

FEDERATED LEARNING WITH HETEROGENEOUS DATA FORMAT

Under the guidance of

Prof Mrinal Kanti Das

Department of Data Science, IIT Palakkad

-Santhosh V(142302020)

-Siva Kumar(142302008)

Quick Outline

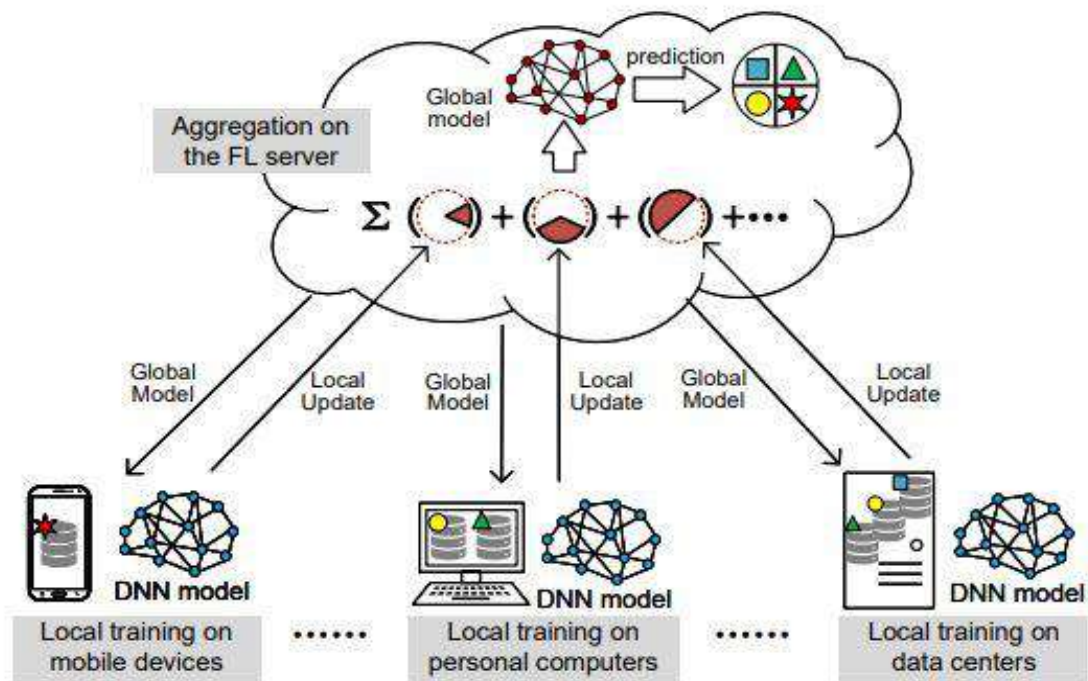
- What is Federated Learning
- Objective
- Applications
- Types Of Federated Learning (Data Partitioning)
- Types Of Federated Learning(Based On Client)
- Aggregation Algorithms
- Results & Observations
- References

Shift from centralized data to decentralized data

The standard ML considers a centralized dataset processed in tightly integrated system.

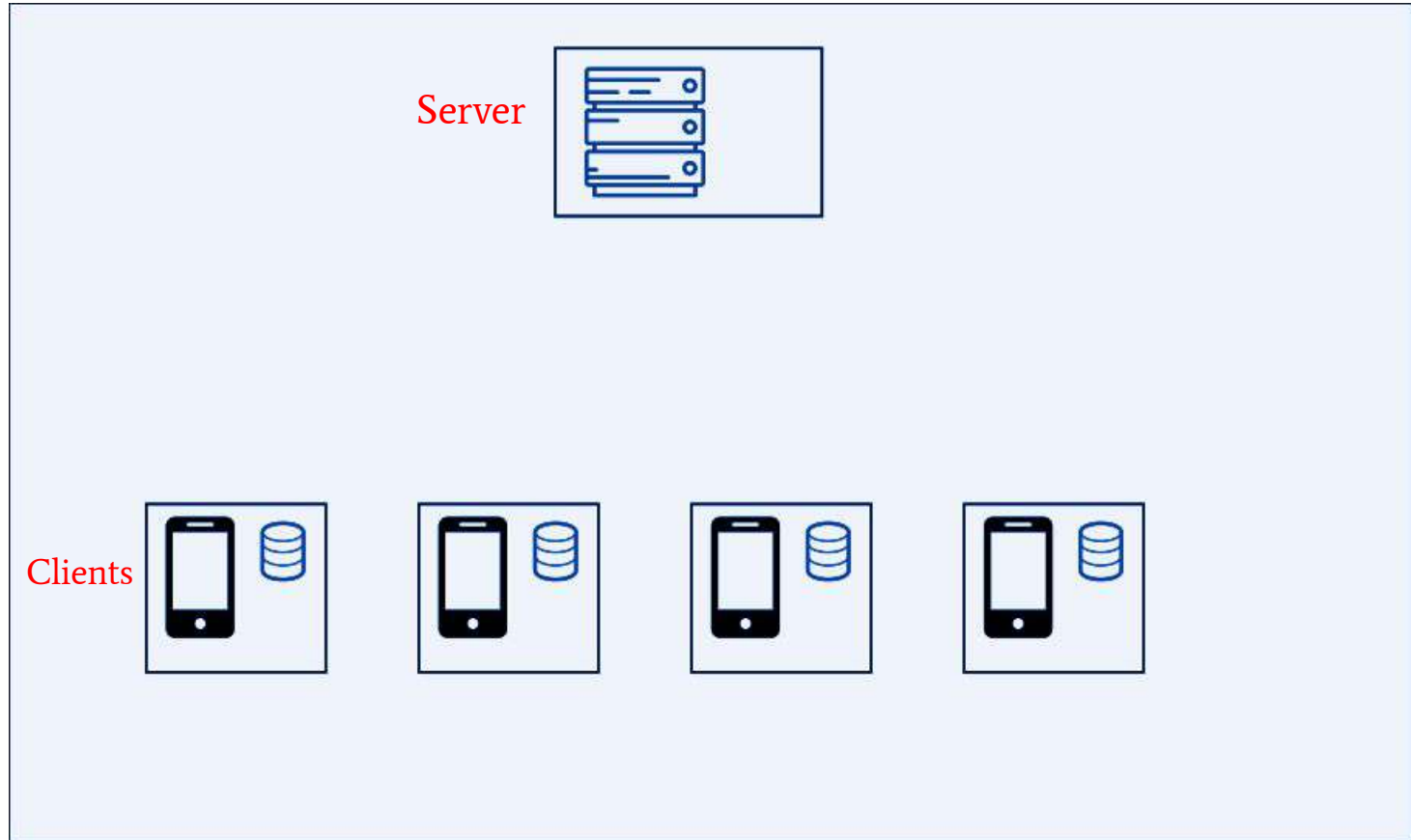
- Sending data to cloud for centralized ML is too costly
 - Self driving cars generate several TBs of data every day
 - Some wireless networks have limited bandwidth
- Data is too sensitive (medical reports)
 - Data privacy
 - Keeping the control over data and give the competitive advantage in business and research

What is Federated Learning?

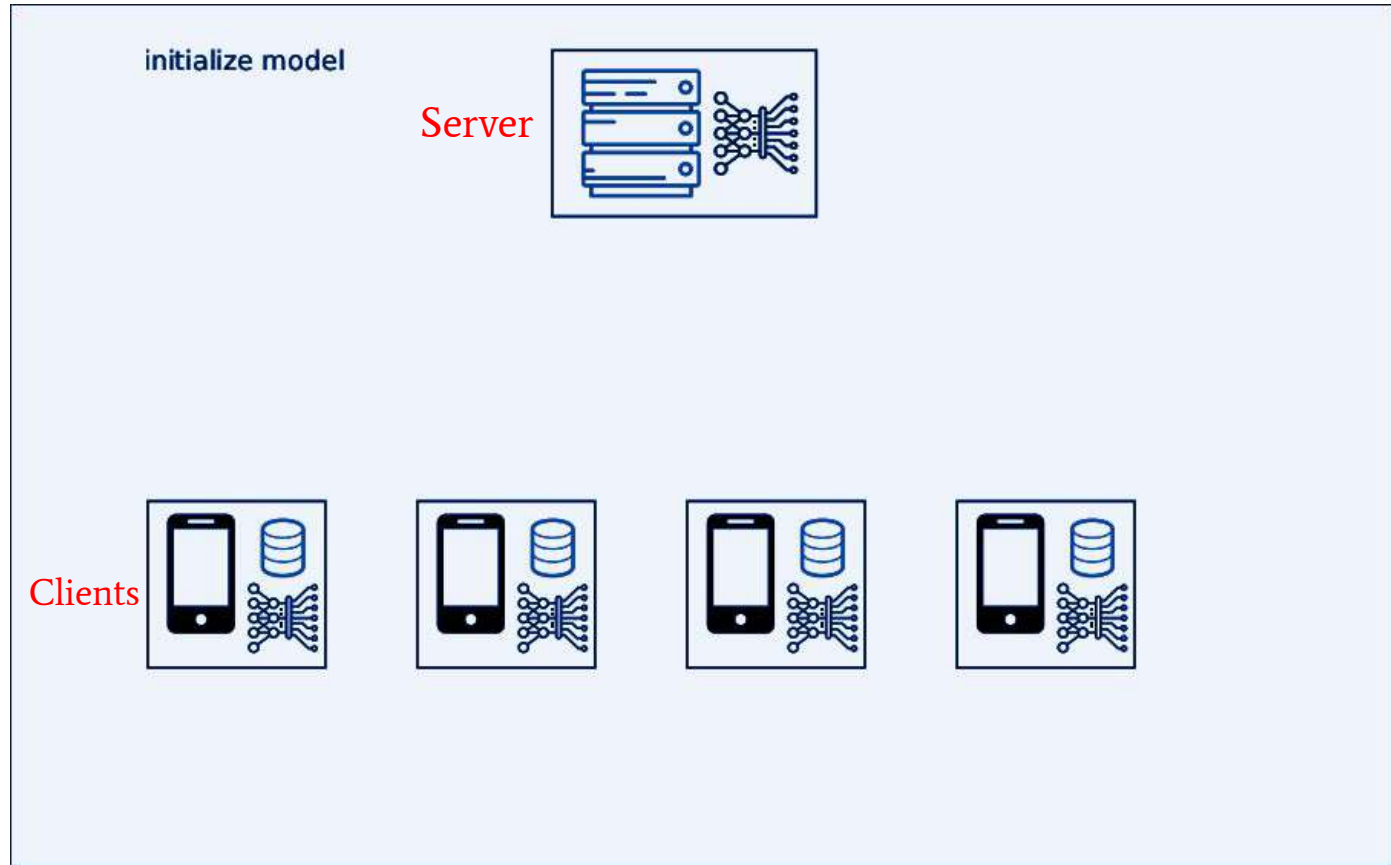


Federated Learning Process

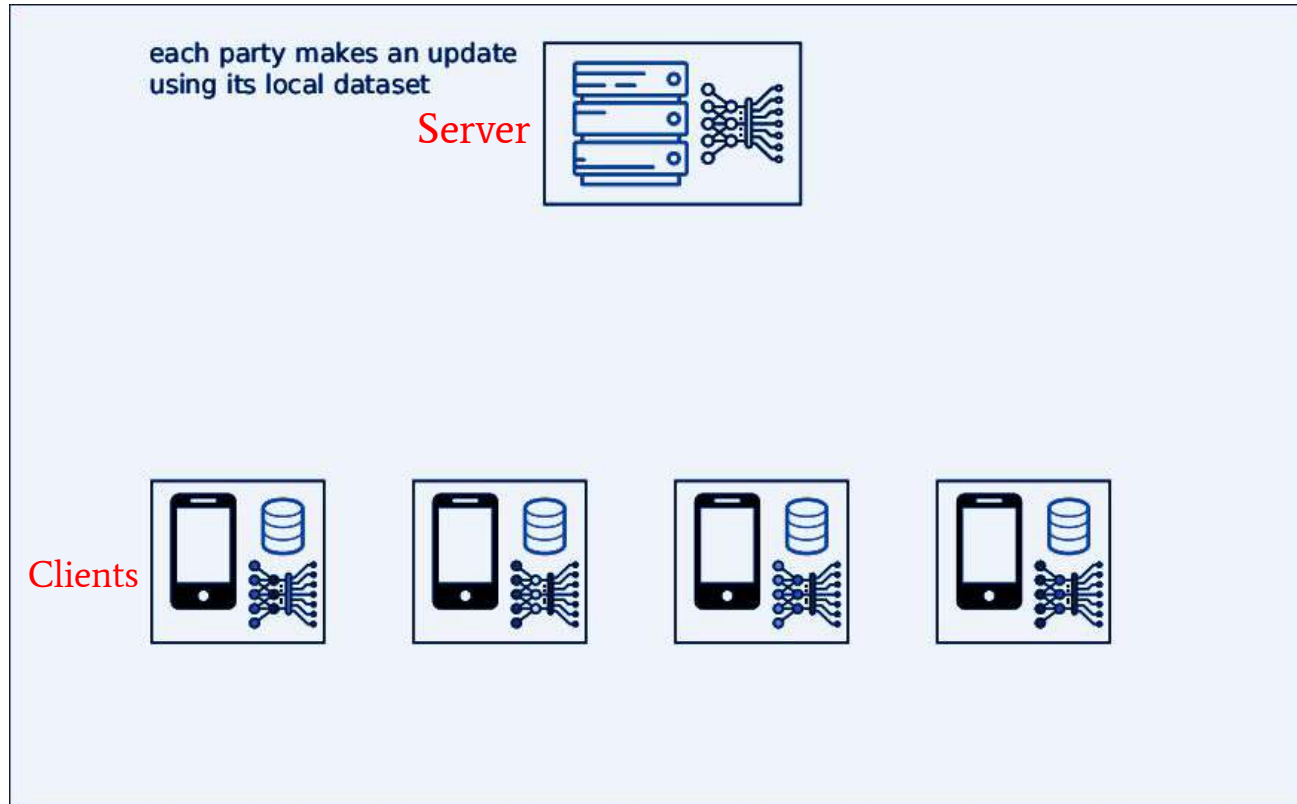
Client - Server Setup



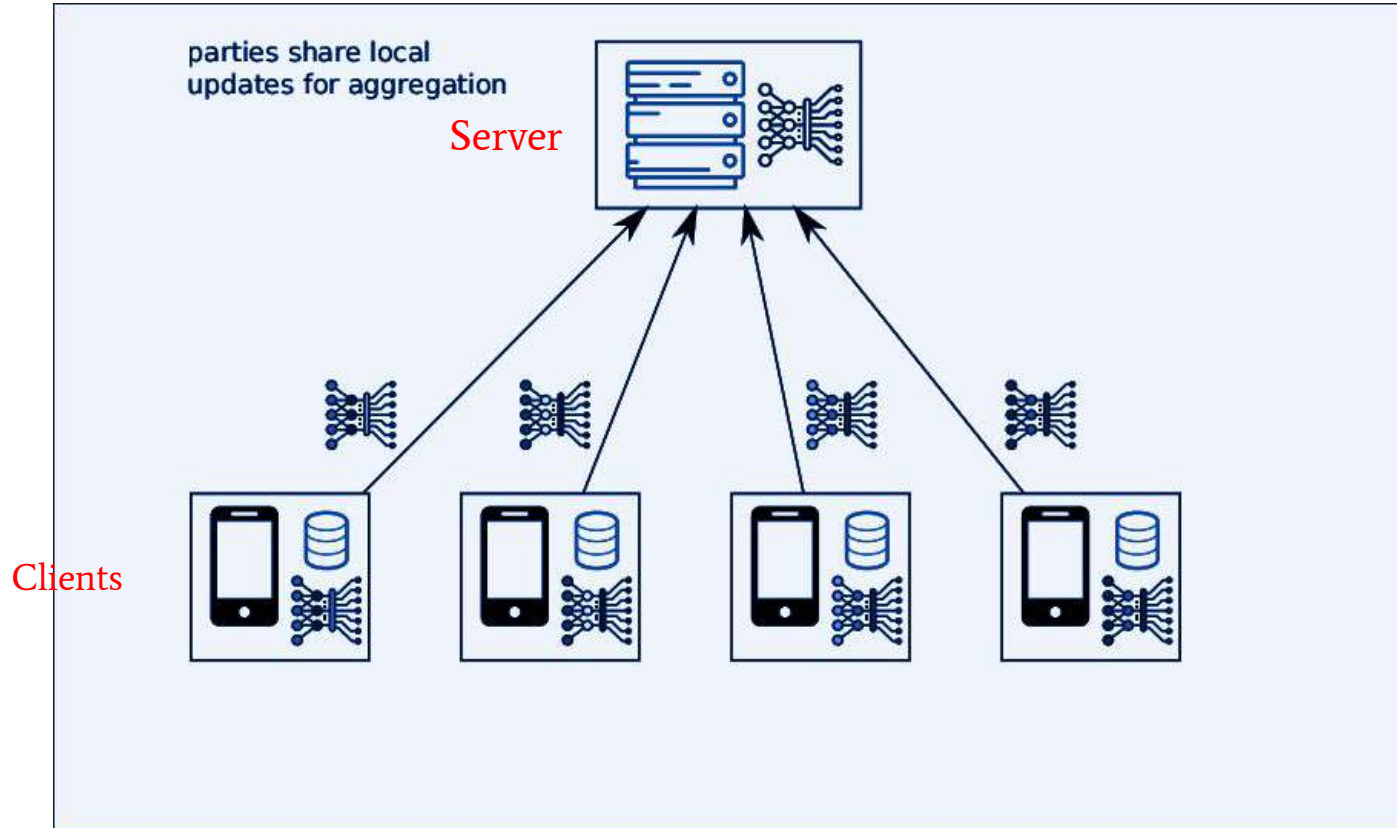
Step1: Initialization of Model



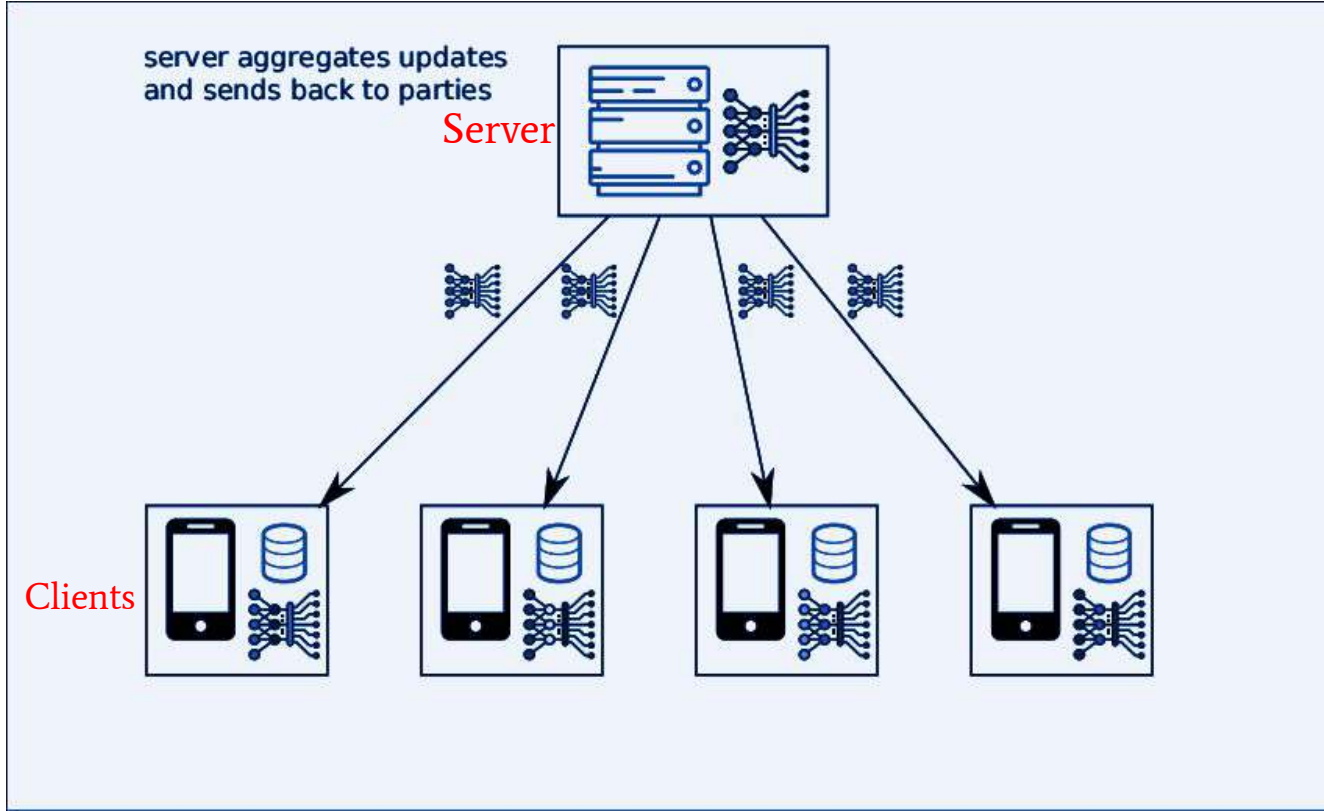
Step2: Training the model on client data



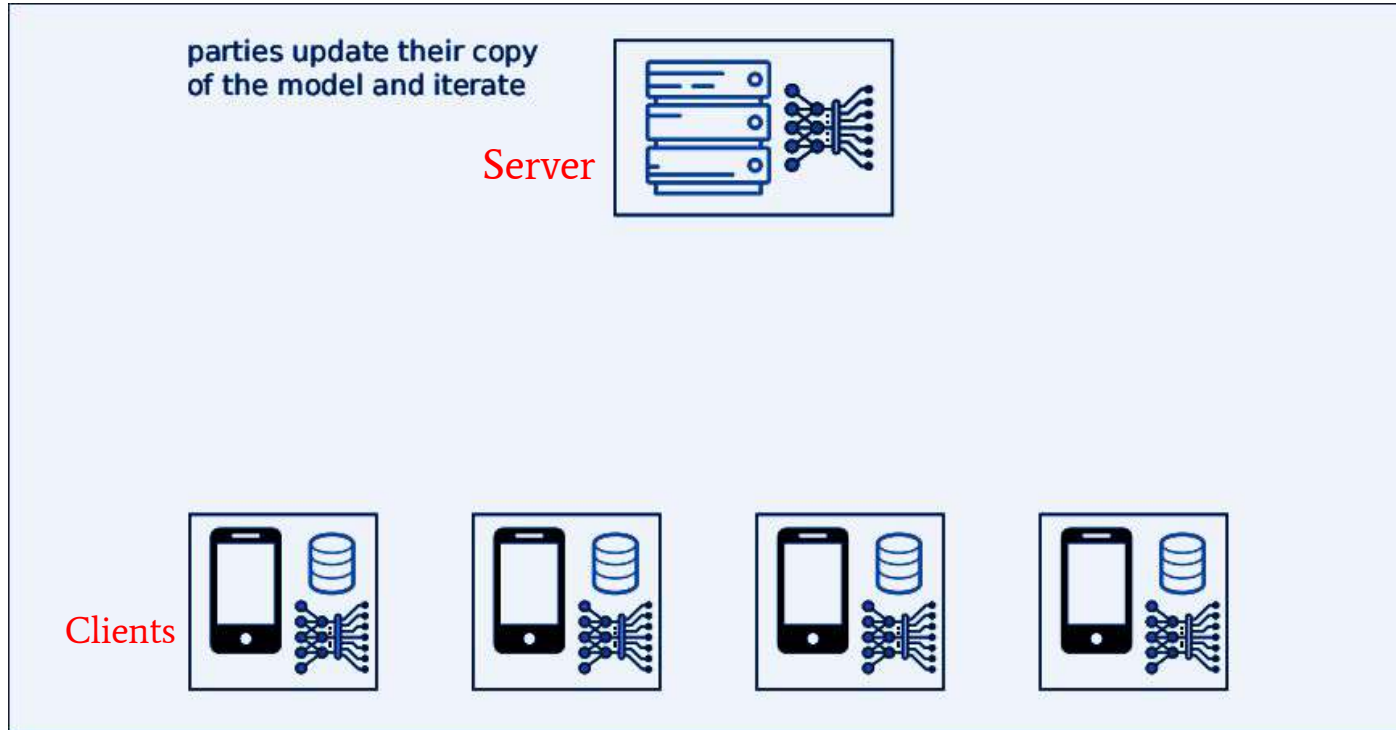
Step3: Sharing the parameters to Global model



Step4: Aggregation of parameters at server

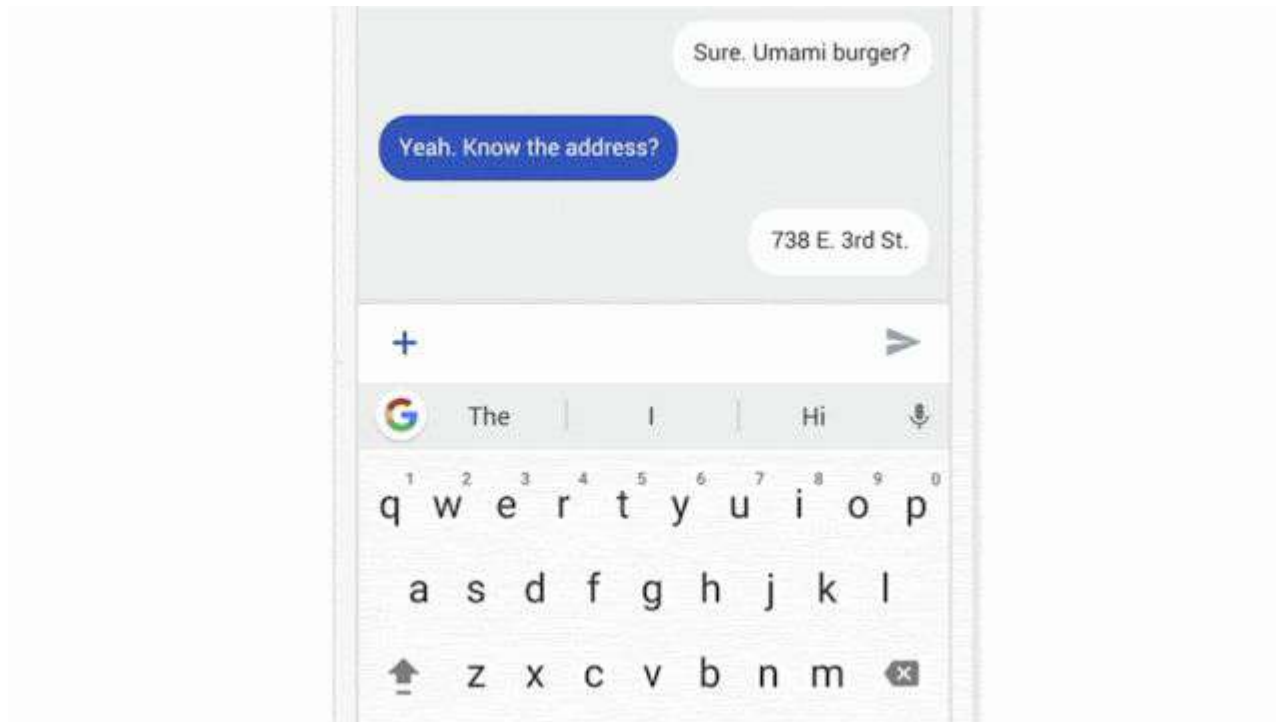


Step5: Sharing the updated parameters to clients



Applications

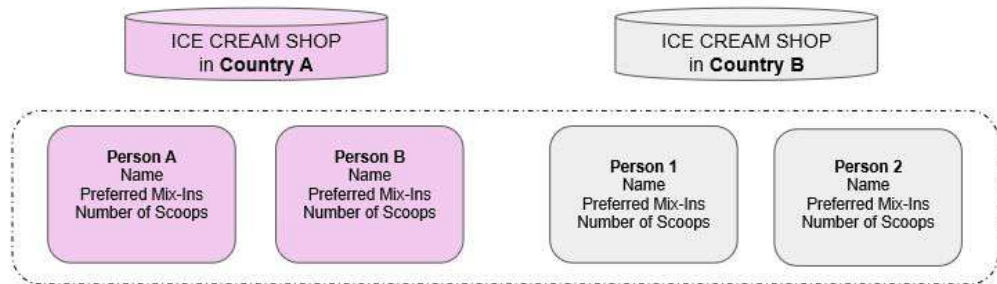
Google G-Board - Federated GRU(Gated Recurrent Unit)



Types of federated learning - Data partitioning

Horizontal Federated Learning

- Sample-based federated learning or homogenous federated learning
- Involves separating the data that has the same features but operates within a different sample space

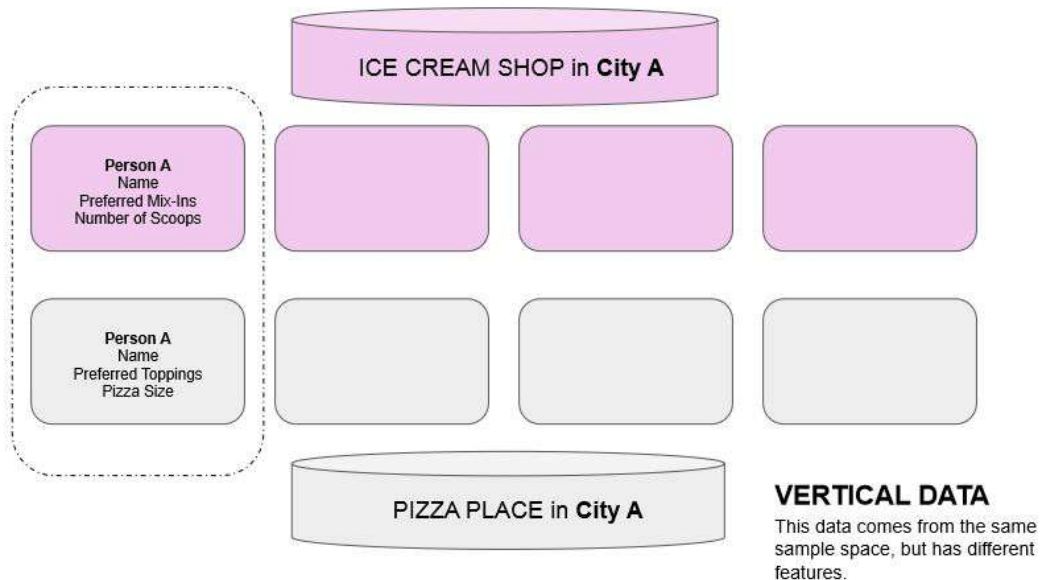


HORIZONTAL DATA

This data comes from different sample spaces, but has very similar features.

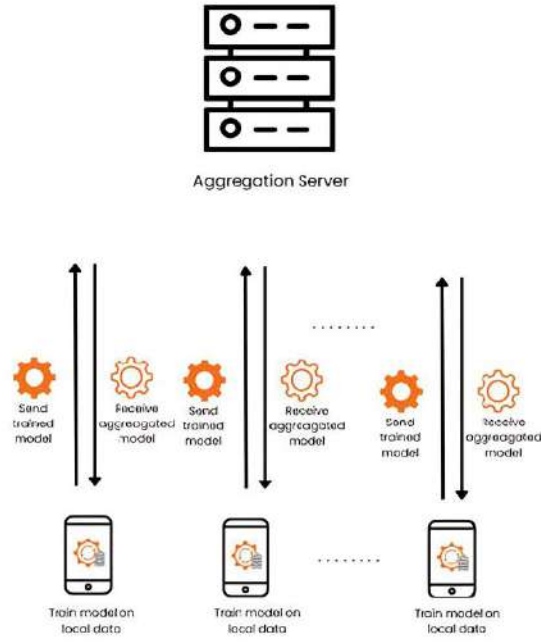
Vertical Federated Learning

- Feature-based federated learning or heterogeneous federated learning.
- Data shares the same sample space, but different feature space



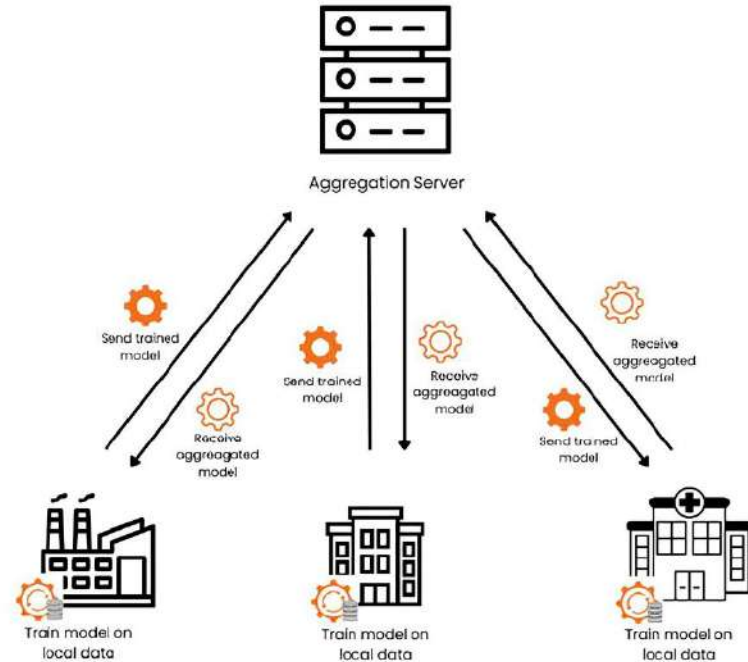
Types of federated learning - Clients

Cross-Device vs Cross-Silo FL



- Small size node (smartphones, edge devices, etc)
- Might not be available at each iteration.

a) Cross-Device Federated Learning



- Large size node (companies organizations)
- Necessary to participate in each iteration

b) Cross-Silo Federated Learning

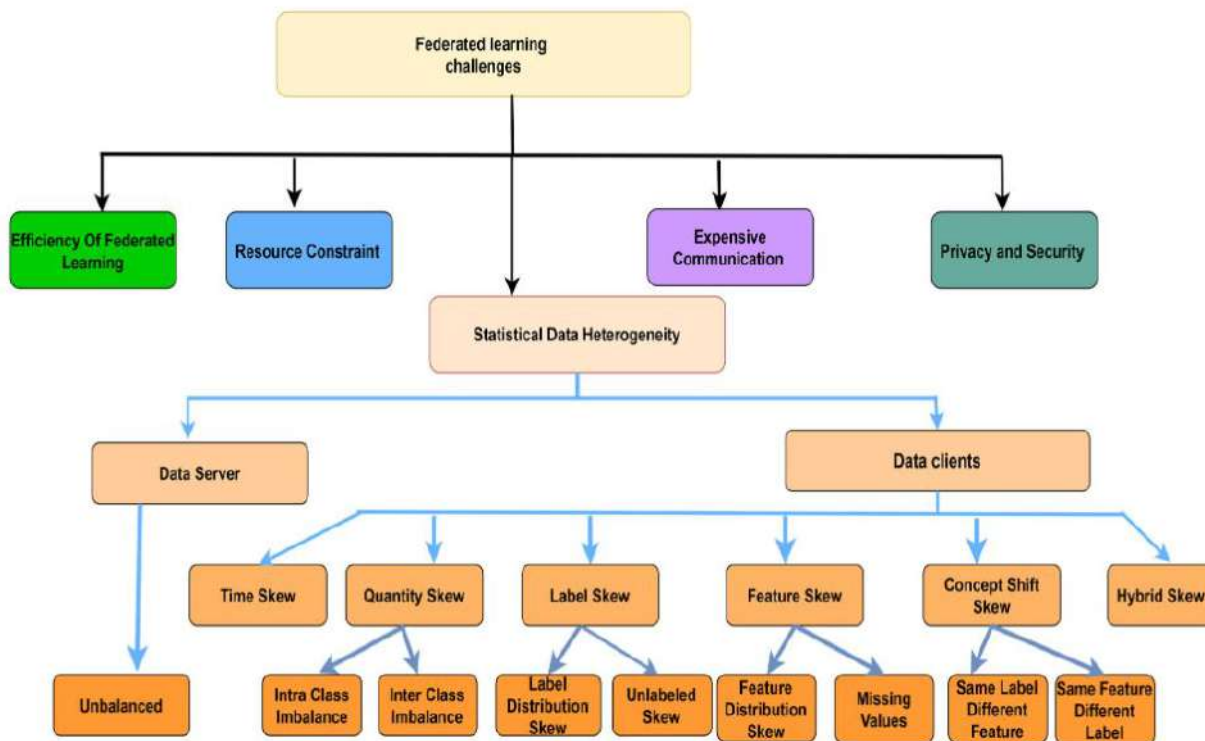
Cross-device FL

- Massive number of parties
- Small dataset per party (could be size 1)
- Limited availability and reliability
- Some parties may be malicious
- Communication is often the primary bottleneck.
- Partitioning by example (horizontal)

Cross-silo FL

- 2-100 parties
- Medium to large dataset per party
- Reliable parties, almost always available
- Parties are typically honest
- Might be computation or communication.
- Example-partitioned (horizontal) or feature-partitioned (vertical).

The main challenges in Federated learning ?



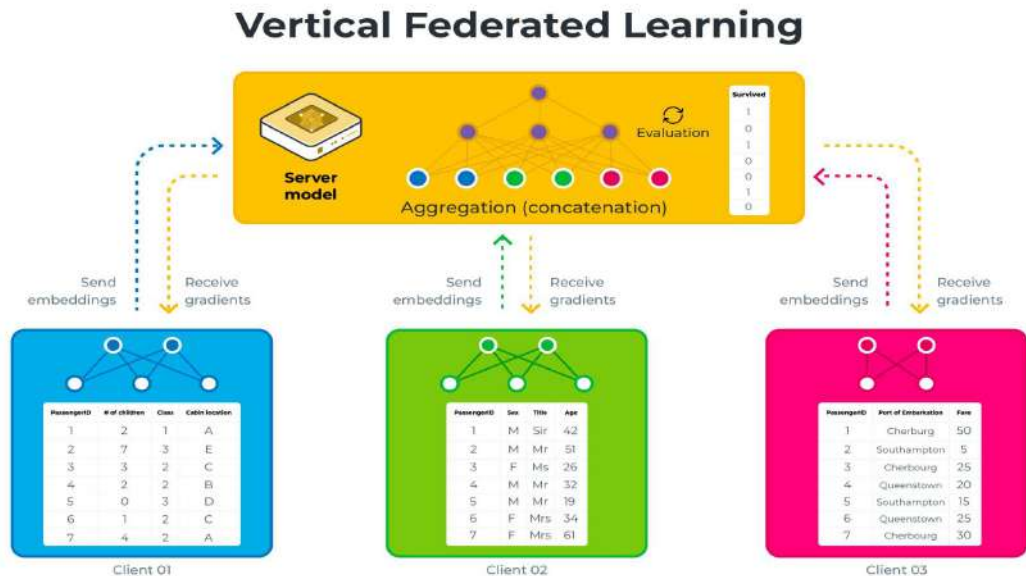
Problem Statement

FL systems gaining wider adoption for **privacy preserving machine learning**

- Heterogeneity is expected to cause **lowered performance of the trained models** with **longer convergence time**.
- Leading to **excessive energy consumption** for both the cloud infrastructure and battery powered devices.
- Find Out Innovative algorithms for FL task with **lower convergence** time with **minimal impact on data privacy**.

The main challenges in Vertical Federated learning

- One of the main challenges in vertical federated learning is **aggregation of the weights**



AGGREGATION ALGORITHMS

FEDAVG- (Communication-Efficient Learning of Deep Networks from Decentralized Data)

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
     $w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t)$ 
   $m_t \leftarrow \sum_{k \in S_t} n_k$ 
   $w_{t+1} \leftarrow \sum_{k \in S_t} \frac{n_k}{m_t} w_{t+1}^k$  // Erratum4
```

ClientUpdate(k, w): // Run on client k
 $B \leftarrow$ (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

FEDPROX- (Federated Optimization in Heterogeneous Networks)

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 γ_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

Algorithm 2 q -FedAvg

- 1: **Input:** $K, E, T, q, 1/L, \eta, w^0, p_k, k = 1, \dots, m$
- 2: **for** $t = 0, \dots, T - 1$ **do**
- 3: Server selects a subset S_t of K devices at random (each device k is chosen with prob. p_k)
- 4: Server sends w^t to all selected devices
- 5: Each selected device k updates w^t for E epochs of SGD on F_k with step-size η to obtain \bar{w}_k^{t+1}
- 6: Each selected device k computes:

$$\Delta w_k^t = L(w^t - \bar{w}_k^{t+1})$$

$$\Delta_k^t = F_k^q(w^t) \Delta w_k^t$$

$$h_k^t = q F_k^{q-1}(w^t) \|\Delta w_k^t\|^2 + L F_k^q(w^t)$$

- 7: Each selected device k sends Δ_k^t and h_k^t back to the server
- 8: Server updates w^{t+1} as:

$$w^{t+1} = w^t - \frac{\sum_{k \in S_t} \Delta_k^t}{\sum_{k \in S_t} h_k^t}$$

- 9: **end for**
-

Open Source Federated Learning Frameworks

FATE


TensorFlow
federated


NVIDIA FLARE

 **Flower**

OpenFL
intel

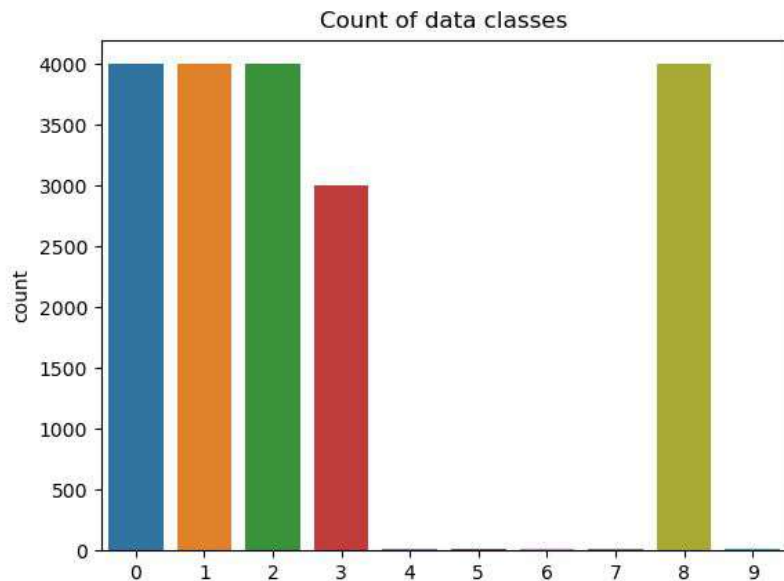
 PySyft

 **FedML**

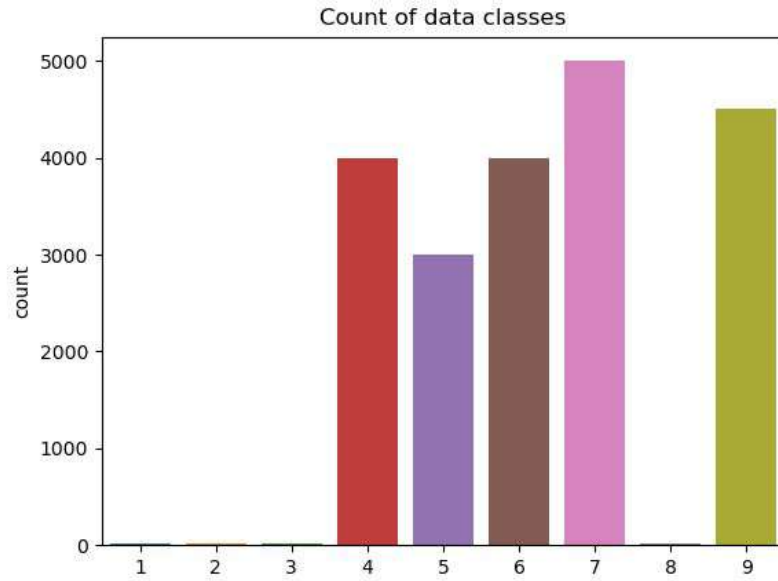
**Federated
Learning** 

Baseline Models - Flower Framework

Horizontal Federated Learning - MNIST dataset(with different class distributions to each client)



Client 1



Client 2

```

53 | server.py:222 | fit_round 10: strategy sampled 2
25 | server.py:236 | fit_round 10 received 2 results
Fit history: {'accuracy': [0.98750050045], 'val_accuracy': [0.725799977793884], 'loss': [0.05435093492269516], 'val_loss': [1.8747718334197998]}
Eval accuracy: 0.9534000158309937
Fit history: {'accuracy': [0.98950129747], 'val_accuracy': [0.676699977793884], 'loss': [0.04499607905745506], 'val_loss': [1.7222992181777954]}
Eval accuracy: 0.9405999779701233
Fit history: {'accuracy': [0.9871470332145691], 'loss': [0.036762069910764694], 'val_accuracy': [0.6499999761581421], 'val_loss': [1.7135902643203735]}
Eval accuracy: 0.9534000158309937
Fit history: {'accuracy': [0.9898734092712402], 'loss': [0.029244087636470795], 'val_accuracy': [0.6620000004768372], 'val_loss': [2.1011500358581543]}
Eval accuracy: 0.9473999738693237
Fit history: {'accuracy': [0.9913339614868164], 'loss': [0.026174629107117653], 'val_accuracy': [0.6406999826431274], 'val_loss': [1.861505389213562]}
Eval accuracy: 0.9520999789237976
Fit history: {'accuracy': [0.9929406046867371], 'loss': [0.020072871819138527], 'val_accuracy': [0.6643000245094299], 'val_loss': [1.7174336910247803]}
Eval accuracy: 0.9546999931335449
Fit history: {'accuracy': [0.992112934589386], 'loss': [0.021396998316049576], 'val_accuracy': [0.6766999959945679], 'val_loss': [1.8145411014556885]}
Eval accuracy: 0.9527999758720398
DEBUG flwr 2024-03-15 01:57:52,398 | connection.py:220 | gRPC channel closed
INFO flwr 2024-03-15 01:57:52,398 | app.py:398 | Disconnect and shut down
siva@siva-Swift-SF314-55G: ~/Desktop/Federated$

```

Client1 accuracy- 0.6766
 Client2 accuracy- 0.7257
 Global accuracy- 0.9527

Horizontal Federated Learning - MNIST fashion dataset(with different class distributions to each client)

```
server.py:173 | evaluate_round 10: strategy
server.py:187 | evaluate_round 10 received
server.py:153 | FL finished in 36.2526219430
siva@siva-Swift-SF314-55G: ~/Desktop/Federated
6 Fit history: {'accuracy': [0.9194741845130, 'val_accuracy': [0.5965999960899353], 'val_loss': [0.76419997215271]
Eval accuracy: 0.76419997215271
Fit history: {'accuracy': [0.9271665215492, 'val_accuracy': [0.6291000247001648], 'val_loss': [0.7864999771118164]
Eval accuracy: 0.7864999771118164
Fit history: {'accuracy': [0.9318403005599, 'val_accuracy': [0.5871000289916992], 'val_loss': [0.7803999781608582]
Eval accuracy: 0.7803999781608582
Fit history: {'accuracy': [0.935589075088501, 'val_accuracy': [0.5968000292778015], 'val_loss': [3.11248779296875]}
Eval accuracy: 0.79830002784729
Fit history: {'accuracy': [0.9352483153343201], 'loss': [0.17073878645896912], 'val_accuracy': [0.6484000086784363], 'val_loss': [2.4795634746551514]}
Eval accuracy: 0.8104000091552734
Fit history: {'accuracy': [0.9387536644935608], 'loss': [0.1595885306596756], 'val_accuracy': [0.5792999863624573], 'val_loss': [3.3508050441741943]}
Eval accuracy: 0.7509999871253967
Fit history: {'accuracy': [0.9414800405502319], 'loss': [0.15539324283599854], 'val_accuracy': [0.5788000226020813], 'val_loss': [3.385594367980957]}
Eval accuracy: 0.8256000280380249
DEBUG flwr 2024-03-15 02:25:47,183 | connection.py:220 | gRPC channel closed
INFO flwr 2024-03-15 02:25:47,183 | app.py:398 | Disconnect and shut down
siva@siva-Swift-SF314-55G: ~/Desktop/Federated$
```

Client1 accuracy- 0.7106
Client2 accuracy- 0.5788
Global accuracy-0.8256

EXPERIMENT SETUP - 2

Centralized VS Decentralized

DATASET : TITANIC DATASET

```
Running centralised training...
Train accuracy: 84.248%
Test accuracy: 82.022%
Running decentralised training...
Iteration 1, loss = 0.64399825
Iteration 2, loss = 0.56226789
Iteration 3, loss = 0.49493030
Iteration 4, loss = 0.45459877
Iteration 5, loss = 0.43634872
Iteration 6, loss = 0.43703215
Iteration 7, loss = 0.43900626
Iteration 8, loss = 0.42954290
Iteration 9, loss = 0.42684795
Iteration 10, loss = 0.42327625
Iteration 11, loss = 0.42247988
Iteration 12, loss = 0.42048271
Iteration 13, loss = 0.41931750
Iteration 14, loss = 0.41798792
Iteration 15, loss = 0.41693441
Client 0 test accuracy: 73.034%
Client 1 test accuracy: 81.461%
Client 2 test accuracy: 71.348%
Combined test accuracy: 80.337%
```

EXPERIMENT SETUP - 2

Centralized VS Decentralized

DATASET : CANCER DATASET

```
Training loss did not improve more than tol=0.000100 for 10 consecutive epochs.  
Stopping.  
client 0 test accuracy: 96.512%  
client 1 test accuracy: 94.186%  
client 2 test accuracy: 94.186%  
combined Test accuracy: 96.512%
```

```
Running centralised training...  
Train accuracy: 98.758%  
Test accuracy: 95.349%
```

Implementation

EXPERIMENT SETUP -3

Experimenting with different aggregating algorithm

1.Fedavg

2.FedProx

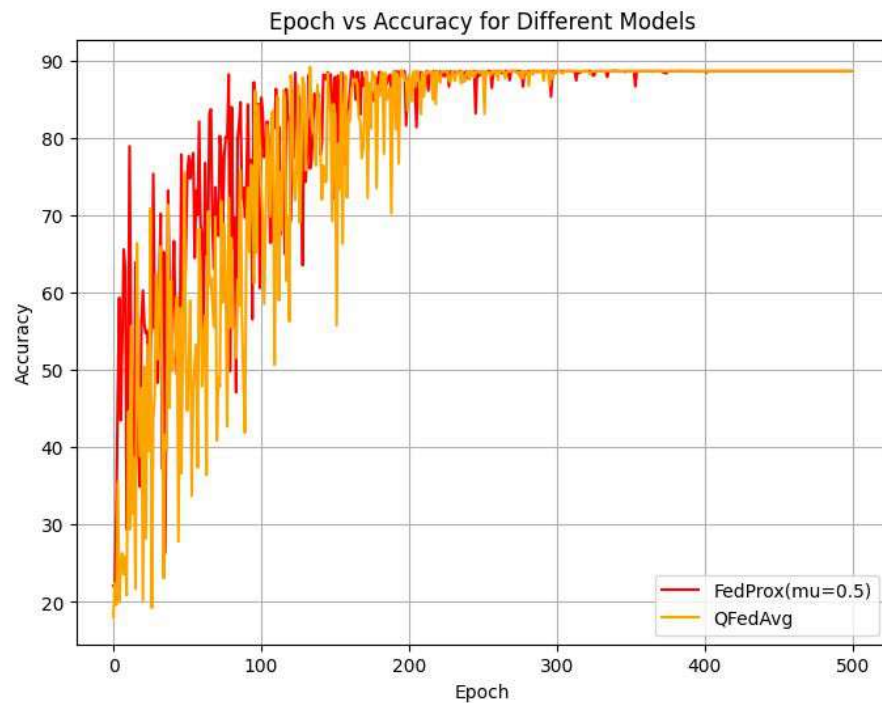
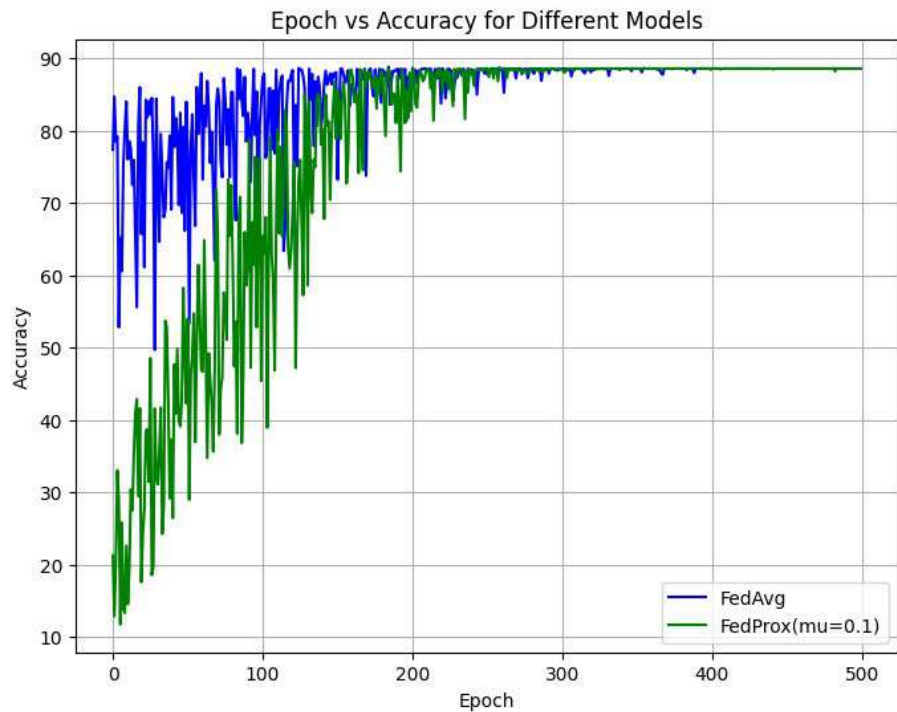
3.QFedavg

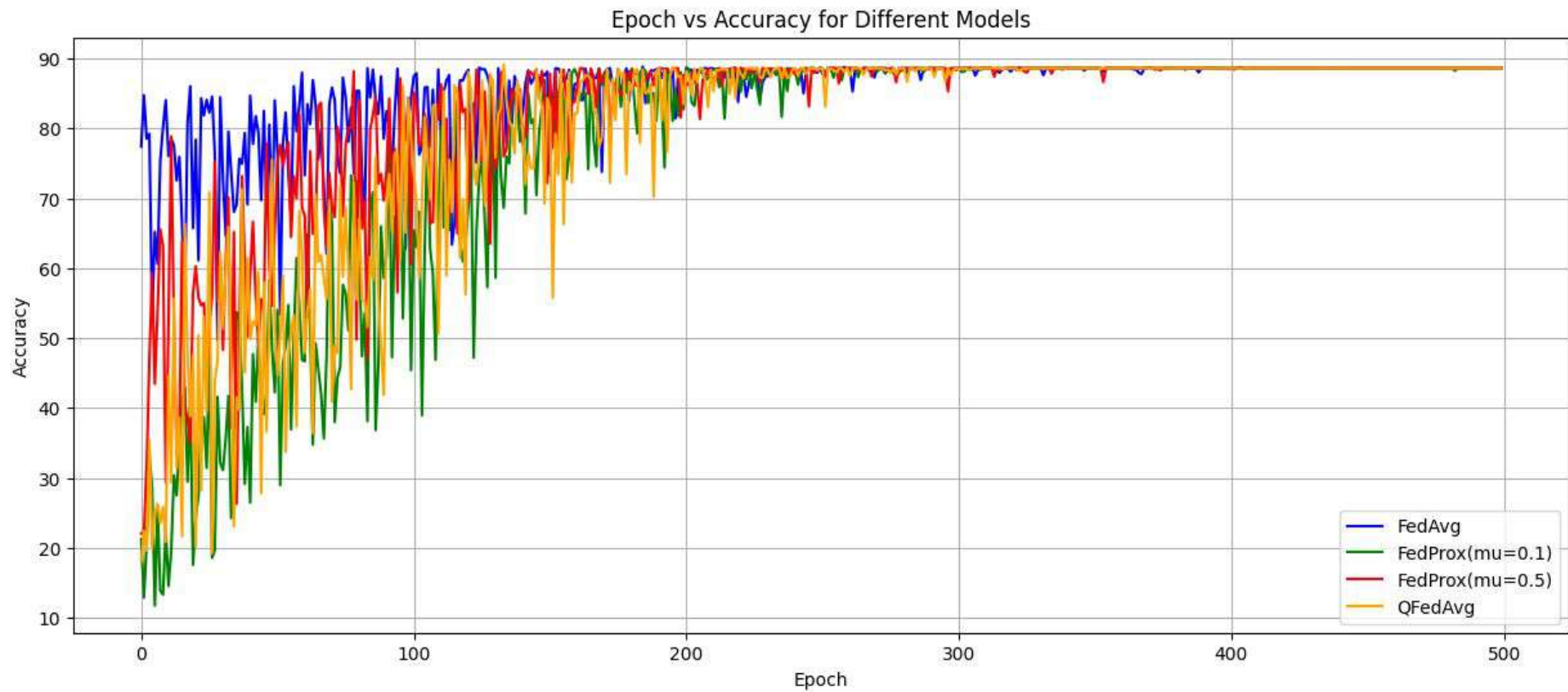
DATASET: TITANIC DATASET, WATER QUALITY DATASET

Water Quality Dataset

- Dataset - Classification Problem(Two class)(20 features)(7999 data points)
- Number of clients : 3
- Client 1 - Features {'aluminium', 'ammonia', 'arsenic', 'barium', 'cadmium', 'chloramine', 'chromium', 'is_safe'}
- Client 2 - Features {'copper', 'flouride', 'bacteria', 'viruses', 'lead', 'nitrates', 'nitrites', 'is_safe'}
- Client 3 - Features{'mercury' 'perchlorate' 'radium' 'selenium' 'silver' 'uranium', 'is_safe'}

Results & Observations

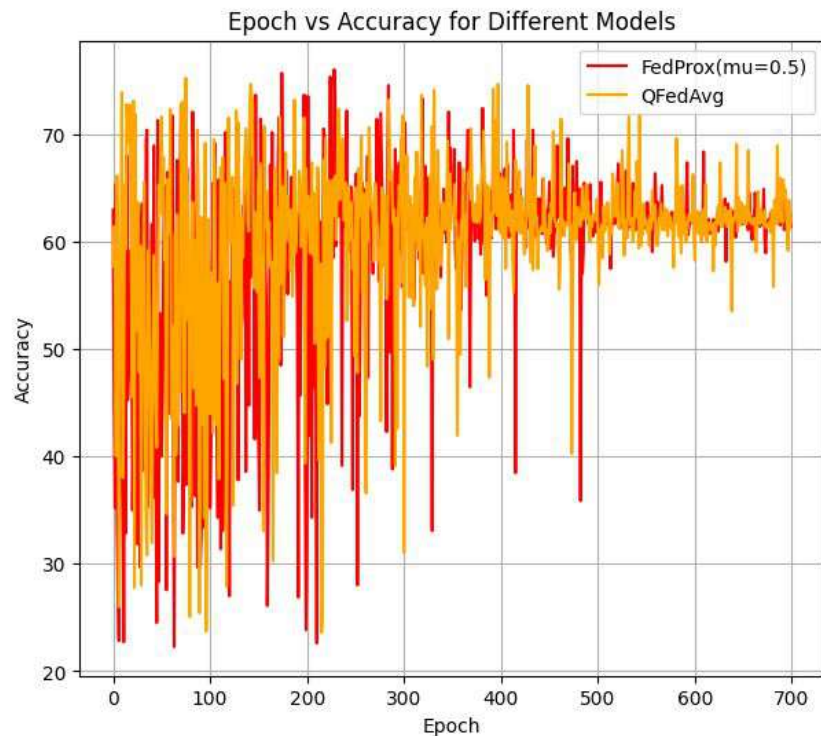
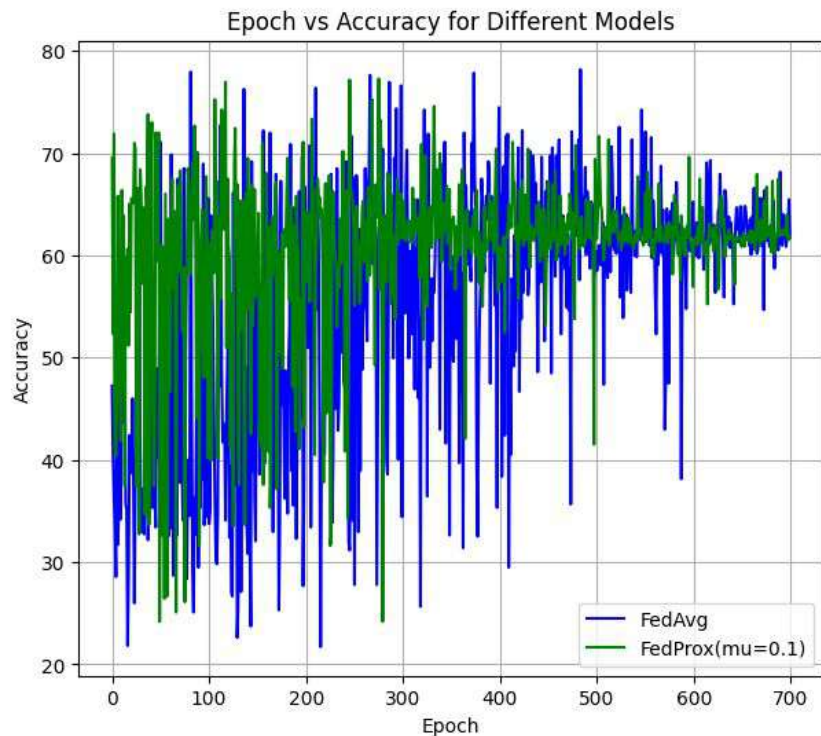




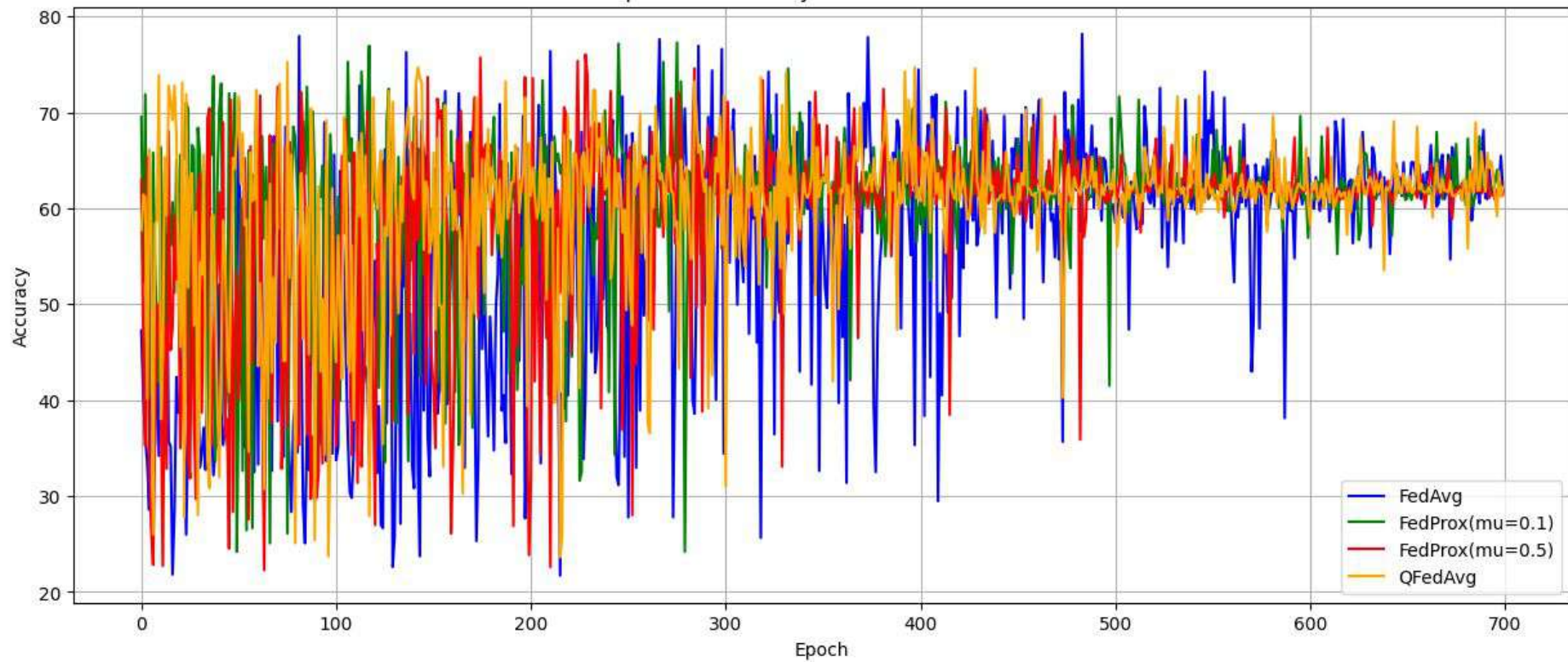
Titanic Dataset

- Dataset - Classification Problem(Two class)(12 features)(891 data points)
- Number of clients : 3
- Client 1 - Features {`"Parch", "Cabin", "Pclass", "Survived"`}
- Client 2 - Features {`"Sex", "Title", "Survived"`}
- Client 3 -Features{`'Age', 'SibSp', 'Embarked', Cabin`}

Results & Observations



Epoch vs Accuracy for Different Models



References

Github link for resources(Research papers and codes for frameworks):

1. <https://github.com/monk1337/Aweome-Heathcare-Federated-Learning?tab=readme-ov-file>
2. <https://github.com/albarqouni/Federated-Learning-In-Healthcare?tab=readme-ov-file>
3. <https://github.com/adap/flower>(flower framework)
4. <https://github.com/FedML-AI/FedML>-(fedml framework)-<https://github.com/FedML-AI/FedML/tree/master/python/fedml>
5. <https://github.com/OpenMined/PySyft>(pysyft framework)

Thank you