	1.0	2488	3906	
FEDAVG	1	50	1.5	1635 (1.5 \times)	549 (7.1 \times)	
FEDAVG	5	∞	5.0	613 (4.1 \times)	597 (6.5 \times)	
FEDAVG	1	10	7.4	460 (5.4 \times)	164 (23.8 \times)	
FEDAVG	5	50	7.4	401 (6.2 \times)	152 (25.7 \times)	
FEDAVG	5	10	37.1	192 (13.0 \times)	41 (95.3 \times)	

FedSGD: $C = 1$ and Full-Batch Optimization

FedAvg: $C = 0.1$

C: Proportion of clients participating in calculation

E: Number of training cycles per client between two communications

B: Mini-batch size of each client. ∞ indicates full-batch

$u_k = E \frac{n_k}{B}$: The total number of updates in each iteration on client k

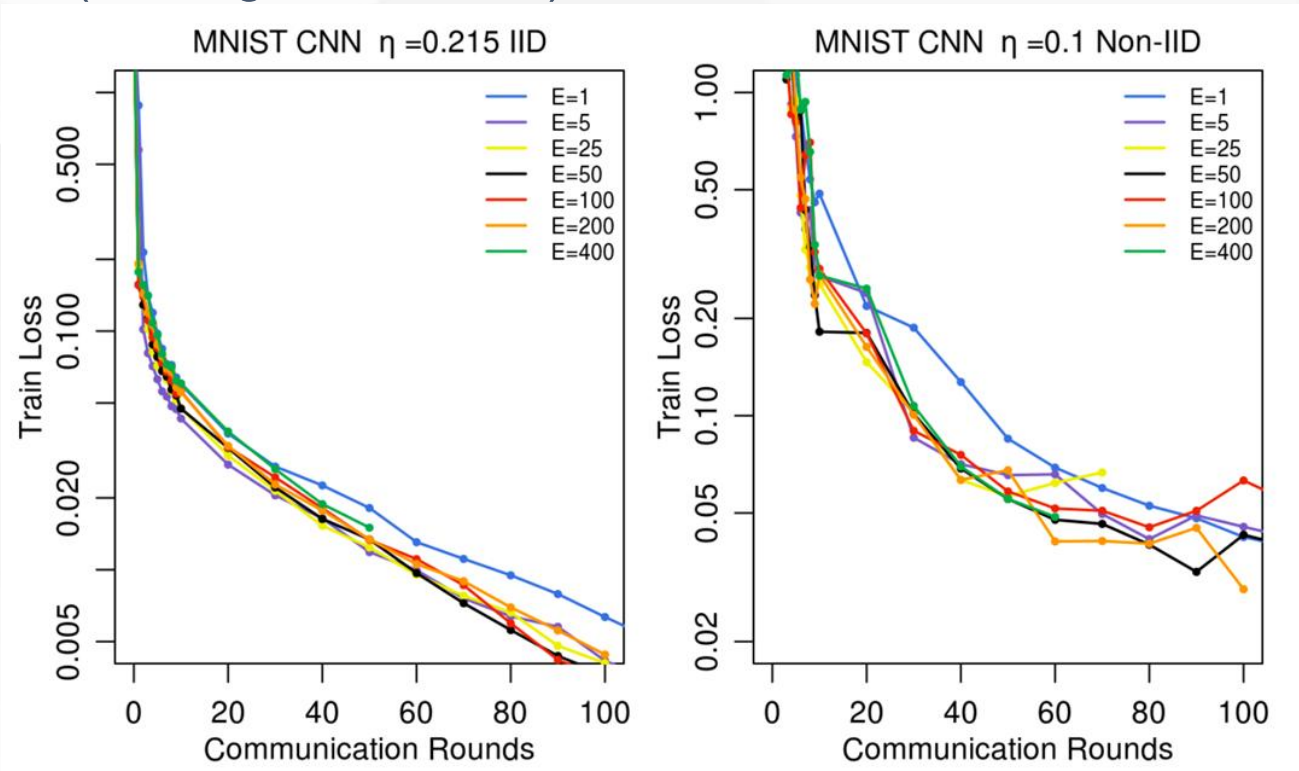




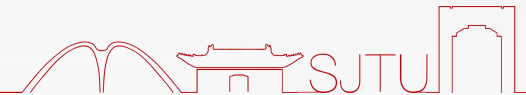
FedAvg algorithm



Experiments (training loss curves)



- C:** Proportion of clients participating in calculation, $C=0.1$
 - E:** Number of training cycles per client between two communications
 - B:** Mini-batch size of each client. $B=10$.
- $u_k = E \frac{n_k}{B}$: The total number of updates in each iteration on client k

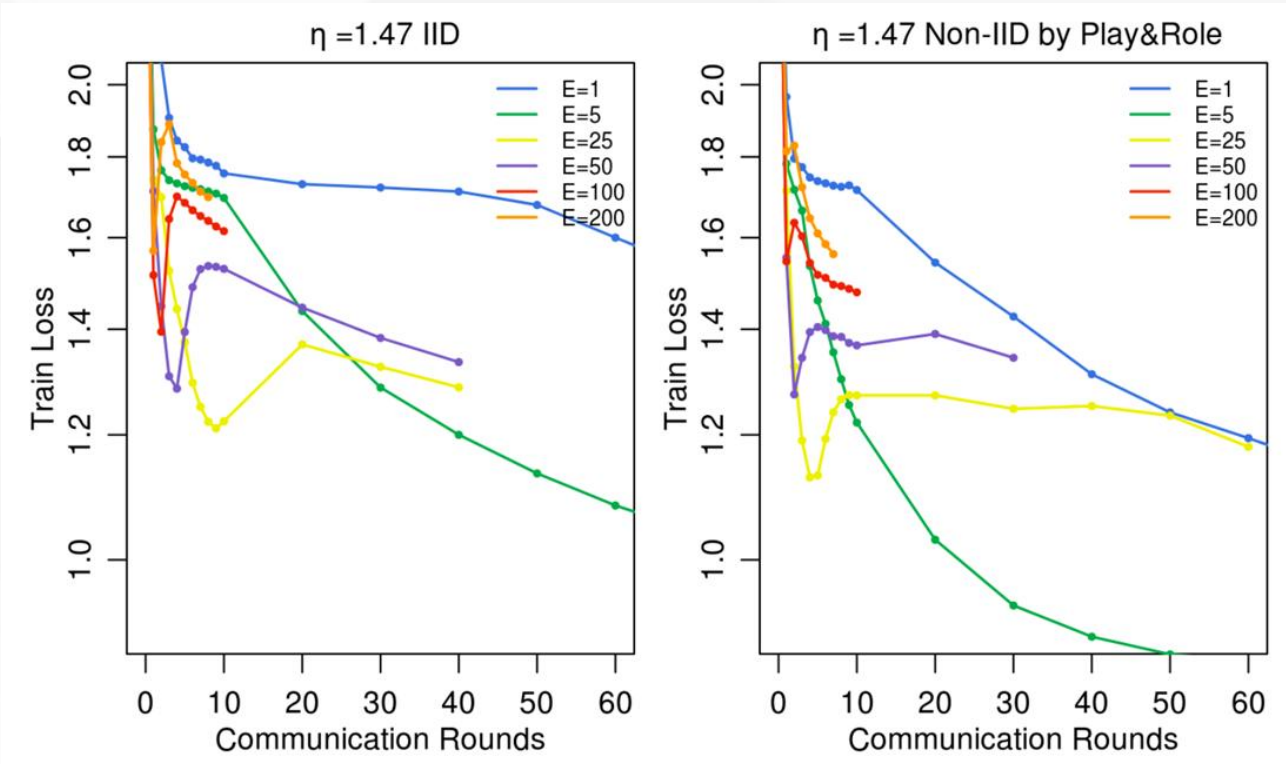




FedAvg algorithm



Experiments (training loss curves)



LSTM

C: Proportion of clients participating in calculation, **C=0.1**

E: Number of training cycles per client between two communications

B: Mini-batch size of each client. **B=10.**

$u_k = E \frac{n_k}{B}$: The total number of updates in each iteration on client k



02

Benchmark algorithms

- MOON
- FedDyn
- KT-pFL
- FedMA

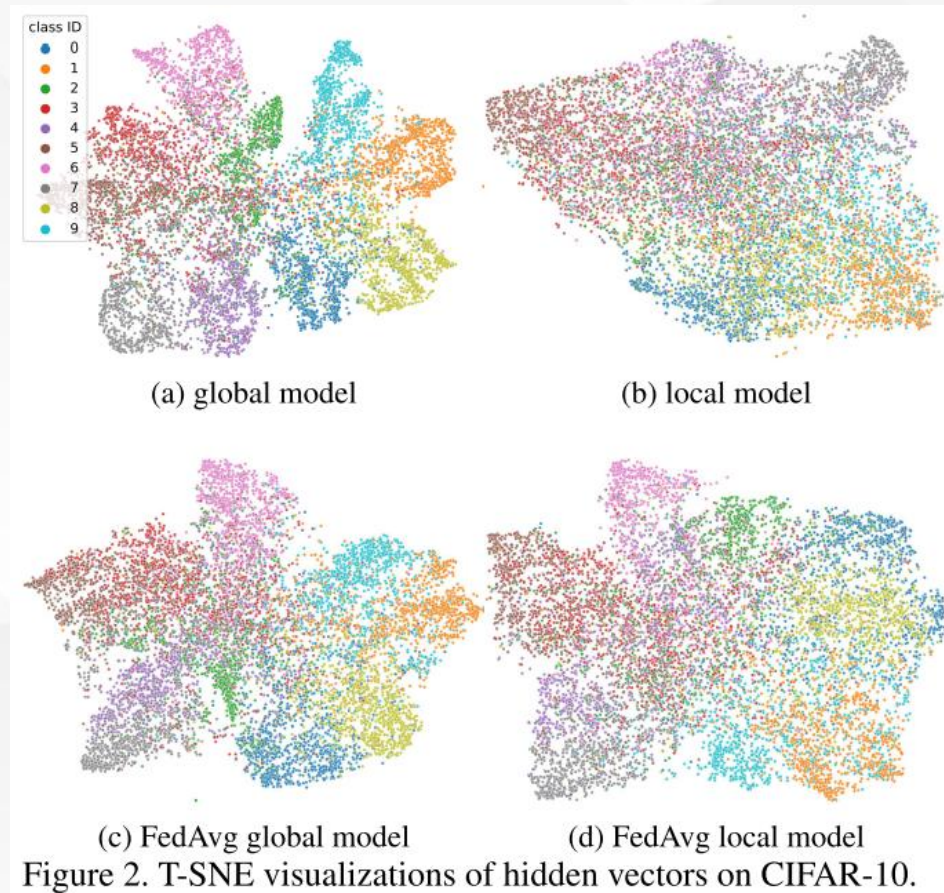


MOON



Observation

- The global model trained on a whole dataset can learn a better representation than the local model trained on a skewed subset.





Observation

- The global model trained on a whole dataset can learn a better representation than the local model trained on a skewed subset.
- Propose model-contrastive learning (MOON), which corrects the local updates by maximizing the agreement of representation learned by the current local model and the representation learned by the global model.



Contrastive learning

- The key idea of contrastive learning is to reduce the distance between the representations of different augmented views of the same image (i.e., positive pairs), and increase the distance between the representations of augmented views of different images (i.e., negative pairs)

$$l_{i,j} = -\log \frac{\exp(\text{sim}(x_i, x_j)/\tau)}{\sum_{k=1}^{2N} \mathbb{I}_{[k \neq i]} \exp(\text{sim}(x_i, x_k)/\tau)}$$



MOON

- Global representation $z_{glob} = R_{w^t}(x)$
- Local representation w_i^{t-1} (i.e., $z_{prev} = R_{w_i^{t-1}}(x)$).
- Current representation $z = R_{w_i^t}(x)$
-

$$\ell_{con} = -\log \frac{\exp(\text{sim}(z, z_{glob})/\tau)}{\exp(\text{sim}(z, z_{glob})/\tau) + \exp(\text{sim}(z, z_{prev})/\tau)}$$

$$\ell = \ell_{sup}(w_i^t; (x, y)) + \mu \ell_{con}(w_i^t; w_i^{t-1}; w_i^t; x),$$

$$\min_{w_i^t} \mathbb{E}_{(x, y) \sim D^i} [\ell_{sup}(w_i^t; (x, y)) + \mu \ell_{con}(w_i^t; w_i^{t-1}; w_i^t; x)].$$



MOON

Algorithm 1: The MOON framework

Input: number of communication rounds T ,
 number of parties N , number of local
 epochs E , temperature τ , learning rate η ,
 hyper-parameter μ

Output: The final model w^T

1 **Server executes:**

2 initialize w^0

3 **for** $t = 0, 1, \dots, T - 1$ **do**

4 **for** $i = 1, 2, \dots, N$ **in parallel do**

5 send the global model w^t to P_i

6 $w_i^t \leftarrow \mathbf{PartyLocalTraining}(i, w^t)$

7 $w^{t+1} \leftarrow \sum_{k=1}^N \frac{|\mathcal{D}^k|}{|\mathcal{D}|} w_k^t$

8 **return** w^T

9 **PartyLocalTraining**(i, w^t):

10 $w_i^t \leftarrow w^t$

11 **for** epoch $i = 1, 2, \dots, E$ **do**

12 **for** each batch $\mathbf{b} = \{x, y\}$ of \mathcal{D}^i **do**

13 $\ell_{sup} \leftarrow \text{CrossEntropyLoss}(F_{w_i^t}(x), y)$

14 $z \leftarrow R_{w_i^t}(x)$

15 $z_{glob} \leftarrow R_{w^t}(x)$

16 $z_{prev} \leftarrow R_{w_i^{t-1}}(x)$

17 $\ell_{con} \leftarrow$

$-\log \frac{\exp(\text{sim}(z, z_{glob})/\tau)}{\exp(\text{sim}(z, z_{glob})/\tau) + \exp(\text{sim}(z, z_{prev})/\tau)}$

18 $\ell \leftarrow \ell_{sup} + \mu \ell_{con}$

19 $w_i^t \leftarrow w_i^t - \eta \nabla \ell$

20 **return** w_i^t to server





FedDyn



Intuition

- Training models on local data that minimize local empirical loss appears to be meaningful, but yet, doing so is fundamentally inconsistent with minimizing the global empirical loss.
- Dynamically modify the device objective with a penalty term so that, in the limit, when model parameters converge, they do so to stationary points of the global empirical loss.



FedDyn

Algorithm 1: Federated Dynamic Regularizer - (FedDyn)

Input: $T, \theta^0, \alpha > 0, \nabla L_k(\theta_k^0) = \mathbf{0}$.

for $t = 1, 2, \dots, T$ **do**

Sample devices $\mathcal{P}_t \subseteq [m]$ and transmit θ^{t-1} to each selected device,

for each device $k \in \mathcal{P}_t$, **and in parallel do**

Set $\theta_k^t = \underset{\theta}{\operatorname{argmin}} L_k(\theta) - \langle \nabla L_k(\theta_k^{t-1}), \theta \rangle + \frac{\alpha}{2} \|\theta - \theta^{t-1}\|^2$,

Set $\nabla L_k(\theta_k^t) = \nabla L_k(\theta_k^{t-1}) - \alpha (\theta_k^t - \theta^{t-1})$,

Transmit device model θ_k^t to server,

end for

for each device $k \notin \mathcal{P}_t$, **and in parallel do**

Set $\theta_k^t = \theta_k^{t-1}, \nabla L_k(\theta_k^t) = \nabla L_k(\theta_k^{t-1})$,

end for

Set $\mathbf{h}^t = \mathbf{h}^{t-1} - \alpha \frac{1}{m} (\sum_{k \in \mathcal{P}_t} \theta_k^t - \theta^{t-1})$,

Set $\theta^t = \left(\frac{1}{|\mathcal{P}_t|} \sum_{k \in \mathcal{P}_t} \theta_k^t \right) - \frac{1}{\alpha} \mathbf{h}^t$

end for



Analysis

$$\boldsymbol{\theta}_k^t = \operatorname{argmin}_{\boldsymbol{\theta}} \left[\mathfrak{R}_k(\boldsymbol{\theta}; \boldsymbol{\theta}_k^{t-1}, \boldsymbol{\theta}^{t-1}) \triangleq L_k(\boldsymbol{\theta}) - \langle \nabla L_k(\boldsymbol{\theta}_k^{t-1}), \boldsymbol{\theta} \rangle + \frac{\alpha}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^{t-1}\|^2 \right]$$

- The first order condition

$$\nabla L_k(\boldsymbol{\theta}_k^t) - \nabla L_k(\boldsymbol{\theta}_k^{t-1}) + \alpha(\boldsymbol{\theta}_k^t - \boldsymbol{\theta}^{t-1}) = \mathbf{0}$$

- If local device models converge, they converge to the server model, and the convergence point is a stationary point of the global loss.

if $\boldsymbol{\theta}_k^t \rightarrow \boldsymbol{\theta}_k^\infty$, it generally follows that, $\nabla L_k(\boldsymbol{\theta}_k^t) \rightarrow \nabla L_k(\boldsymbol{\theta}_k^\infty)$, and as a consequence, we have $\boldsymbol{\theta}^t \rightarrow \boldsymbol{\theta}_k^\infty$. In turn this implies that $\boldsymbol{\theta}_k^\infty \rightarrow \boldsymbol{\theta}^\infty$, i.e., is independent of k .



KT-pFL

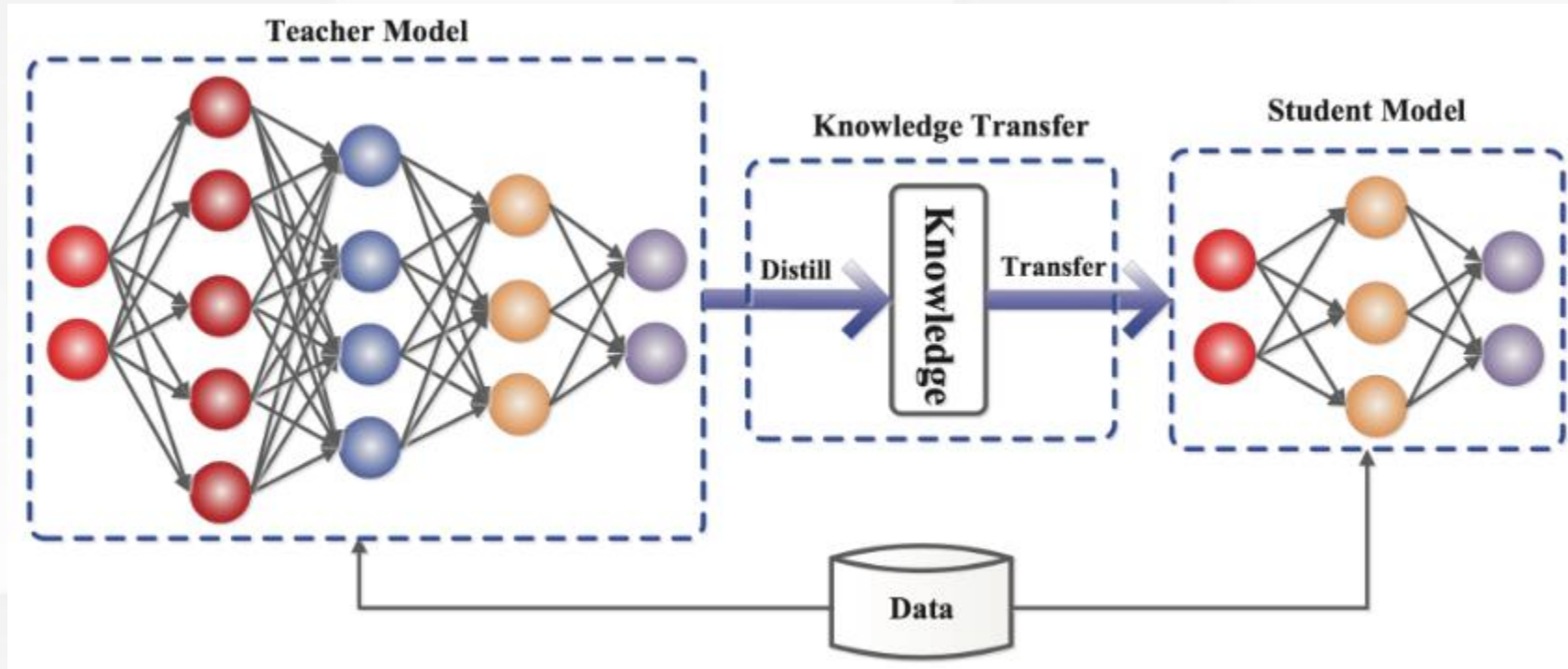


Intuition

- Main idea is to allow each client to maintain a personalized soft prediction at the server that can be updated by a linear combination of all clients' local soft predictions using a knowledge coefficient matrix.
- Regardless of model structures

Knowledge Distillation (KD)

- Transfer knowledge from well-learned teacher model to student model



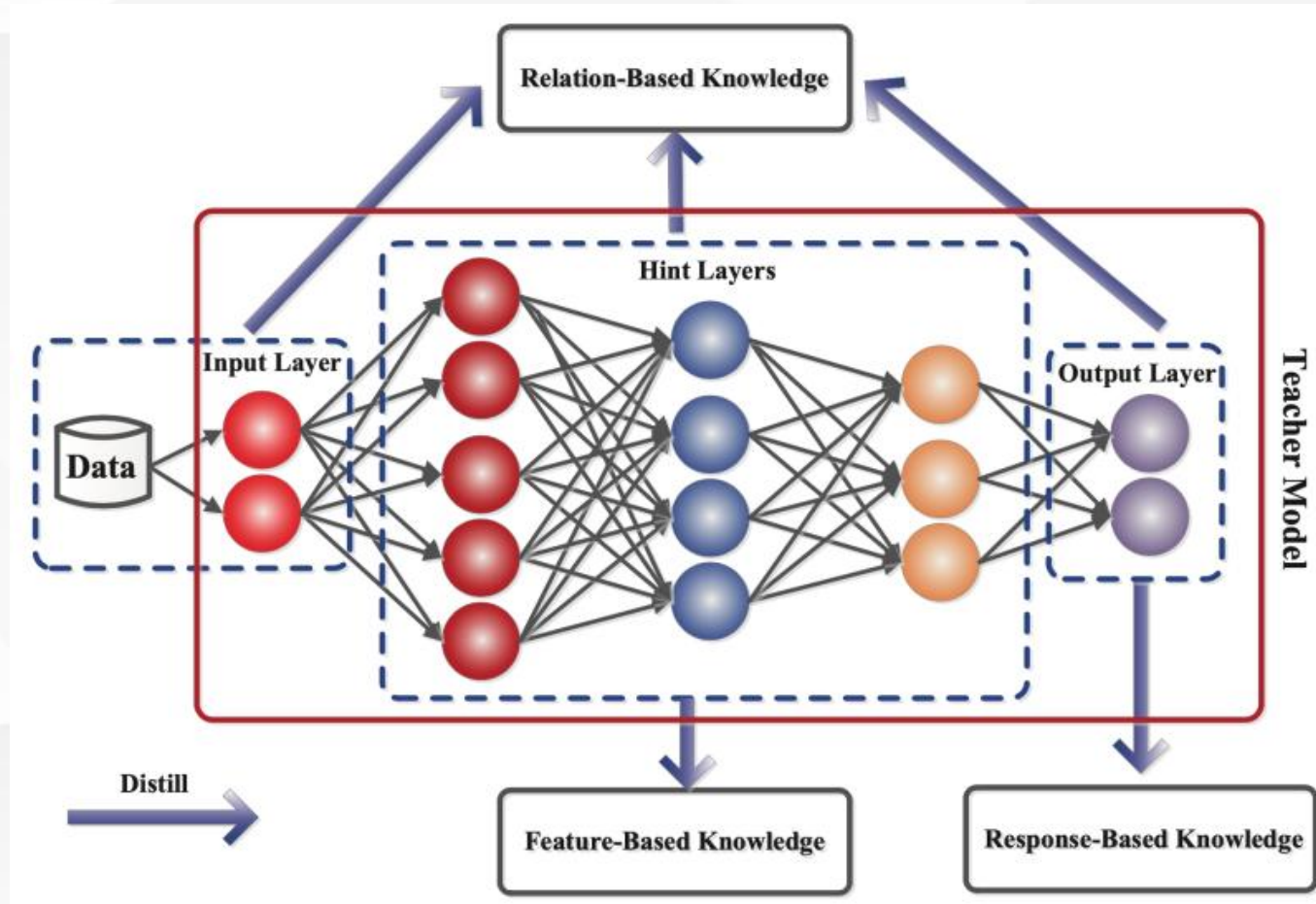


KT-pFL algorithm



Knowledge Distillation (KD)

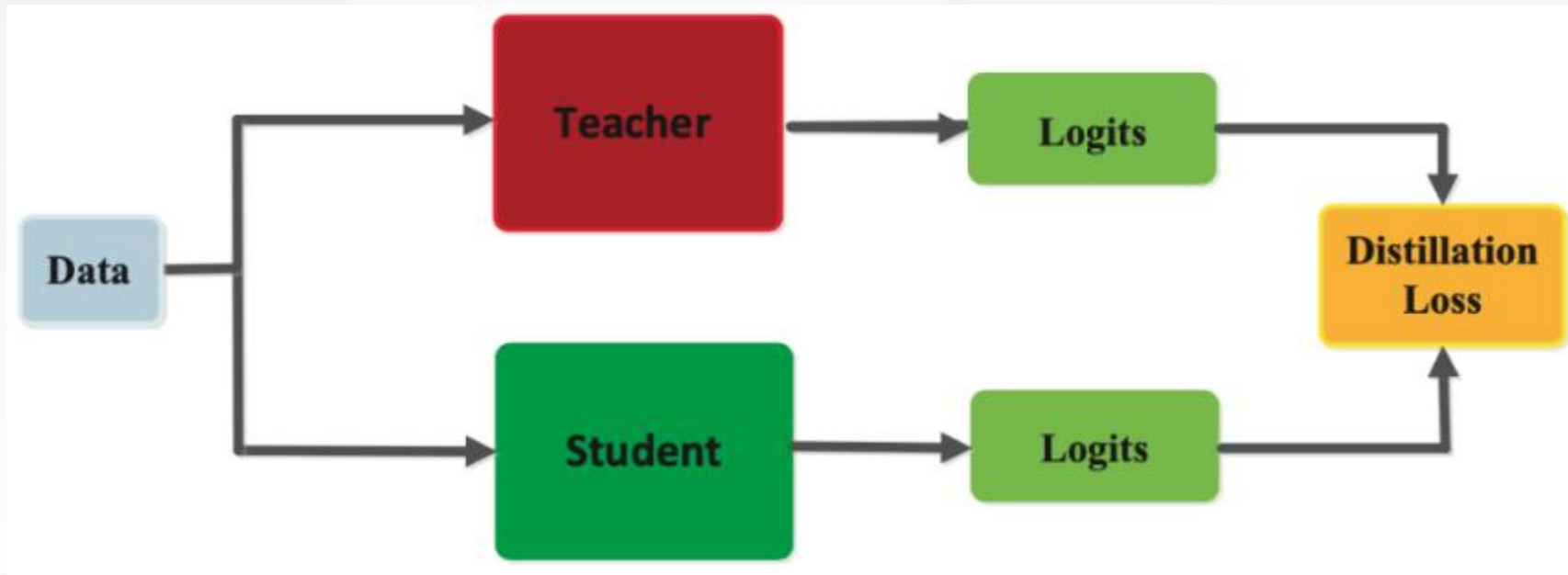
- Classification





Knowledge Distillation (KD)

- Response-based KD





KT-pFL

- Objective

$$\min_{\mathbf{w}} \mathcal{L}(\mathbf{w}) := \sum_{n=1}^N \frac{D_n}{D} \mathcal{L}_n(\mathbf{w}), \text{ where } \mathcal{L}_n(\mathbf{w}) = \frac{1}{D_n} \sum_{i=1}^{D_n} \mathcal{L}_{CE}(\mathbf{w}; x_i, y_i).$$

$$\min_{\mathbf{w}^1, \dots, \mathbf{w}^N} \mathcal{L}(\mathbf{w}^1, \dots, \mathbf{w}^N) := \sum_{n=1}^N \frac{D_n}{D} L_n(\mathbf{w}^n)$$



KT-pFL

- Personalized loss function
 - Kullback–Leibler (KL) Divergence
 - c_{mn} is the knowledge coefficient which is used to estimate the contribution from client m to n .
 - $s(\mathbf{w}^n, \hat{x})$ can be deemed to be a soft prediction of the client n

$$\mathcal{L}_{per,n}(\mathbf{w}^n) := \mathcal{L}_n(\mathbf{w}^n) + \lambda \sum_{\hat{x} \in \mathbb{D}_r} \mathcal{L}_{KL} \left(\sum_{m=1}^N c_{mn} \cdot s(\mathbf{w}^m, \hat{x}), s(\mathbf{w}^n, \hat{x}) \right)$$

$$s(\mathbf{w}^n, \hat{x}) = \frac{\exp(z_c^n / T)}{\sum_{c=1}^C \exp(z_c^n / T)},$$



KT-pFL

- Knowledge coefficient matrix

$$\mathbf{c} = \begin{Bmatrix} c_{11} & c_{12} & \cdots & c_{1N} \\ c_{21} & c_{22} & \cdots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{N1} & c_{N2} & \cdots & c_{NN} \end{Bmatrix}.$$



KT-pFL

- Objective

$$\min_{\mathbf{w}, \mathbf{c}} \mathcal{L}(\mathbf{w}, \mathbf{c}) := \sum_{n=1}^N \frac{D_n}{D} \mathcal{L}_{per,n}(\mathbf{w}^n) + \rho \left\| \mathbf{c} - \frac{\mathbf{1}}{N} \right\|^2.$$

$$\mathbf{w} = [\mathbf{w}^1, \dots, \mathbf{w}^N] \in \mathbb{R}^{\sum_{n=1}^N d_n}$$



KT-pFL

- Training

- Update w

- Local Training
- Distillation

- Update c

$$\mathbf{w}^n \leftarrow \mathbf{w}^n - \eta_1 \nabla_{\mathbf{w}^n} \mathcal{L}_n(\mathbf{w}^n; \xi_n),$$

$$\mathbf{w}^n \leftarrow \mathbf{w}^n - \eta_2 \nabla_{\mathbf{w}^n} \mathcal{L}_{KL} \left(\sum_{m=1}^N \mathbf{c}_m^{*,T} \cdot s(\mathbf{w}^m, \xi_r), s(\mathbf{w}^n, \xi_r) \right)$$

$$\mathbf{c} \leftarrow \mathbf{c} - \eta_3 \lambda \sum_{n=1}^N \frac{D_n}{D} \nabla_{\mathbf{c}} \mathcal{L}_{KL} \left(\sum_{m=1}^N \mathbf{c}_m \cdot s(\mathbf{w}^{m,*}, \xi_r), s(\mathbf{w}^{n,*}, \xi_r) \right) - 2\eta_3 \rho \left(\mathbf{c} - \frac{\mathbf{1}}{N} \right)$$



KT-pFL

Algorithm 1 KT-pFL Algorithm

Input: $\mathbb{D}, \mathbb{D}_r, \eta_1, \eta_2, \eta_3$ and T

Output: $\mathbf{w} = [\mathbf{w}^1, \dots, \mathbf{w}^N]$

- 1: Initialize \mathbf{w}_0 and \mathbf{c}_0
- 2: **procedure** SERVER-SIDE OPTIMIZATION
- 3: Distribute \mathbf{w}_0 and \mathbf{c}_0 to each client
- 4: **for** each communication round $t \in \{1, 2, \dots, T\}$ **do**
- 5: **for** each client n **in parallel do**
- 6: $\mathbf{w}_{t+1}^n \leftarrow \text{ClientLocalUpdate}(n, \mathbf{w}_t^n, \mathbf{c}_{t,n})$
- 7: Update knowledge coefficient matrix \mathbf{c} via (7)
- 8: Distribute \mathbf{c}_{t+1} to all clients
- 9: **procedure** CLIENTLOCALUPDATE($n, \mathbf{w}_t^n, \mathbf{c}_{t,n}$)
- 10: Client n receives \mathbf{w}_t^n and \mathbf{c}_n from the server
- 11: **for** each local epoch i from 1 to E **do**
- 12: **for** mini-batch $\xi_t \subseteq \mathbb{D}_n$ **do**
- 13: **Local Training:** update model parameters on private data via (5)
- 14: **for** each distillation step j from 1 to R **do**
- 15: **for** mini-batch $\xi_{r,t} \subseteq \mathbb{D}_r$ **do**
- 16: **Distillation:** update model parameters on public data via (6)
- return** local parameters \mathbf{w}_{t+1}^n

KT-pFL

- Illustration

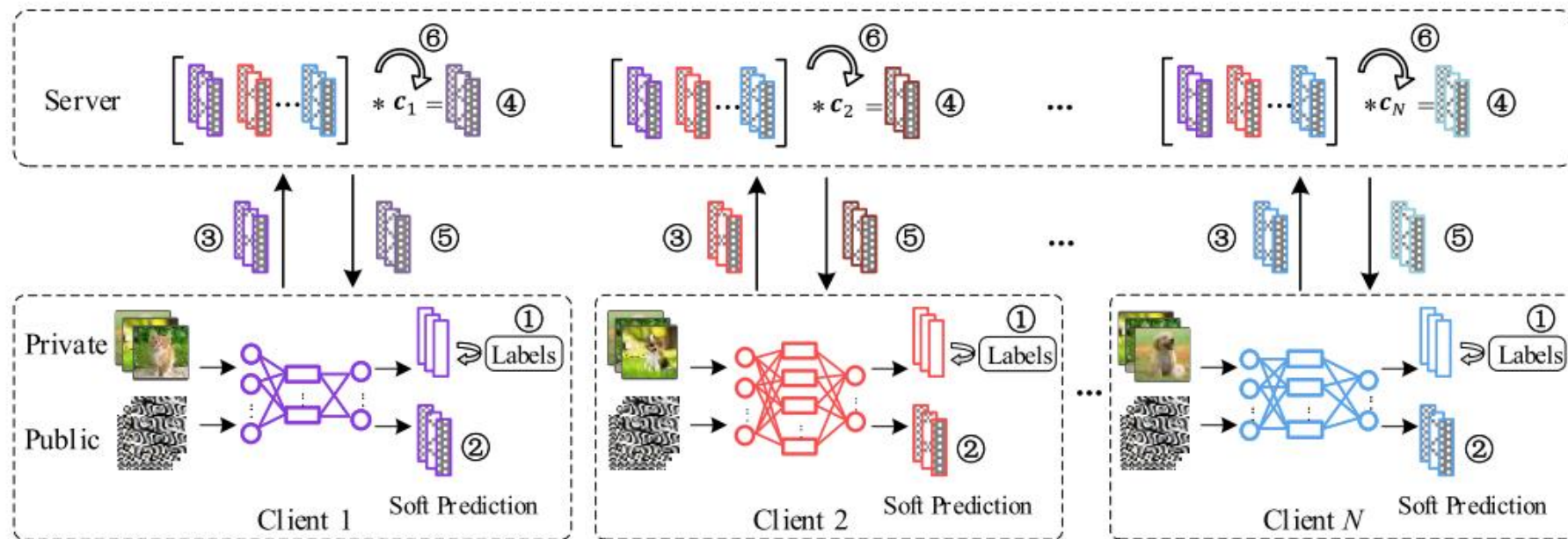


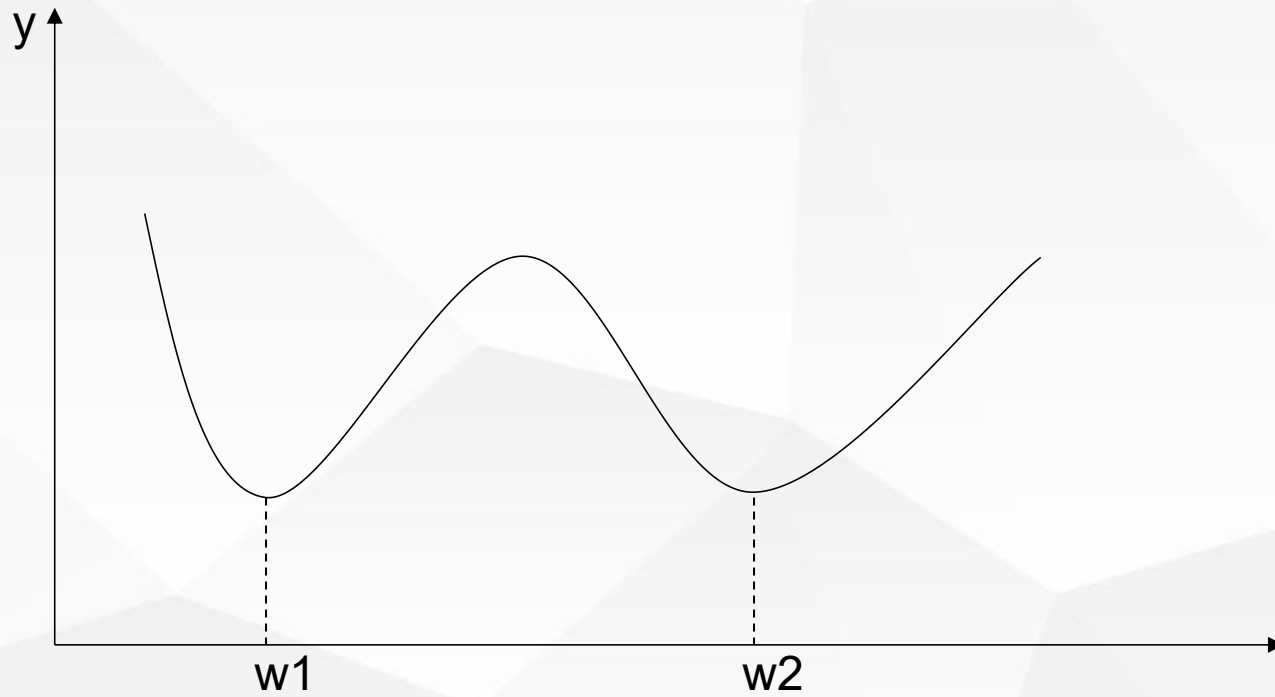
Figure 1: Illustration of the KT-pFL framework. The workflow includes 6 steps: ① local training on private data; ②, ③ each client outputs the local soft prediction on public data and sends it to the server; ④ the server calculates each client's personalized soft prediction via a linear combination of local soft predictions and knowledge coefficient matrix; ⑤ each client downloads the personalized soft prediction to perform distillation phase; ⑥ the server updates the knowledge coefficient matrix.



FedMA



Permutation invariance

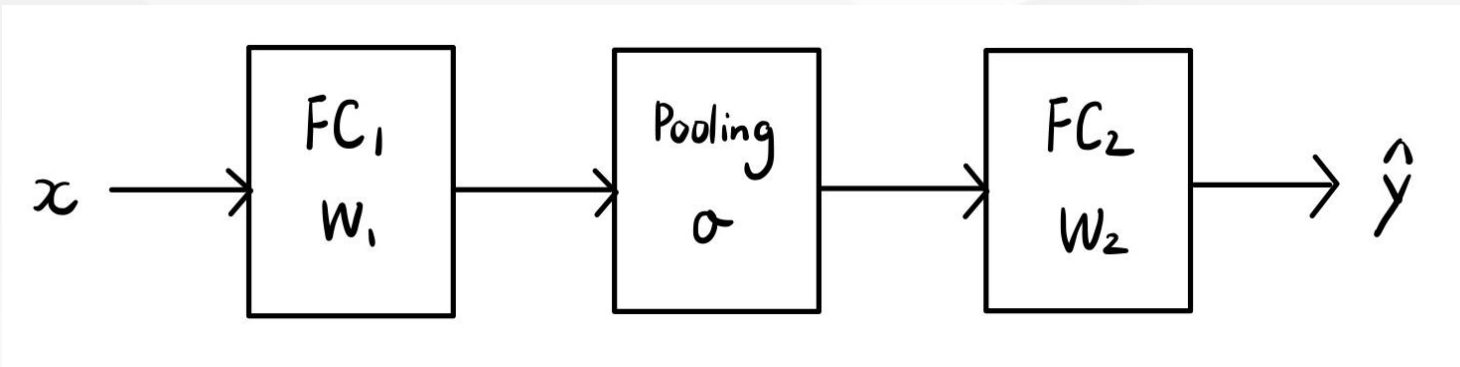




FedMA algorithm



⊙ Permutation invariance (fully-connected (FC) layer)



$$\hat{y} = \sigma(xW_1)W_2$$

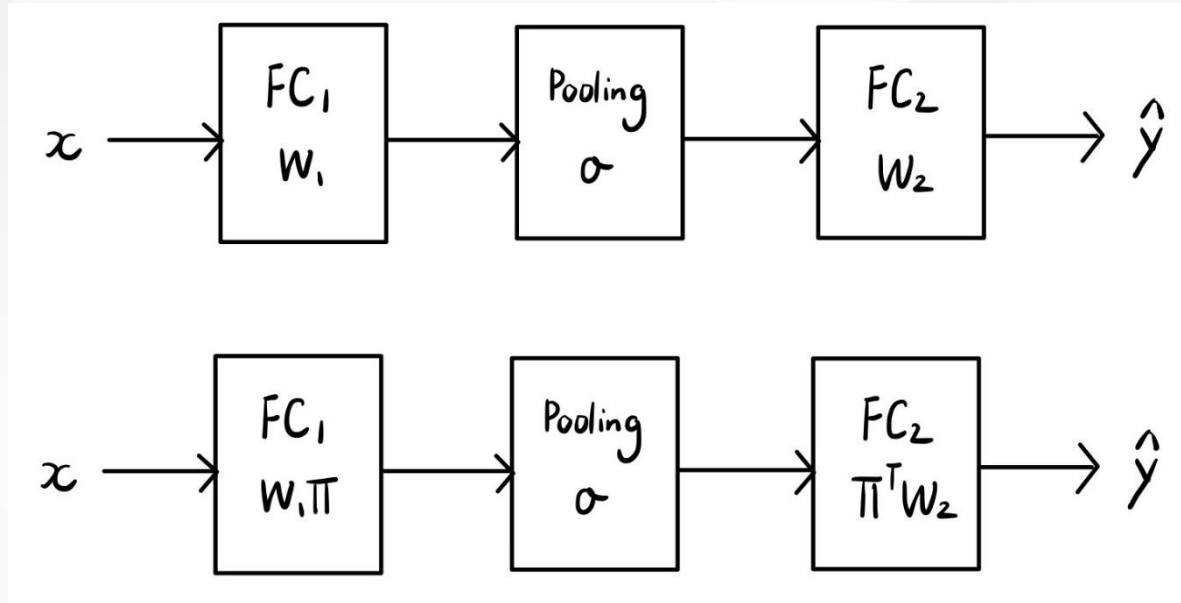




FedMA algorithm



④ Permutation invariance (fully-connected (FC) layer)



$$\hat{y} = \sigma(xW_1\Pi)\Pi^T W_2, \text{ where } \Pi \text{ is any } L \times L \text{ permutation matrix.}$$



FedMA algorithm



⊗ Permutation invariance (fully-connected (FC) layer)

Client A:

$$\{W_1 \Pi_j, \Pi_j^T W_2\}$$

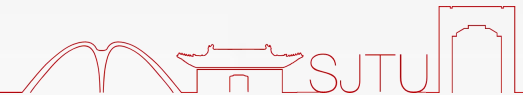
Client B:

$$\{W_1 \Pi_{j'}, \Pi_{j'}^T W_2\}$$

Aggregation

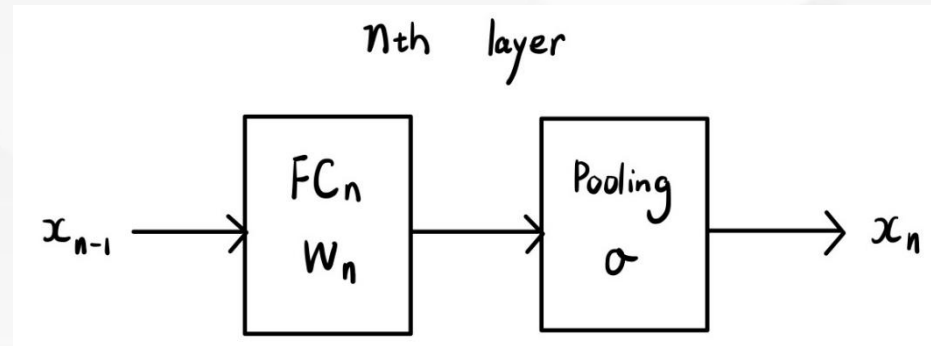
$$(W_1 \Pi_j + W_1 \Pi_{j'}) / 2 \neq W_1 \Pi \text{ for any } \Pi$$

$$\text{Solution: } (W_1 \Pi_j \Pi_j^T + W_1 \Pi_{j'} \Pi_{j'}^T) / 2 = W_1$$





Permutation invariance (FCs)



Simple FCs: $\hat{y} = \sigma(xW_1\Pi)\Pi^T W_2$

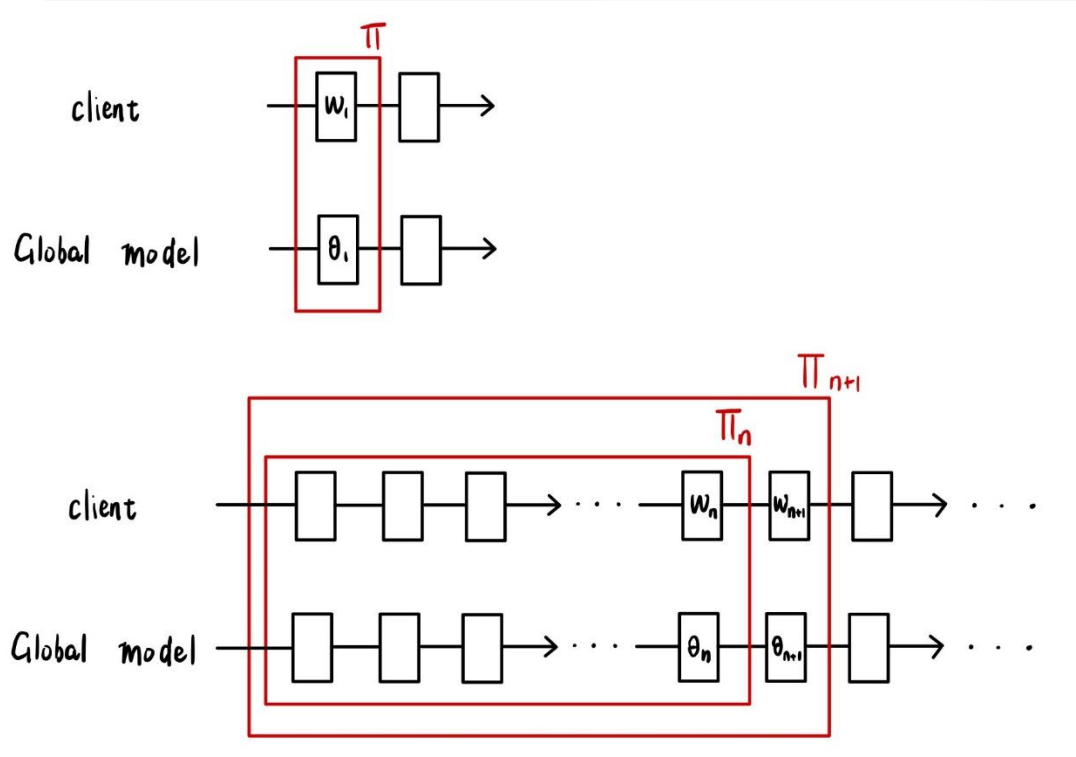
Deep FCs: $x_n = \sigma(x_{n-1}\Pi_{n-1}^T W_n \Pi_n)$



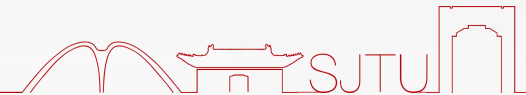
FedMA algorithm



Permutation invariance (FCs)



$$x_n = \sigma(x_{n-1} \Pi_{n-1}^T W_n \Pi_n)$$





Permutation invariance (CNN)

FC : $x_n = \sigma(x_{n-1} \Pi_{n-1}^T W_n \Pi_n)$

CNN : $x_n = \sigma(\text{Conv}(x_{n-1}, \Pi_{n-1}^T W_n \Pi_n))$



Permutation invariance (recall)

Client A:

$$\{W_1 \Pi_j, \Pi_j^T W_2\}$$

Client B:

$$\{W_1 \Pi_{j'}, \Pi_{j'}^T W_2\}$$

Aggregation

$$(W_1 \Pi_j + W_1 \Pi_{j'}) / 2 \neq W_1 \Pi \text{ for any } \Pi$$

Solution: $(W_1 \Pi_j \Pi_j^T + W_1 \Pi_{j'} \Pi_{j'}^T) / 2 = W_1$



FedMA algorithm



Matched averaging formulation

Client Model (trained on dataset j)

Global Model

$$\text{layer } W_j \xleftarrow{W_j \Pi = \theta} \text{layer } \theta$$

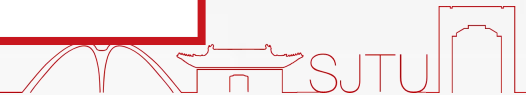
$$[w_{j1} | w_{j2} | w_{j3}]$$

$$[\theta_1 | \theta_2 | \theta_3]$$

which ?

$$\min_{i \in \{1,2,3\}} C(w_{j1}, \theta_i)$$

$$\min_{\{\pi_{li}^j\}} \sum_{i=1}^L \sum_{j,l} \min_{\theta_i} \pi_{li}^j c(w_{jl}, \theta_i) \text{ s.t. } \sum_i \pi_{li}^j = 1 \forall j, l; \sum_l \pi_{li}^j = 1 \forall i, j.$$

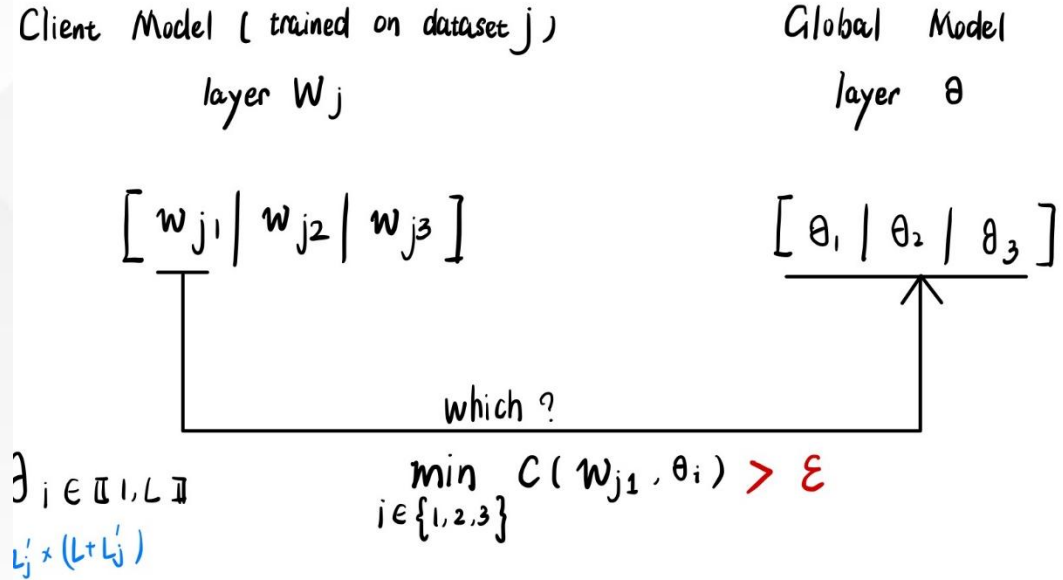




FedMA algorithm



Matched averaging formulation



After aggregation, $[\theta_1 | \theta_2 | \theta_3 | w_{j1}]$

$$\min_{\{\pi_{li}^{j'}\}_{l,i}} \sum_{i=1}^{L+L_{j'}} \sum_{j=1}^{L_{j'}} \pi_{li}^{j'} C_{li}^{j'} \text{ s.t. } \sum_i \pi_{li}^{j'} = 1 \forall l; \sum_l \pi_{li}^j \in \{0, 1\} \forall i, \text{ where}$$

$$C_{li}^{j'} = \begin{cases} c(w_{j'l}, \theta_i), & i \leq L \\ \epsilon + f(i), & L < i \leq L + L_{j'}. \end{cases}$$

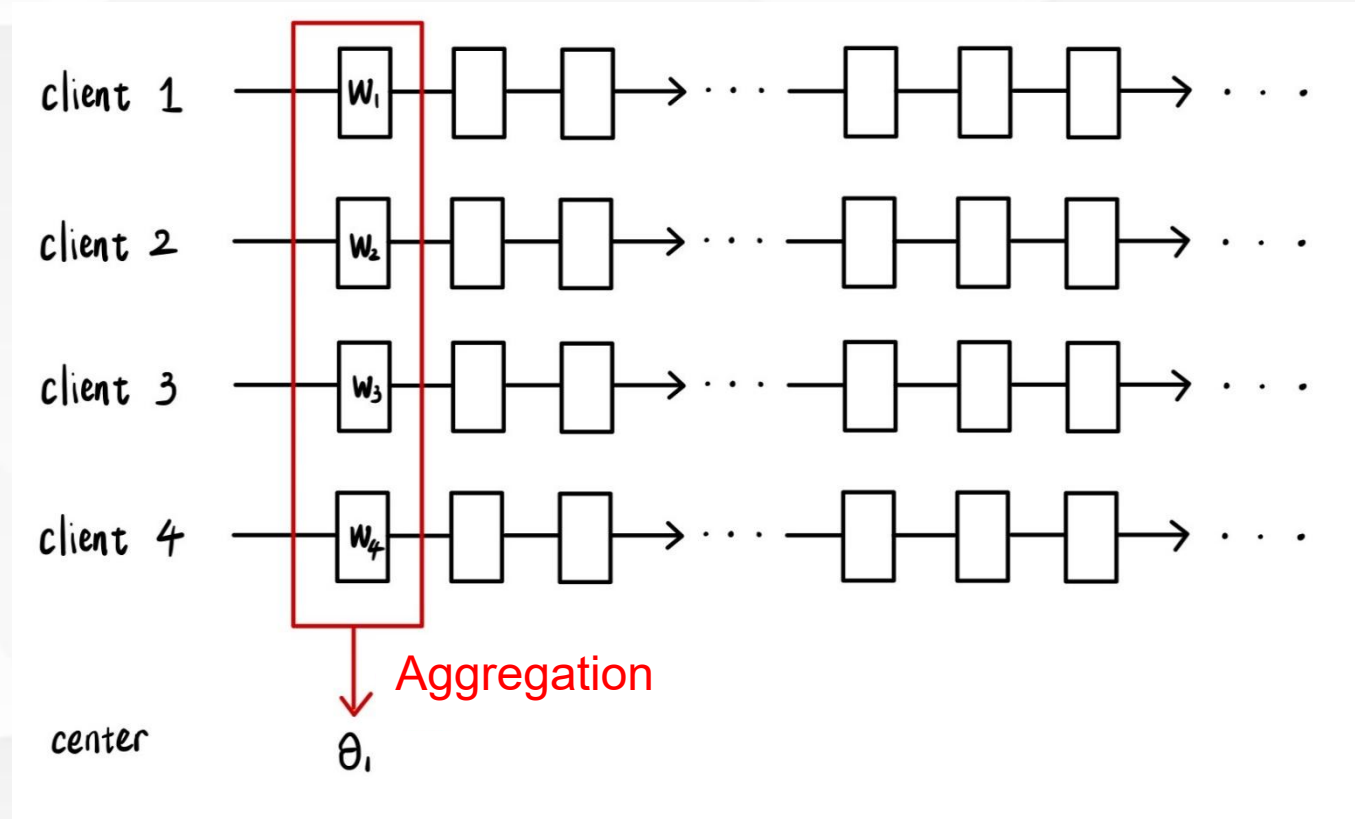




FedMA algorithm



FedMA (<https://github.com/IBM/FedMA>)

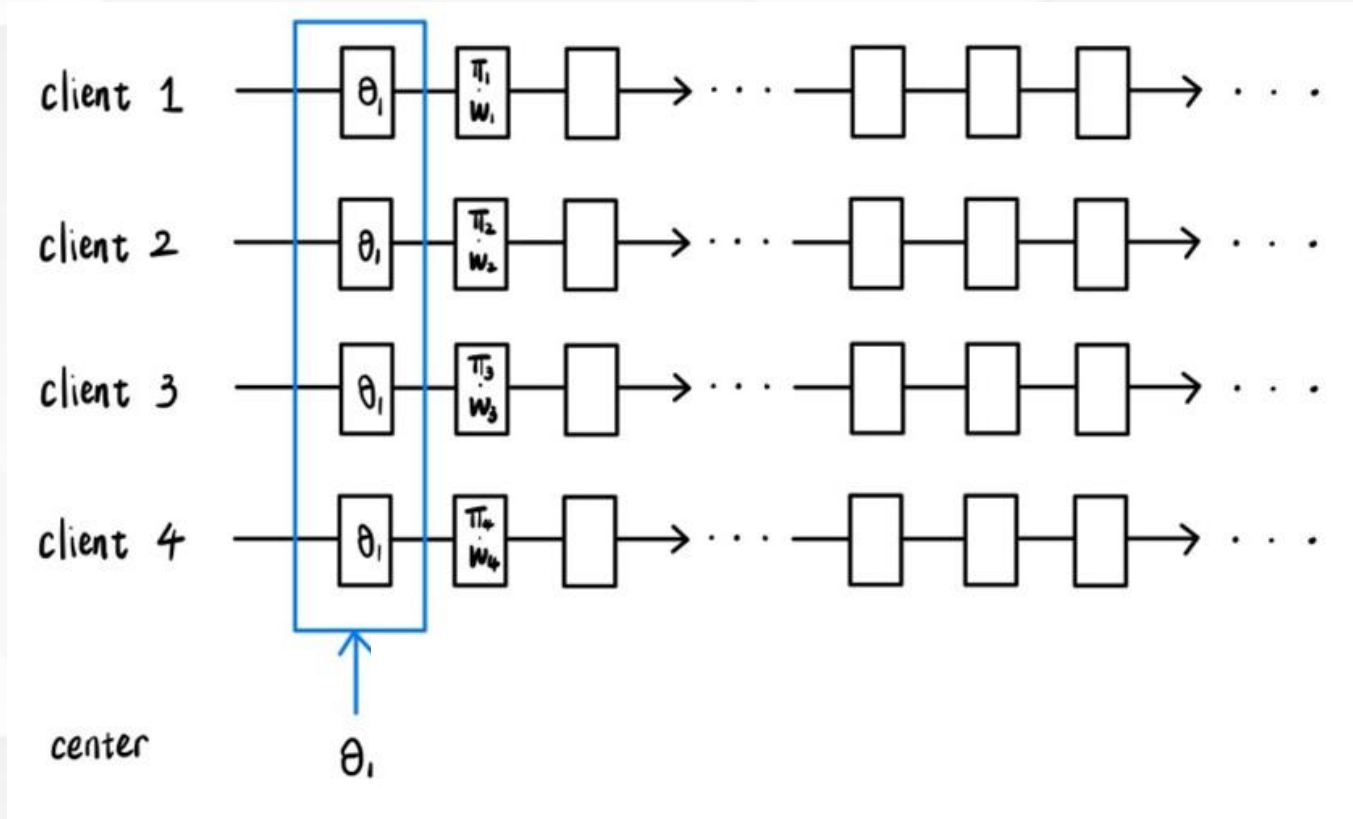




FedMA algorithm



FedMA (<https://github.com/IBM/FedMA>)

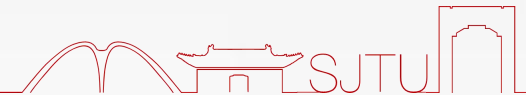
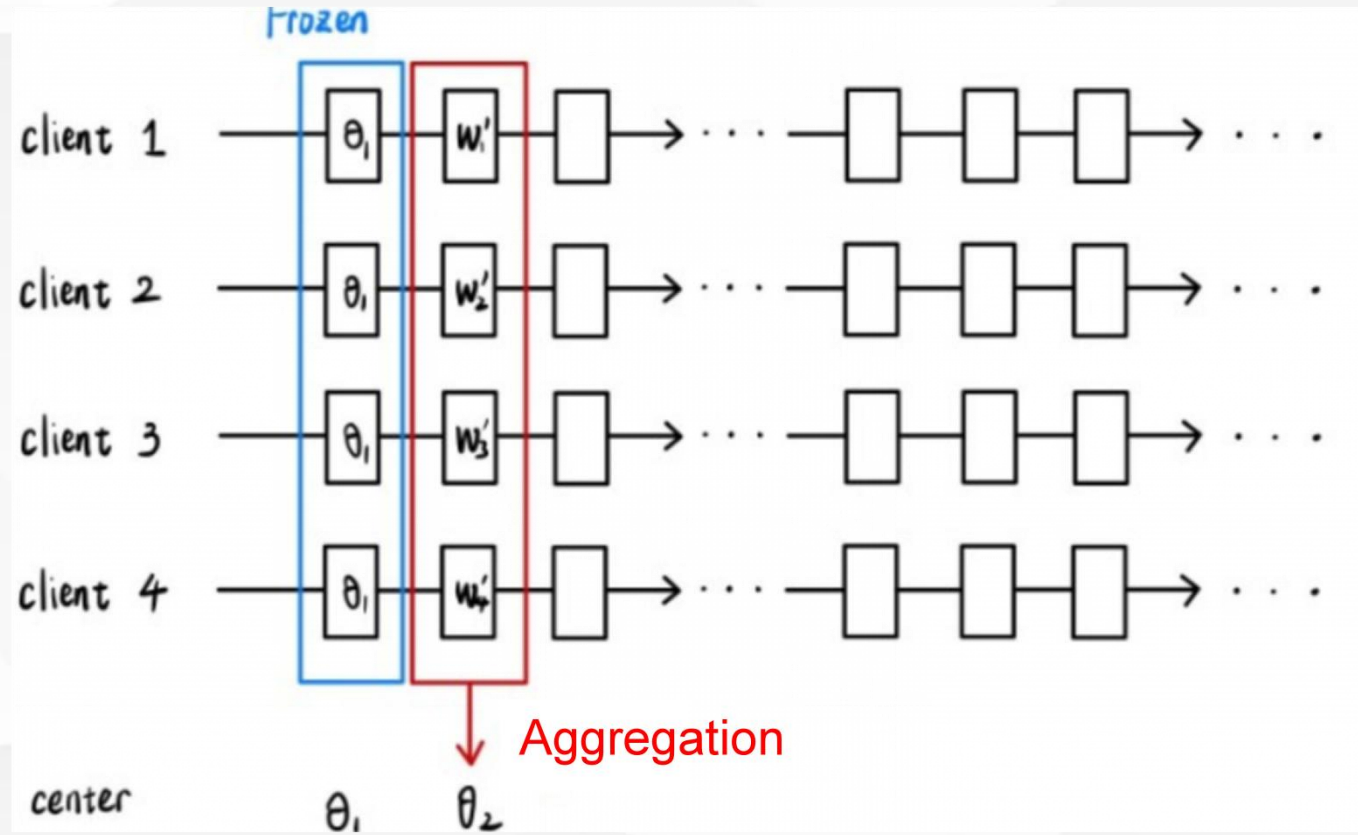




FedMA algorithm



FedMA (<https://github.com/IBM/FedMA>)

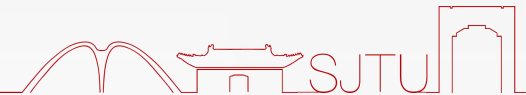
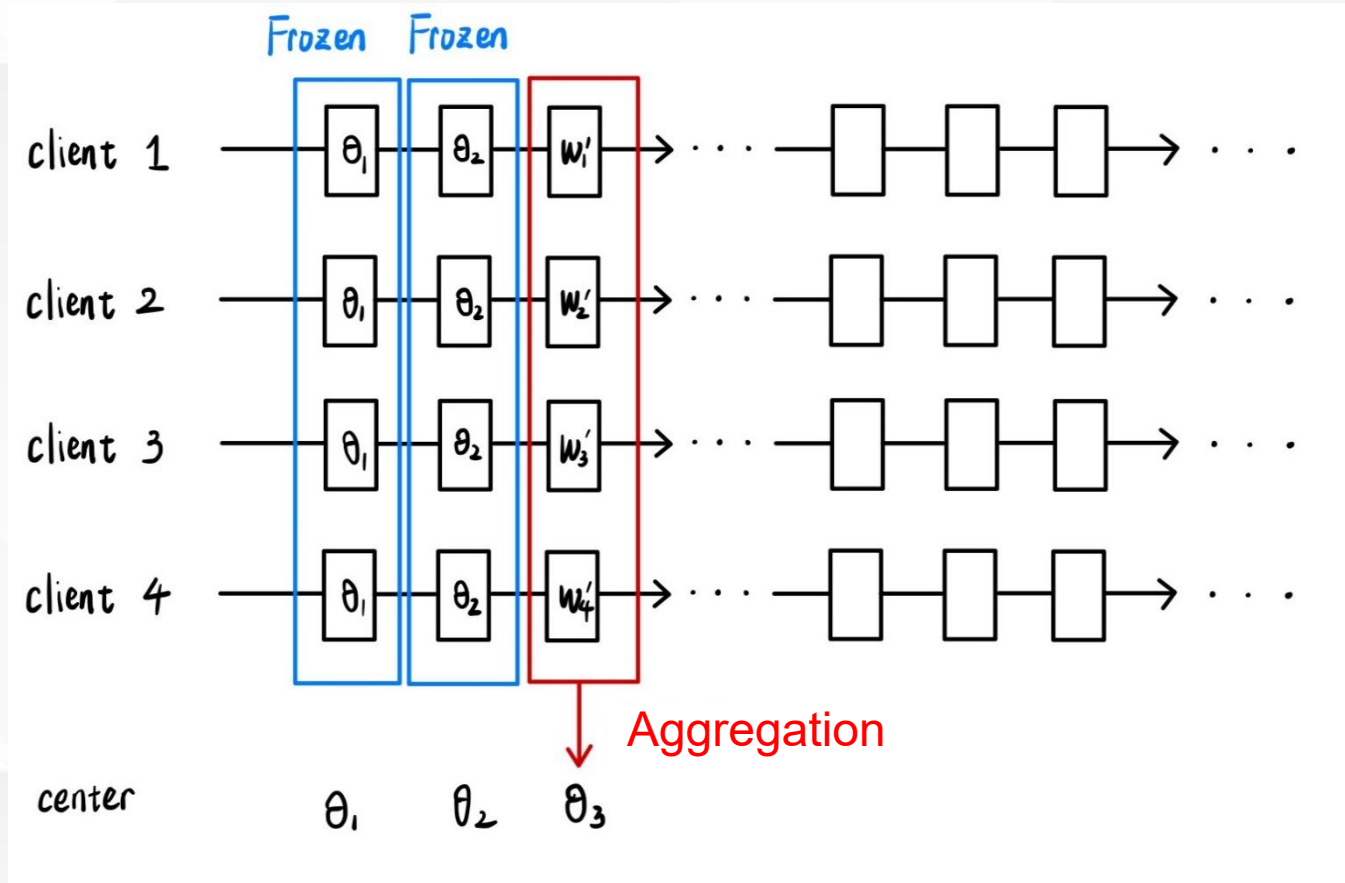




FedMA algorithm



FedMA (<https://github.com/IBM/FedMA>)





FedMA (<https://github.com/IBM/FedMA>)

Algorithm 1: Federated Matched Averaging (FedMA)

Input : local weights of N -layer architectures $\{W_{j,1}, \dots, W_{j,N}\}_{j=1}^J$ from J clients

Output: global weights $\{W_1, \dots, W_N\}$

$n = 1;$

while $n \leq N$ **do**

if $n < N$ **then**

$\{\Pi_j\}_{j=1}^J = \text{BBP-MAP}(\{W_{j,n}\}_{j=1}^J);$ // call BBP-MAP to solve Eq. 2

$W_n = \frac{1}{J} \sum_j W_{j,n} \Pi_j^T;$

else

$W_n = \sum_{k=1}^K \sum_j p_{jk} W_{j,n}$ where p_k is fraction of data points with label k on worker j ;

end

for $j \in \{1, \dots, J\}$ **do**

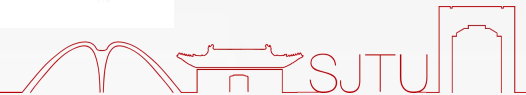
$W_{j,n+1} \leftarrow \Pi_j W_{j,n+1};$ // permute the next-layer weights

 Train $\{W_{j,n+1}, \dots, W_{j,L}\}$ with W_n frozen;

end

$n = n + 1;$

end



A photograph of a modern building with a white, faceted facade and large glass windows, set against a blue sky with light clouds. The building is the background for the slide.

03

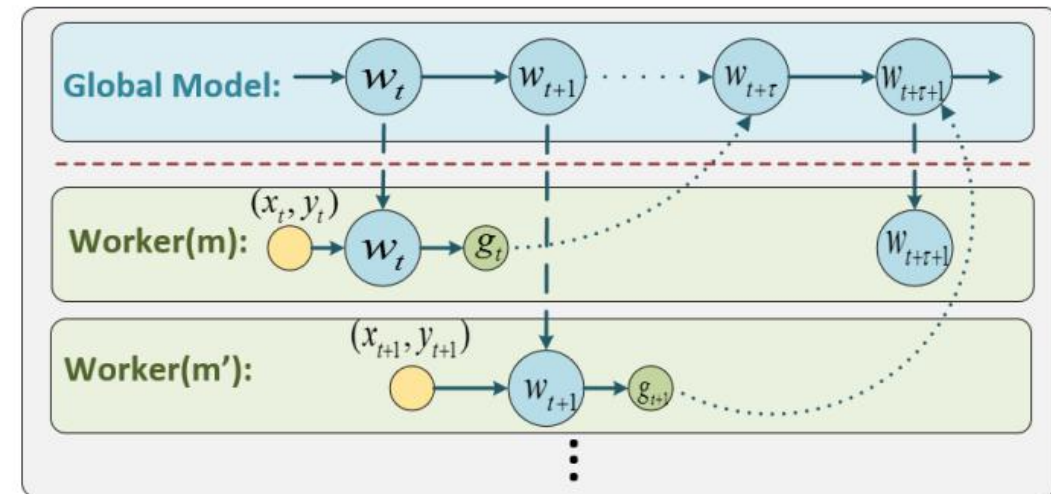
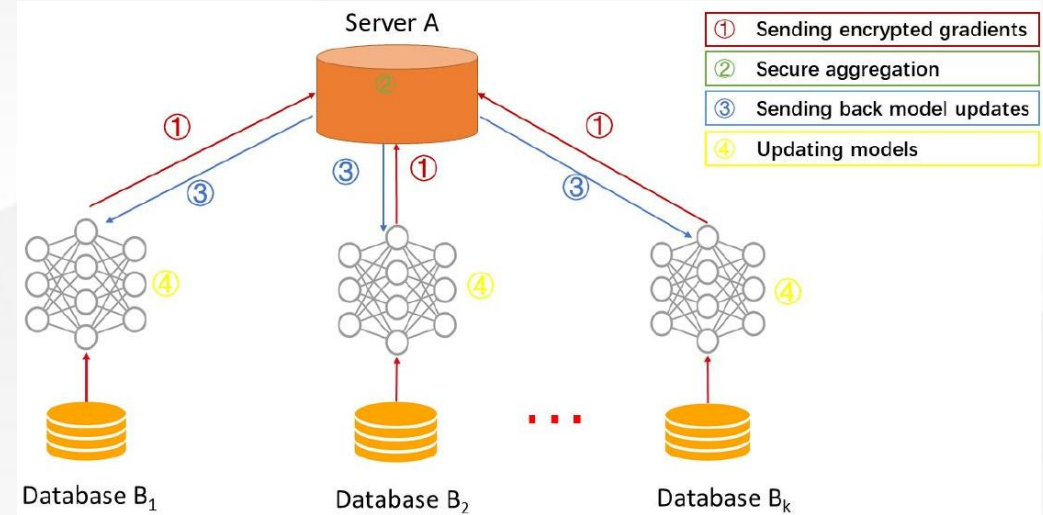
Advanced algorithms

- SageFlow



SageFlow

- ① Stragglers: slow devices
 - **Keep waiting:** slow down the overall process
 - **Drop out:** important data missing
 - **Asynchronous(staleness):** + adversaries?
- ② Attackers: malicious attacks launched by adversaries
 - **untargeted attacks:** model poisoning, data poisoning
 - **targeted/backdoor attacks:** misclassify the targeted subtasks
 - Robust Federated Averaging & Multi-Krum
- Large portion of adversaries
 - Straggler: increase attack ratio

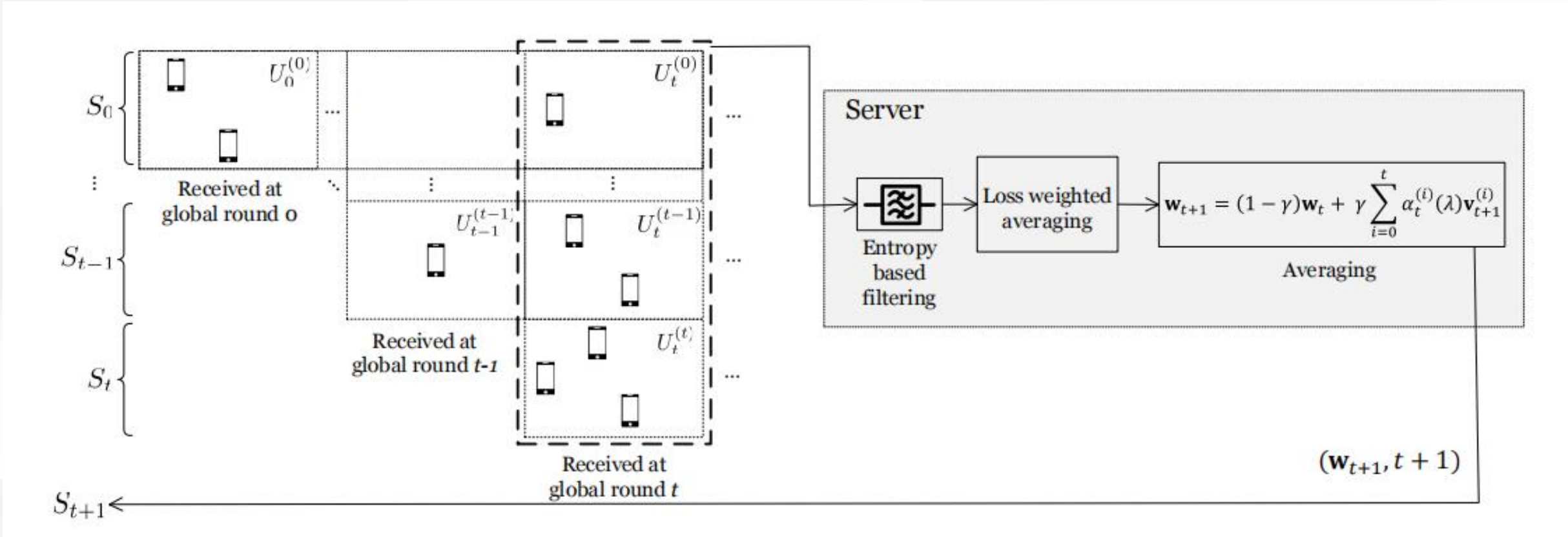




Sageflow algorithm



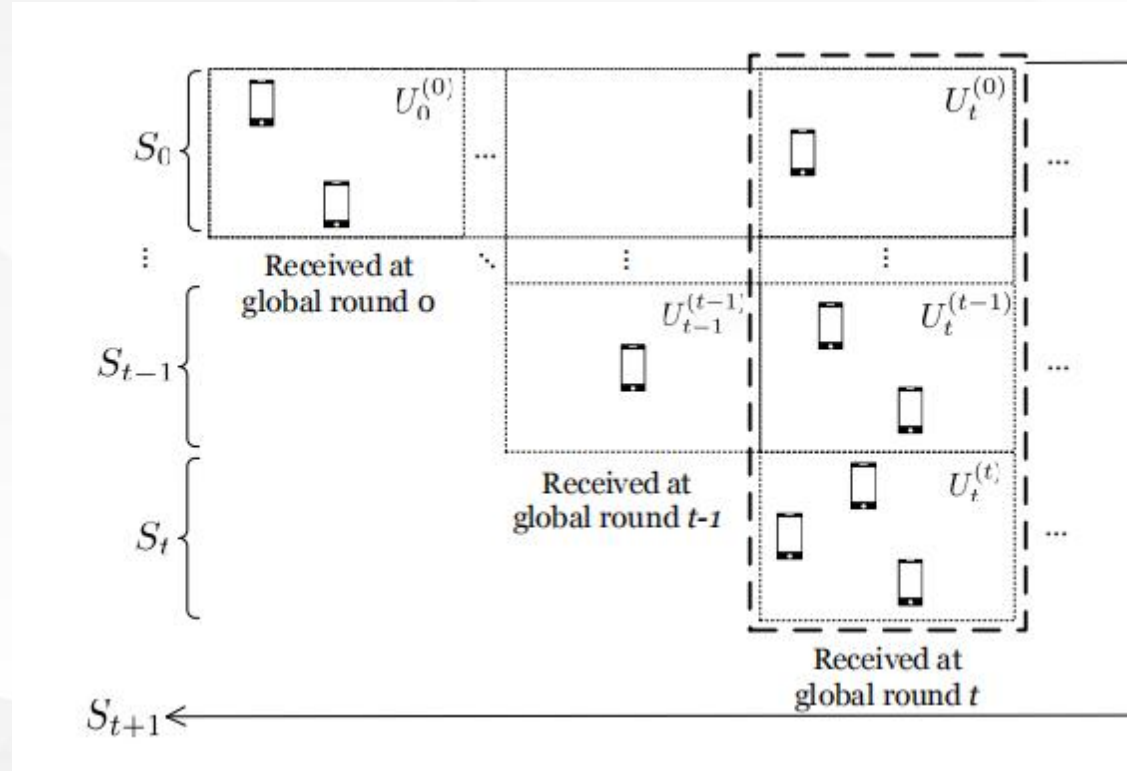
- ⊗ Staleness-aware grouping
- ⊗ Entropy-based filtering + Loss-weighted averaging





Staleness-aware grouping

- perform **periodic** global aggregation (fixed time deadline)
- allow stragglers to be aggregated in **later rounds**
- group with **same staleness** -> group representative model
- aggregate according to **staleness**



Staleness group

Number of data samples

$$\mathbf{v}_{t+1}^{(i)} = \sum_{k \in U_t^{(i)}} \frac{m_k}{\sum_{k \in U_t^{(i)}} m_k} \mathbf{w}_i(k)$$

$$\mathbf{w}_{t+1} = (1 - \gamma) \mathbf{w}_t + \gamma \sum_{i=0}^t \alpha_t^{(i)}(\lambda) \mathbf{v}_{t+1}^{(i)}$$

$$\alpha_t^{(i)}(\lambda) \propto \frac{\sum_{k \in U_t^{(i)}} m_k}{(t-i+1)^\lambda}$$

Staleness function

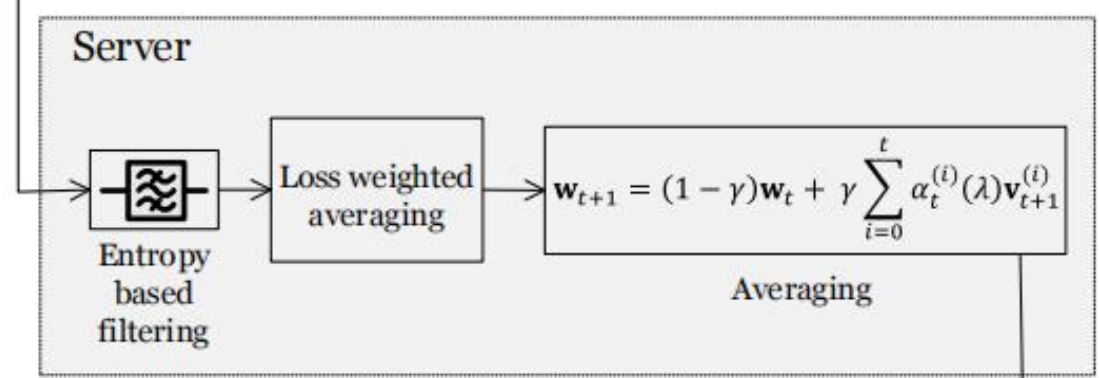




Sageflow algorithm

Entropy-based filtering

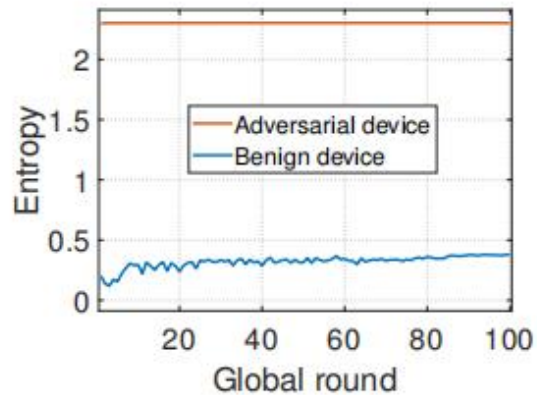
- **Public data** on server
- **Filter out high entropy** models (loss)
- For model poisoning



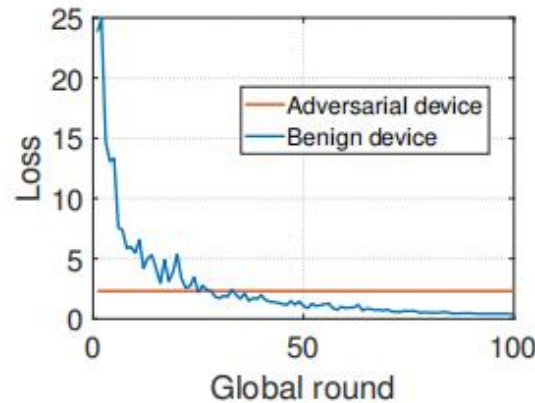
$$E_{x_{pub,j}}(k) = - \sum_{q=1}^Q P_{x_{pub,j}}^{(q)}(k) \log P_{x_{pub,j}}^{(q)}(k)$$

$$E(k) = \frac{1}{n_{pub}} \sum_{j=1}^{n_{pub}} E_{x_{pub,j}}(k)$$

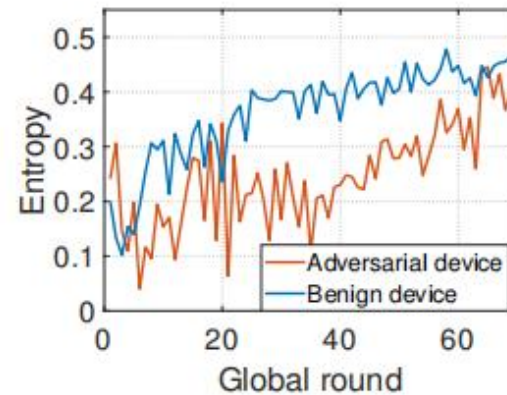
Shannon entropy



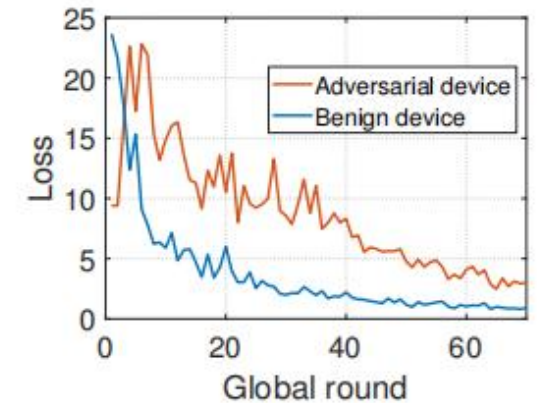
(a) Model poison, entropy



(b) Model poison, loss



(c) Data poison, entropy



(d) Data poison, loss



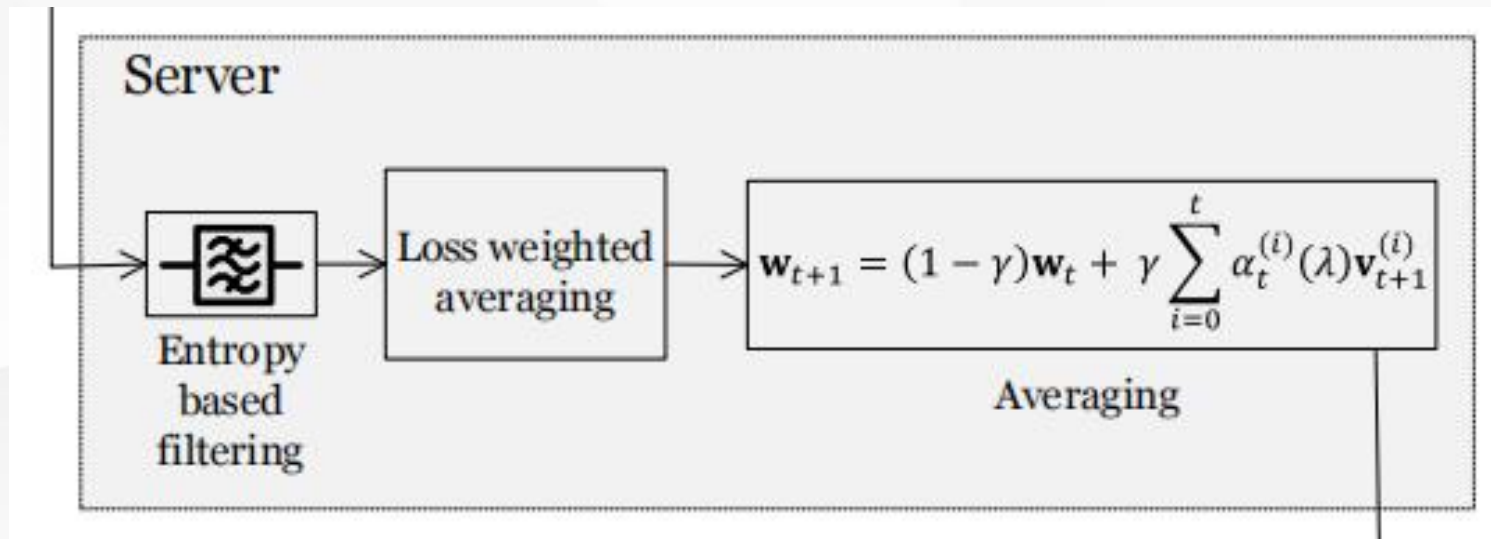
Sageflow algorithm



- Loss-weighted averaging
 - Aggregation weight according to local models' **measured qualities**
 - Measure by **loss** on public data
 - data-poisoned model -> small weight + less impact
 - For data poisoning & scaled backdoor

$$\beta_t^{(k)}(\delta) \propto \frac{m_k}{\{F_{pub}(\mathbf{w}_t(k))\}^\delta} \quad \text{and} \quad \sum_{k \in S_t} \beta_t^{(k)}(\delta) = 1.$$

$$\mathbf{w}_{t+1} = \sum_{k \in S_t} \beta_t^{(k)}(\delta) \mathbf{w}_t(k)$$





Sageflow algorithm



Time complexity

Public data Model parameters Clients number

$$\mathcal{O}(n_{pub}|w|K)$$

Algorithm 1 Proposed Sageflow Algorithm

Input: Initialized model \mathbf{w}_0 , **Output:** Final global model \mathbf{w}_T

Process at the Server

- 1: **for** each global round $t = 0, 1, \dots, T - 1$ **do**
- 2: Choose S_t and send the current model and the global round (\mathbf{w}_t, t) to the devices
- 3: Wait for T_d and then:
- 4: **for** $i = 0, 1, \dots, t$ **do**
- 5: $U_t^{(i)}(E_{th}) = \{k \in U_t^{(i)} | E(k) < E_{th}\}$ // Entropy-based filtering in each group
- 6: $\mathbf{v}_{t+1}^{(i)} = \sum_{k \in U_t^{(i)}(E_{th})} \beta_t^{(k)}(\delta) \mathbf{w}_i(k)$ // Loss-weighted averaging in each group (with same staleness)
- 7: **end for**
- 8: $\mathbf{w}_{t+1} = (1 - \gamma) \mathbf{w}_t + \gamma \sum_{i=0}^t \alpha_t^{(i)}(\lambda) \mathbf{v}_{t+1}^{(i)}$ // Averaging of representative models (with different staleness)
- 9: **end for**

Process at the Device: Device k receives (\mathbf{w}_t, t) from the server and performs local updates to obtain $\mathbf{w}_t(k)$. Then each benign device k sends $(\mathbf{w}_t(k), t)$ to the server, while a malicious adversary sends a poisoned model depending on the type of attack.





Sageflow algorithm



Theoretical Analysis

Convergence analysis

- Assumption 1: μ -strongly convex + L-smooth
- Assumption 2: unbiased estimation

Theoretical bound

$$F(x) \leq F(y) + \nabla F(x)^T(x-y) - \frac{\mu}{2}\|x-y\|^2$$

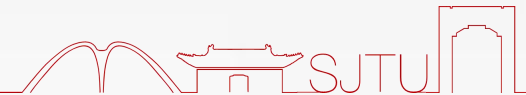
$$F(x) \geq F(y) + \nabla F(x)^T(x-y) - \frac{L}{2}\|x-y\|^2$$

$$\mathbb{E}\|\nabla \hat{F}_k(\mathbf{w}_t^i(k), \xi_t^i(k)) - \nabla F(\mathbf{w}_t^i(k))\|^2 \leq \rho_1$$

$$\mathbb{E}[F(\mathbf{w}_T) - F(\mathbf{w}^*)] \leq \nu^T [F(\mathbf{w}_0) - F(\mathbf{w}^*)] + (1 - \nu^T) Z(\lambda, E_{th}, \delta)$$

Convergence speed

Error





Sageflow algorithm



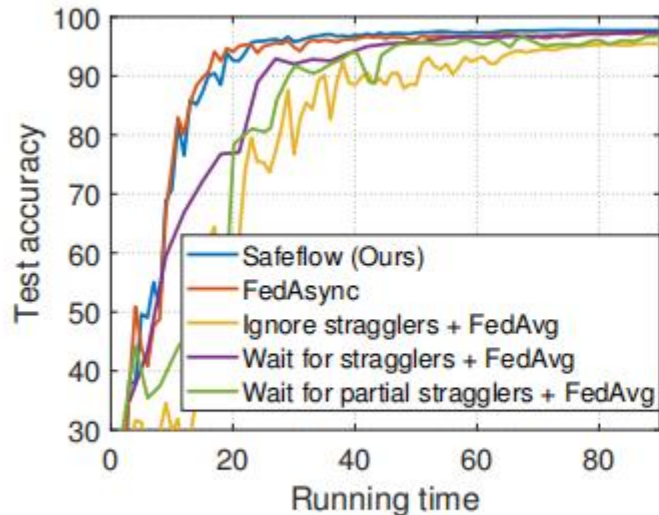
- ① Datasets: MNIST, FMNIST, CIFAR10
 - 2% as public data
- ① Models: CNN(2conv+2fc), CNN(2conv+1fc), CNN(VGG-11)
 - ignore batchnorm
- ① FL setting: 100 clients, two classes for each client, 5 local epochs, batch size of 10



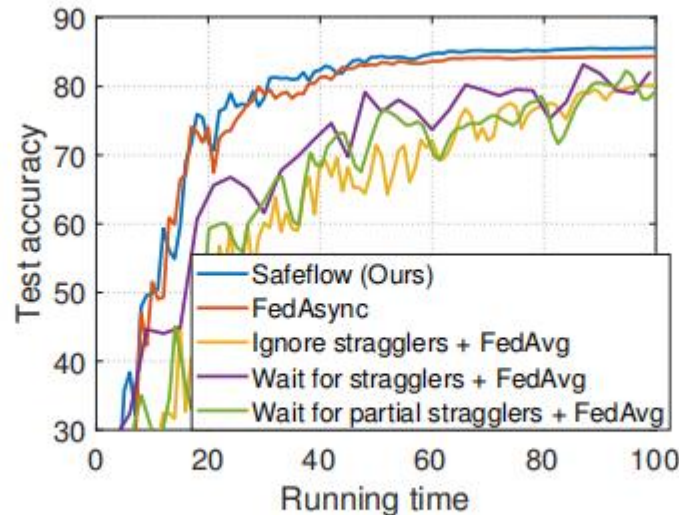
Sageflow algorithm



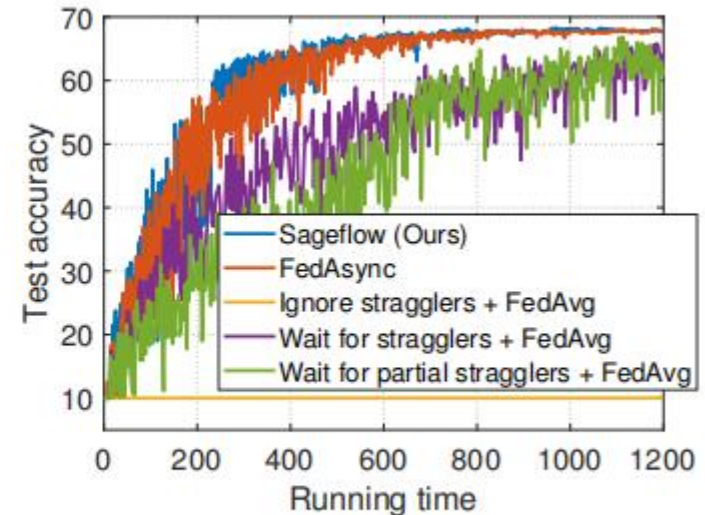
- Only stragglers: 10% participants
- Baselines: FedAvg(waiting, ignoring, waiting 50%), FedAsync
- Settings: uniform delay of [0,1,2] global rounds
- Ignoring lose significant data converges to a **suboptimal point**
- Waiting(all, 50%) requires the **largest running time** until convergence



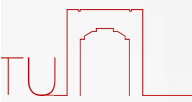
(a) MNIST



(b) FMNIST



(c) CIFAR10

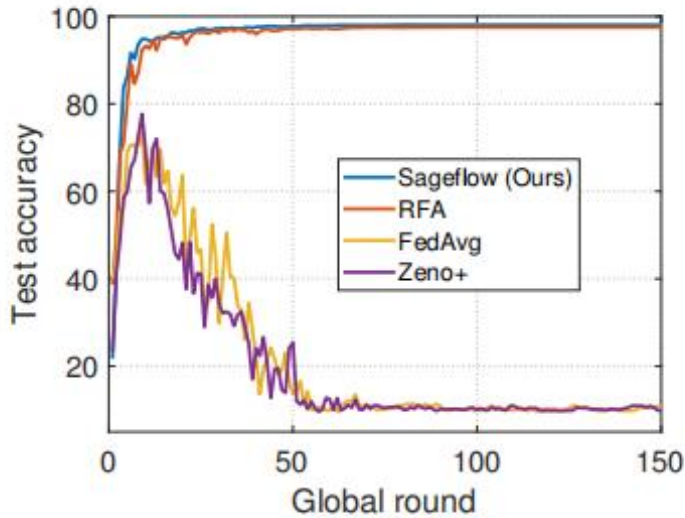




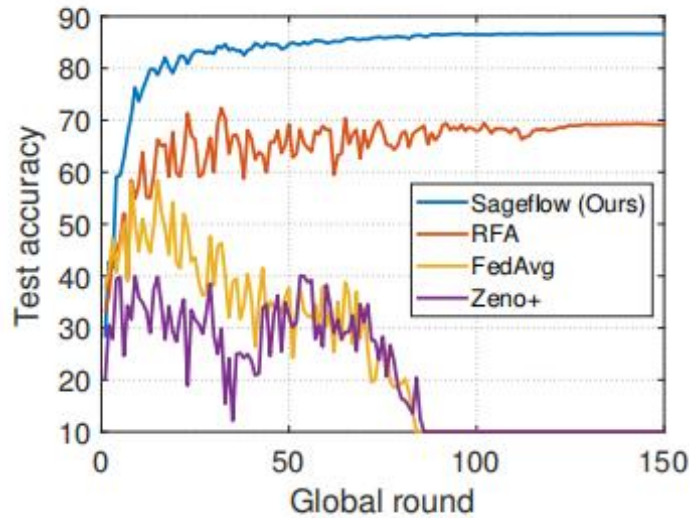
Sageflow algorithm



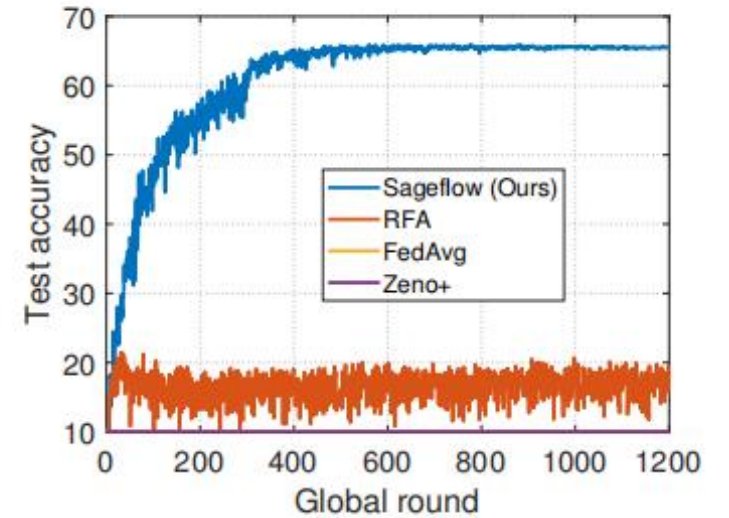
- Only adversaries: 20% participants
- Baselines: RFA, FedAvg, synchronized Zeno+, Multi-Krum
- Attacks: model(-0.1w), data(label-flipping), backdoor(model replacement, pixel-pattern attack)
- FedAvg **does not work well** on all datasets
- Zeno+: bad on poisoning **but** good for backdoor
- RFA: **complex model** led to worse performance
- Sageflow: **slow down** poisoning



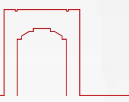
(a) MNIST



(b) FMNIST



(c) CIFAR10

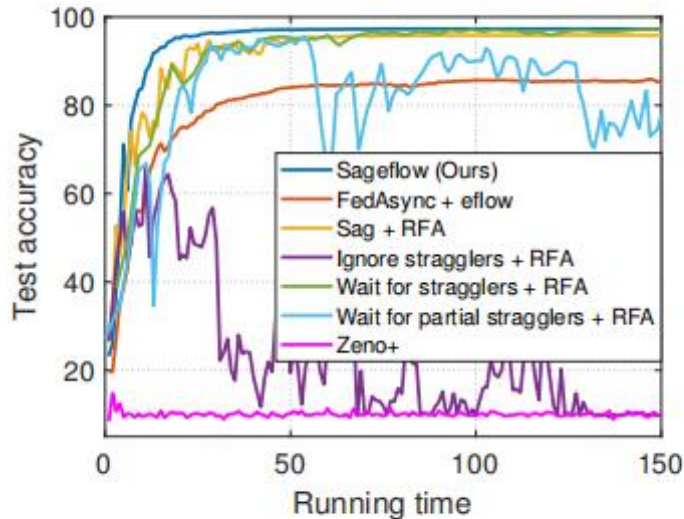




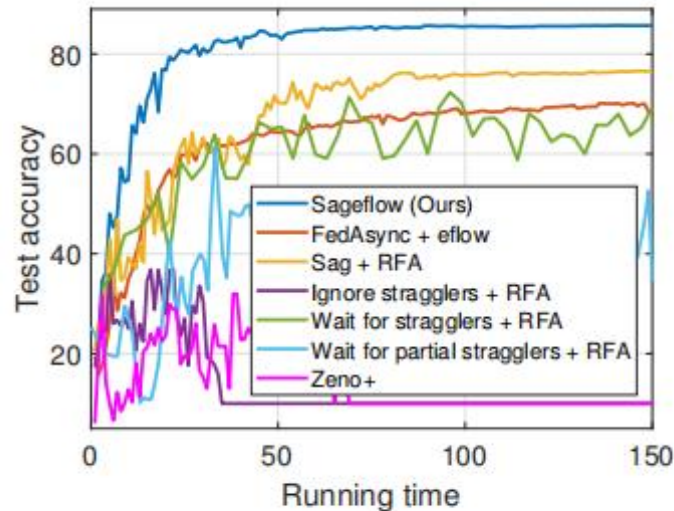
Sageflow algorithm



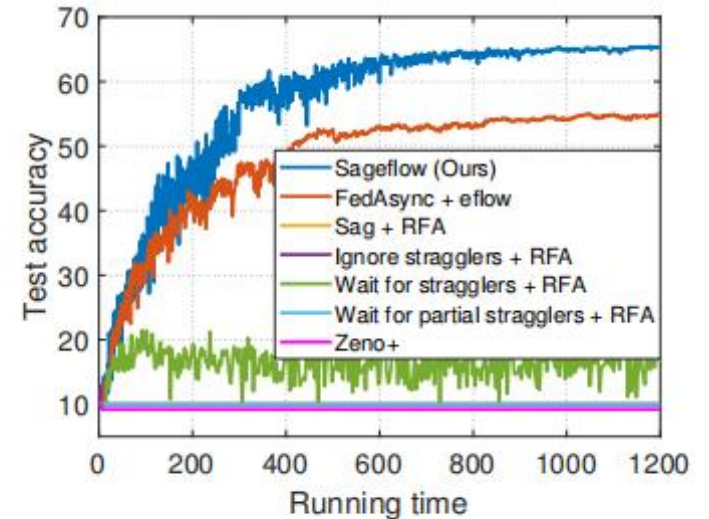
- ⊗ Stragglers + adversaries: 20%(model/data), 10%(backdoor) participants
- ⊗ Baselines: asynchronized Zeno+, Multi-Krum
- ⊗ Zeno+: **does not perform well**(ignore staleness & entropy)
- ⊗ Waiting + RFA: suffer from **straggler**
- ⊗ Ignoring/Sag + RFA: poor(**high attack ratio**) ⊗ eflow/RFA + FedAsync: poor(**one-by-one arrivals**)



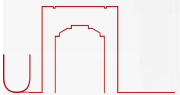
(a) MNIST



(b) FMNIST



(c) CIFAR10





Sageflow algorithm



- ① Sageflow: robust FL scheme handle both **stragglers and adversaries**
 - staleness-aware grouping: stragglers
 - entropy-based filtering: model poisoning
 - loss-weighted averaging: data poisoning + backdoor
- ① Theoretical convergence analysis
- ① Extensive experimental results
- ① Future issues: Sageflow + secure aggregation



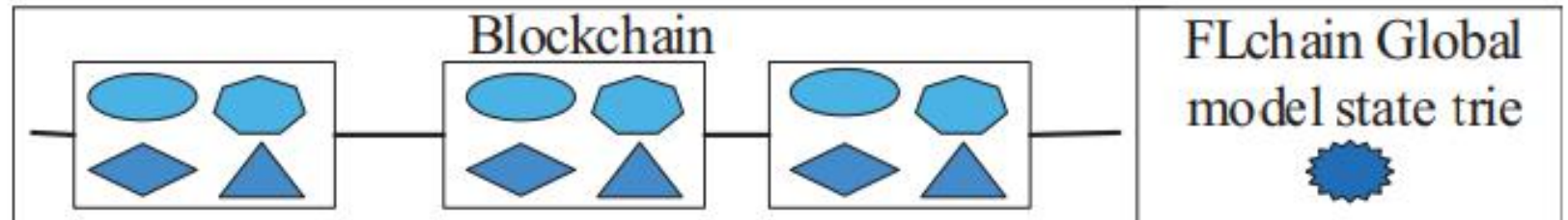
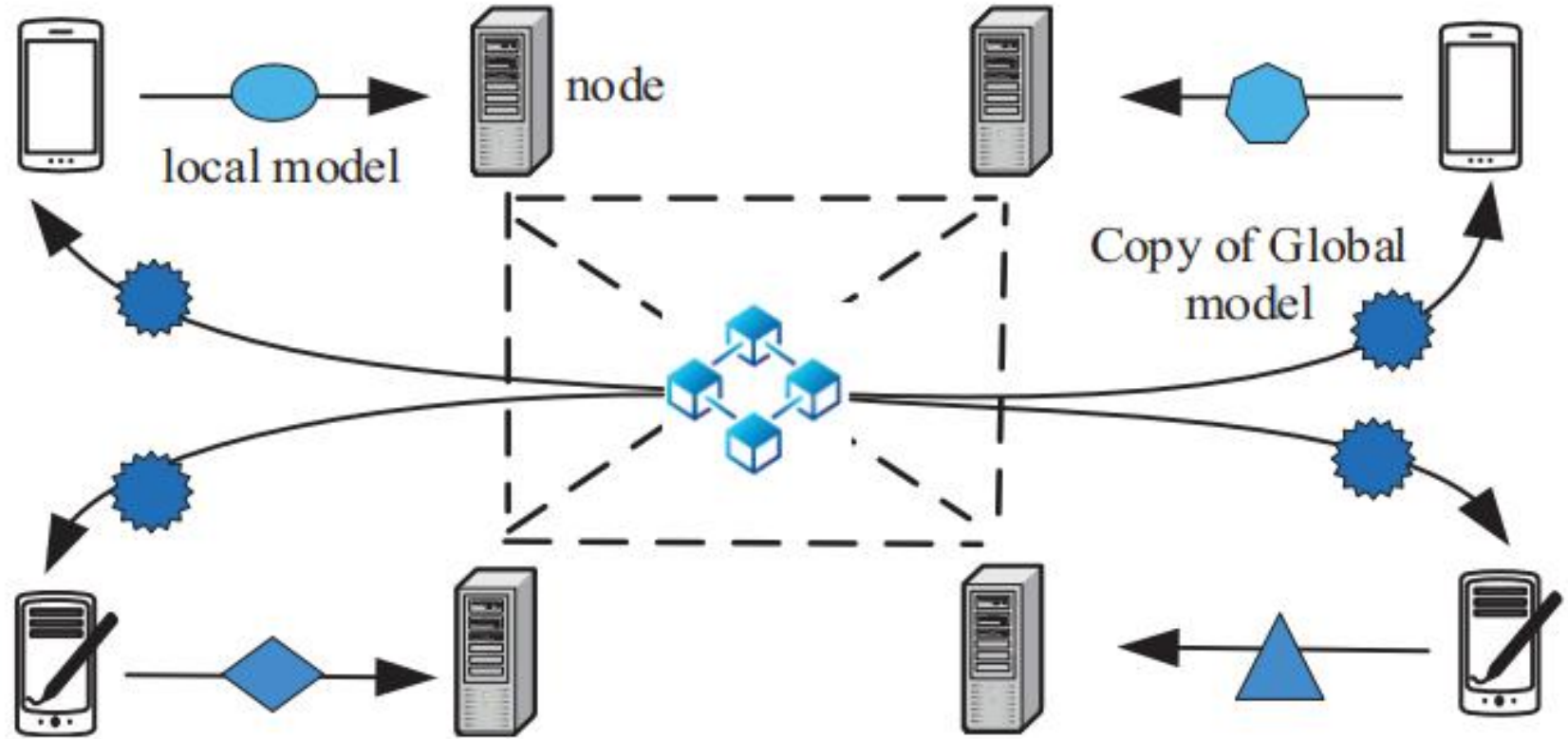
FLChain



FL + Blockchain



- ④ complete dependency on the reliability of a central server for storage and computation of the global model update
- ④ Any malicious activity leads to flawed global model update which is detrimental for accuracy of subsequent local model updates

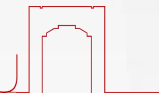


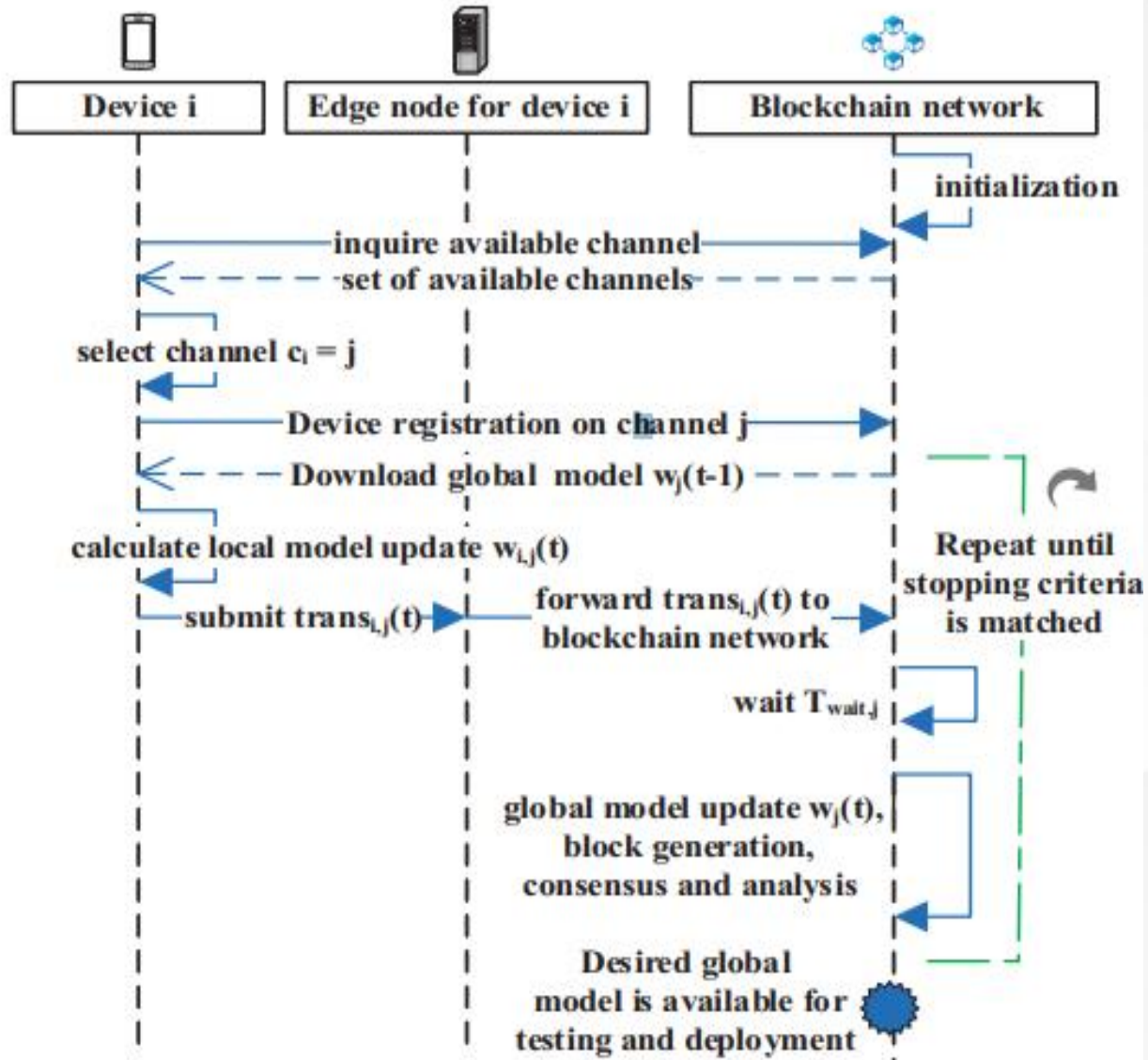


- ① Initialization
- ② Channel Inquiry
- ③ Channel Selection
- ④ Device Registration
- ⑤ Local Model Update
- ⑥ Transaction Pool
- ⑦ Global Model Update
- ⑧ Consensus Protocol
- ⑨ Analysis

Algorithm 1 : FLchain operation for channel j

- 1: Setup channel j for global model M_j
 - 2: initialization: $t = 0$; $w_j(0), w_{i,j}(0) \in [0, w_{j,max}]$;
 - 3: **for all** $i \in D_j$ **do in parallel**
 - 4: Inquire available channels
 - 5: Select channel $c_i = j \in C \triangleq \{1, 2, 3, \dots, C_n\}$
 - 6: Register device i to channel j
 - 7: **end for**
 - 8: **while** $\|w_j(t) - w_j(t-1)\|_2 \leq \varepsilon_{threshold,j}$ **do**
 - 9: **for all** $i \in D_j$ **do in parallel**
 - 10: Download $w_j(t)$ from channel j to device i
 - 11: $t \leftarrow t + 1$; $w_{i,j}^0(t) = w_j(t)$;
 - 12: **for** $v = 1, \dots, V$ **do**
 - 13: $w_{i,j}^v(t) = w_{i,j}^{v-1}(t) - \eta \nabla \Phi$, and Eq. (6)
 - 14: **end for**
 - 15: $w_{i,j}(t) = w_{i,j}^V(t)$, Generate $trans_{i,j}(t)$ and forward to blockchain network
 - 16: Wait for notification from channel j
 - 17: **end for**
 - 18: Calculate $w_j(t)$ using Eq. (7)
 - 19: global model state trie updation, block generation and consensus
 - 20: **end while**
-



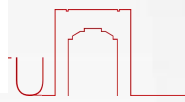
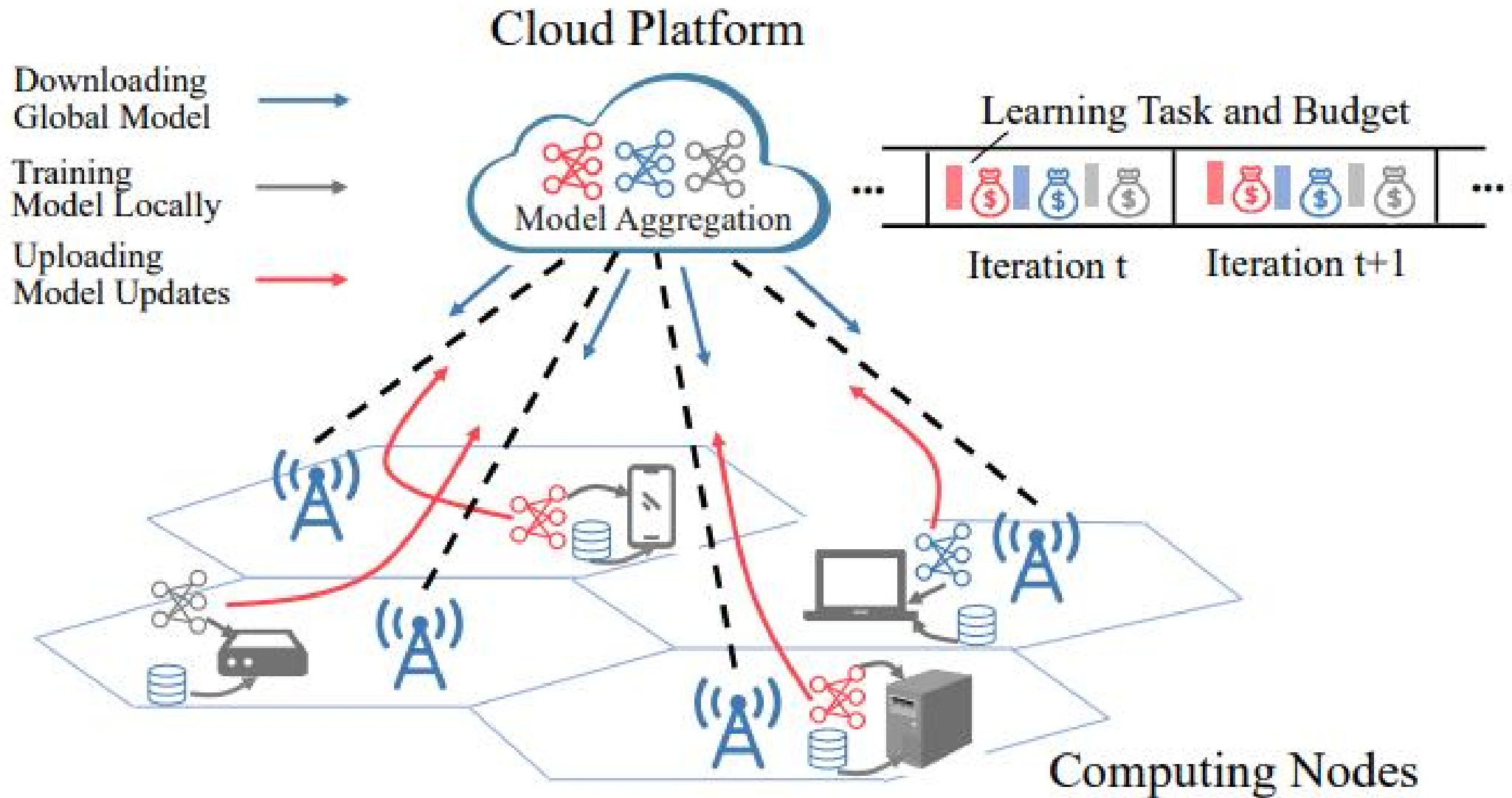




FAIR



Federated Learning





⊗ Problems

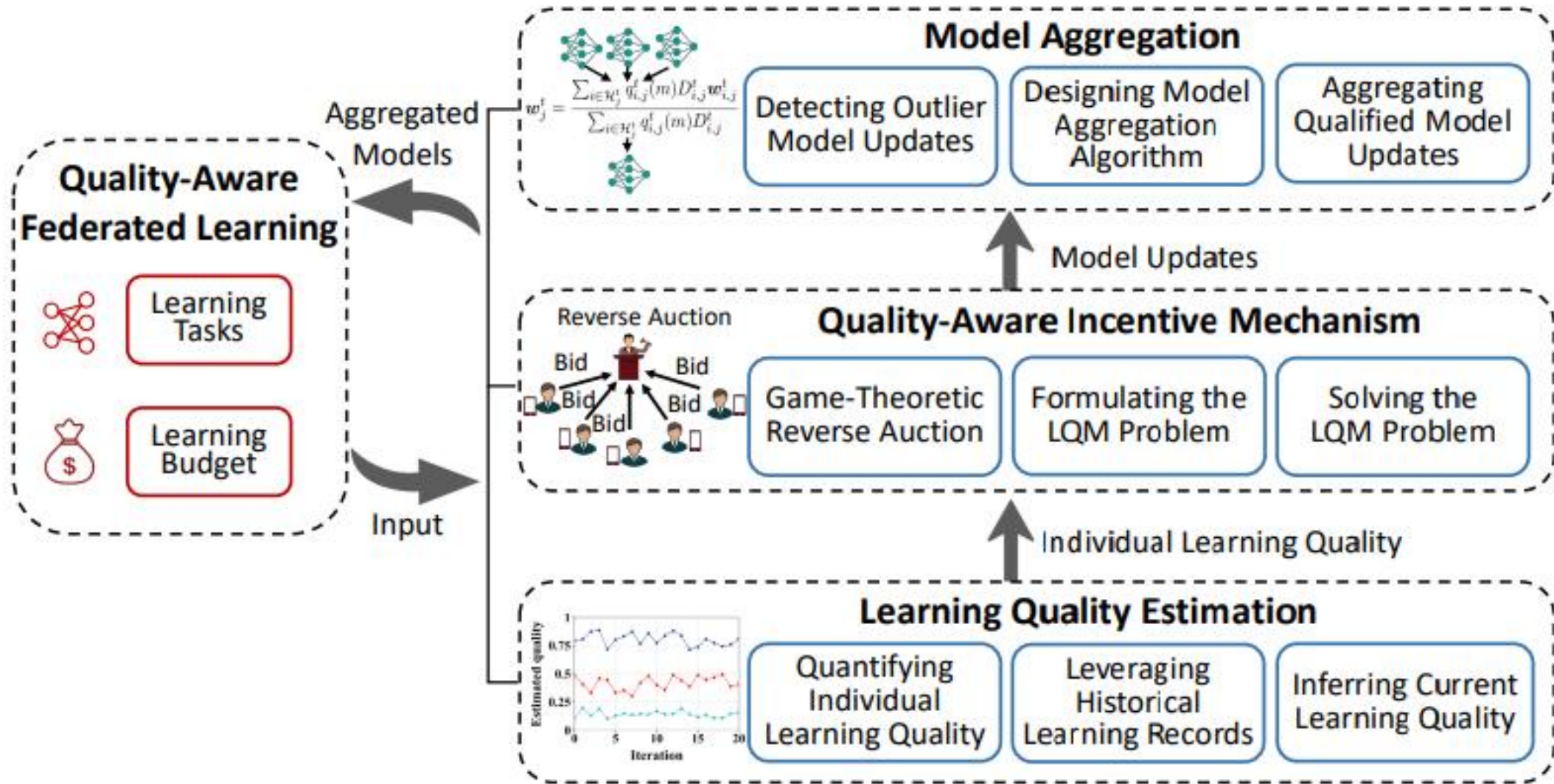
- mobile users do not participate in learning due to computation and energy consumption
- model update quality of mobile devices can vary dramatically, inclusively aggregating low-quality model updates can deteriorate the global model quality

⊗ Solution: Federated Learning with quality awareness

- learning quality estimation: leverage historical learning records to estimate the user learning quality, the exponential forgetting function is utilized for weight assignment;
- quality-aware incentive mechanism: within the recruiting budget, model a reverse auction problem to encourage the participation of high-quality learning users;
- model aggregation: integrates the model quality into aggregation and filters out non-ideal model updates, to further optimize the global learning model



FAIR Architecture





Estimating Learning Quality

$$m_{i,j}^t = \text{loss}_j(t_s) - \text{loss}_{i,j}(t_e).$$

$$q_{i,j}^t = \begin{cases} m_{i,j}^t D_{i,j}^t, & \text{if } t_s < t_{i,j}^u \leq t_e; \\ 0, & \text{otherwise.} \end{cases}$$

$$\hat{q}_{i,j}^t = \frac{\sum_{k=0}^r \rho^{t_r - t_k} q_{i,j}^{t_k}}{\sum_{k=0}^r \rho^{t_r - t_k}}.$$

Quality-Aware Incentive Mechanism

Algorithm 1: Solving the LQM problem.

Input : (1) bid price $b_{i,j}^t$; (2) budget B_j^t ; (3) task set \mathcal{L}_i^t ;
(4) quality estimation values $\hat{q}_{i,j}^t$.

Output : (1) the task allocation results $s_{i,j}^t$; (2) the payments $r_{i,j}^t$.

```

1 Initialize  $\mathcal{N}_j^t \leftarrow \emptyset$ ,  $p_j^t \leftarrow 0$  for each  $l_j^t \in \mathcal{L}^t$ ;
2 Initialize  $x_i^t \leftarrow 1$  for each  $i \in \mathcal{N}$ ;
3 Initialize  $r_{i,j}^t \leftarrow 0$ ,  $s_{i,j}^t \leftarrow 0$  for each  $i \in \mathcal{N}$ ,  $l_j^t \in \mathcal{L}^t$ ;
4 foreach  $i \in \mathcal{N}$  do
5   foreach  $l_j^t \in \mathcal{L}_i^t$  do
6      $\mathcal{N}_j^t \leftarrow \mathcal{N}_j^t + \{i\}$ ;
7   end
8 end
9 while  $\exists x_i^t = 0$  and  $\exists p_j^t = 0$  do
10  Initialize  $\mathcal{M}_j^t \leftarrow \emptyset$  for each  $l_j^t \in \mathcal{L}^t$ ;
11  foreach  $l_j^t \in \mathcal{L}^t$  do
12    if  $p_j^t = 0$  then
13      Sort all  $i \in \mathcal{N}_j^t$  in descending order of  $\frac{q_{i,j}^t}{b_{i,j}^t}$ ;
14      Find the smallest  $k$  such that
15         $\sum_{i=1}^k \frac{b_{i,j}^t}{q_{i,j}^t} q_{i,j}^t x_i^t > B_j^t$ ;
16      for  $i \leftarrow 1$  to  $k-1$  do
17         $\mathcal{M}_j^t \leftarrow \mathcal{M}_j^t + \{i\}$ ;
18         $r_{i,j}^t \leftarrow \frac{b_{i,j}^t}{q_{i,j}^t} q_{i,j}^t$ ;
19      end
20    end
21  end
22  Find the task  $l_k^t$  with maximum  $\sum_{i \in \mathcal{M}_k^t} q_{i,k}^t x_i^t$ ;
23  Set  $p_k^t \leftarrow 1$ ;
24  foreach  $i \in \mathcal{M}_k^t$  do
25    if  $x_i^t = 1$  then
26       $s_{i,k}^t \leftarrow 1$ ;
27       $x_i^t \leftarrow 0$ ;
28    end
29  end
30 return  $(s_{i,j}^t, r_{i,j}^t)$ ;

```

The defined LQM problem can be formulated as

$$\max_{\mathcal{M}_j^t, \mathcal{R}_j^t} \sum_{l_j^t \in \mathcal{L}^t} \sum_{i \in \mathcal{M}_j^t} \hat{q}_{i,j}^t, \quad (10)$$

$$s.t. \sum_{i \in \mathcal{N}} r_{i,j}^t s_{i,j}^t \leq B_j^t, \quad \forall l_j^t \in \mathcal{L}^t, \quad (11)$$

$$r_{i,j}^t \geq b_{i,j}^t, \quad \forall l_j^t \in \mathcal{L}^t, \forall i \in \mathcal{M}_j^t, \quad (12)$$

$$s_{i,j}^t \in \{0, 1\}, \quad \forall i \in \mathcal{N}, \forall l_j^t \in \mathcal{L}^t, \quad (13)$$

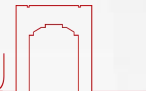
$$s_{i,j}^t = 0, \quad \forall l_j^t \notin \mathcal{L}_i^t, \forall i \in \mathcal{N}, \quad (14)$$

$$\sum_{l_j^t \in \mathcal{L}^t} s_{i,j}^t \leq 1, \quad \forall i \in \mathcal{N}. \quad (15)$$

```

19   | end
20   | end
21   Find the task  $l_k^t$  with maximum  $\sum_{i \in \mathcal{M}_k^t} q_{i,k}^t x_i^t$ ;
22   Set  $p_k^t \leftarrow 1$ ;
23   foreach  $i \in \mathcal{M}_k^t$  do
24     | if  $x_i^t = 1$  then
25       |  $s_{i,k}^t \leftarrow 1$ ;
26       |  $x_i^t \leftarrow 0$ ;
27     | end
28   | end
29 end
30 return  $(s_{i,j}^t, r_{i,j}^t)$ ;

```





Model Aggregation

$$w_j^t = \frac{\sum_{i=1}^N D_{i,j}^t w_{i,j}^t}{\sum_{i=1}^N D_{i,j}^t},$$

$$w_j^t = \frac{\sum_{i \in \mathcal{M}_j^t} m_{i,j}^t D_{i,j}^t w_{i,j}^t}{\sum_{i \in \mathcal{M}_j^t} m_{i,j}^t D_{i,j}^t},$$

Algorithm 2: Model aggregation algorithm.

Input : (1) winner node set \mathcal{M}_j^t ; (2) local model parameters $w_{i,j}^t$; (3) data size used for training $D_{i,j}^t$; (4) data quality $m_{i,j}^t$.

Output : aggregated model parameters w_j^t .

- 1 Initialize high-quality node set $\mathcal{H}_j^t \leftarrow \{i : i \in \mathcal{M}_j^t\}$;
- 2 Set $D \leftarrow \sum_{i \in \mathcal{M}_j^t} m_{i,j}^t D_{i,j}^t$;
- 3 Set $w_j^t \leftarrow \sum_{i \in \mathcal{M}_j^t} \frac{m_{i,j}^t D_{i,j}^t}{D} w_{i,j}^t$;
- 4 **foreach** $i \in \mathcal{M}_j^t$ **do**
- 5 | Compute $d_{i,j}^t \leftarrow \text{similarity}(w_j^t, w_{i,j}^t)$;
- 6 **end**
- 7 Compute $\bar{\mu}_d, \hat{\mu}_d, \sigma_d$;
- 8 **if** $\bar{\mu}_d > \hat{\mu}_d$ **then**
- 9 | **foreach** $i \in \mathcal{H}_j^t$ **do**
- 10 | | **if** $d_{i,j}^t > \hat{\mu}_d + \eta \sigma_d$ **then**
- 11 | | | $\mathcal{H}_j^t \leftarrow \mathcal{H}_j^t - \{i\}$;
- 12 | | **end**
- 13 | **end**
- 14 **else**
- 15 | **foreach** $i \in \mathcal{H}_j^t$ **do**
- 16 | | **if** $d_{i,j}^t < \hat{\mu}_d - \eta \sigma_d$ **then**
- 17 | | | $\mathcal{H}_j^t \leftarrow \mathcal{H}_j^t - \{i\}$;
- 18 | | **end**
- 19 | **end**
- 20 **end**
- 21 Set $D \leftarrow \sum_{i \in \mathcal{H}_j^t} m_{i,j}^t D_{i,j}^t$;
- 22 Set $w_j^t \leftarrow \sum_{i \in \mathcal{H}_j^t} \frac{m_{i,j}^t D_{i,j}^t}{D} w_{i,j}^t$;
- 23 **return** w_j^t ;



Theoretical

- Truthful
- Individually rational
- Computational Efficiency

Truthfulness 每一个报价都必须都是真实的，等于自身的cost，如果不等于，则utility会变小

Individual Rationality 每一轮迭代的utility都大于等于0

Computational Efficiency 所有算法的执行时间为多项式时间



Settings

- Models: MLP, LeNet, MobileNet, VGG-11, EfficientNet-B0, ResNet-18
- Dataset: MNIST, Fashion-MNIST, CIFAR-10, Street View House Numbers
- Baselines: Theoretically optimal mechanism, Knapsack greedy mechanism, Bid price first mechanism
- Hyperparameters

Settings	$ \mathcal{N} $	$ \mathcal{L}^t $	$b_{i,j}^t$	$D_{i,j}^t$	N_e	R_e	B_j^t
I	20	1	[1,3]	[1000,3000]	5	[0.05,0.5]	5
II	20	1	[1,3]	5000	[0,20]	0.5	10
III	100	4	[1,10]	[1000,10000]	30	[0,1]	[10,30]

- **Theoretically optimal mechanism:** It adopts the depth-first search approach to find the theoretically optimal solution for the LQM problem, which however cannot guarantee the truthfulness of nodes.
- **Knapsack greedy mechanism:** It greedily selects winner nodes based on the amount of data used for training divided by the bid price, i.e., $D_{i,j}^t/b_{i,j}^t$, where the data quality and truthfulness of nodes are not considered.
- **Bid price first mechanism:** It preferentially selects nodes with the lowest bid price without guaranteeing the truthfulness of nodes.



④ Different incentive mechanisms

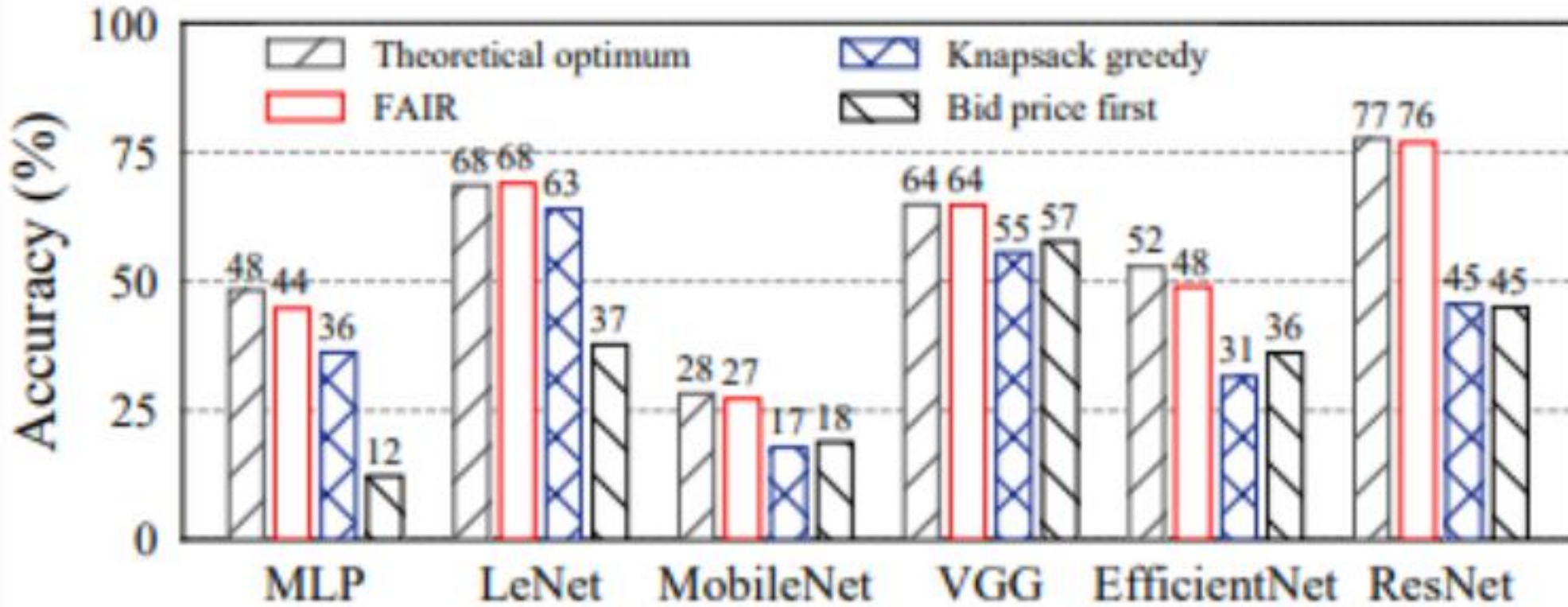


Fig. 4. The average model accuracy with different incentive mechanisms.



Attack scenarios

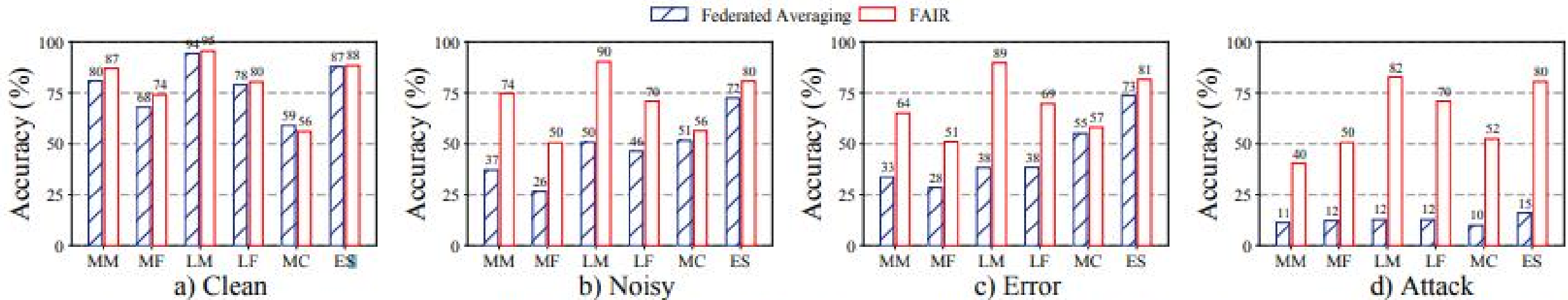
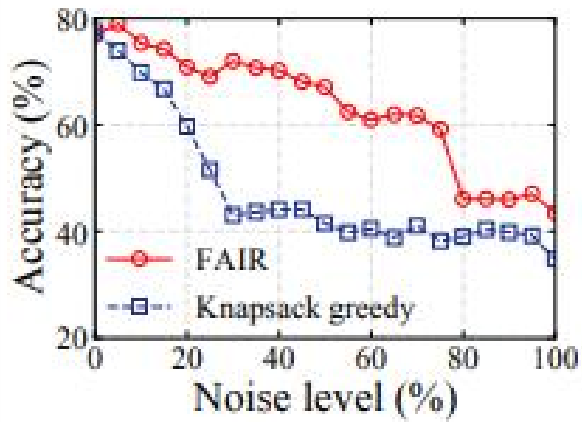


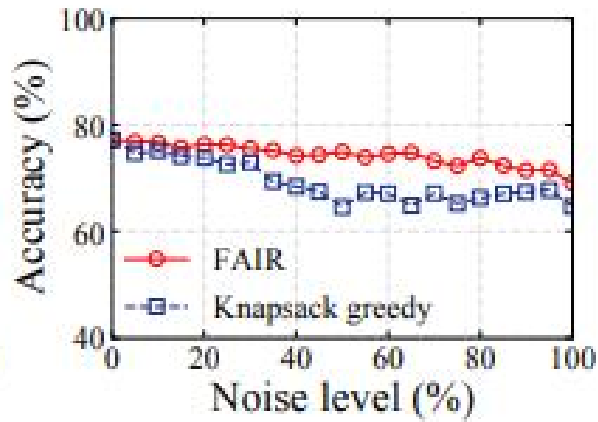
Fig. 3. The model aggregation performance of MLP-MNIST (MM), MLP-FMNIST (MF), LeNet-MNIST (LM), LeNet-FMNIST (LF), MobileNet-CIFAR10 (MC), EfficientNet-SVHN (ES) under different scenarios: a) Clean datasets; b) Noisy label datasets; c) Error label datasets; d) Attack datasets.



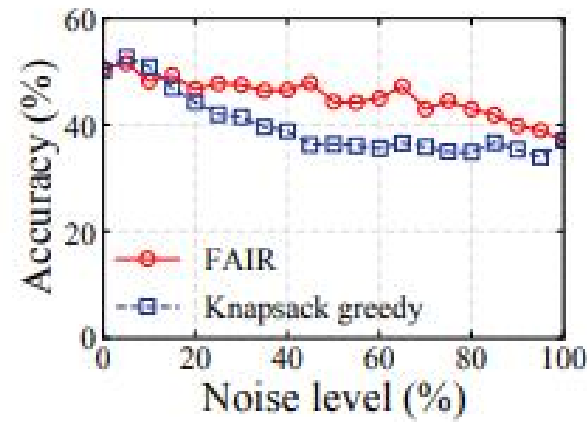
⊗ Noise levels



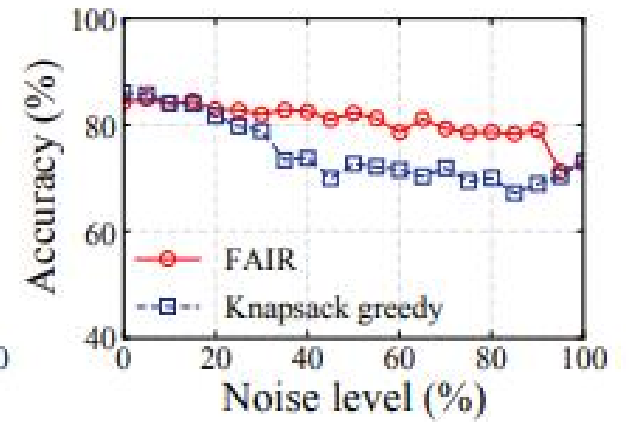
a) MLP MNIST



b) LeNet FMNIST



c) MobileNet CIFAR10



d) EfficientNet SVHN

Fig. 5. The model accuracy vs. noise levels.



Learning budgets

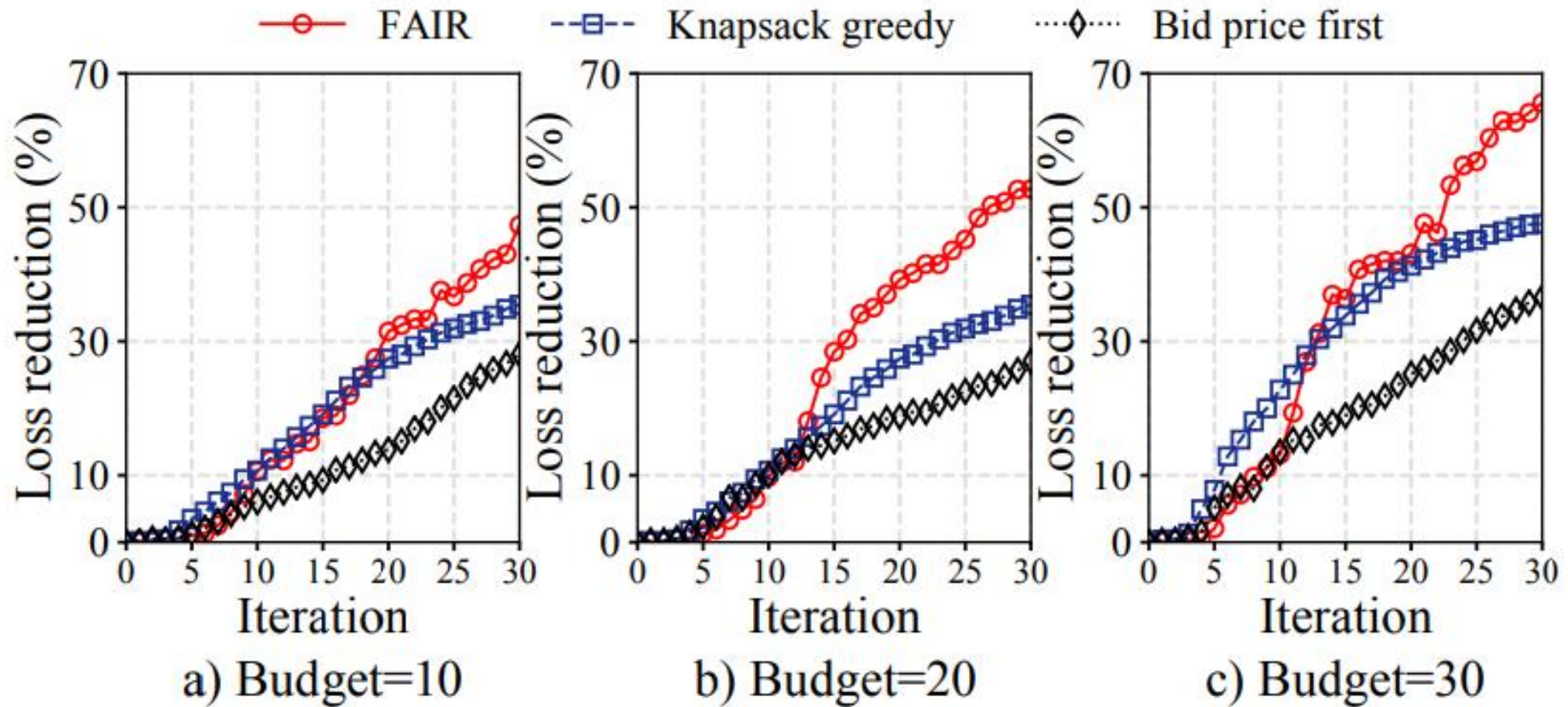


Fig. 6. The average test loss reduction of learning models vs. learning budgets.



Q&A

饮水思源 爱国荣校