

Privacy Enhancing Technologies for Gigapixel Medical Image Analysis

Jesse Cresswell, PhD

Senior Machine Learning Scientist
Layer 6 AI at TD

*Mar. 17 2022 - Endless Summer School:
Healthcare Roundup - Vector Institute*

Agenda

- Obstacles with Medical Data
- Privacy Enhancing Technologies
- Gigapixel Image Analysis
- Decentralized Private Learning

Obstacles with Medical Data

Obstacles

An individual's right to privacy is of utmost importance with medical data.

Privacy laws and regulations have hindered the development of ML in this space by limiting access to data.

ML advances occur with the release of large, diverse datasets.

But the healthcare data available to institutions is often

- Small scale
- Restricted to a single source (not multi-centric)
- Specialized (not diverse)
- Unlabeled (not annotated with regions of interest)
- Subject to legal/ethical review



Obstacles

Typical solutions to these obstacles may not be possible due to privacy concerns:

A. Data can't be pooled across institutions.

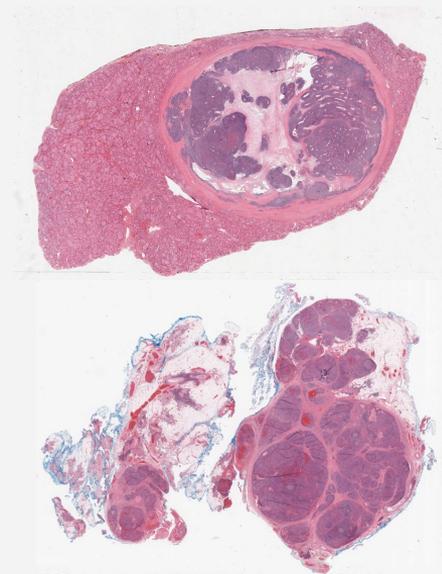
In addition, image sizes are increasing. In digital histopathology, images can easily have 50,000x50,000 pixel.

B. Large images can't be directly fed into NNs.

Obstacle B sounds like an ML problem we could solve.

Obstacle A sounds like a regulatory problem out of our scope.

Privacy Enhancing Technologies may provide a ML-type solution to the regulatory problem.





Privacy Enhancing Technologies

Privacy Enhancing Technologies

The aim of PETs are to **minimize the risk** to individuals that their personally identifiable data will be exposed, while **maximizing the utility** of that data for analysis.

Open data is useful, but not private. Siloed data is safe, but not useful.

Four emerging PETs actively being researched in ML:

Federated Learning	Differential Privacy
Secure Multi-Party Computation	Homomorphic Encryption

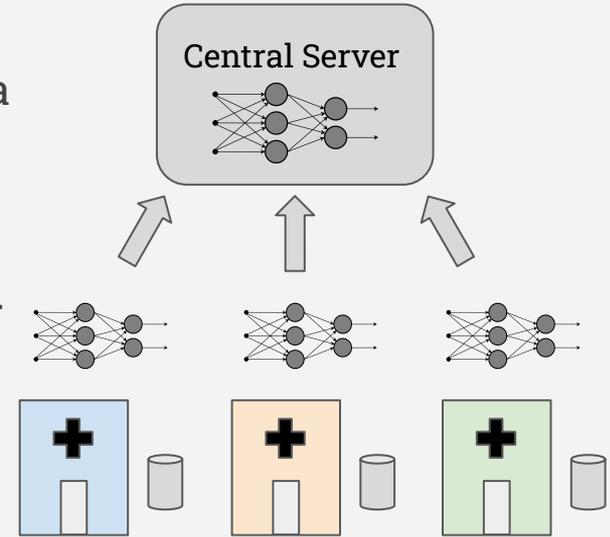
Federated Learning

Federated Learning is a distributed ML approach where data is not pooled together on a centralized server.

Models are trained at the institution where data is collected.

Only gradient updates are shared with the central server.

Intuitively this seems “more private”, since the raw data never leaves the institution where it was generated.

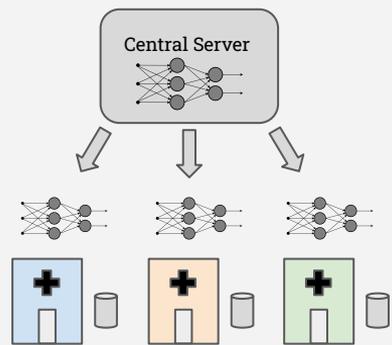


Federated Learning - *FedAvg*

[\[McMahan et al. 2017\]](#)

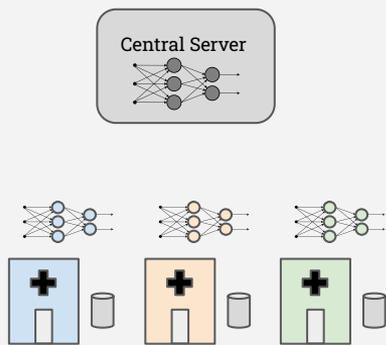
Repeat until convergence:

1.



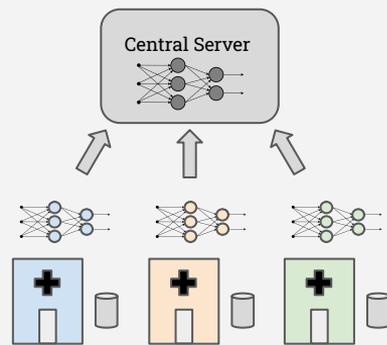
The central model is shared to each institution.

2.



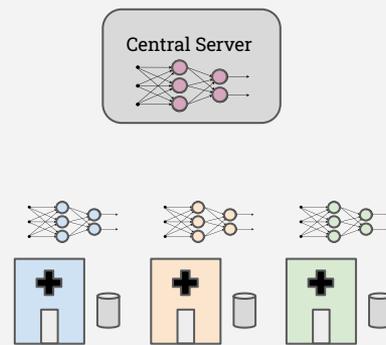
Institutions locally train the model on their data.

3.



Model updates (gradients) are sent to the server.

4.



The server averages the updates, and applies them to the central model

Federated Learning - Weaknesses

Memorization and reconstruction attacks

There are many examples of neural networks memorizing individual training datapoints, which can be recovered by analyzing the model. [\[Fredrikson, Jha, Ristenpart 2015\]](#)



Image seen by model during training



Reconstructed image inferred from trained model

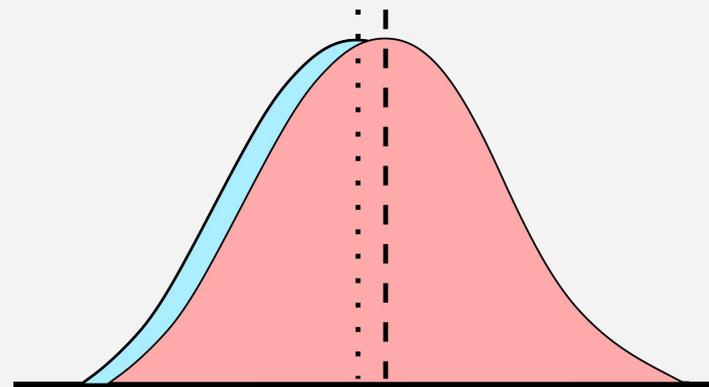
Differential Privacy

The previous example makes us think “is sharing model updates more private than sharing raw data?”

Key point: privacy is not binary. It is a resource.

Differential privacy (DP) is a mathematical framework for **quantifying how private** some analysis is, and providing rigorous guarantees to individuals.

DP bounds how much outcomes can change with the addition of one individual's data.





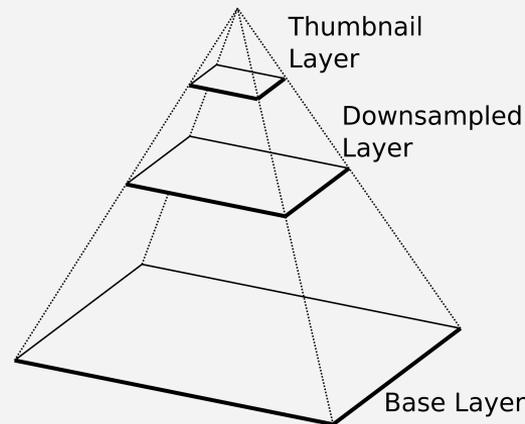
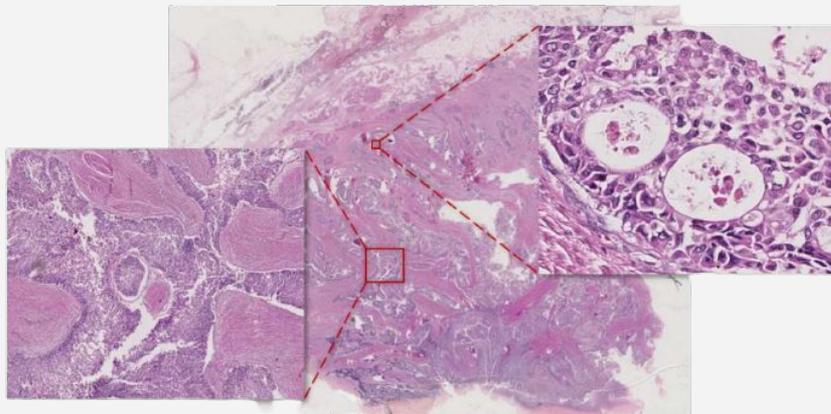
Gigapixel Histopathology

Gigapixel Digital Histopathology

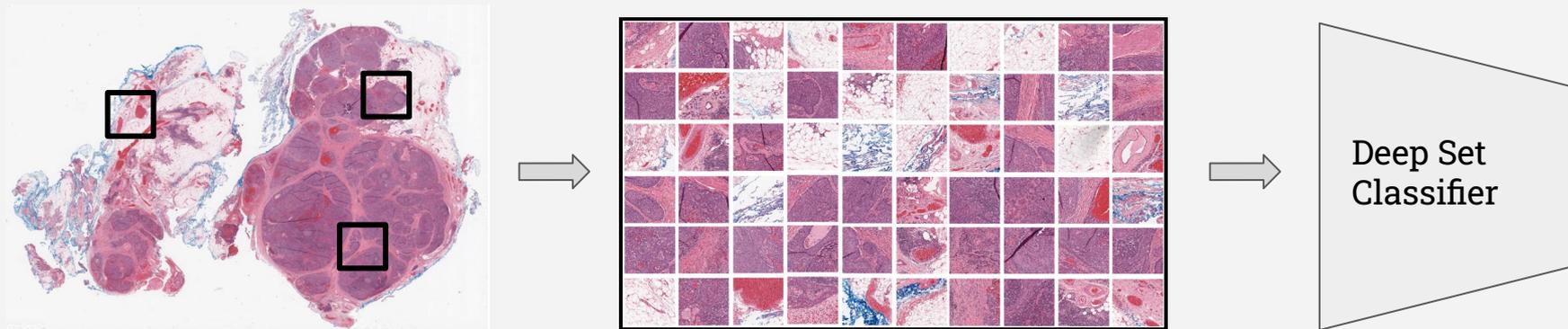
Digital histopathology produces ultra-high resolution Whole Slide Images (WSI) which are commonly larger than 50,000x50,000 pixels.

But NNs can't currently handle such large images.

Downsampling removes crucial information [\[Tizhoosh, Pantanowitz 2018\]](#)



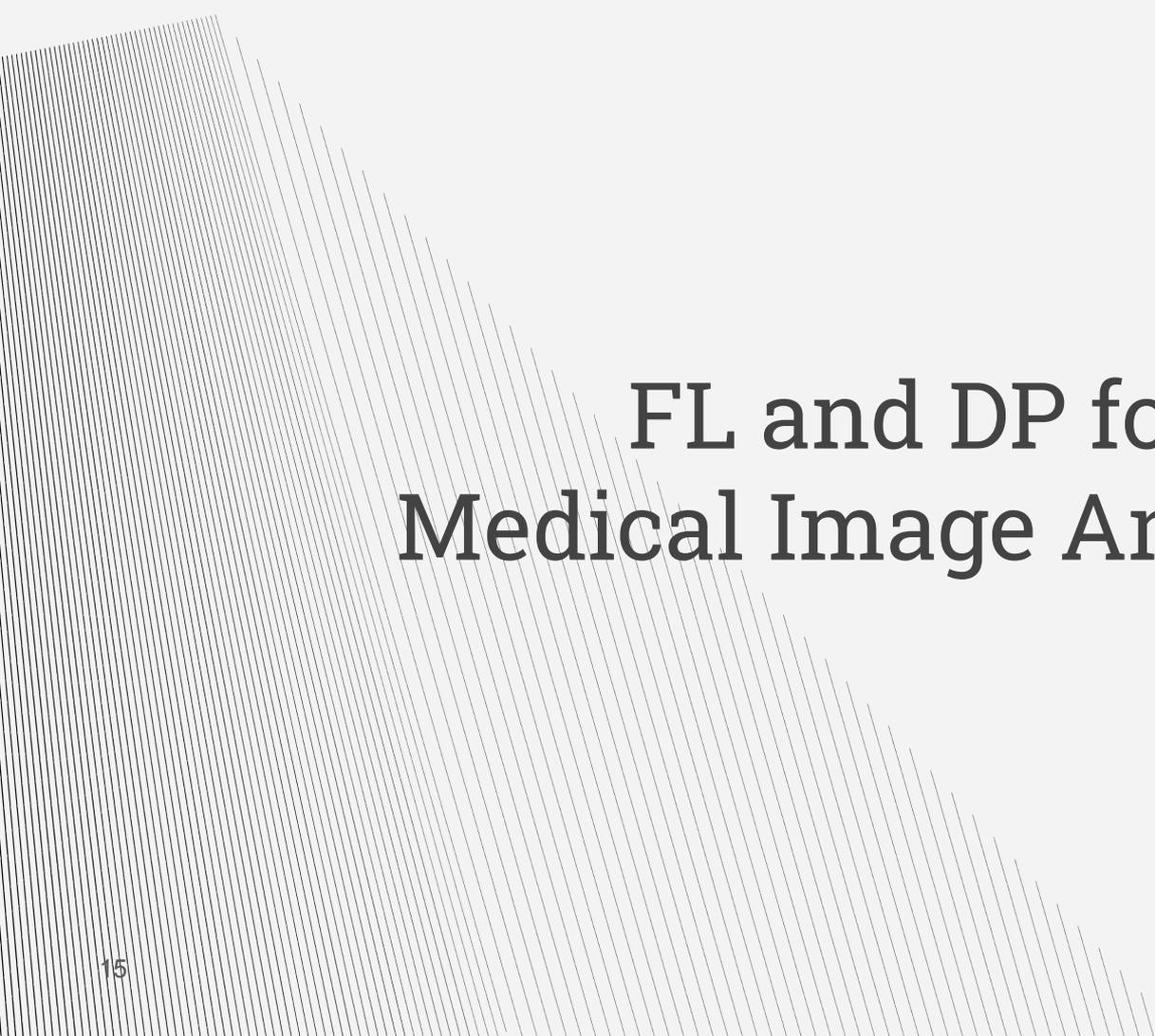
Gigapixel Digital Histopathology



Instead, break image into patches which are classified as a set [\[Kalra et al. 2020\]](#).

We used The Cancer Genome Atlas [\[Weinstein et al. 2013\]](#), a dataset of WSIs which originate from different institutions.

Aim is to classify the cancer subtype given the anatomical site.



FL and DP for Medical Image Analysis

FL and DP for Medical Image Analysis

[\[Adnan, Kalra, Cresswell, Taylor, Tizhoosh 2022\]](#)

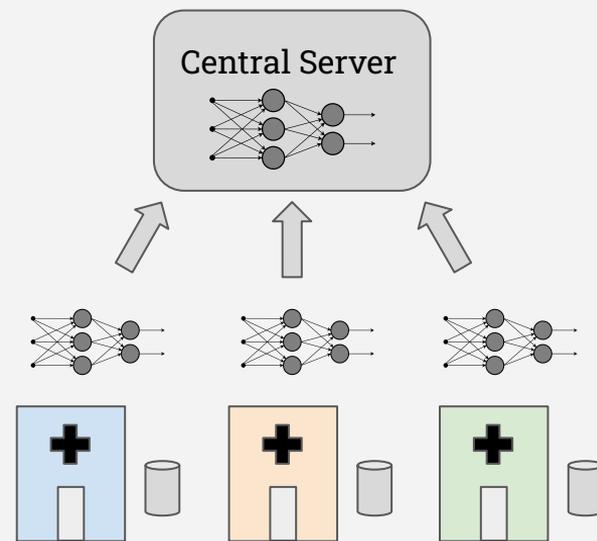
When institutions cannot pool their raw data, it may be possible to increase effective dataset size through FL.

Institutions will have diverse imaging methods.

FL must contend with

- non-IID datasets for each client
- generalization to previously unseen institutions

We conducted a case study to determine the viability of FL with DP for digital histopathology images.



FL and DP for Medical Image Analysis

- A. Is FL for histopathology viable when institutions have varying amounts of data, and varying distributions?

We selected >2500 WSIs from TCGA containing either Lung Adenocarcinoma, or Lung Squamous Cell Carcinoma (binary classification).

The data is divided into varying numbers of clients, and can be IID or non-IID.

Data distribution	Number of clients n	Accuracy		
		Without FL	With FL	Centralized
IID	4	0.731 ± 0.03	0.824 ± 0.02	0.848 ± 0.02
	8	0.620 ± 0.06	0.780 ± 0.05	
	16	0.570 ± 0.03	0.726 ± 0.06	
	32	0.527 ± 0.02	0.641 ± 0.09	
Non IID	4	0.682 ± 0.10	0.824 ± 0.01	0.848 ± 0.02
	8	0.561 ± 0.08	0.823 ± 0.05	
	16	0.524 ± 0.03	0.750 ± 0.06	
	32	0.520 ± 0.03	0.550 ± 0.20	

FL and DP for Medical Image Analysis

B. Is FL with DP viable for institutions to provide privacy guarantees?

We divided the WSIs by their source hospital. Testing was done on local, in-distribution data, as well as external data from unseen hospitals.

Applying DP negatively affects performance, but provides rigorous guarantees on how much individuals could be affected ($\epsilon=2.90$, $\delta=1e-4$).

Source hospital	Non-collaborative training		DP-FL training		FL training		Combined training	
	Test	External	Test	External	Test	External	Test	External
International Genomics Consortium	0.654	0.631	0.823 ± 0.01	0.707 ± 0.01	0.823 ± 0.01	0.741 ± 0.01	0.839 ± 0.01	0.768 ± 0.003
Indivumed	0.648	0.556						
Asterand	0.709	0.701						
John Hopkins	0.681	0.600						



Decentralized Private Learning

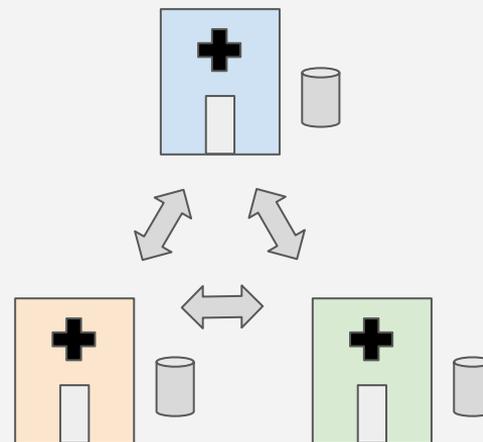
Decentralized Private Learning

In FL examples thus far, a central model is trained and sent back to clients.
At inference time, one shared model is used by all clients.

Clients can't develop their own model architectures, or personalize to their local data.

More preferable would be

- Direct collaboration without intermediaries
- Model heterogeneity and secrecy
- Personalization to local data
- Differential privacy guarantees



ProxyFL: Decentralized FL through Proxy Model Sharing

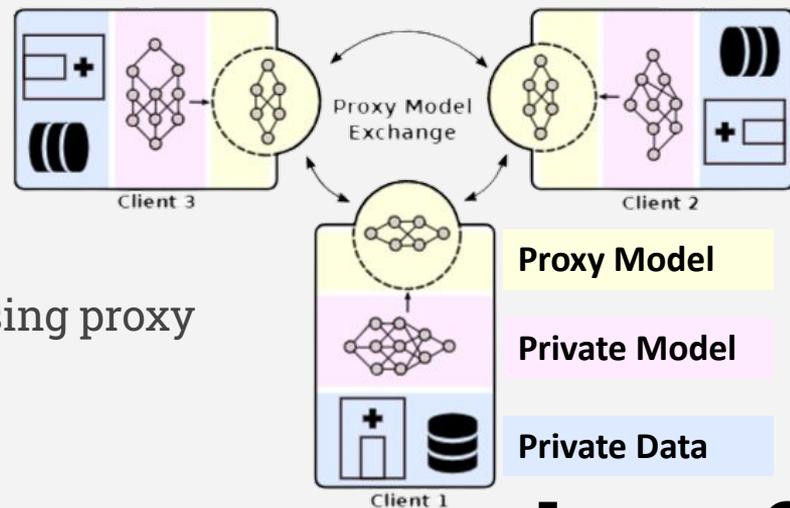
[\[Kalra, Wen, Cresswell, Volkovs, Tizhoosh 2021\]](#)

Handle all peer-to-peer communication through proxy models.

Each client maintains two models:

1. Private model - any architecture, not shared.
2. Proxy model - common architecture, small scale, shared to other clients.

Efficient information transfer is facilitated by passing proxy models, and mutual learning.



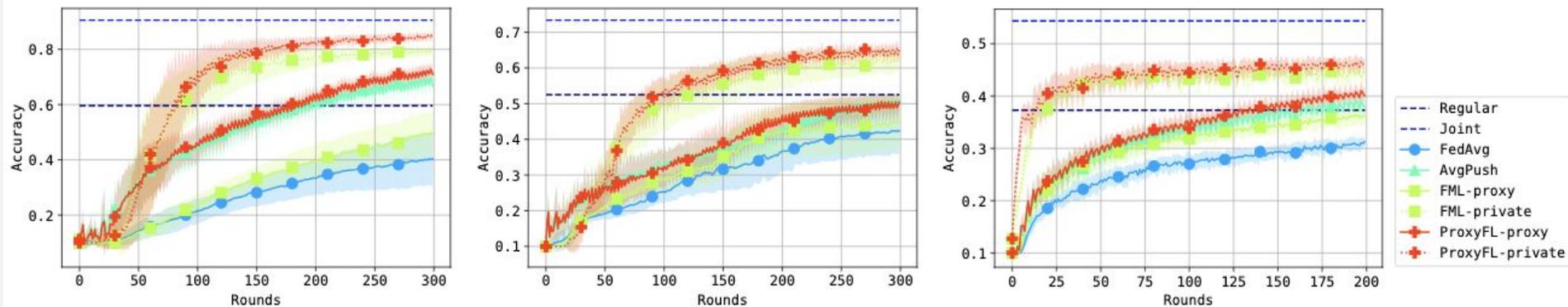
ProxyFL: Decentralized FL through Proxy Model Sharing

On benchmarks datasets, ProxyFL outperforms other FL schemes for the same privacy consumption.

MNIST

Fashion-MNIST

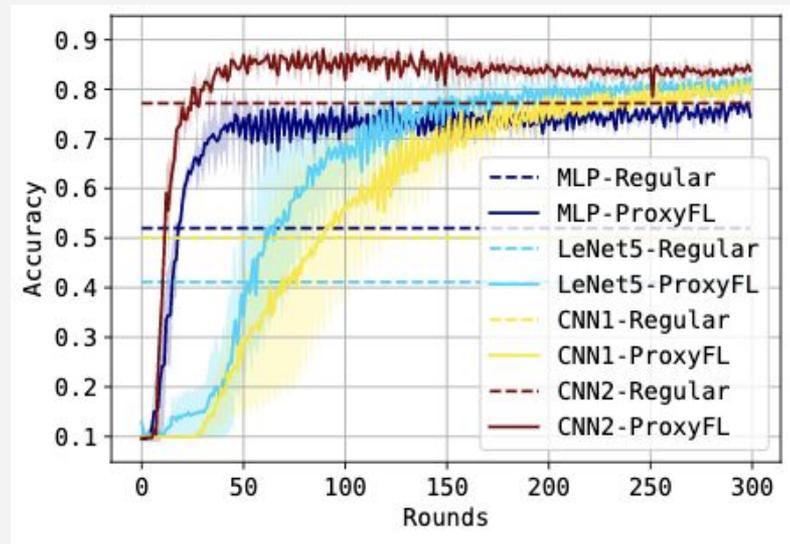
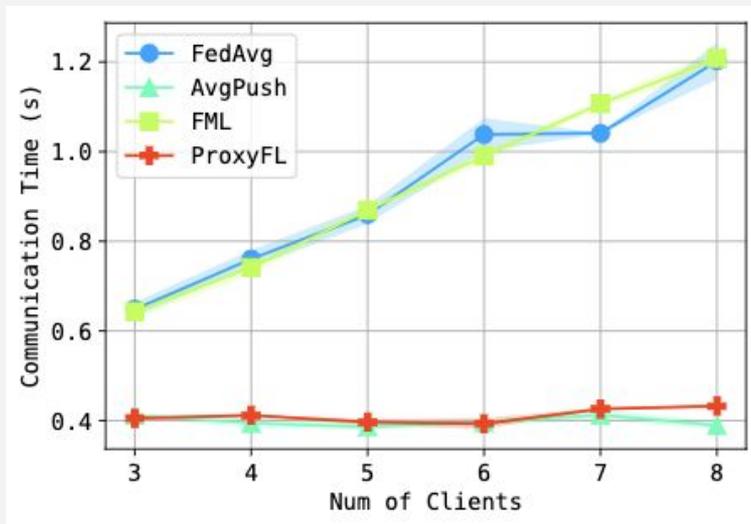
CIFAR-10



Eight clients each with 1k images, and non-IID distributions - 80% of training data from a single class. (3k images and 30% for CIFAR-10)

ProxyFL: Decentralized FL through Proxy Model Sharing

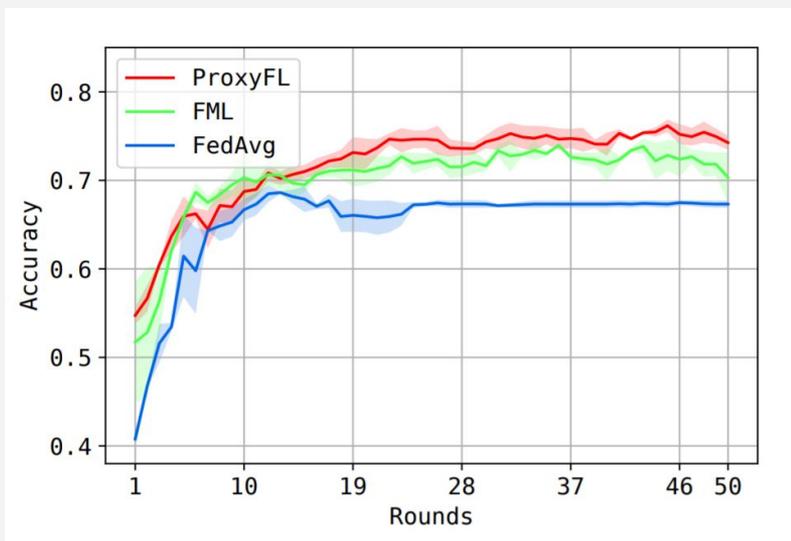
Decentralization leads to improved communication efficiency



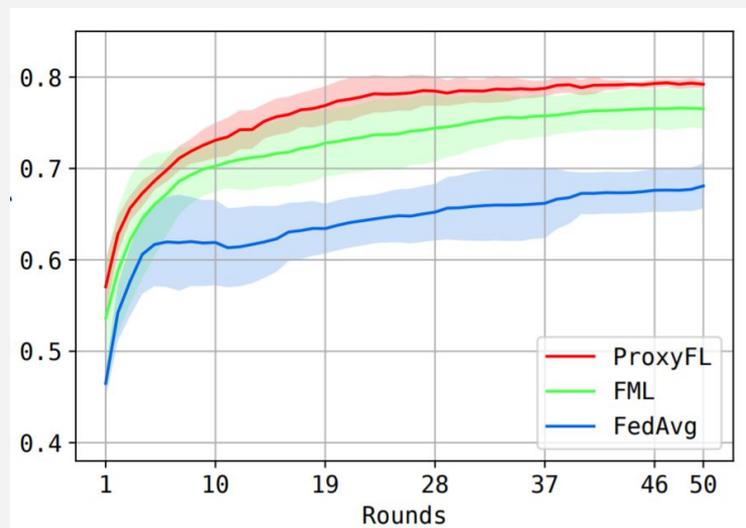
And ProxyFL allows clients to use different private model architectures

ProxyFL: Decentralized FL through Proxy Model Sharing

We also tested on histopathology data, using >5,600 WSIs from four institutions.



Strong privacy
per-client ($\epsilon=2.1, \delta=1e-3$)



Relatively weaker privacy
per-client ($\epsilon=6.2, \delta=1e-3$)



Conclusions

Conclusions

Federated learning may enable healthcare institutions to collaborate with one another while respecting patient's rights to privacy.

Collaboration effectively increases dataset size, diversity, and generality, leading to more accurate ML models and better clinician information.

Differential privacy can supplement FL to provide rigorous privacy guarantees.

We are currently researching:

Fairness of PETs - can exacerbate negative impact on underrepresented groups

Personalization for FL - learn from others, but tune outcomes to local data

Representation learning - patient representations for cross-task knowledge transfer

References

Adnan, Kalra, Cresswell, Taylor, Tizhoosh. ***Federated learning and differential privacy for medical image analysis***. Nature Scientific Reports 12 (1953), 2022

Kalra, Wen, Cresswell, Volkovs, Tizhoosh. ***ProxyFL: Decentralized Federated Learning through Proxy Model Sharing***. arXiv:2111.11343

Dwork, McSherry, Nissim. ***Calibrating Noise to Sensitivity in Private Data Analysis***. Journal of Privacy and Confidentiality 7 (3), 2017

Fredrickson, Jha, Ristenpart. ***Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures***. 22nd ACM SIGSAC, 2015

McMahan, Moore, Ramage, Hampson, Aguera y Arcas. ***Communication-Efficient Learning of Deep Networks from Decentralized Data***. 20th ICAIS, 2017

Tizhoosh, Pantanowitz. ***Artificial intelligence and digital pathology: Challenges and opportunities***. Journal of Pathology Informatics 9:39, 2018

Zhang, Xiang, Hospedales, Lu. ***Deep Mutual Learning***. CVPR 2018

We are open to collaboration!

jesse@layer6.ai

