

Federated Deep Learning for Healthcare Data

Erik Babu

23 June 2020

Supervisor: Prof. Daniel Rueckert

Motivation

Motivation

- › Study was conducted on the accuracy of cancer diagnoses in 2015
- › 16 observers had to decide whether images of breast tissue were cancerous
- › Observers achieved 99% classification accuracy (majority voting)

Motivation

› The identity of the observers?



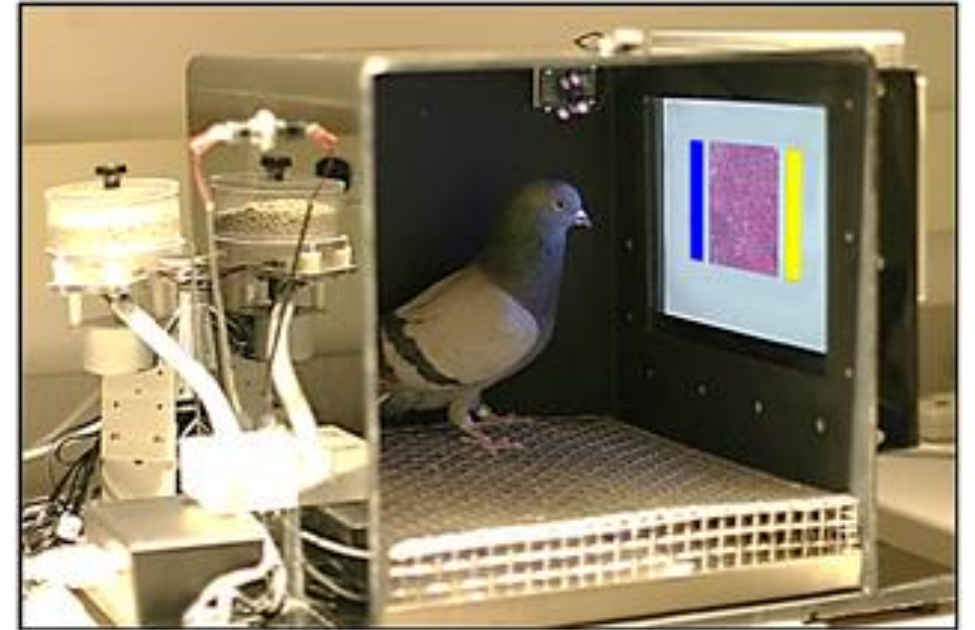
Motivation

› The identity of the observers?



Motivation

› The identity of the observers?



Identifying patterns in medical data is **not a uniquely human skill!**

Motivation

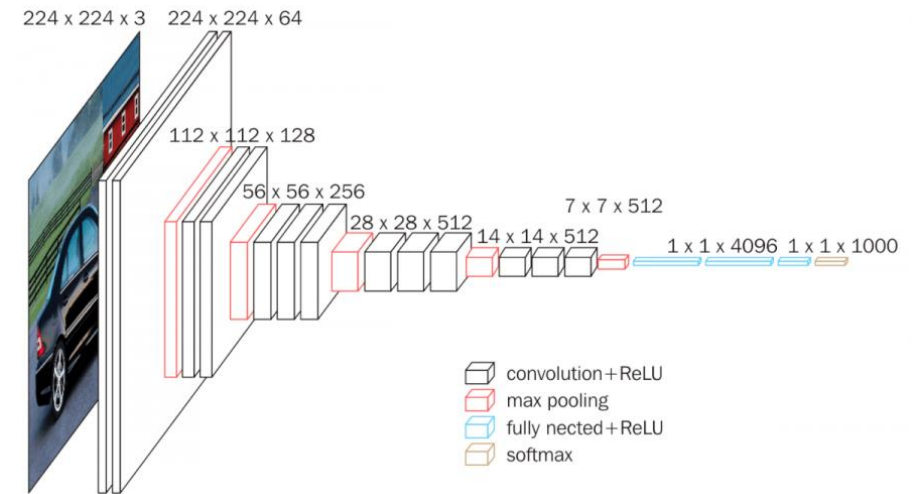
Identifying patterns in medical data is **not** a uniquely human skill!

- › DL algorithms excel at uncovering patterns from data
- › They have the potential to address medical diagnosis problems
 - But they require **large training sets** to achieve **high performance**

Motivation

SOTA DL classification models are data-hungry

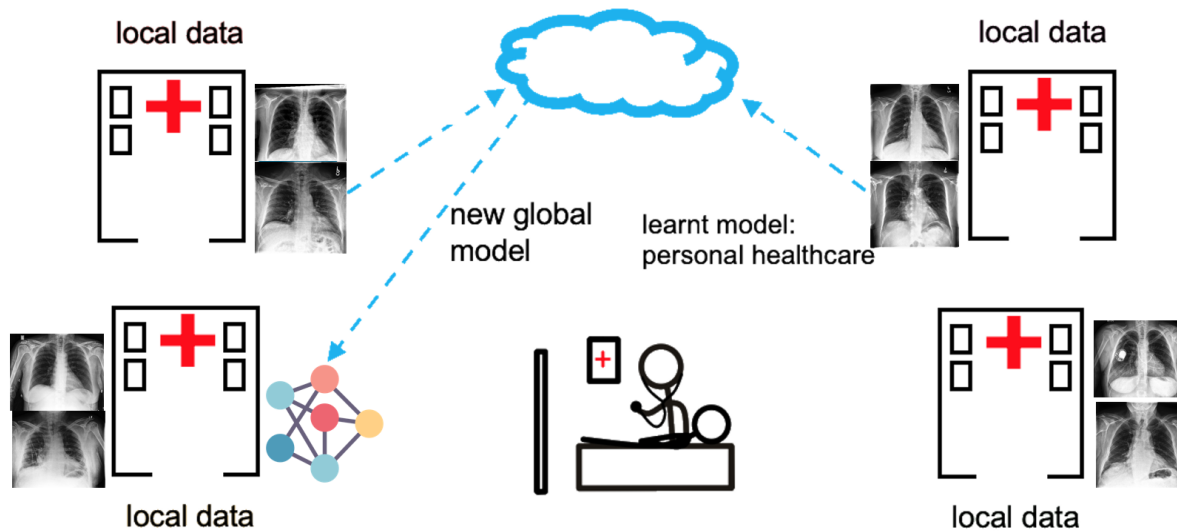
- › E.g. VGG16
- › More than 138M trainable parameters
- › All parameters tuned for model to converge to a solution
- › Insufficient training data leads to poor approximation



VGG16

Possible Solution

Aggregate the raw data centrally



Enter your message: [expand message input area](#)

Send

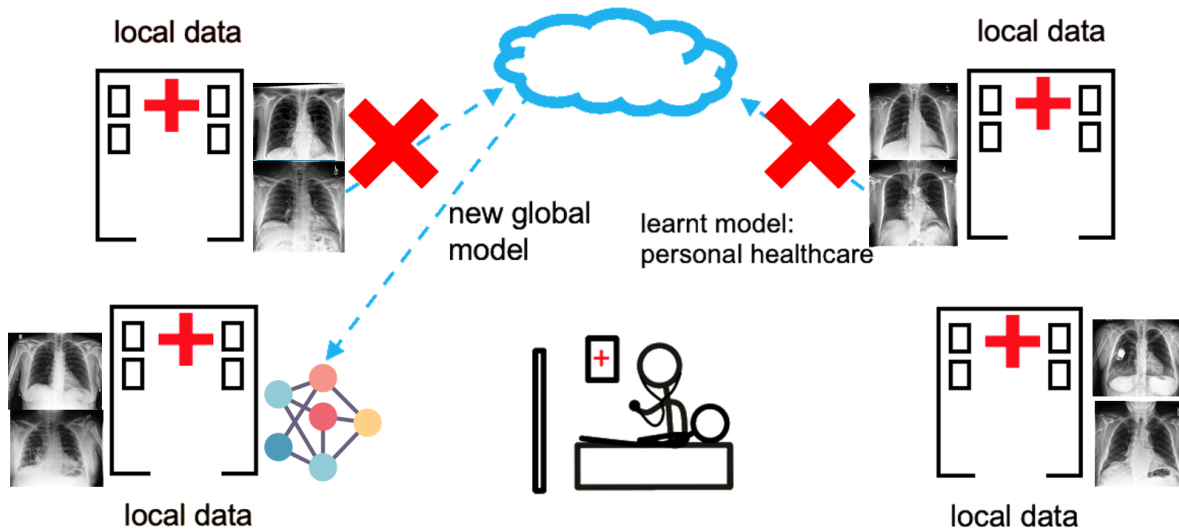
Server: Send all your images pls
Hospital 1: Ok
Hospital 2: Sure
Hospital 3: Coming up

Live Chat

On new message: ☒ Bring to top

Impossible Solution.

Aggregate the raw data centrally



Enter your message: [expand message input area](#)

Send

Server: Send all your images pls
Hospital 1: Ok
Hospital 2: Sure
Hospital 3: Coming up

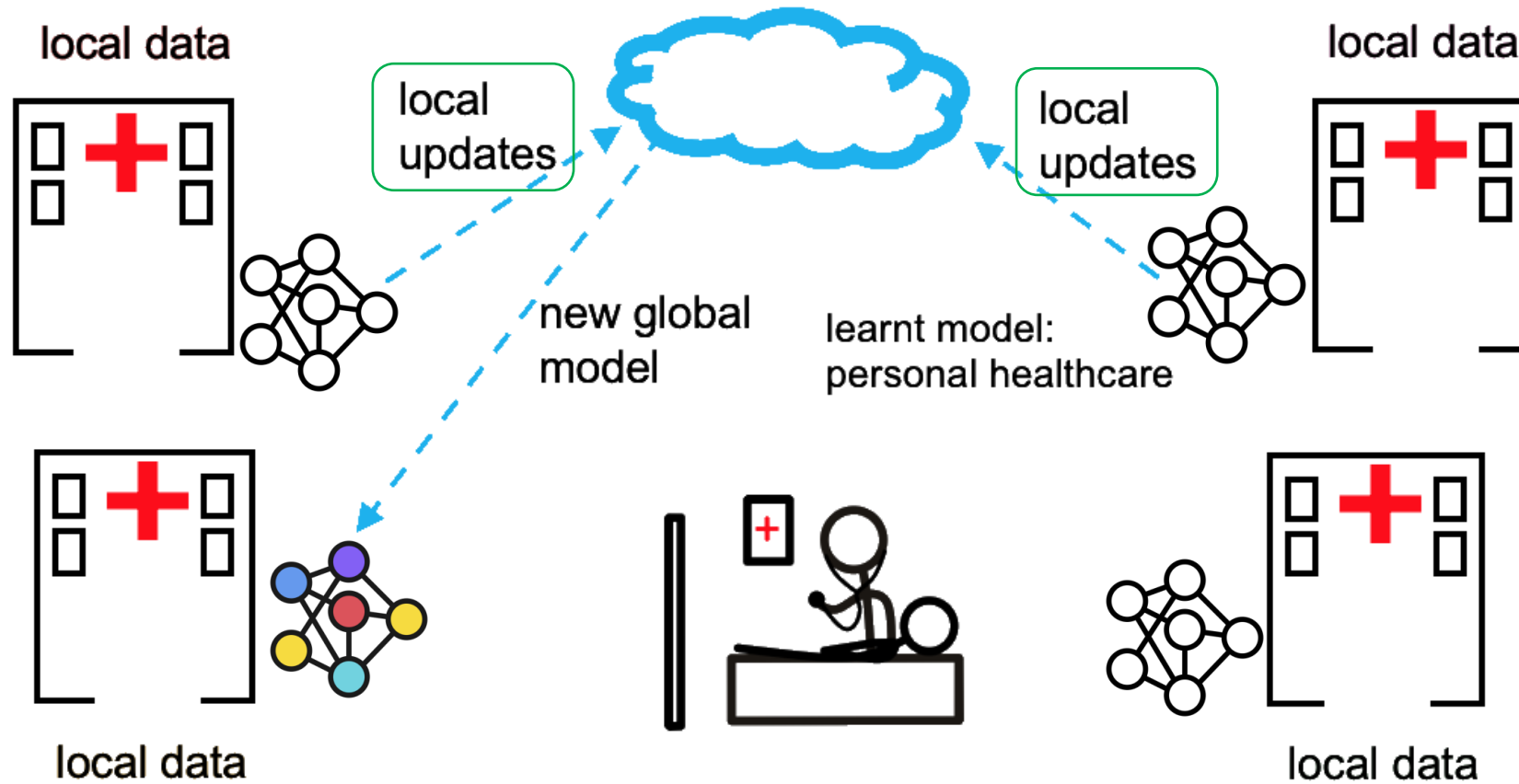
*** GDPR has entered the chat ***

Live Chat

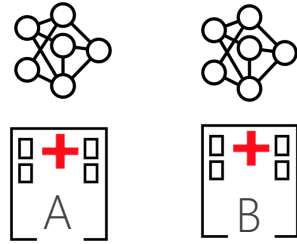
On new message: ☒ Bring to top

Our Solution

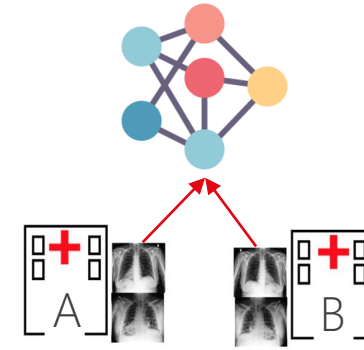
Federated Learning (FL)



Terminology



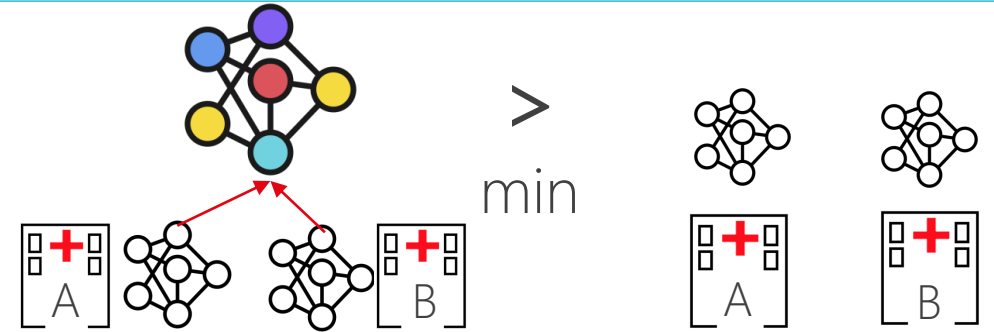
- › Institutional/Baseline models
- › Expected **lower bound**



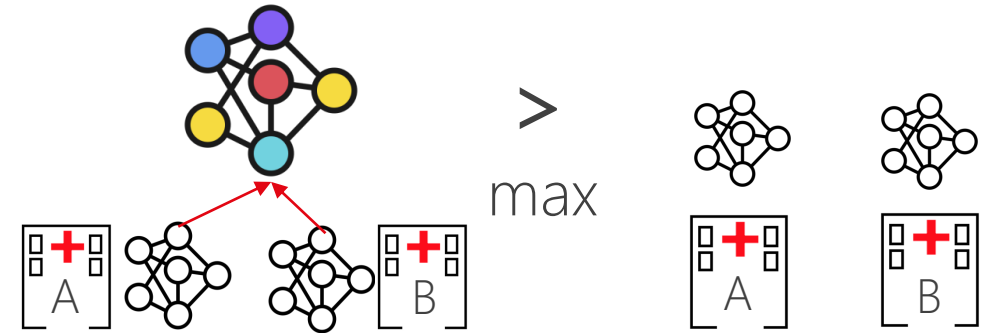
- › Centralised/Benchmark model
- › Expected **upper bound**

Investigation

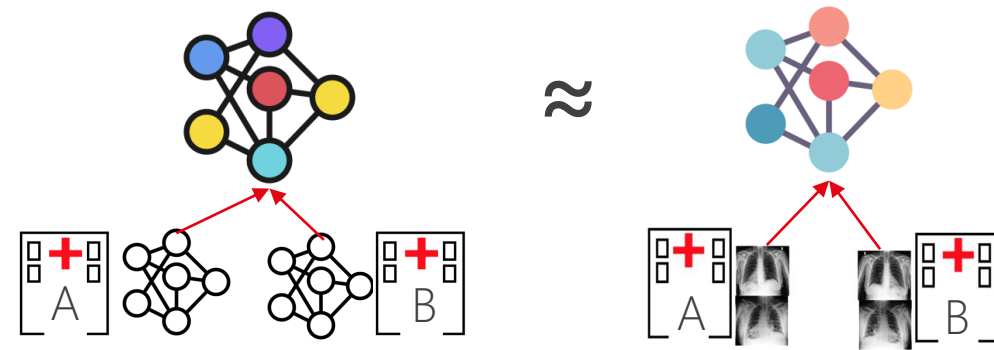
› Can the FL models outperform **any** baseline models?



› Can the FL models outperform **all** baseline models?

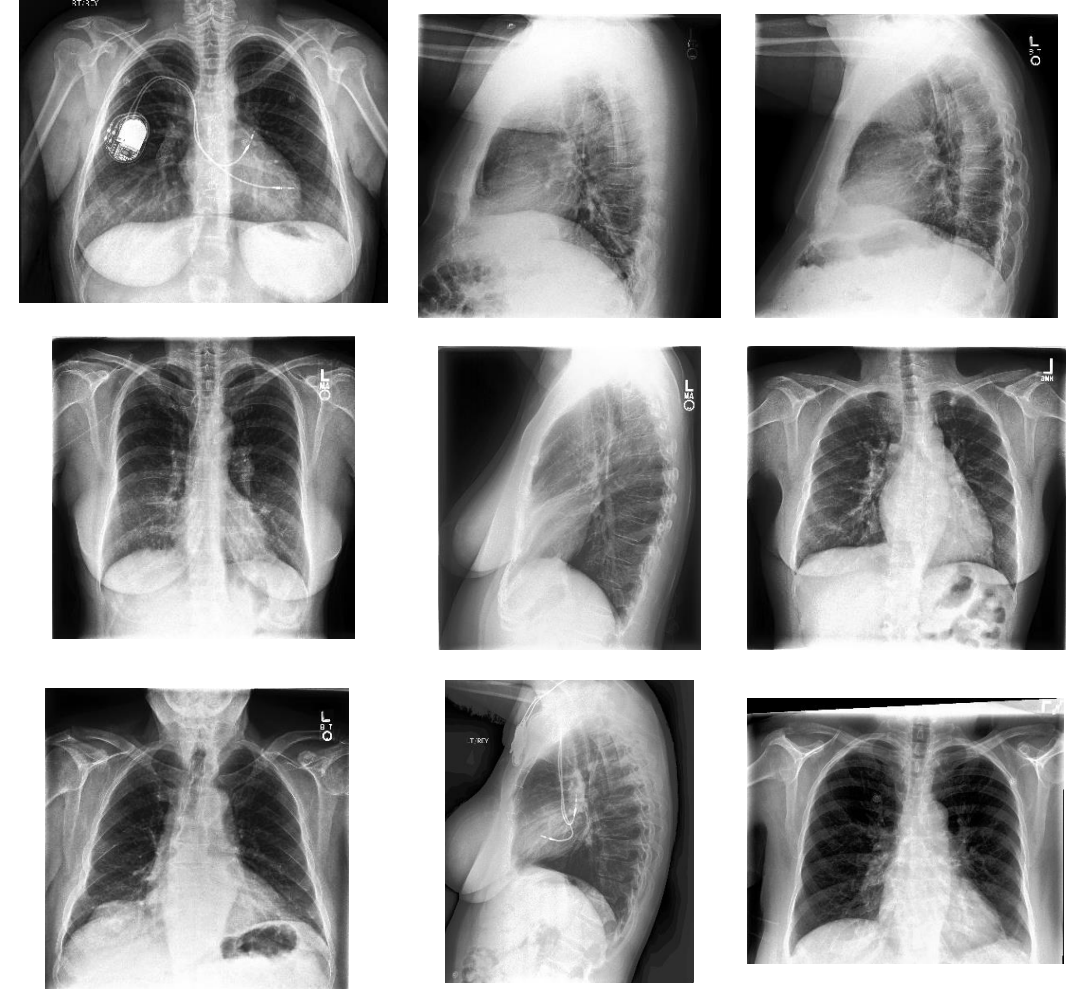


› Can the FL models perform **similarly** to the benchmark model?



Dataset

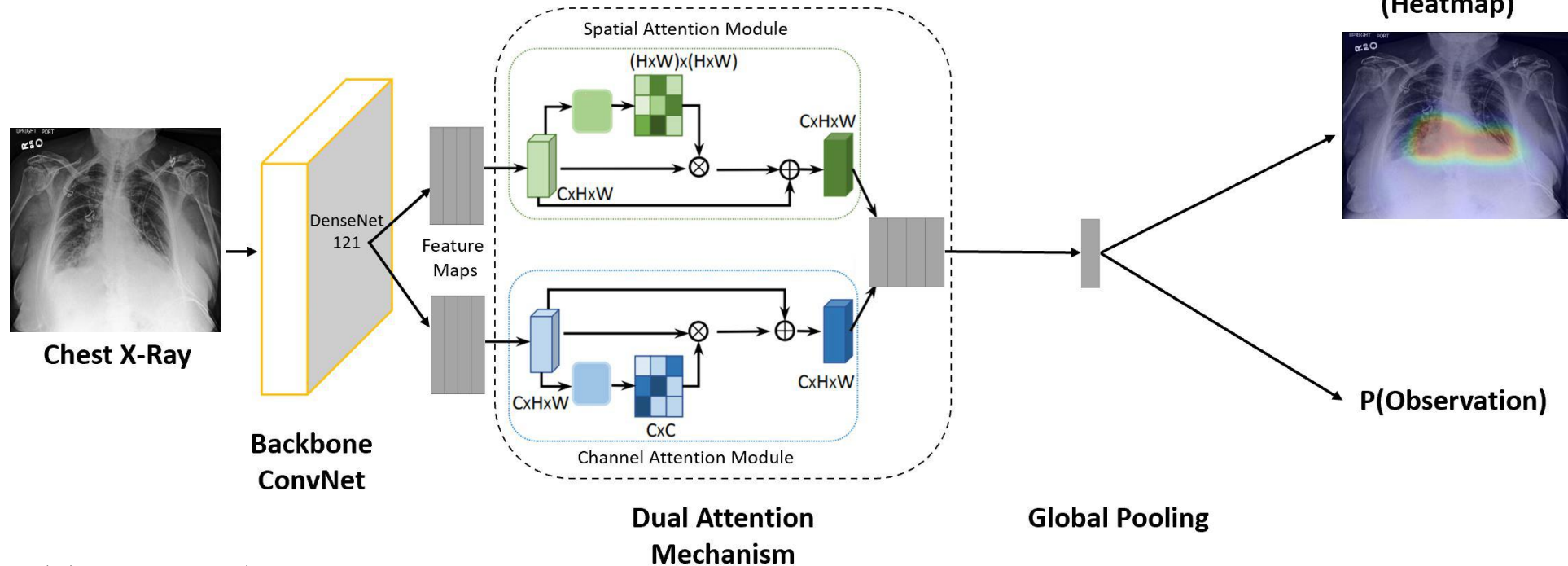
- › CheXpert (downsampled)
- › 224,316 multi-view radiographs
- › 14 typical chest observations (focus on 5)
- › Uses ROC-AUC as evaluation metric
- › Report Section 3.2



Implementation

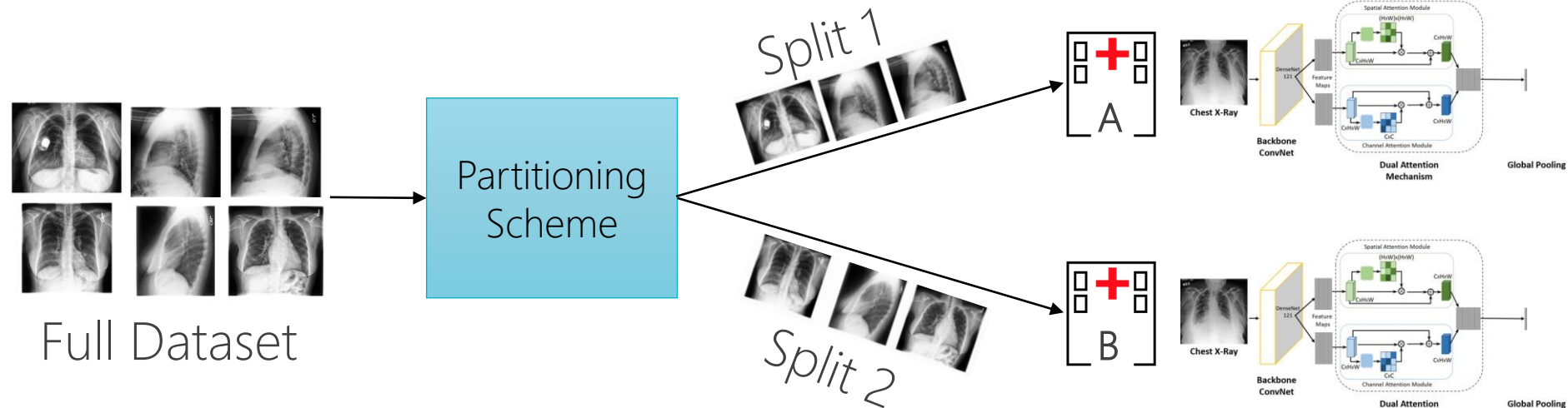
Benchmark

- › Off-the-shelf implementation (4th highest on CheXpert leaderboard)
- › Trained on full dataset -> **0.901** mean AUC
- › Report Section 3.3



Baseline

- › Same underlying model as benchmark
- › Trained on partition of dataset
- › No “collaboration” between institutions



Baseline

Data Partitioning Scheme

- › Equal & unequal amounts between 2 institutions (50/50 & 75/25)
- › Equal amounts between 5 and 10 institutions (stretch goal)
- › Non-i.i.d. partitioning:
 - › Iteratively sample from full dataset without replacement
 - › Assign different weights to the observations to skew label distribution
- › Report Appendix A

- › Various FL frameworks considered (Report Section 4.1)
- › Initially proceeded with PySyft
 - ✓ Similar API to PyTorch
 - ✓ Multiple available tutorials
 - ✓ Abstractions to help with implementing ϵ -DP and SMPC

- › Eventually abandoned PySyft
- ✗ Tutorials not always up-to-date with API
- ✗ Cannot assign custom / uneven data partitions to clients
- ✗ Momentum not yet implemented
- › No FL framework used

FL

- › Manually implement training loops
 - › Wrap main training procedure in external loops

```
for e in range(num_epochs):  
    train(model, data)
```

FL

- › Inner-loop -> local training at an institution
- › First outer-loop -> local training at all institutions

```
for client in clients:  
    model = initialise_model(global_weights)  
    data = client.data  
  
    for e in range(num_epochs):  
        train(model, data)  
  
    updates.append(model.weights)
```

FL

- › Outer-loop -> communication rounds (server-aggregation)

```
for cr in range(num_comm_rounds):  
    updates = []  
  
    for client in clients:  
        model = initialise_model(global_weights)  
        data = client.data  
  
        for e in range(num_epochs):  
            train(model, data)  
  
        updates.append(model.weights)  
  
    global_weights = server.aggregate(updates)
```

- › Manually implement training loops (Report Section 4.3)
 - › Wrap main training procedure in external loops
 - › Inner-loop: local rounds of training at all institutions
 - › Outer-loops: communication rounds (server-aggregation)
- › Manually implement aggregation code
- › Institutions use data partitions from baseline approach

Approach #1: FedAvg (server-side)

- › Equally weight updates according to number of participants

$$w_{t+1} \leftarrow \boxed{\frac{1}{K}} \sum_{k=1}^K w_{t+1}^k$$

- › Weight updates according to proportion of data contributed by each participant

$$w_{t+1} \leftarrow \sum_{k=1}^K \boxed{\frac{n_k}{n}} w_{t+1}^k$$

Approach #1: FedAvg (server-side)

$$w_{t+1} \leftarrow \frac{1}{K} \sum_{k=1}^K w_{t+1}^k$$

```
1 def fed_avg(weights):
9     n_clients = len(weights)
10    factor = float(1 / n_clients)
11
12    # Perform averaging
13    avg = weights[0]
14    for k in avg.keys():      # iterate through all weight parameters
15        for i in range(1, n_clients):
16            avg[k] += weights[i][k]
17            avg[k] = torch.mul(avg[k], factor)
18
19    # Averaged weights to be sent back to participants
20    return avg
```

Approach #2: FedProx (client-side)

- › Used in conjunction with FedAvg on server-side
- › Convergence guarantees in presence of non-i.i.d. data

$$\text{Loss} = \text{Loss} + \frac{\mu}{2} ||w - w_k||^2$$

- › Adds regularization term to client training loss, based on difference between global and local weights
- › μ hyperparameter determines aggressiveness of regularisation

Approach #2: FedProx (client-side)

$$\text{Loss} = \text{Loss} + \frac{\mu}{2} ||w - w_k||^2$$

```
1 # Get loss as usual
2 loss = get_loss(output, target)
3
4 # Compute FedProx regularisation term and update loss value
5 reg = 0.0
6 for param_index, param in enumerate(model.parameters()):
7     reg += ((mu / 2) * torch.norm((param - global_weights[param_index]))**2)
8
9 loss += reg
10
11 # Compute gradient of loss as usual
12 loss.backward()
```

Measuring Overhead

Why track communication and computational overhead?

- › FL is a distributed learning problem
- › Existing research tends to focus solely on model performance
- › Serves as baseline for future optimisations to be evaluated against

Measuring Overhead

Computational Overhead

- › THOP: Pytorch-OpCounter library
- › Institutions in centralised approach incur 0 computational overhead
- › $\text{FLOPs} = \text{flops_in_batch} * \text{batches_in_epoch} * \text{total_epochs}$
- › $\text{fl_total_epochs} = \text{local_epochs} * \text{comm_rounds}$

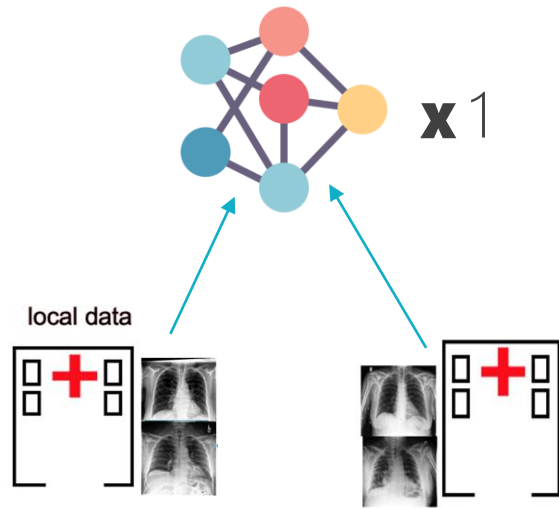
Measuring Overhead

Communication Overhead

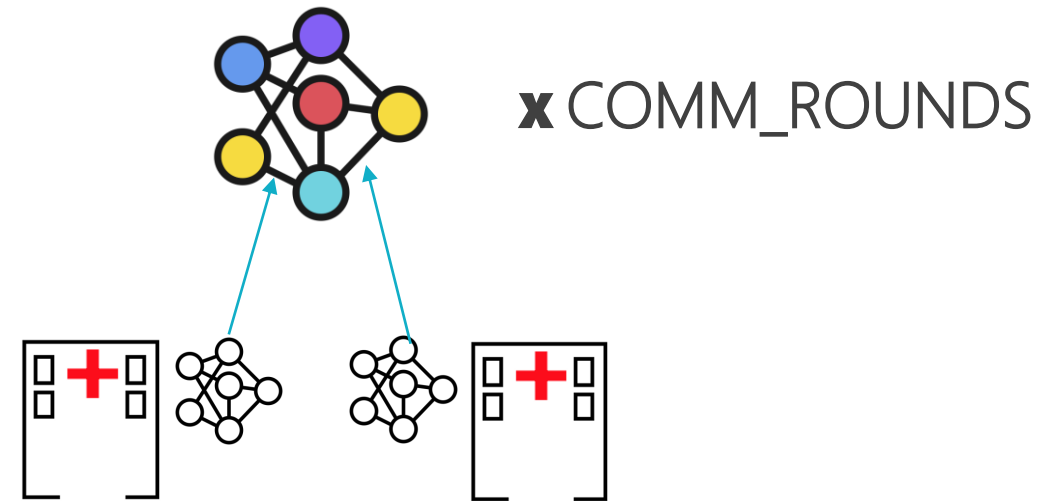
- › Implemented using our own modelling assumptions
- › Institutions in baseline approach incur 0 communication overhead
- › $\text{comm_overhead} = \text{total_bytes_uploaded} + \text{total_bytes_downloaded}$

Measuring Overhead

Communication Overhead (Institution Upload)



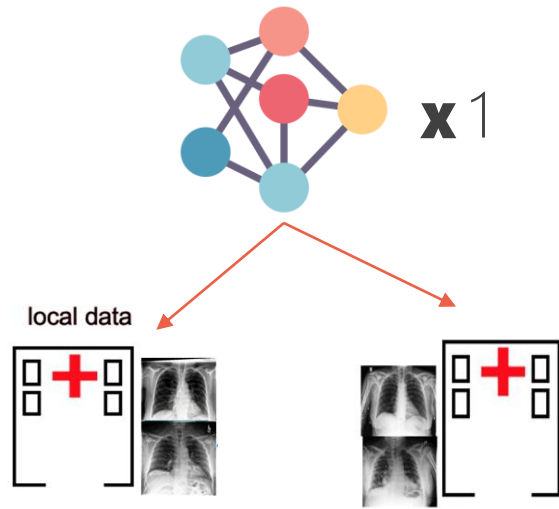
Centralised Approach



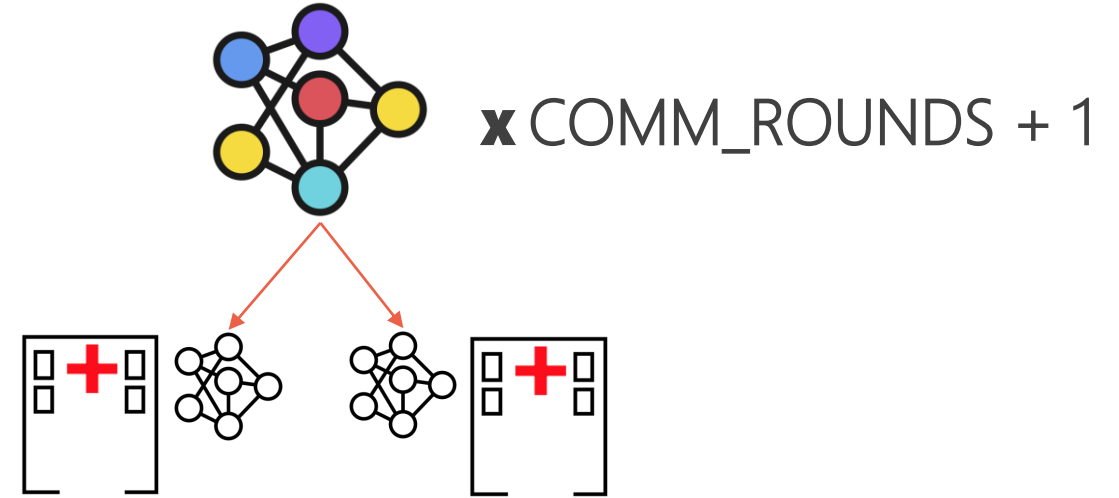
FL Approach

Measuring Overhead

Communication Overhead (Institution Download)



Centralised Approach

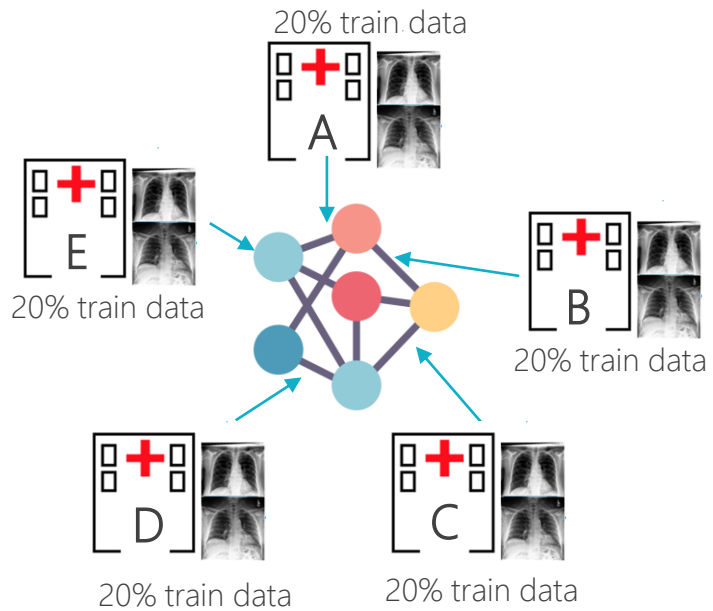


FL Approach

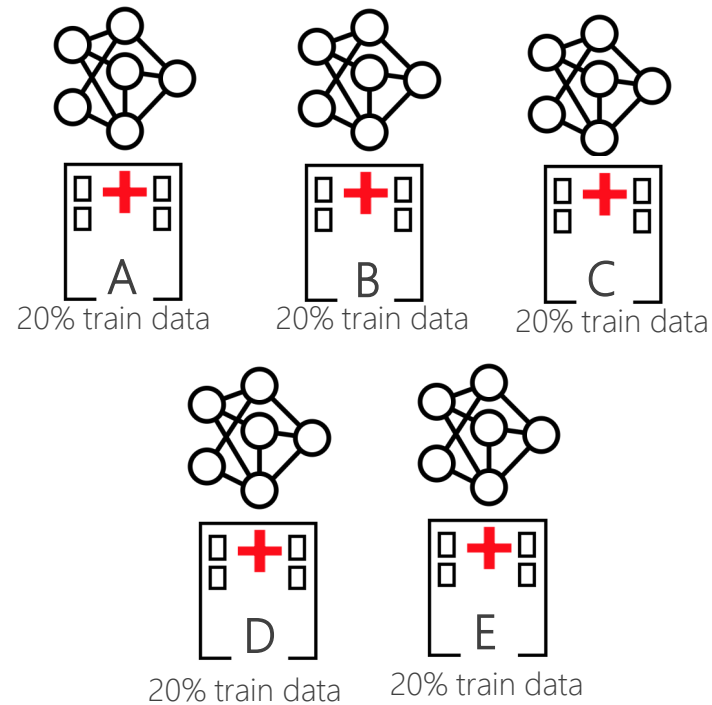
Evaluation

Five Institutions

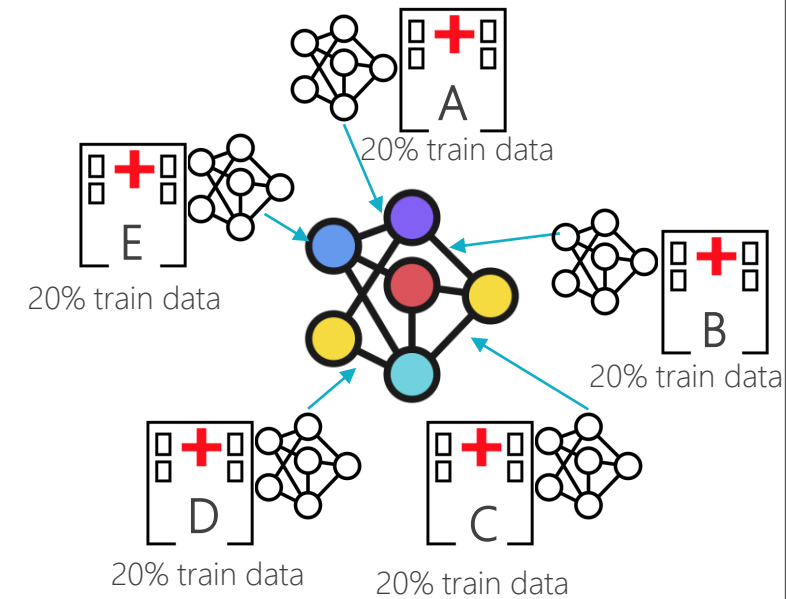
Centralised (Benchmark)



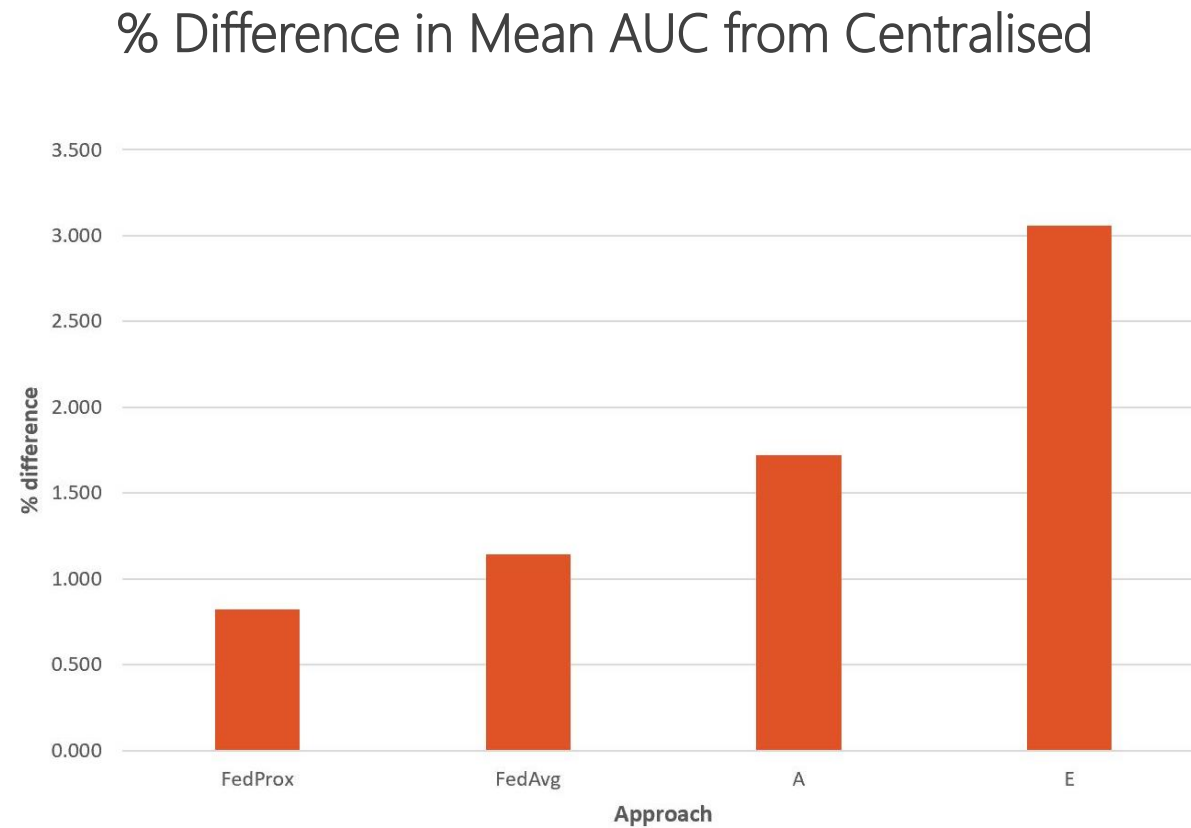
Institutional (Baseline)



FL



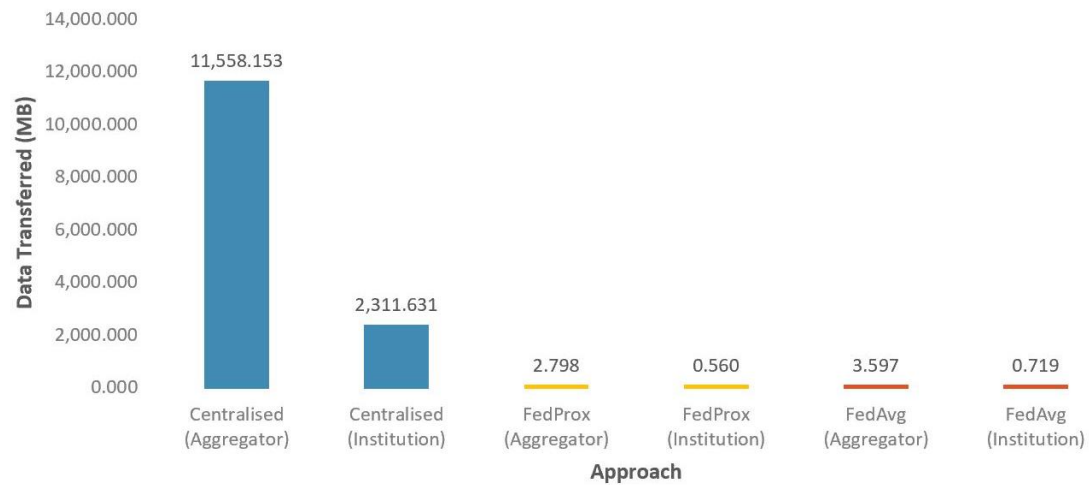
• Five Institutions •



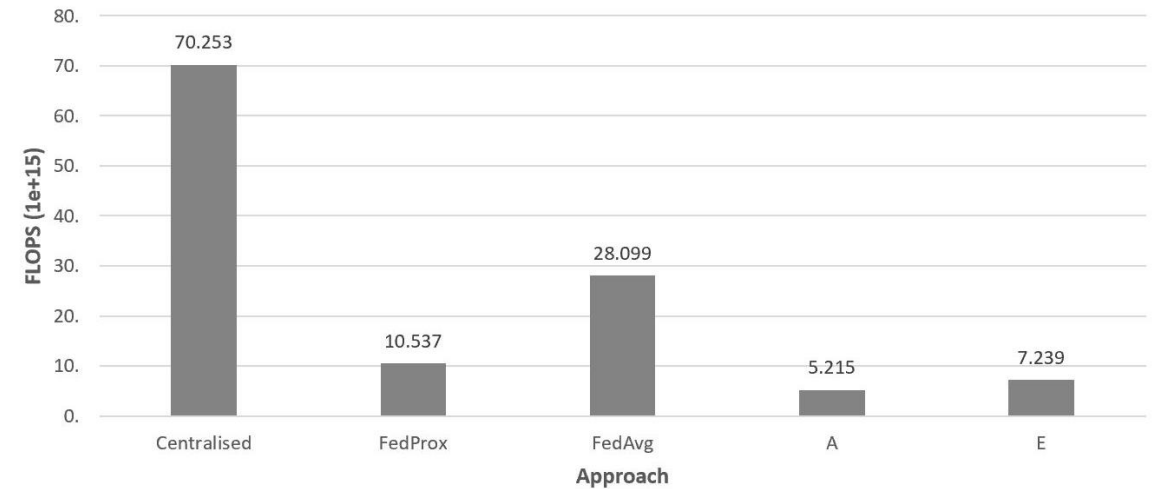
Model Performance

Five Institutions

Communication Overhead



Computational Overhead



Training Overhead

• Ten Institutions •

- › Shows similar trends to five institution split:
 - › Avg baseline model performance decreases (6.7% off benchmark)
 - › FL model performance remains similar (1.4% off benchmark)
 - › FedProx outperforms FedAvg (model performance & overhead)
- › Report Section 5

Summary

Summary

- › Implement 2 different FL training techniques (+ variant)
- › Develop method to track systems overhead
- › Compare model & systems performance between centralised, institutional, and FL approaches on 2 / 5 / 10 institutions
- › Best FL model achieves similar performance to benchmark (SOTA) in all experiments

THANKS FOR LISTENING

Questions?