# ADVANCES IN PRIVACY PRESERVING FEDERATED LEARNING TO REALIZE A TRULY LEARNING HEALTHCARE SYSTEM

**RAVI MADDURI**
Senior Scientist
Data Science and Learning Division
Argonne National Laboratory
madduri@anl.gov

**ZILINGHAN LI**
Machine Learning Research Engineer
Data Science and Learning Division
Argonne National Laboratory
zilinghan.li@anl.gov

Washington D.C.
Oct 29., 2024

# FUNDING ACKNOWLEDGEMENTS

Argonne
NATIONAL LABORATORY

# KEY FEATURES OF A LEARNING HEALTH SYSTEM (LHS)

## From Institute of Medicine Report

- Data Integration and Interoperability
  - Diverse Data Sources: Combines EHRs, genomics, imaging, and patient-reported outcomes.

- Collaborative Culture
  - Stakeholder Collaboration: Involves clinicians, patients, researchers, and policymakers.
  - Shared Goals: Focuses on improving outcomes collectively.

- Ethical Data Use and Privacy
  - Data Security: Implements strong protections for patient information.
  - Governance: Follows ethical guidelines for data access and use.

Argonne
NATIONAL LABORATORY

# KEY FEATURES OF A LEARNING HEALTH SYSTEM (LHS)
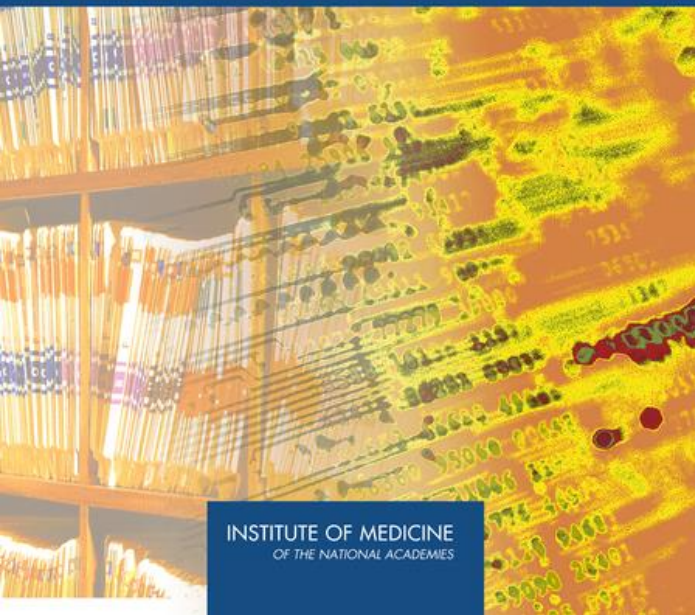
## From Institute of Medicine Report

- Continuous Learning Cycle
  - Data-Driven Improvement: Routinely collects and analyzes clinical data.
  - Rapid Feedback: Insights are quickly integrated back into practice.

- Feedback Mechanisms
  - Clinician Support: Offers decision support tools and alerts.
  - Patient Feedback: Incorporates patient experiences to refine care delivery.

- Evidence Generation at Point of Care
  - Embedded Research: Research activities are part of clinical workflows.
  - Real-Time Analytics: Provides immediate evidence to inform decisions.

- Patient-Centered Care
  - Active Engagement: Patients participate in their own care and data sharing.
  - Personalization: Care plans reflect individual preferences and needs.

Argonne
NATIONAL LABORATORY

Multimodal machine learning with continuous learning, enhanced collaboration, and patient privacy preservation capabilities can enable LHS.
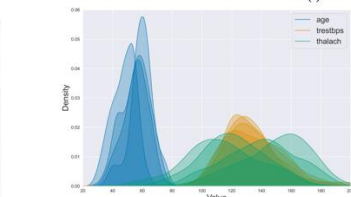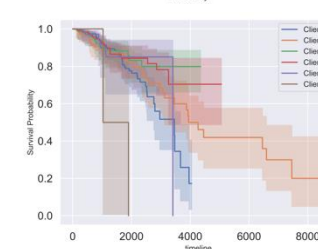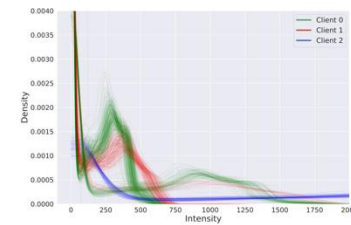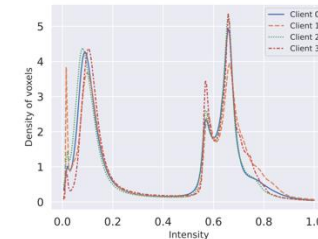
# CHALLENGES TOWARDS LHS

## Multimodal and Heterogeneous Medical Data in Distributed Silos

### Biomedical Health Data are Multimodal

Biomedical health data are multimodal as they include different types like images, texts, signals, and structured data from various sources.

### Biomedical Health Data are Heterogeneously Distributed

Biomedical Health data are heterogeneous across distributed data silos due to variations in patient populations, medical practices, and data collection processes among hospitals.

Argonne
NATIONAL LABORATORY

# CHALLENGES TOWARDS LHS

## Training Multimodal Medical AI models is Essential But Challenging

**Training a multimodal medical AI model is essential**

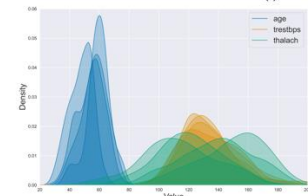Training multimodal medical AI models is essential because it can leverage diverse data types to provide more comprehensive and accurate insights, improving diagnosis and treatment outcomes.

**Training a robust multimodal AI model within single data silo is difficult**

Training robust multimodal AI models within a single data silo is difficult due to the homogeneous patient population and limited data modalities, hindering generalization to real-world scenarios.

**Collecting medical data centrally is challenging**

Collecting the distributed medical data centrally is challenging due to privacy concerns and varying regulations (e.g. HIPAA).

# MULTIMODAL FL APPROACHES TO LHS

## Federated Learning Offers a Viable Solution

- Federated learning (FL) is a distributed learning paradigm with multiple data silos as clients and a central server for orchestration.

- Each FL client has its own computing facilities and trains model using their *private* local data.

- Each FL client only shares the *model* trained on their private local data to the FL server for aggregation.

- A global model is obtained by aggregating models from different FL clients, thus implicitly leveraging private data from various data silos.
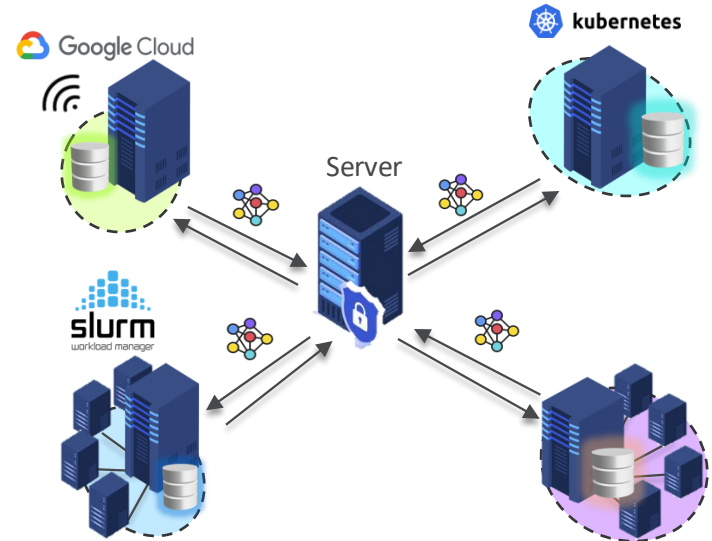


**Fig.** Federated Learning Illustration.

# CASE STUDIES IN FEDERATED LEARNING FOR BIOMEDICINE

## Federated Learning Shows Effectiveness in Uni-modal Medical Models

### Challenge: Predicting age from ECG signals.

Two Clients

- ECG-ANL: the dataset is aggregated from multiple open-source datasets.

- ECG-Broad: private dataset collected by Broad Institute.

FL can obtain a model that performs relatively well on both datasets.

| Dataset | Train | Val | Test | Total |
|---------|-------|-----|------|-------|
| ECG-ANL | 64518 | 7905 | 7905 | 80328 |
| ECG-Broad | 33140 | 4143 | 4143 | 41426 |

| Training Dataset | Testing Set | | |
|------------------|-------------|---------|---------|
| | ECG-ANL | ECG-Broad | Average |
| ECG-ANL *(local training)* | 109.95 | 224.48 | 149.33 |
| ECG-Broad *(local training)* | 225.41 | 38.93 | 161.28 |
| ECG-ANL+Broad - FedAvg[1] | 125.00 | 41.70 | 96.35 |

Mean Squared Error (the lower, the better)

Argonne
NATIONAL LABORATORY

# CASE STUDIES IN FEDERATED LEARNING FOR BIOMEDICINE

## Federated Learning Shows Effectiveness in Uni-modal Medical Models

### *Challenge: COVID-19 prediction from Chest X-Rays.*

Two Clients

- ANL-COVID: the dataset is aggregated from multiple open-source datasets.

- UChicago-COVID: private dataset collected by UChicago.





Better Classification

Local Model 1 (AUC = 0.75)
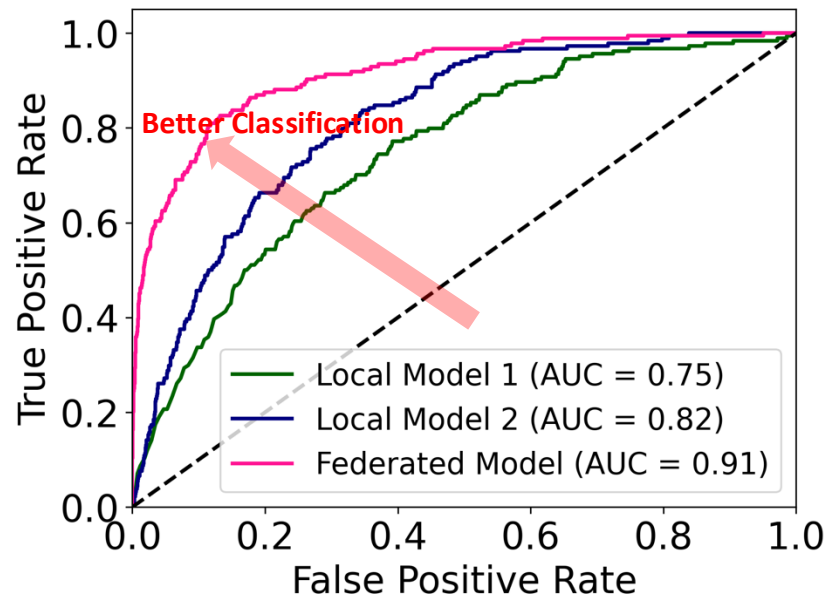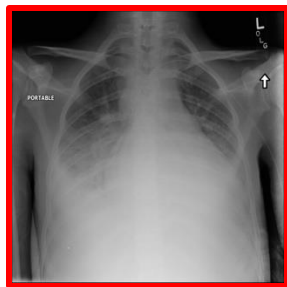Local Model 2 (AUC = 0.82)
Federated Model (AUC = 0.91)

# CASE STUDIES IN FEDERATED LEARNING FOR BIOMEDICINE

## Federated Learning Shows Effectiveness in Uni-modal Medical Models

*Challenge: Preventing the reconstruction of Chest X-Rays from model gradients*

- Federated learning itself is not privacy preserving – as the training data can be reversely constructed from model gradients.

- Differential privacy (DP) techniques, can significantly increase the difficulty of reconstruction by adding noises to model parameters



Increasing Differential Privacy $(\epsilon, c)$

Baseline (a) | $(\epsilon = 0.1, c = 1)$ (b) | $(\epsilon = 0.05, c = 1)$ (c) | $(\epsilon = 0.01, c = 1)$ (d)

# VISION

## From Uni-modal to Multimodal Privacy-Preserving Federated Learning

The examples only train uni-modal biomedical models using privacy-preserving federated learning, which limits the range of applications.
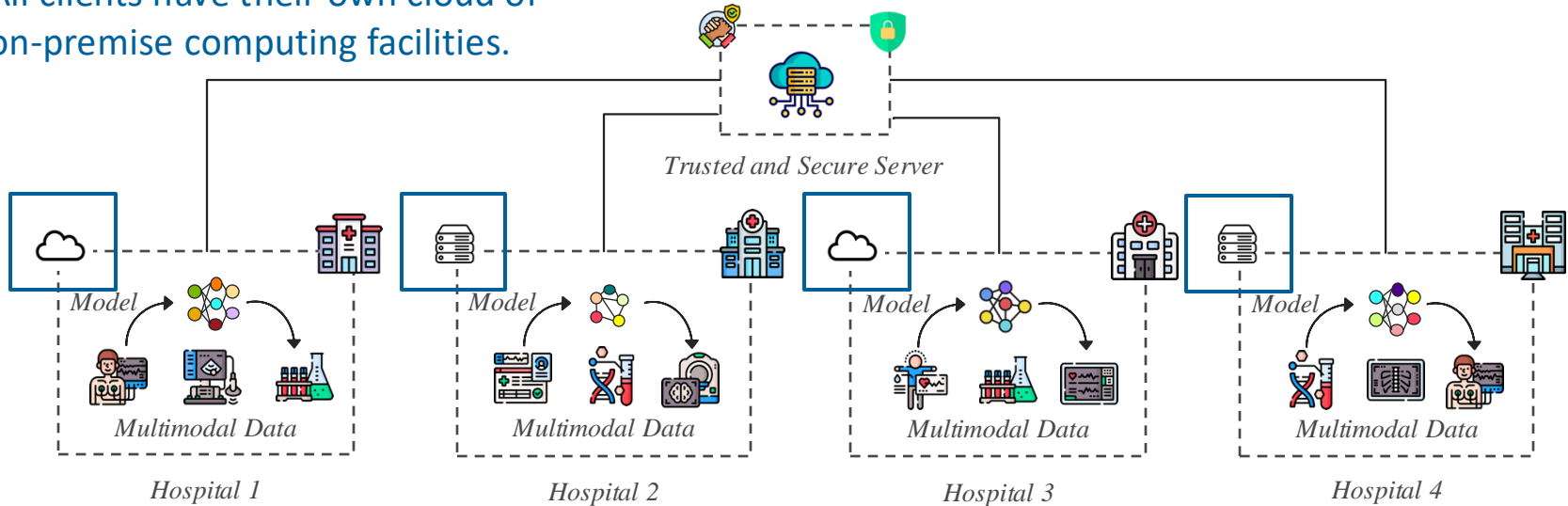
➜ Here we envision an ideal federated learning framework to:

❖ Enable secure and privacy-preserved training of multimodal biomedical models.

❖ Continuously update the trained models as new data accumulate.

❖ Help realize a truly learning healthcare system.

Argonne
NATIONAL LABORATORY

# VISION

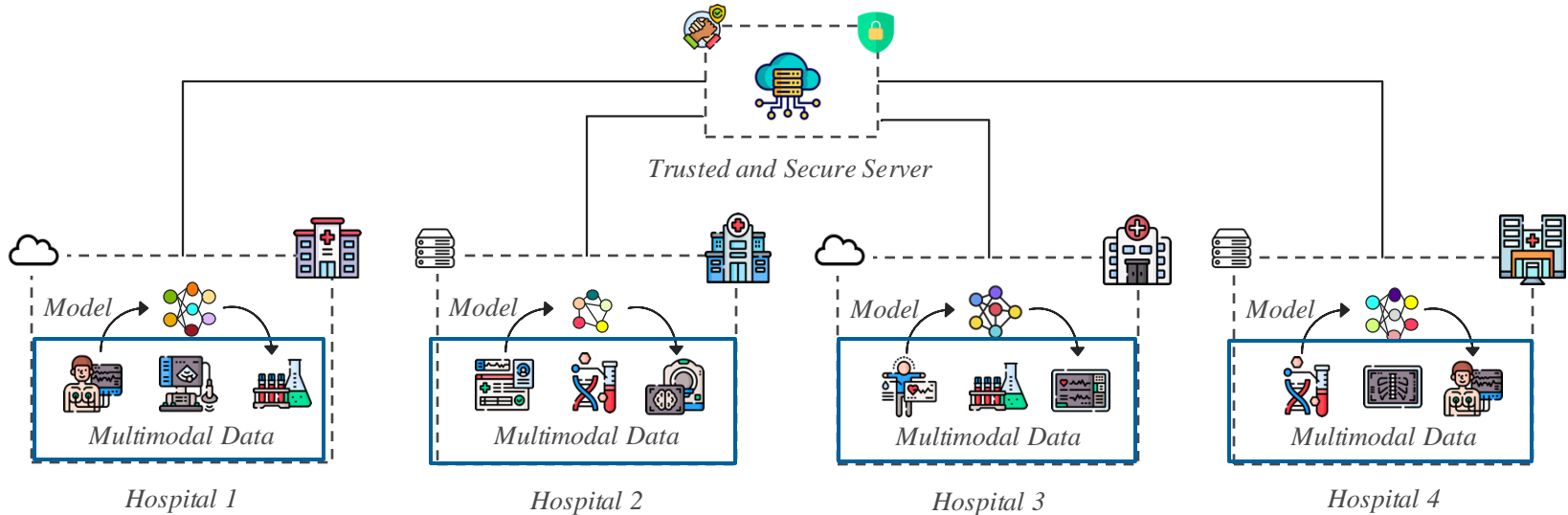## From Uni-modal to Multimodal Privacy-Preserving Federated Learning

All clients have their own cloud or
on-premise computing facilities.



*Trusted and Secure Server*

*Model*

*Multimodal Data*

*Hospital 1*

*Model*

*Multimodal Data*

*Hospital 2*

*Model*

*Multimodal Data*

*Hospital 3*

*Model*

*Multimodal Data*

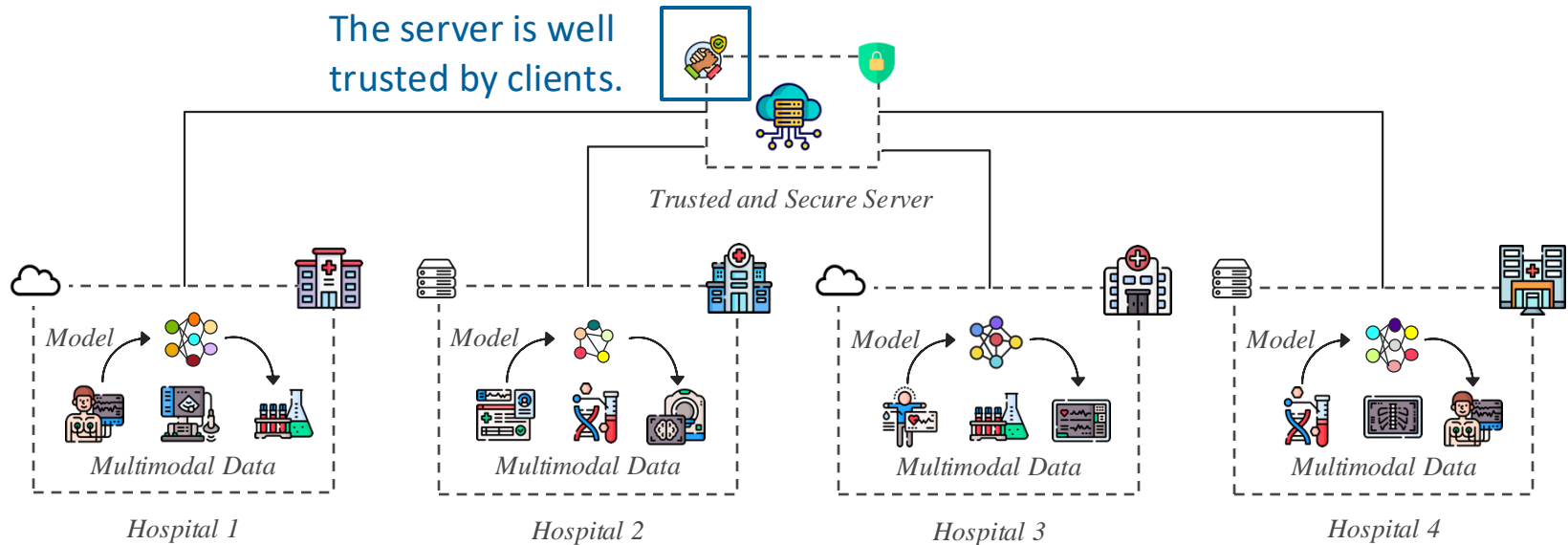*Hospital 4*

Argonne
NATIONAL LABORATORY

# VISION

## From Uni-modal to Multimodal Privacy-Preserving Federated Learning



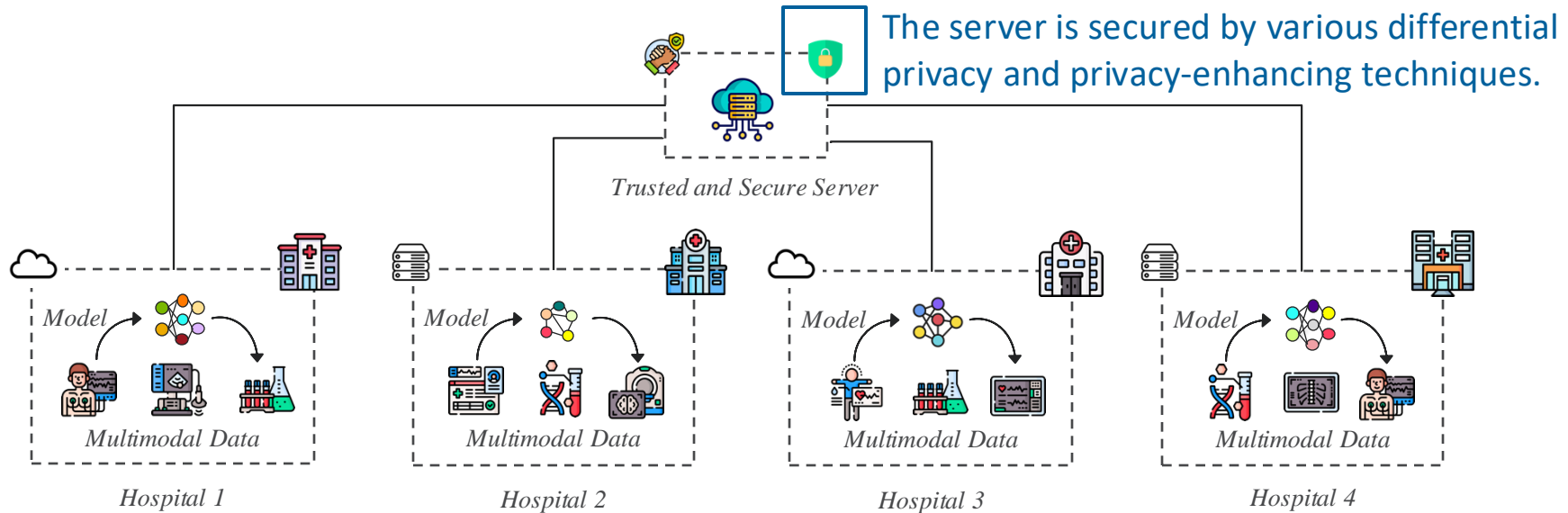Different clients have various data modalities for model training.

# VISION

## From Uni-modal to Multimodal Privacy-Preserving Federated Learning



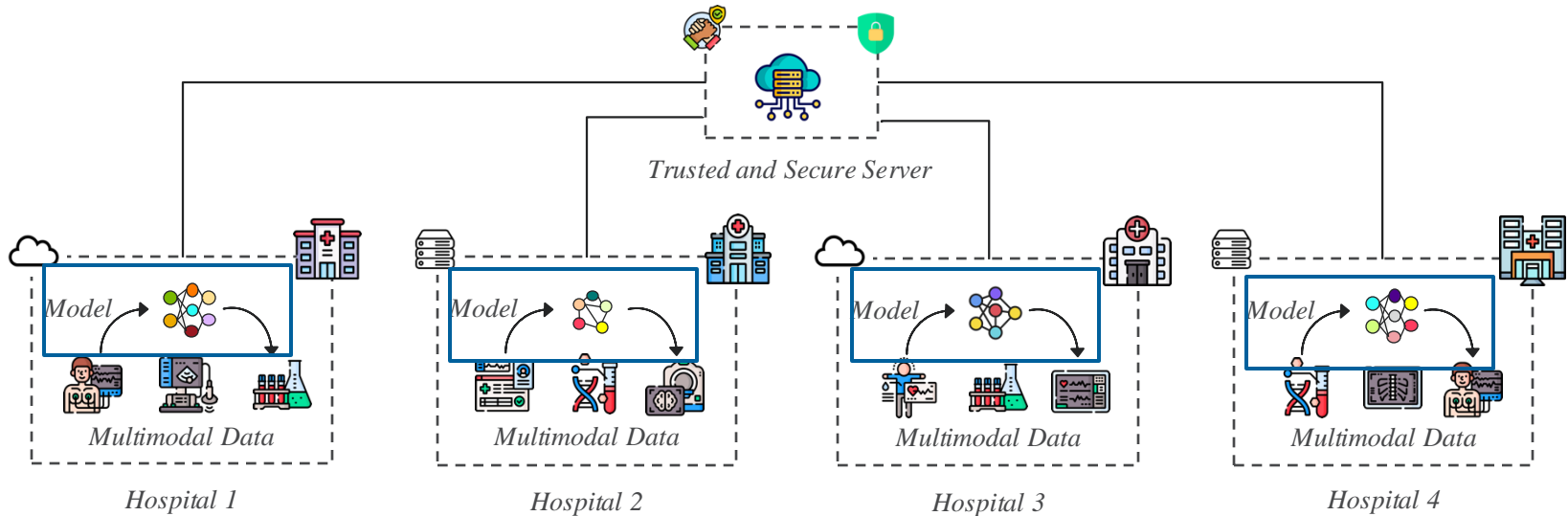The server is well trusted by clients.

*Trusted and Secure Server*

*Model* — *Multimodal Data* — *Hospital 1*

*Model* — *Multimodal Data* — *Hospital 2*

*Model* — *Multimodal Data* — *Hospital 3*

*Model* — *Multimodal Data* — *Hospital 4*

# VISION

## From Uni-modal to Multimodal Privacy-Preserving Federated Learning



The server is secured by various differential privacy and privacy-enhancing techniques.

*Trusted and Secure Server*

*Model*

*Multimodal Data*

*Hospital 1*

*Model*

*Multimodal Data*

*Hospital 2*

*Model*

*Multimodal Data*

*Hospital 3*

*Model*

*Multimodal Data*
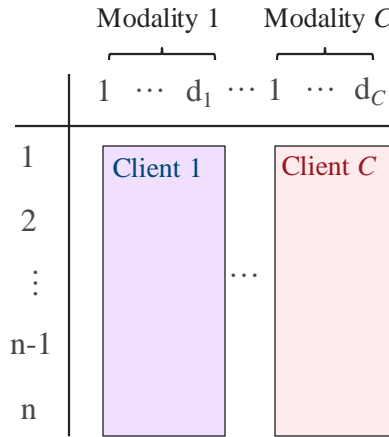
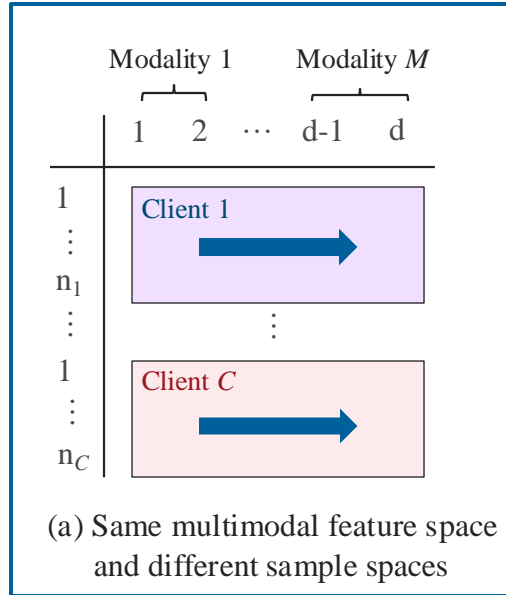*Hospital 4*

Argonne
NATIONAL LABORATORY

# VISION
## From Uni-modal to Multimodal Privacy-Preserving Federated Learning
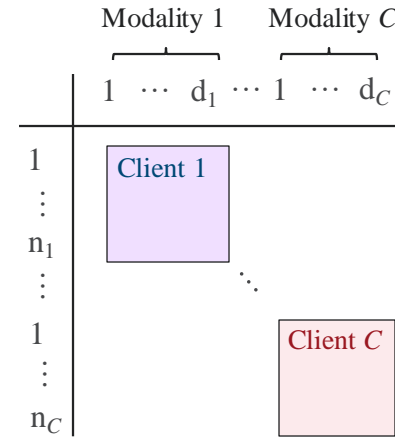


The models can be continuously updated to adapt to any shifts in data distributions, availability of new data, and evolving health trends in real-time.

# BUILDING BLOCKS

## Multimodal Federated Learning



(a) Same multimodal feature space and different sample spaces

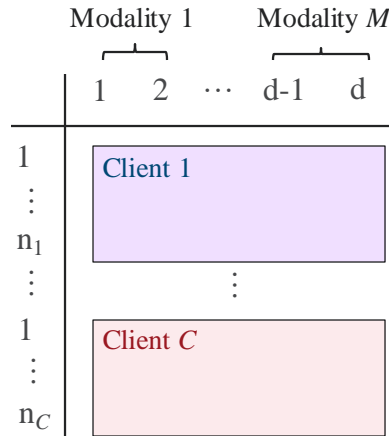(b) Same sample space and different multimodal feature spaces

(c) Different sample spaces and multimodal feature spaces
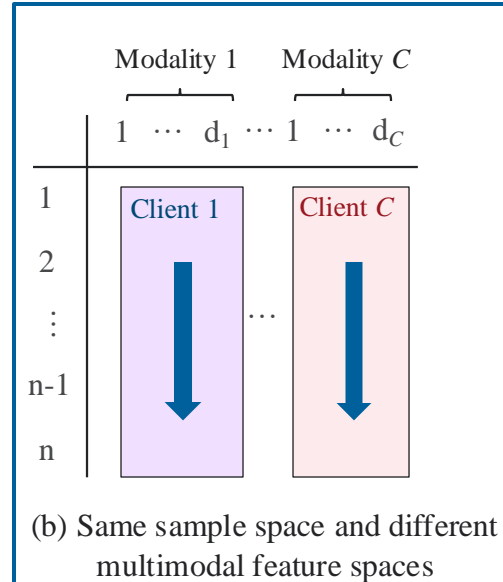
All clients have the same data modalities: All clients train the same multimodal modal using traditional federated learning. ➔ Horizontal Federated Learning.

# BUILDING BLOCKS

## Multimodal Federated Learning



(a) Same multimodal feature space and different sample spaces

(b) Same sample space and different multimodal feature spaces

VFL Server

Back Propagation

Sample Labels

1
0
1
1
0
1

①: Data Embeddings
②: Embedding Gradients

VFL Client 1    VFL Client 2    VFL Client 3

All clients have different modalities but the same sample space (rare in the real-world): Different clients train different embedding models for different modalities by sharing the data embeddings to the server and updating the local embedding models using embedding gradients sent back from the server. ➔ Vertical Federated Learning.
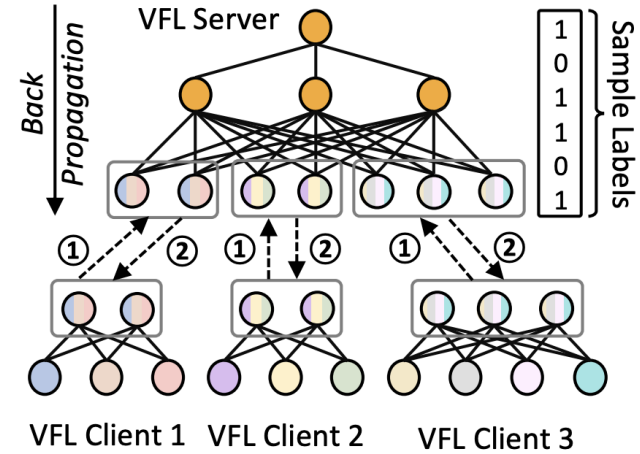
# BUILDING BLOCKS

## Multimodal Federated Learning



(a) Same multimodal feature space and different sample spaces

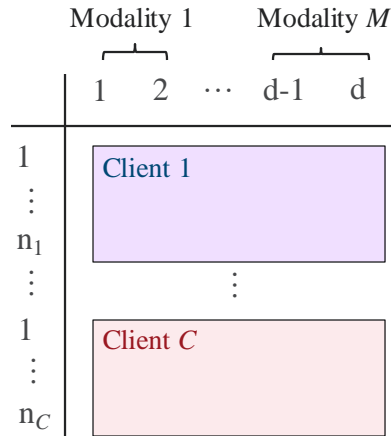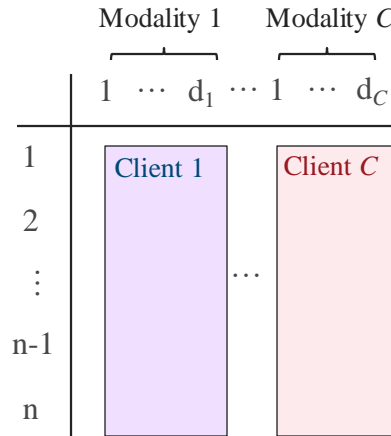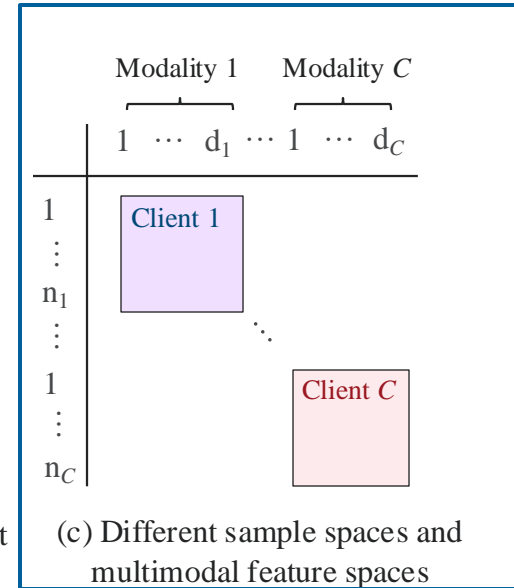(b) Same sample space and different multimodal feature spaces

(c) Different sample spaces and multimodal feature spaces

All clients have different modalities and different sample spaces: Common in real-world scenarios, but very challenging ➜ Potential solutions involve training small embedding models locally, and leveraging synthetic multimodal data or publicly available multimodal data for sample alignment and further tuning.

# BUILDING BLOCKS
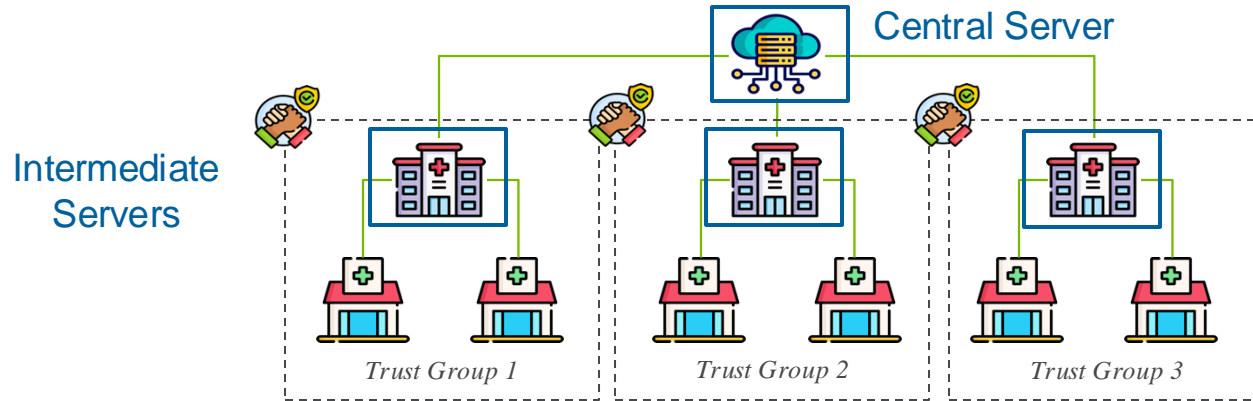## Hierarchical Federated Learning

As training data are vulnerable to data reconstruction, especially without privacy-enhancing techniques:

- Some clients (e.g. small clinics) might not be willing to share their local models with the central server.

- Those small clinics might only be willing to share the local model with some large trusted local medical institutions/university ➔ this prevents the creation of a large federation.

- Hierarchical federated learning enables collaborative training beyond the trust boundary.

Argonne
NATIONAL LABORATORY

# BUILDING BLOCKS

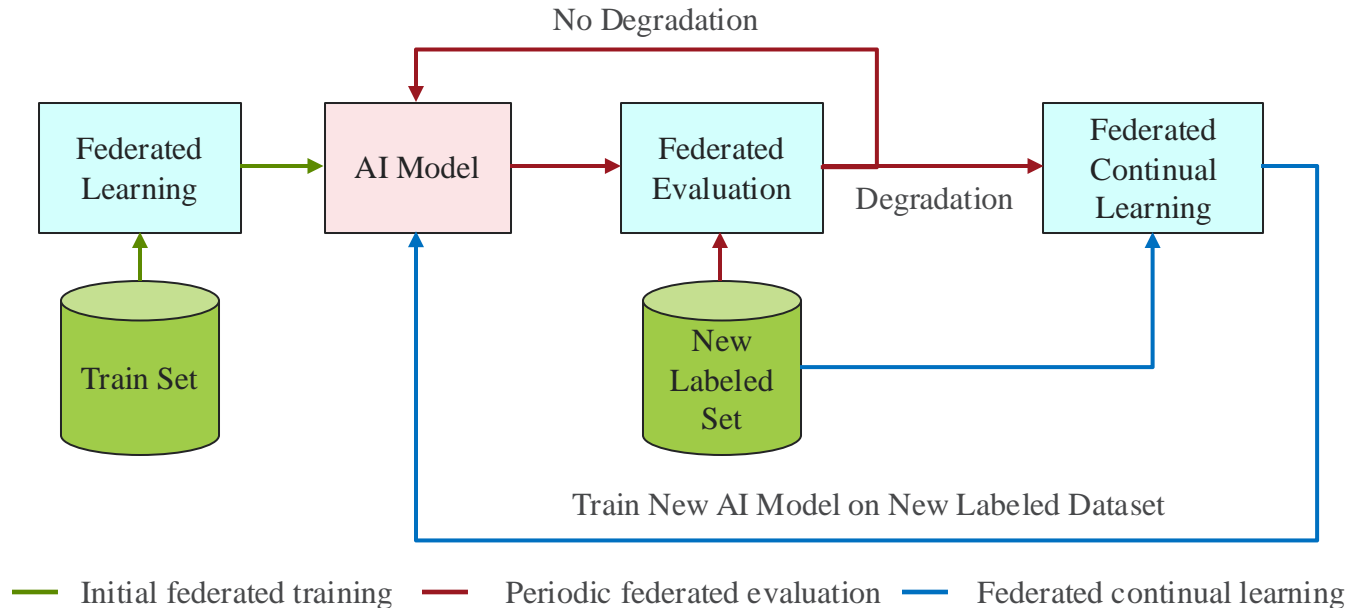## Hierarchical Federated Learning

- Each client first sends its local model to a trusted intermediate server.

- The intermediate server then sends a model aggregated from several clients to the central server.

- As the aggregated model contains the information of several clients, it is hard to reconstruct the training data of any single client.

# BUILDING BLOCKS
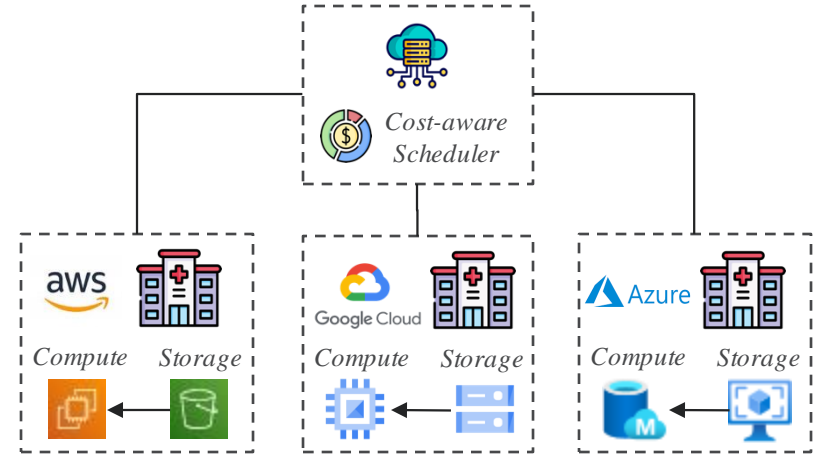
## Federated Continuous Learning

Federated continuous learning ensures that the models evolve to reflect the diversity of healthcare settings and populations, reducing biases and enabling better generalization.

# BUILDING BLOCKS

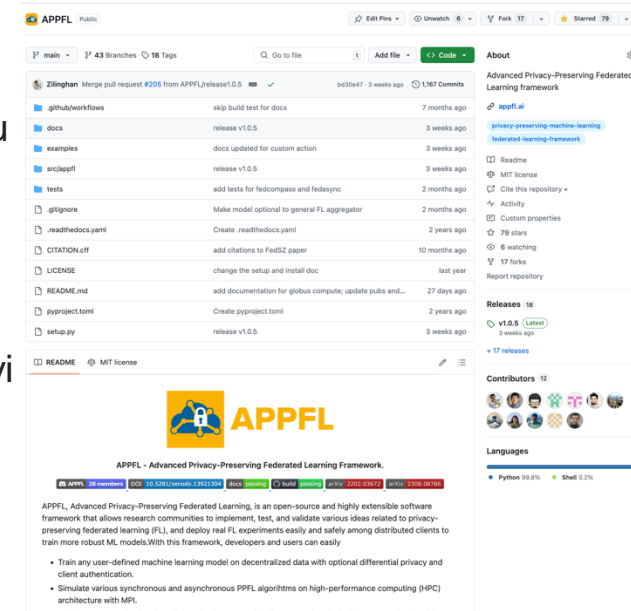## Cost-Aware FL on the Cloud



- Many hospitals have their private data on Cloud Storage (S3, Globus Cloud Storage, etc.) and have their computing on the Cloud as well.

- Training on GPU cloud instances can be costly.

- AWS, Google, and Azure all have "spot computing" – AWS Spot Instances, Google Cloud Preemptable VMs, and Azure Spot VMs, which provide a low-cost computing option, but can be killed at any time with short notice.

- It would be greatly beneficial to have a cost-aware scheduler for the FL server to reduce the cost for FL experiments among heterogeneous cloud computing providers using their spot instances, while maintaining robust and tolerant to potential client failures.

Argonne
NATIONAL LABORATORY

# REFERENCES

- Hoang, Trung-Hieu, Jordan Fuhrman, Ravi Madduri, Miao Li, Pranshu Chaturvedi, Zilinghan Li, Kibaek Kim et al. "Enabling end-to-end secure federated learning in biomedical research on heterogeneous computing environments with APPFLx." *arXiv preprint arXiv:2312.08701* (2023).

- Li, Zilinghan, Shilan He, Ze Yang, Minseok Ryu, Kibaek Kim, and Ravi Madduri. "Advances in APPFL: A Comprehensive and Extensible Federated Learning Framework." *arXiv preprint arXiv:2409.11585* (2024).

- Madduri, Ravi, Zilinghan Li, Tarak Nandi, Kibaek Kim, Minseok Ryu, and Alex Rodriguez. "Advances in Privacy Preserving Federated Learning to Realize a Truly Learning Healthcare System." *arXiv preprint arXiv:2409.19756* (2024).

- https://github.com/APPFL/APPFL

- https://appfl.ai

# THANK YOU!

Argonne
NATIONAL LABORATORY