



Introduction to NVIDIA Federated Learning: Concepts, Technology, and Use Cases

Holger Roth, Principal Federated Learning Scientist, NVIDIA

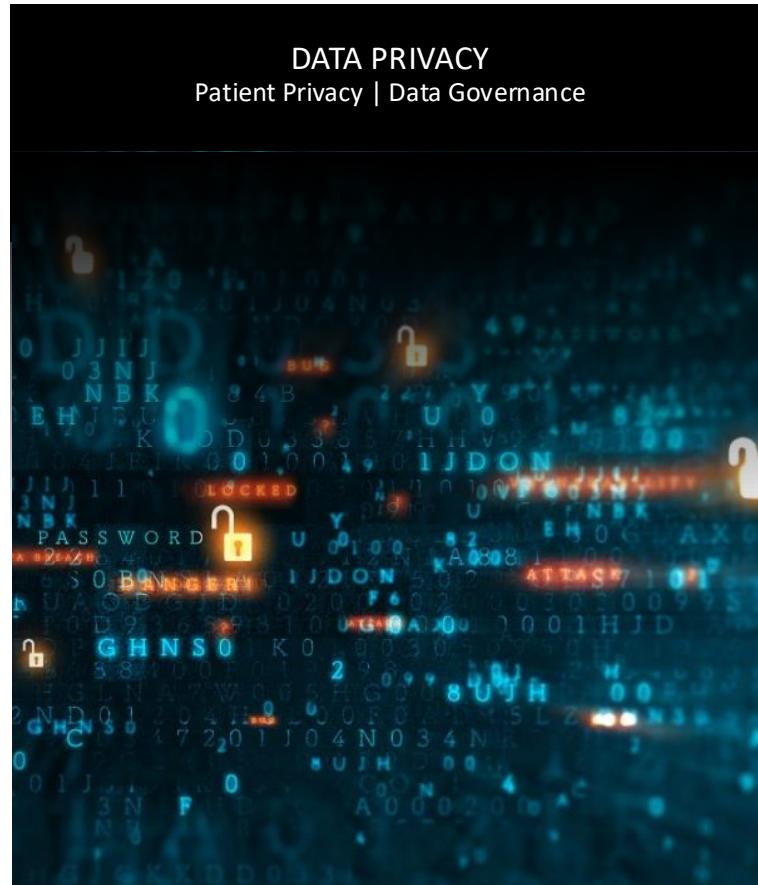
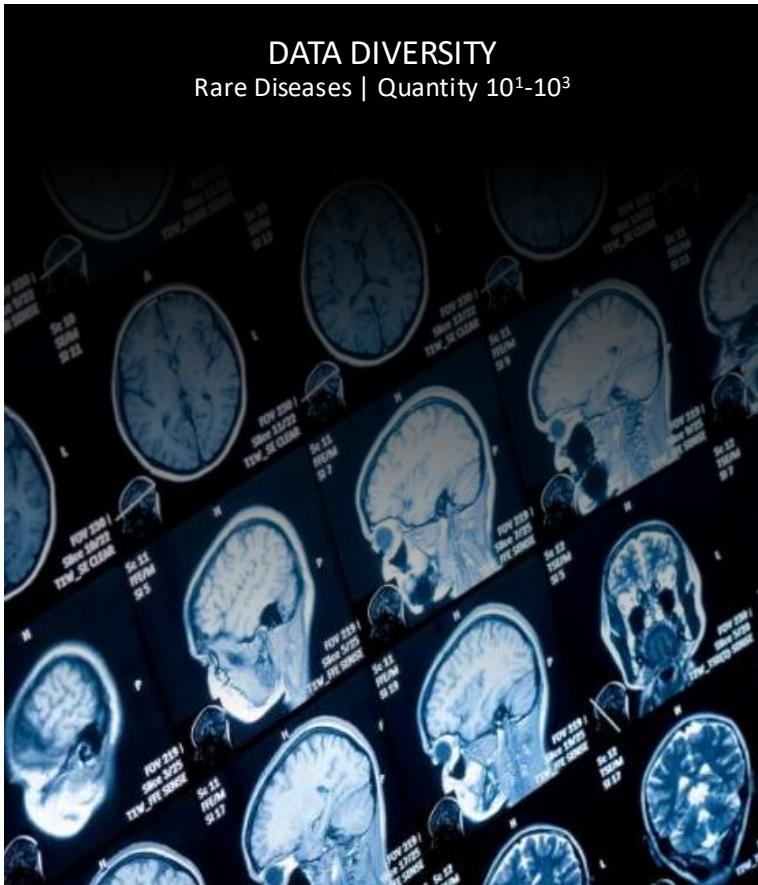
NVIDIA FLARE WEBINAR – Q1, FEB 19TH, 2025



Agenda

- **What is Federated Learning?**
 - NVIDIA Key Technologies for Federated Learning
 - Real-world Use cases of Federated Learning
 - Getting Started with NVIDIA FLARE
 - Research: Addressing Key challenges in Federated Learning
 - Summary & Announcements
-
-
-
-
-
-
-

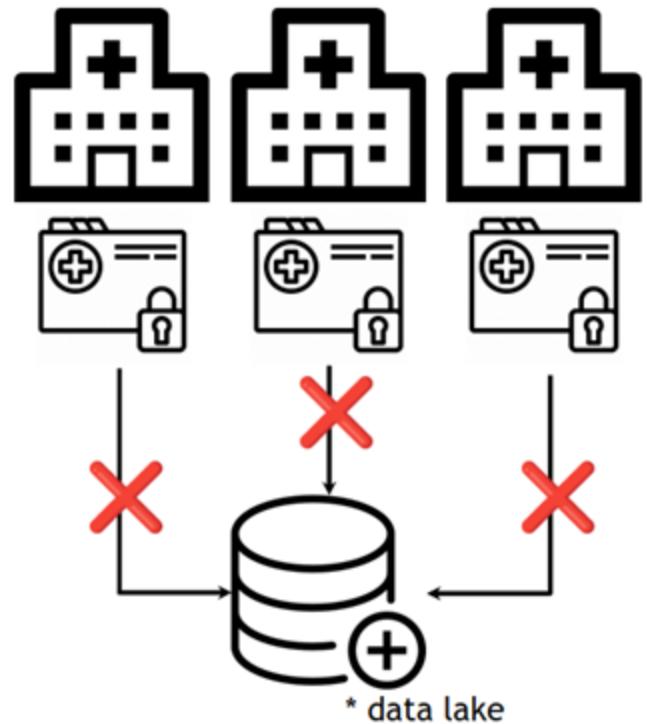
BUILDING ROBUST, GENERALIZABLE AI MODELS IS HARD



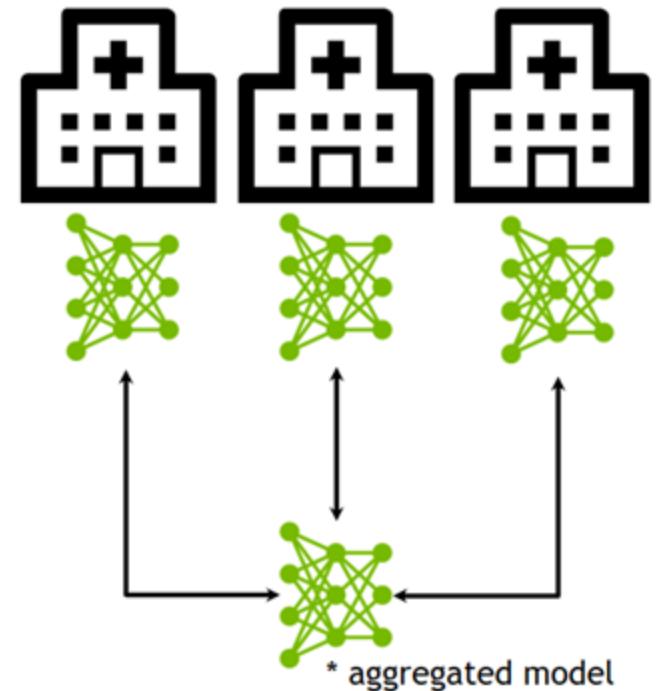
However

- Private data can often not be shared (regulations, complex bureaucracy)
- Data annotation is costly; Data is an asset

FEDERATED LEARNING

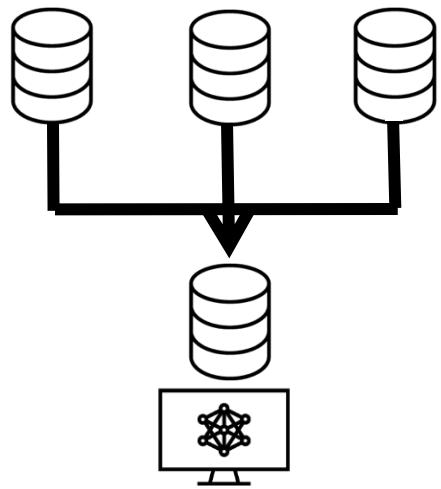


- Data cannot be centralized in many tasks
- **Share models, not data!**



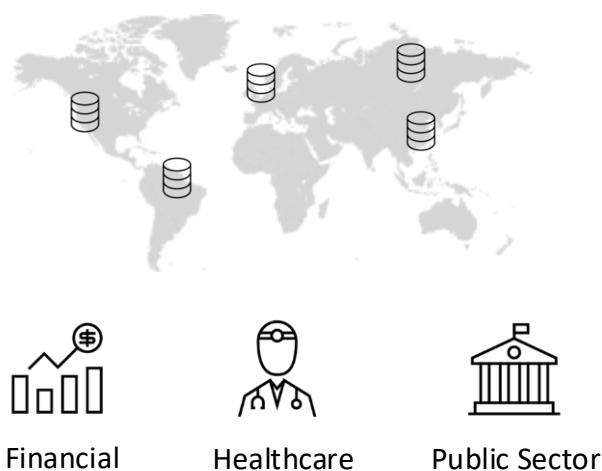
Federated Computing – Removing Data Silos

Avoid Data Copy | Regulatory Compliance | Prevent Private Data Leak



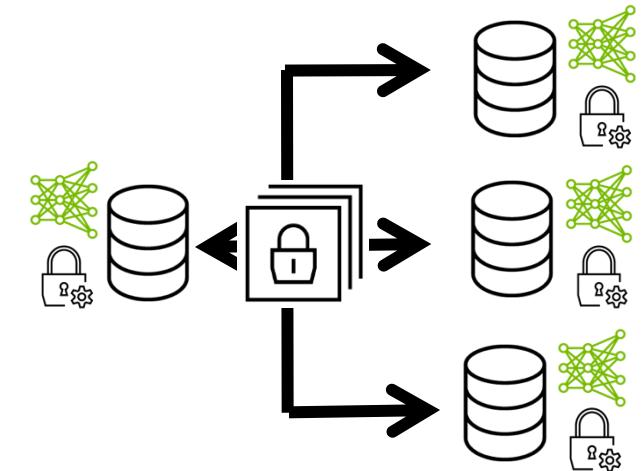
No Data Copy

Centrally aggregating the data is not possible or practical



Compliance

Data sovereignty restrictions and industry regulations



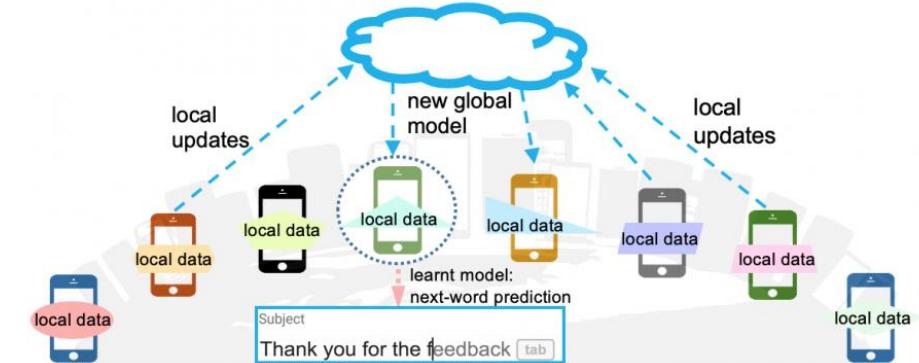
Privacy Enhancing Technology

Multiple layers of security features incl. Homomorphic Encryption, Differential Privacy, & Confidential Computing

CROSS-DEVICE VS CROSS-SILO

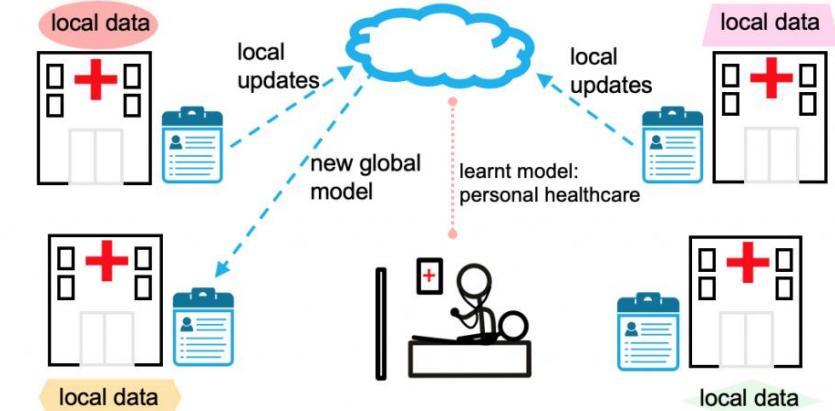
Cross-device:

- Large number of user (millions), unknown to each other
- Intermittent user/data availability (select from a pool) - user/data for each FL round could be different
- Massively distributed
- Lower computation power
- Examples: mobile phones, autonomous cars, robots, etc.



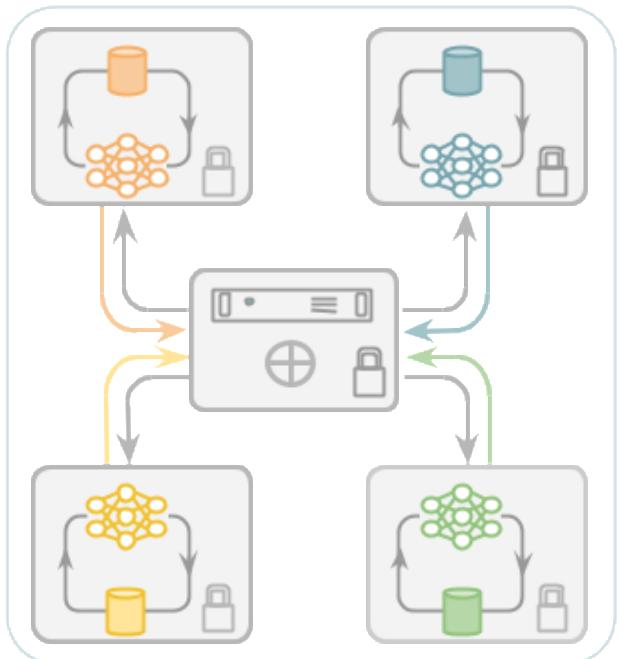
Cross-silo:

- Usually comes from a trusted collaboration, small number (10s-100s)
- Relatively stable user/data and connection
- High computation power, usually dedicated
- Regional/global scale
- Examples: hospitals, financial institutions, enterprises, data centers, etc.

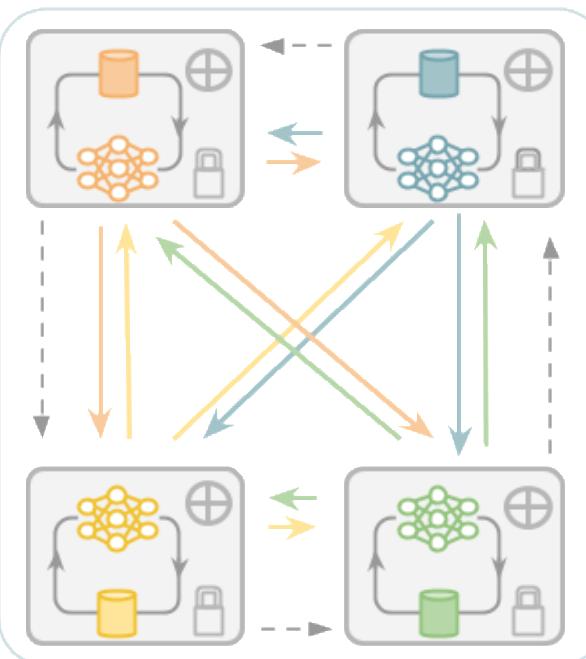


FEDERATED LEARNING ARCHITECTURES

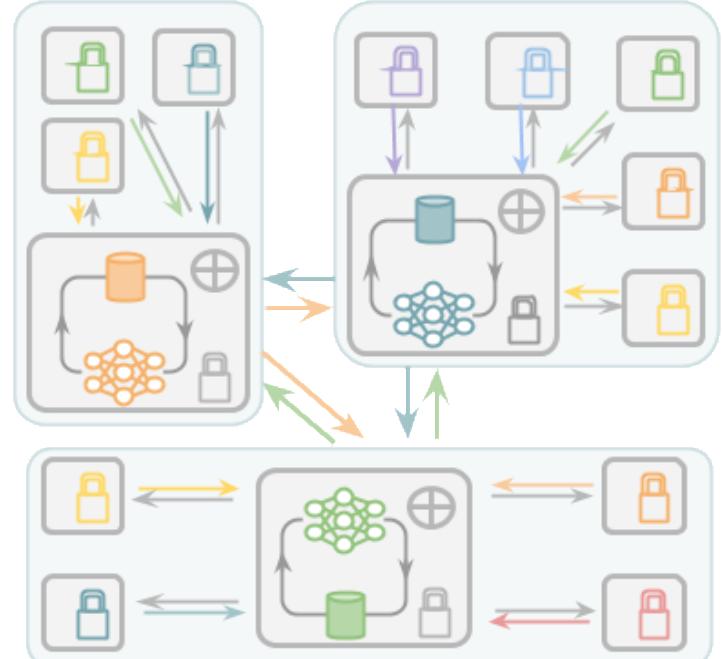
a) Client/Server Architecture



b) Peer to Peer Architecture



c) Hybrid Architecture



Federation of Nodes

Medical Database

→ Consensus Model
Redistribution

Model Aggregation

Locally Trained Model

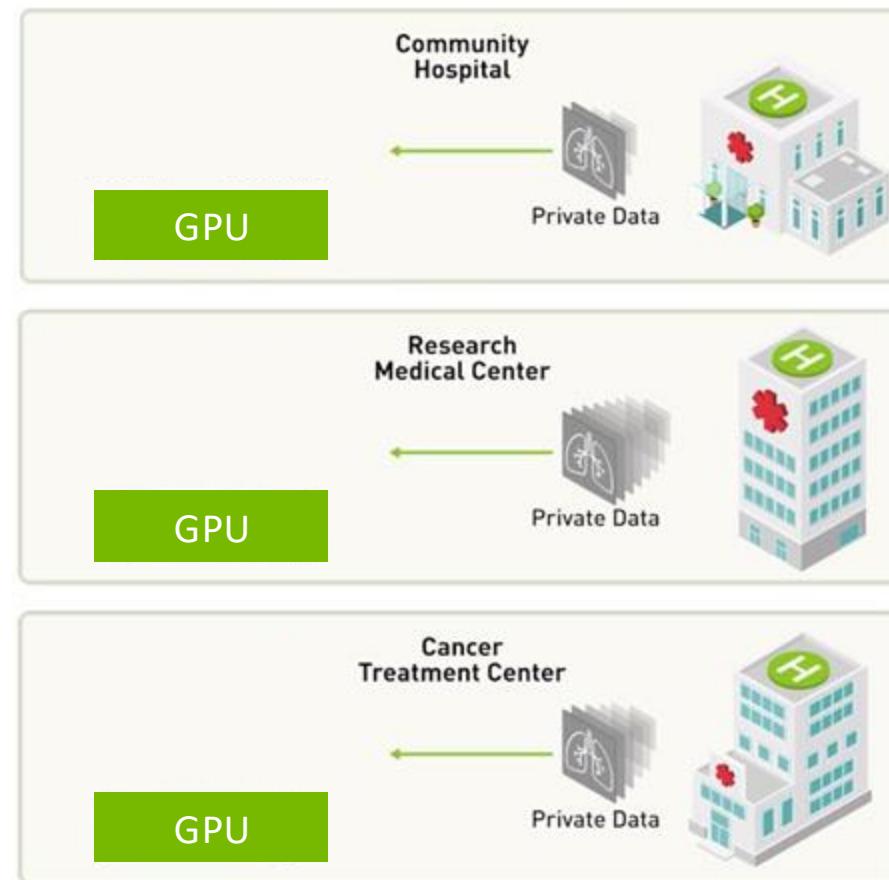
Secure Compute Node

→ Model Forwarding
Cyclic Learning

<https://arxiv.org/abs/2003.08119>

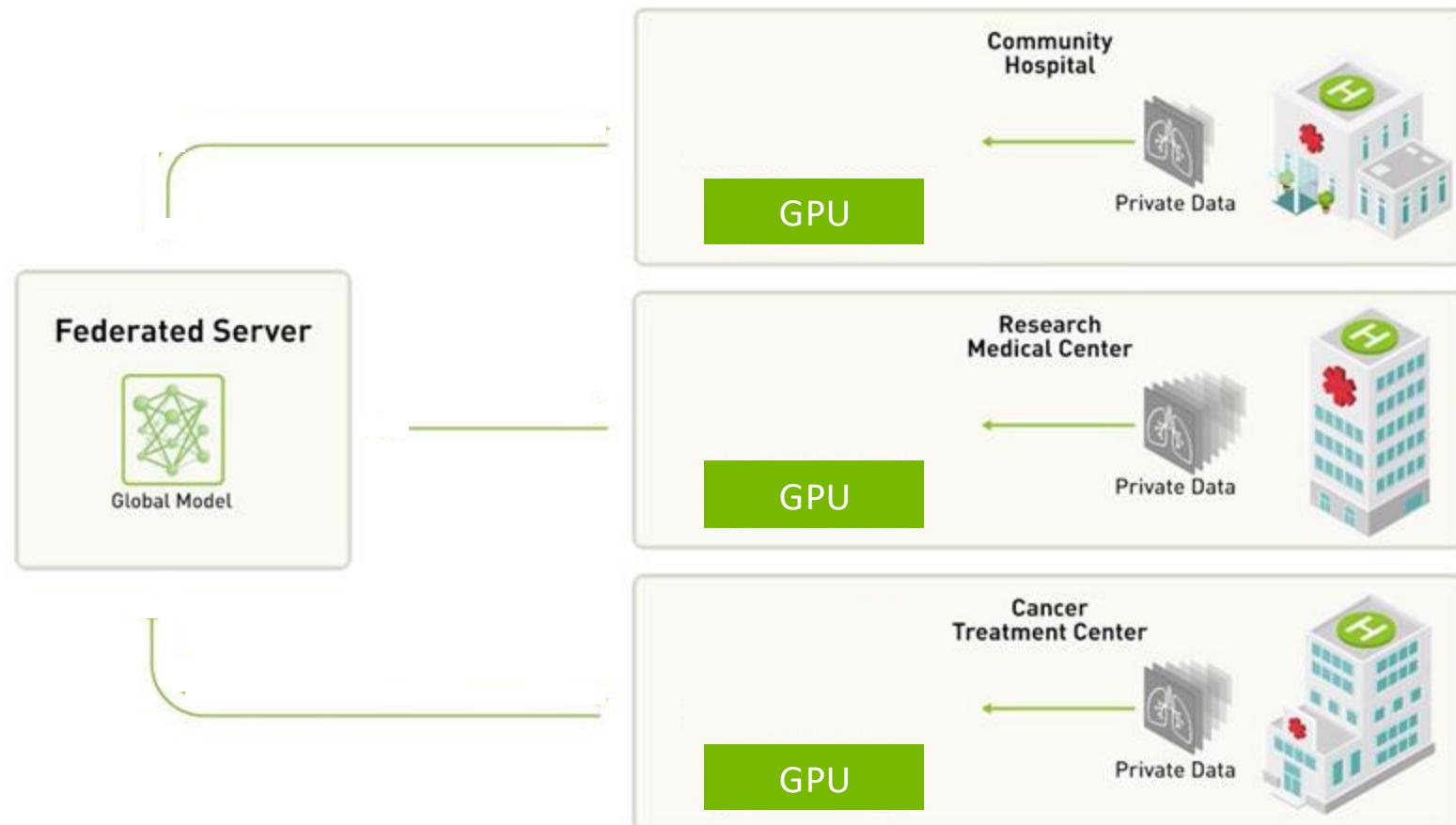
SERVER-CLIENT FEDERATED LEARNING

Changing the way AI algorithms are trained



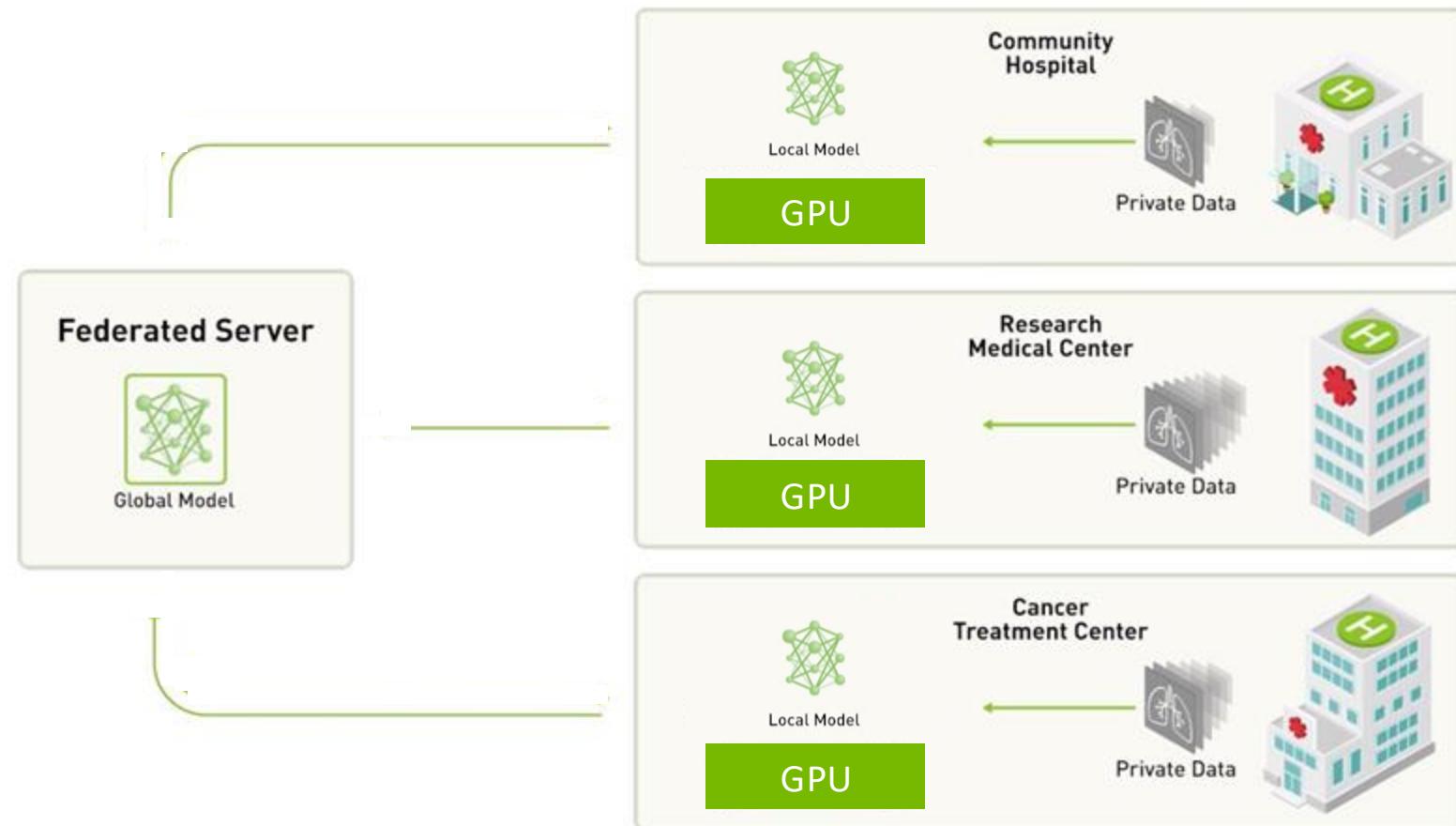
SERVER-CLIENT FEDERATED LEARNING

Changing the way AI algorithms are trained



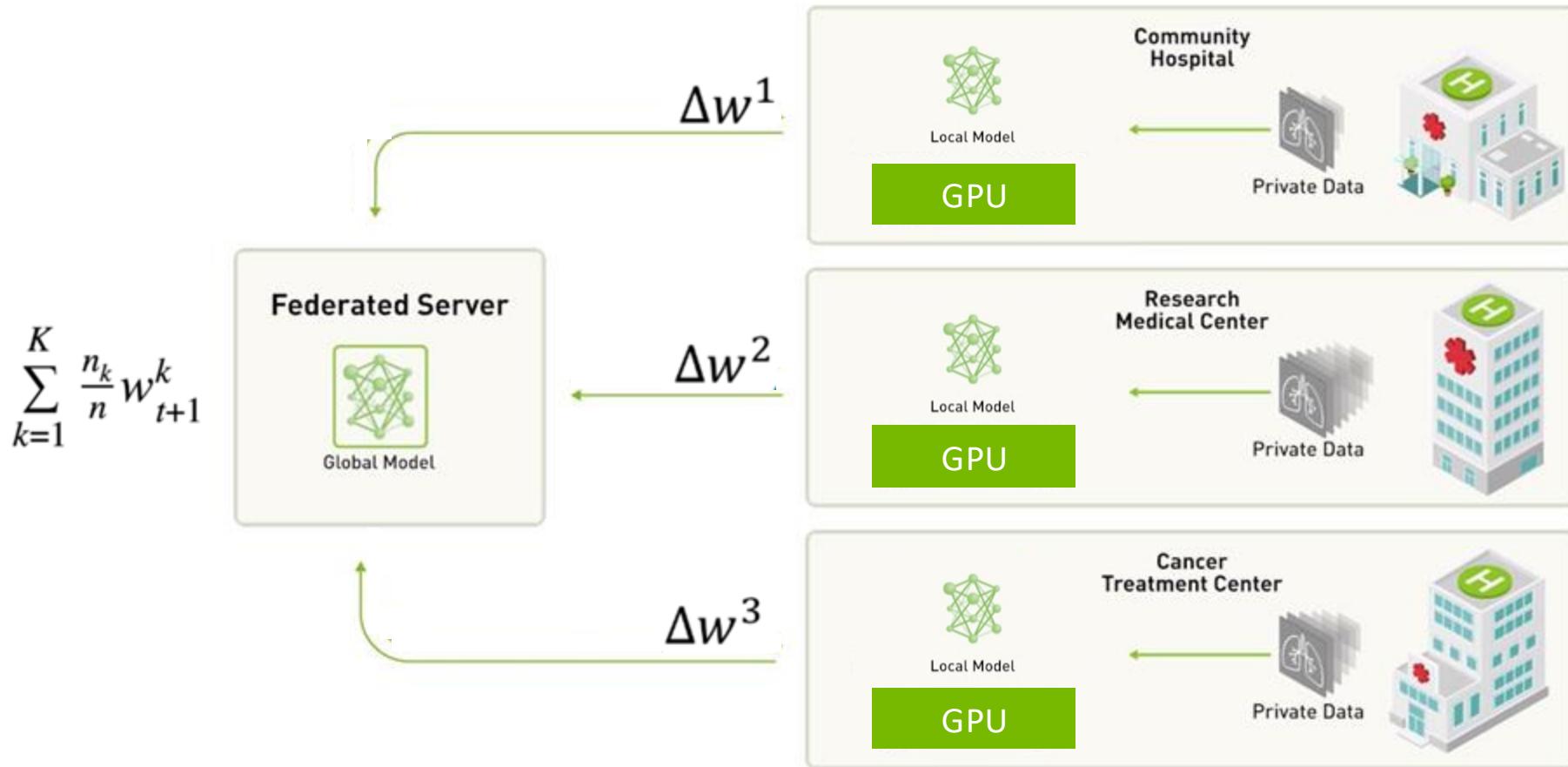
SERVER-CLIENT FEDERATED LEARNING

Changing the way AI algorithms are trained



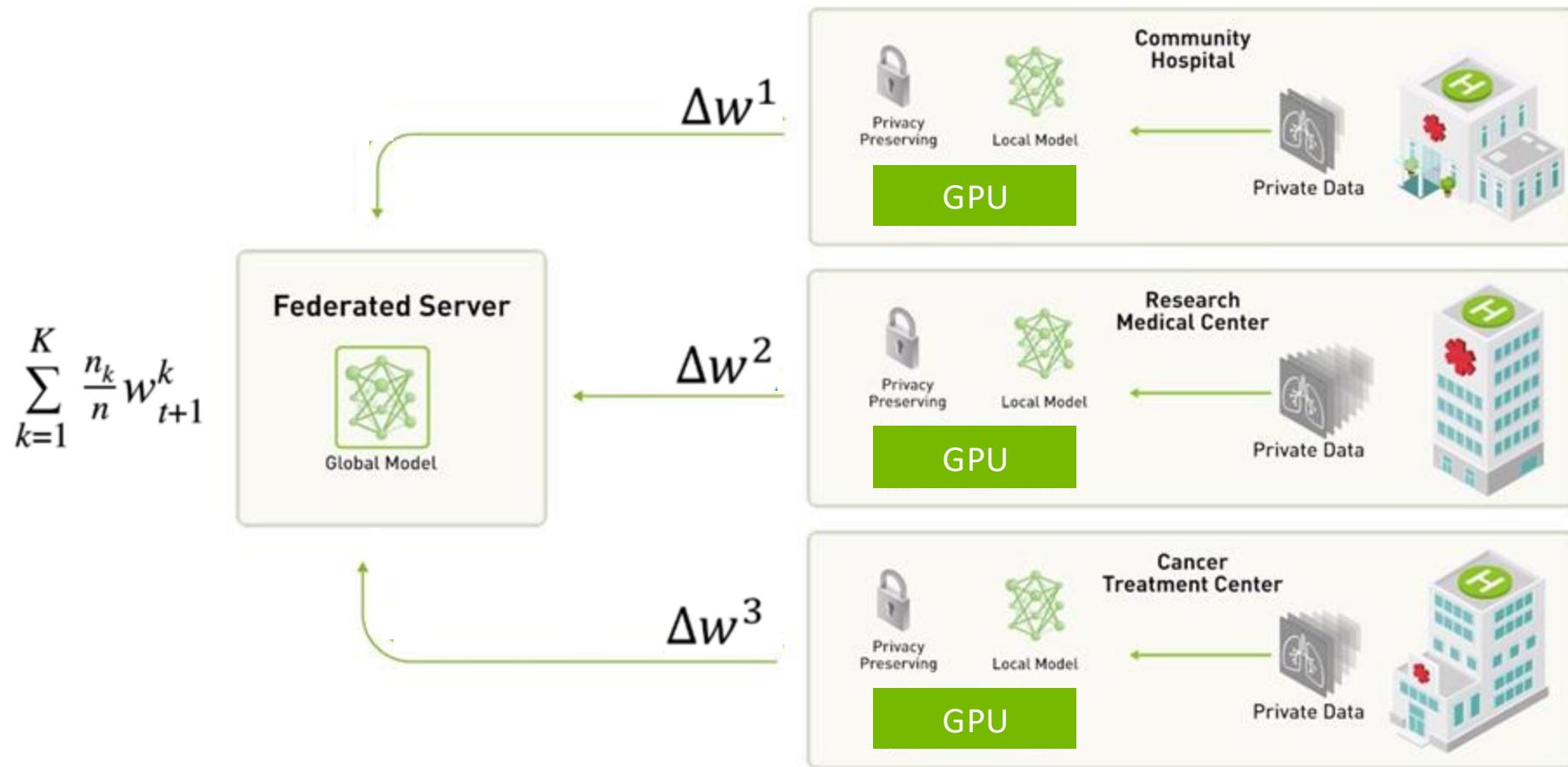
SERVER-CLIENT FEDERATED LEARNING

Changing the way AI algorithms are trained



SERVER-CLIENT FEDERATED LEARNING

Changing the way AI algorithms are trained





Agenda

- What is Federated Learning?
 - **NVIDIA Key Technologies for Federated Learning**
 - Real-world Use cases of Federated Learning
 - Getting Started with NVIDIA FLARE
 - Research: Addressing Key challenges in Federated Learning
 - Summary & Announcements
-

NVIDIA FLARE is Enabling FL in Many Industries

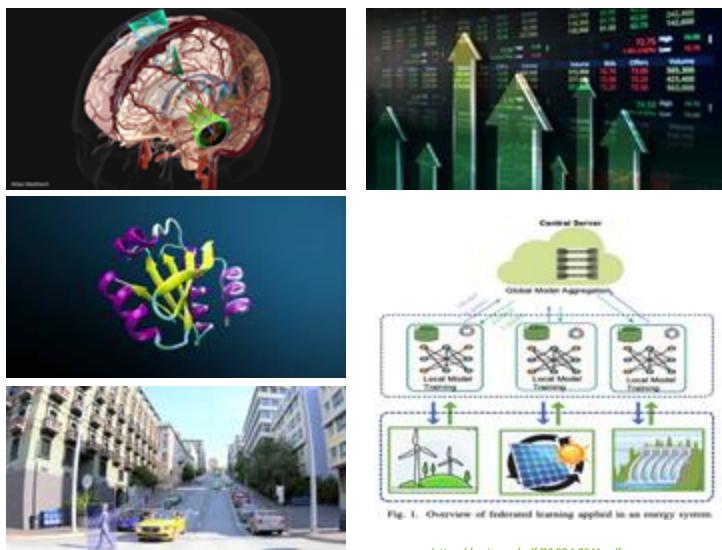
Healthcare – medical imaging, oncology

Financial Service – fraud detection

Pharmaceutical Industry – drug discovery, predictive model

Energy System – distributed energy sources

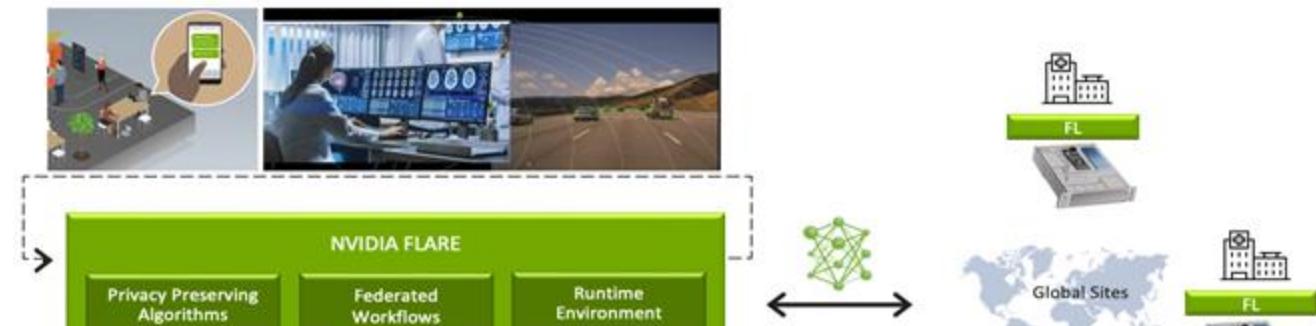
Automotive Industry – autonomous vehicle, object detection & classification



source: <https://arxiv.org/pdf/2208.10941.pdf>

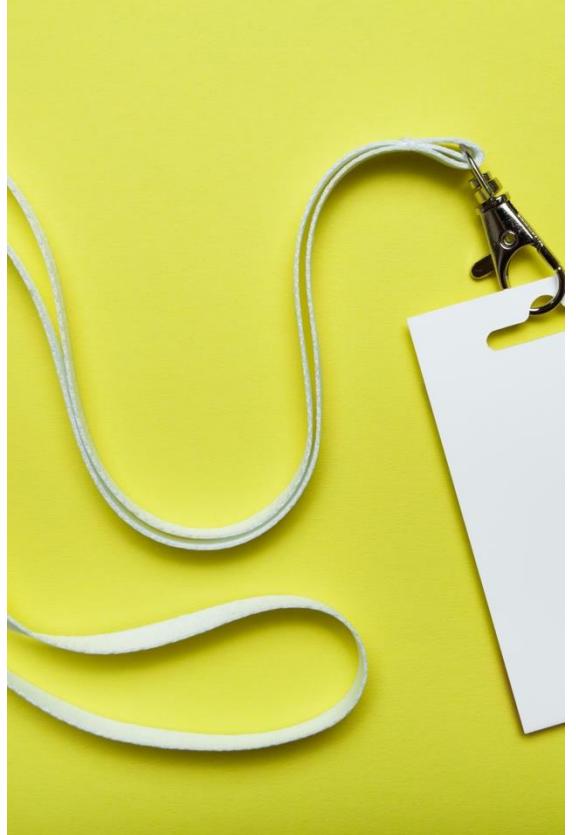
From Research To Production

simulation | easy transition from ML to FL | enterprise security | cross-cloud deployment



NVIDIA FLARE Security and Data Privacy

Defense in Depth approach to protecting data privacy and model IP



User Identity Verification

Certificate and derived token authentication



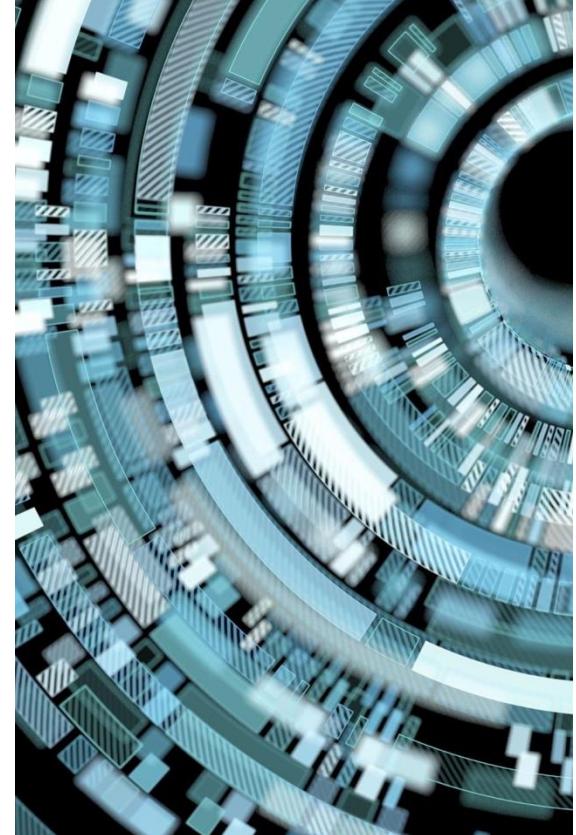
Data Encryption in Transit

Server-Client communication encrypted



User Defined Security Policy

Site-Specific authentication



Privacy Preserving Algorithm

Differential Privacy

Homomorphic Encryption

Confidential Computing

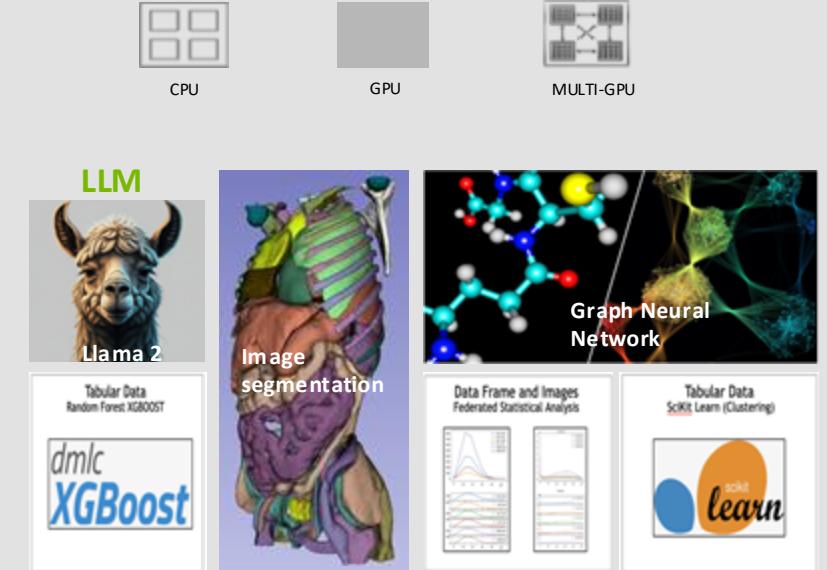
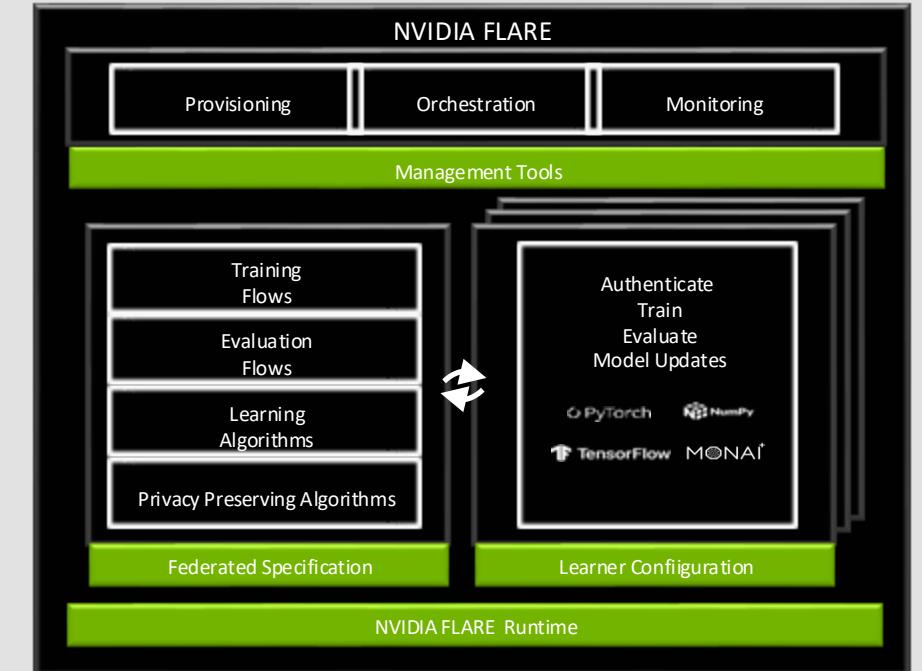
NVIDIA FLARE

Open-Source, Enterprise Federated Learning & Compute Framework

- Apache License 2.0 to catalyze FL research & development
- Designed for production, not just for research
- Enables cross-internet, distributed, multi-party collaborative Learning
- Production scalability with high availability and multi-task execution
- Easy to convert existing ML/DL workflows to a Federated paradigm with few lines of code changes
- LLM streaming, LLM fine tuning
- Framework, model, domain and task agnostic
- Privacy Preserving Technologies
 - Homomorphic Encryption (HE), Differential Privacy (DP)
 - Multi-party computing (Private Set intersection, PSI)
 - Confidential Computing (CC)
- Flexible communication patterns: server-centric, p2p, federated averaging, split learning, swarm learning, etc.
- Layered, pluggable, customizable federated compute architecture
- Secure Provisioning, Orchestration & Monitoring

GitHub: <https://github.com/nvidia/nvFlare>

Web: <https://nvidia.github.io/NVFlare>

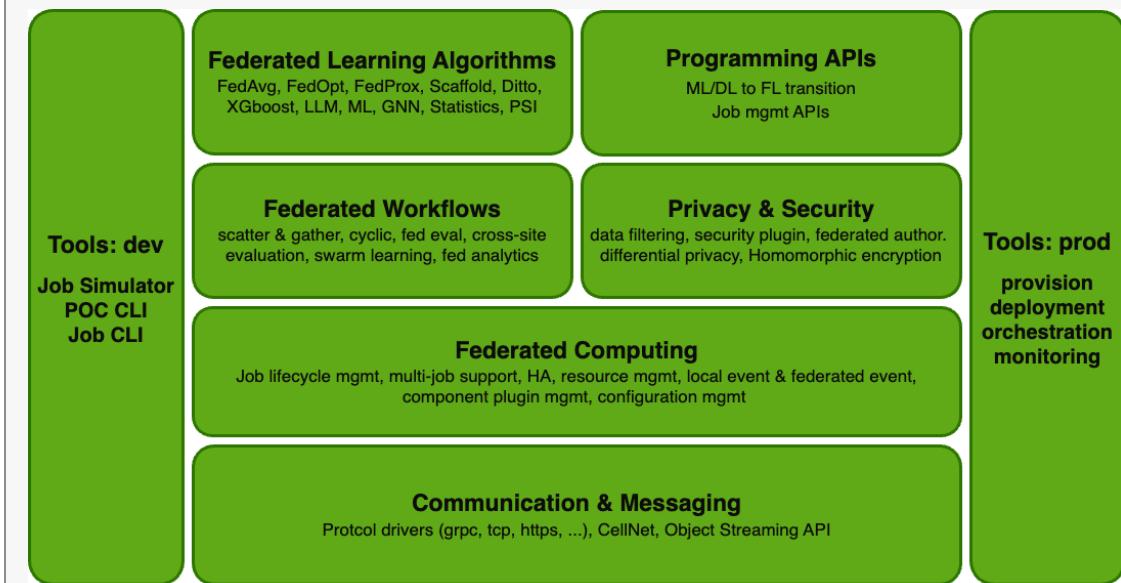


Framework agnostic | Model agnostic | Domain agnostic | Task agnostic

NVIDIA FLARE Architecture

Federated Computing Engine

- **Layered, Pluggable Open Architecture**
 - Each layer's component are composable and pluggable
- **Network: Communication & Messaging layer**
 - Drivers → gRPC, http + websocket, TCP, any plugin driver
 - CellNet: logical end point-to-point (cell to cell) network
 - Message: reliable streaming message
- **Federated Computing Layer**
 - Resource-based job scheduling, job monitoring, concurrent job lifecycle management, High-availability management
 - Plugin component management
 - Configuration management
 - Local event and federated event handling
- **Federated Workflow**
 - SAG, Cyclic, Cross-site Evaluation, split learning, Swarm Learning, Federated Analytics
- **Federated Learning Algorithms**
 - FedAvg, FedOpt, FedProx, Scaffold, Ditto, XGBoost, GNN, PSI, LLM (p-tuning, SFT, PEFT), KM, Scikit-Learn
- **Pythonic Programming APIs**
 - Client API, Controller API, Job Construction API, Job Monitoring API
- **Productivity & Deployment Tools:**
 - Simulator, provision, POC, Cloud deployment, preflight check, more



Use Cases and Integrations

Applications and contributions

- Use cases & contributions

- Cancer Detection with FLARE + MONAI - SIIM (society for imaging informatics and Medicine)
- Renal cell carcinoma: multi-institute collaboration including Case Western, Mayo Clinic, UFL, Vanderbilt, etc.
- Oncology for cancer patients: EU OPTIMA consortium
- Cancer research: OHSU and NCI
- Digital Pathology: Roche
- Federated Learning Interoperability Platform: Kings College London
- Cancer detection / classification: GE Healthcare
- Supercomputing: Oak Ridge National Lab
- Transportation: Cummins
- LLM: Deloitte
- NLP research: UMN
- FL algorithm: UBC

- Integrations

- Medical image analysis framework: MONAI
- Generative AI platform: NeMo
- FL framework: Flower



Use Cases and Integrations

Applications and contributions

- Use cases & contributions

- Cancer Detection with FLARE + MONAI - SIIM (society for imaging informatics and Medicine)
- Renal cell carcinoma: multi-institute collaboration including Case Western, Mayo Clinic, UFL, Vanderbilt, etc.
- Oncology for cancer patients: EU OPTIMA consortium
- Cancer research: OHSU and NCI
- Digital Pathology: Roche
- Federated Learning Interoperability Platform: Kings College London
- Cancer detection / classification: GE Healthcare
- Supercomputing: Oak Ridge National Lab
- Transportation: Cummins
- LLM: Deloitte
- NLP research: UMN
- FL algorithm: UBC

- Integrations

- Medical image analysis framework: MONAI
- Generative AI platform: NeMo
- FL framework: Flower

Deloitte.



Siim

Roche

University of Minnesota



University of Minnesota

**School of Medicine
and Public Health**

UNIVERSITY OF WISCONSIN-MADISON



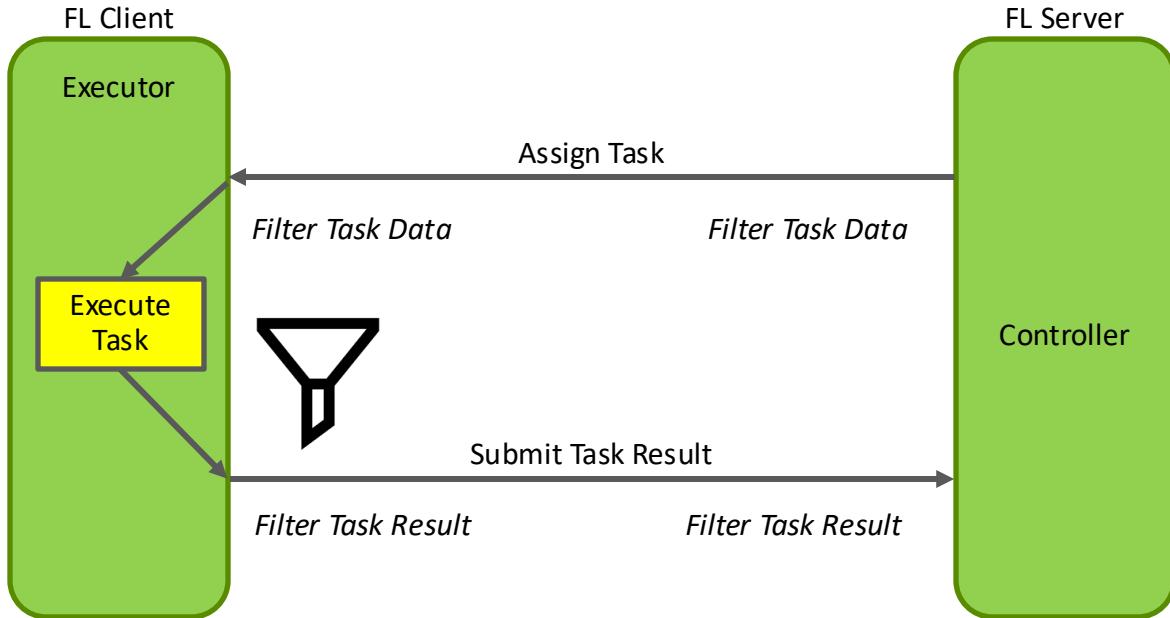
...many more



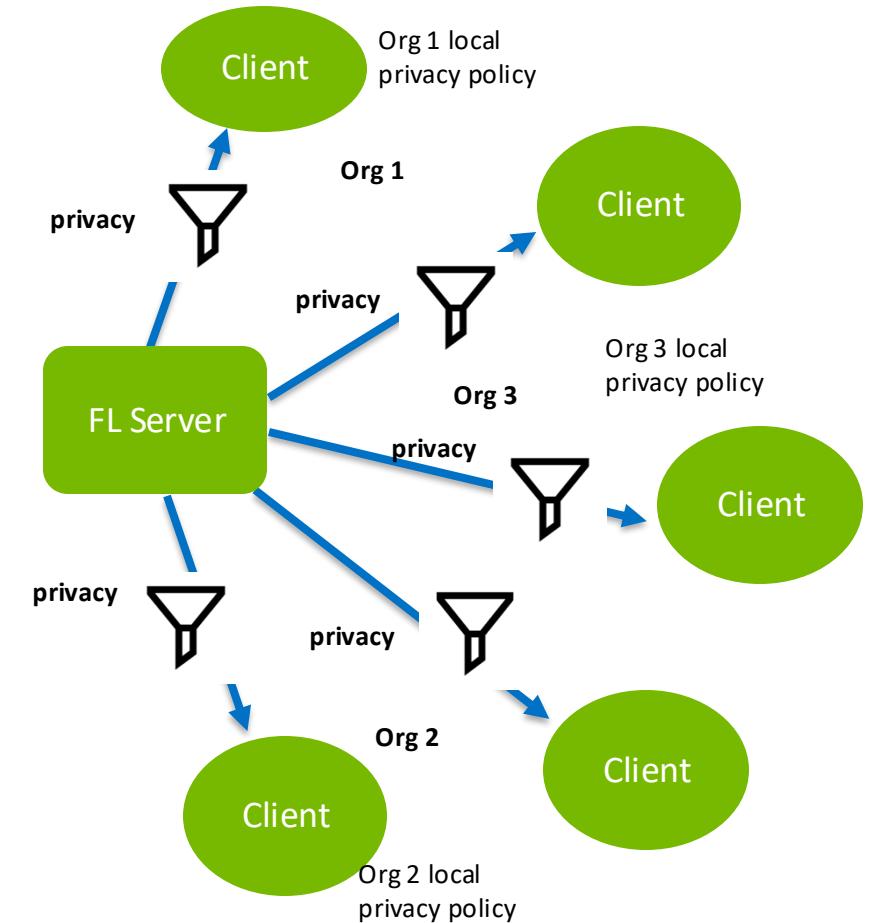
**CASE
WESTERN
RESERVE
UNIVERSITY**

HIGH LEVEL ARCHITECTURE

Data privacy architecture



- Privacy filter can depend on:**
 - Scope: any key-value pair such as datasets
 - Data kind: Weights, Weights Diff or Analytics data
 - Or any other data
- Research develop privacy filter**
- Organization set privacy policy:**
 - privacy budget, noise level as data privacy policy



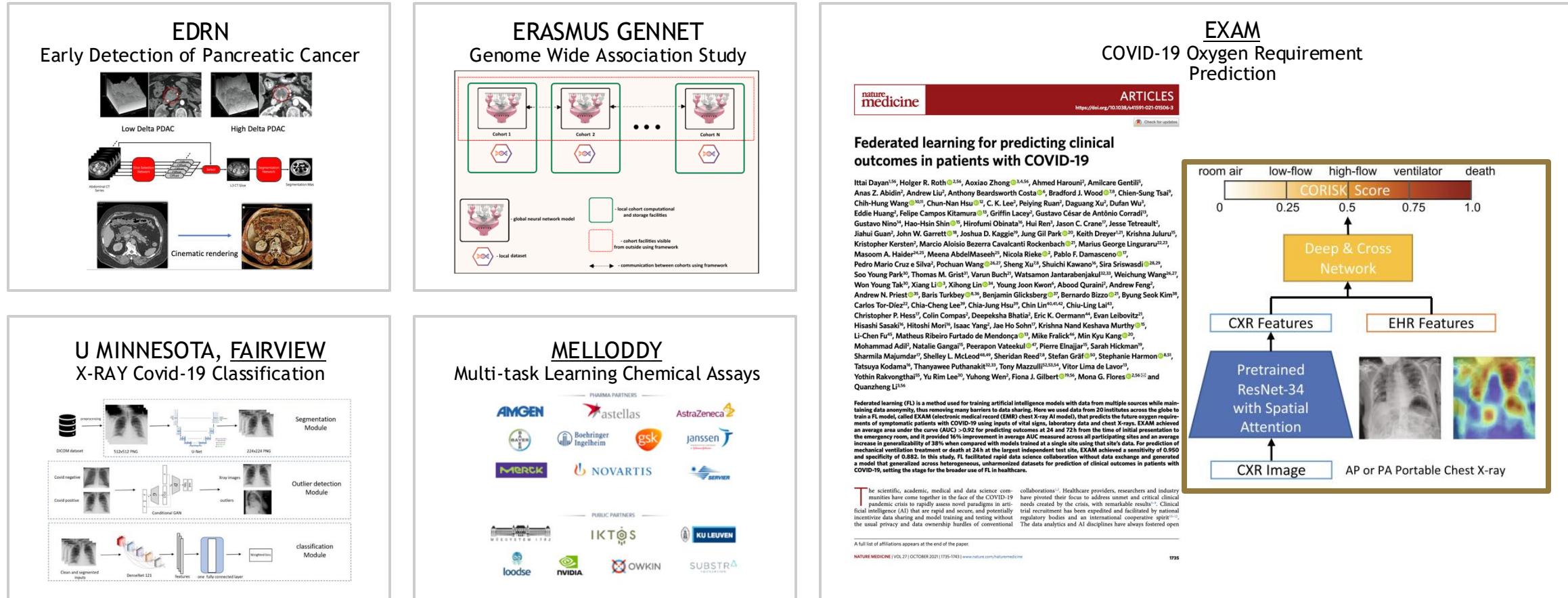


Agenda

- What is Federated Learning?
 - NVIDIA Key Technologies for Federated Learning
 - **Real-world Use cases of Federated Learning**
 - Getting Started with NVIDIA FLARE
 - Research: Addressing Key challenges in Federated Learning
 - Summary & Announcements
-
-
-
-
-
-

Federated learning In HEALTHCARE

Breaking Healthcare Data Siloes



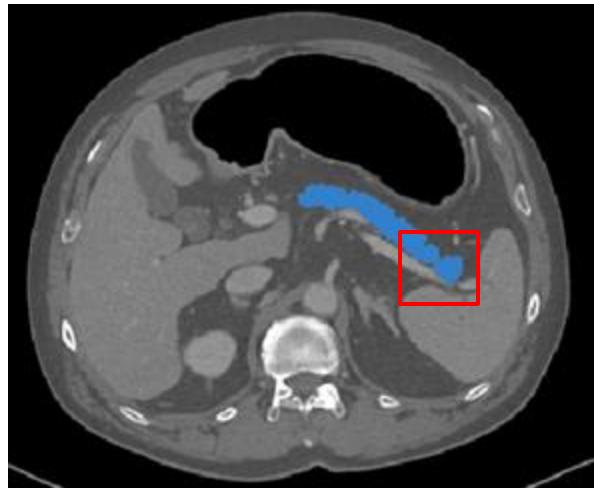
EXAM: <https://www.nature.com/articles/s41591-021-01506-3>

MELLODDY: <https://chemrxiv.org/engage/chemrxiv/article-details/6345c0f91f323d61d7567624>

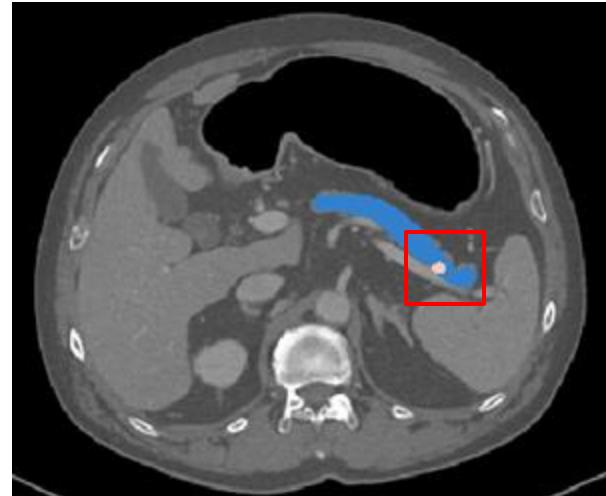
<https://developer.download.nvidia.com/CLARA/Federated-Learning-Training-for-Healthcare-Using-NVIDIA-Clara.pdf>

FEDERATED LEARNING

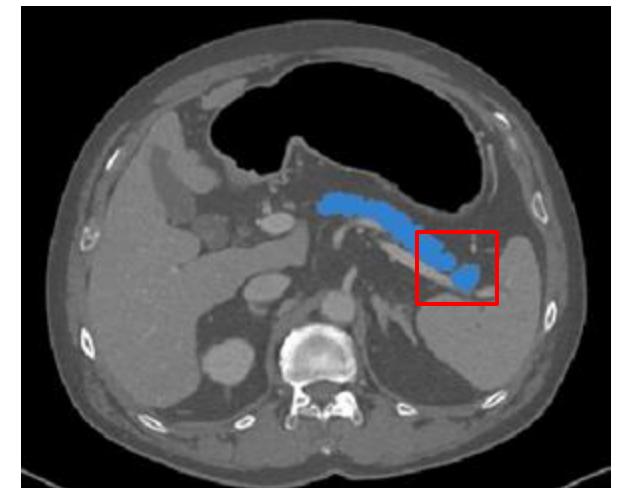
Global vs. local



Ground truth



C2 (local)
Trained on pancreas
& tumor

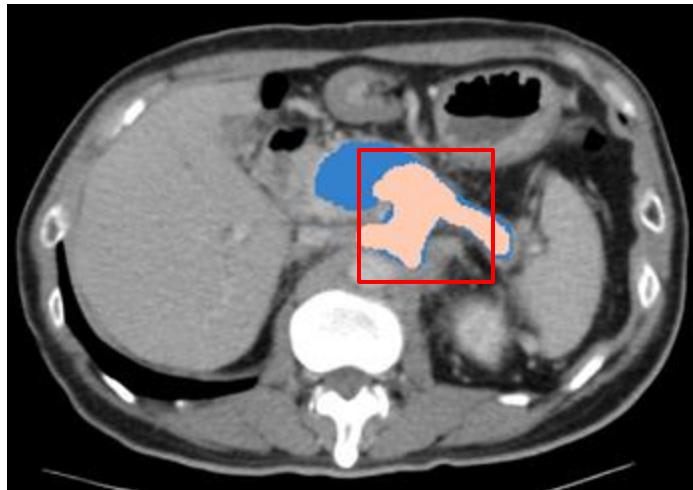


FL (global)

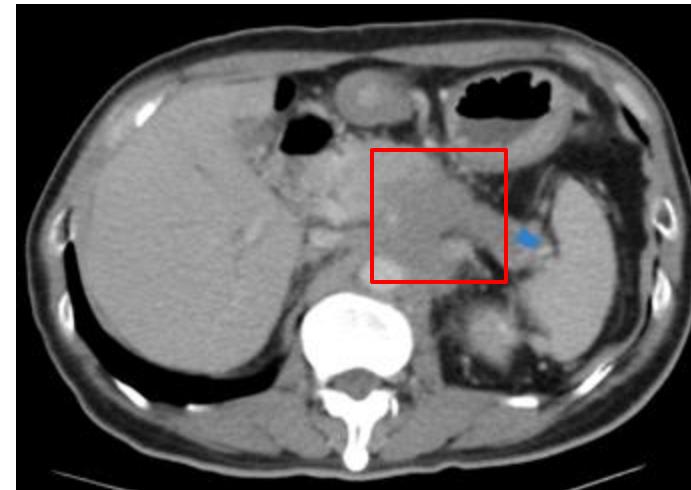
Dataset	Client 1	Client 2
Cases	420	486
Label	2 class (background, pancreas)	3 classes (background, pancreas, tumor)

FEDERATED LEARNING

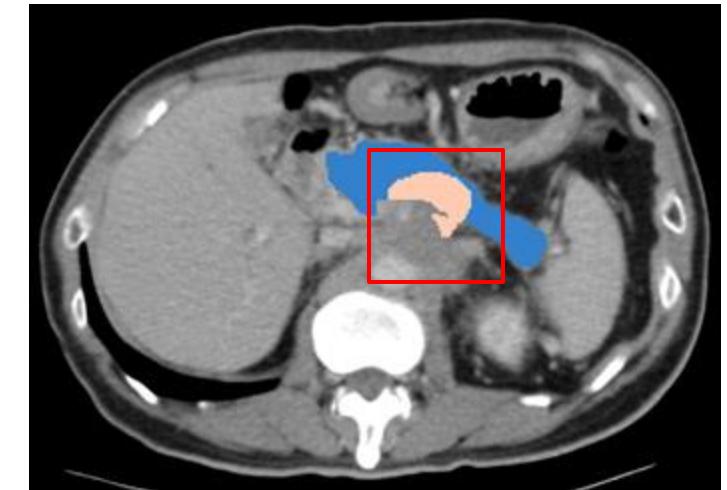
Global vs. local



Ground truth



C1 (local)
Trained on healthy
pancreas



FL (global)

Dataset	Client 1	Client 2
Cases	420	486
Label	2 class (background, pancreas)	3 classes (background, pancreas, tumor)

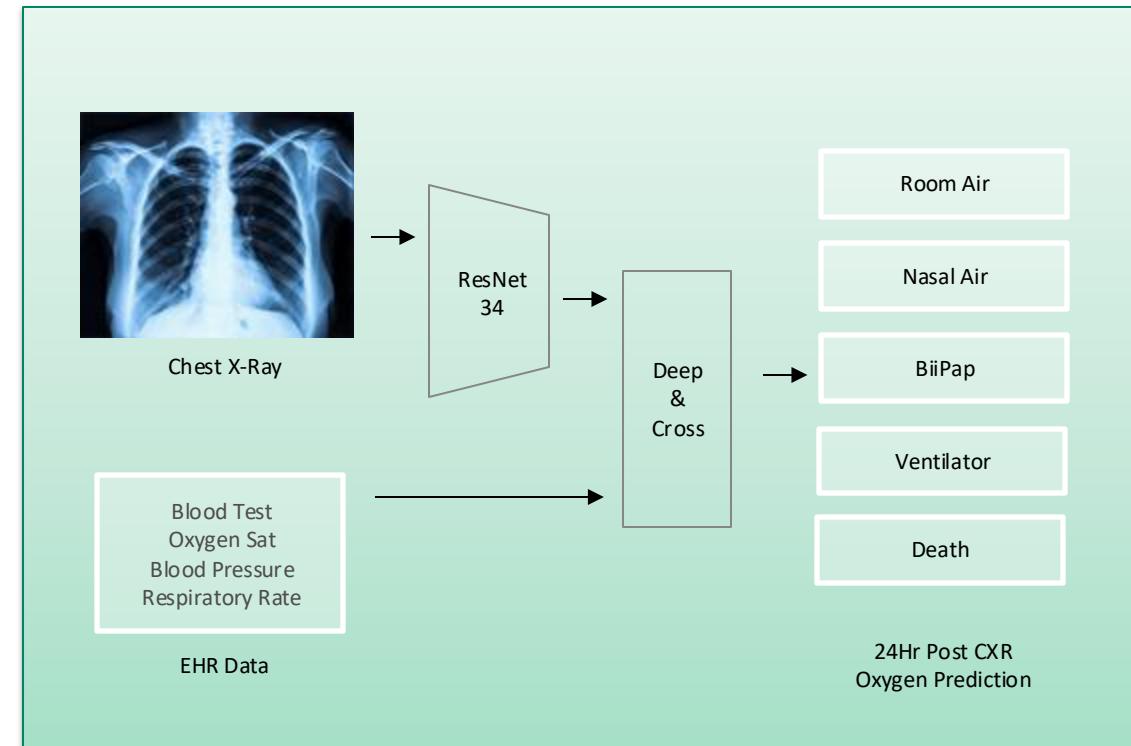
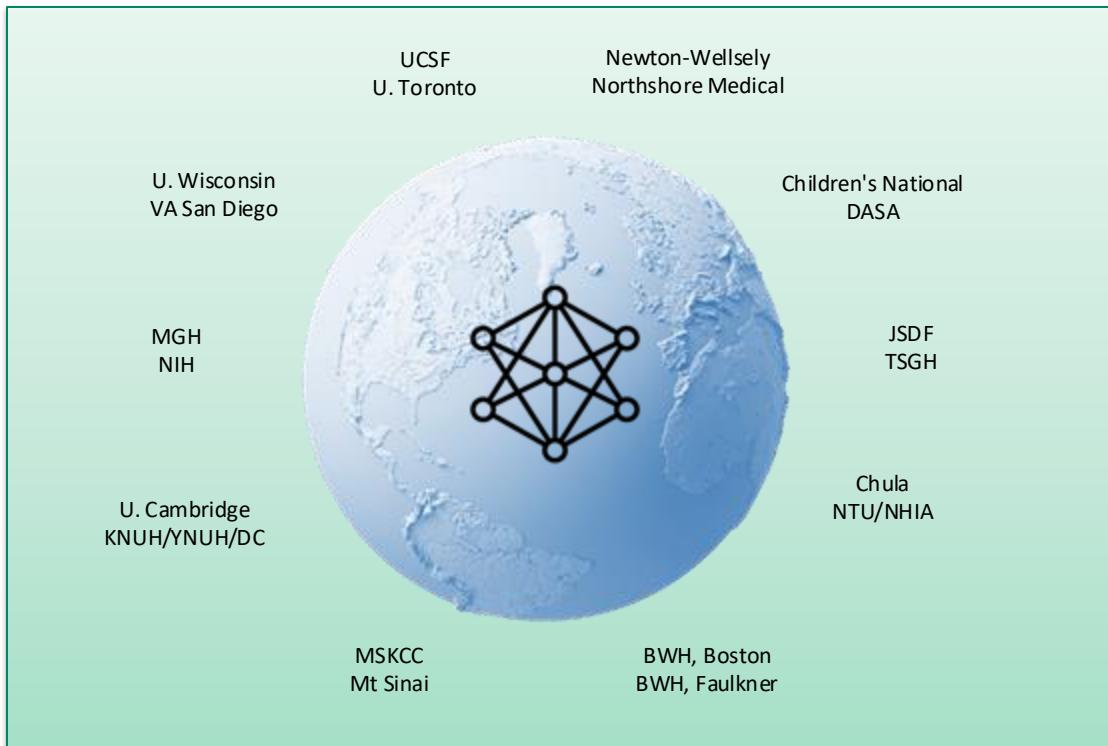


Federated learning for predicting clinical outcomes in patients with COVID-19

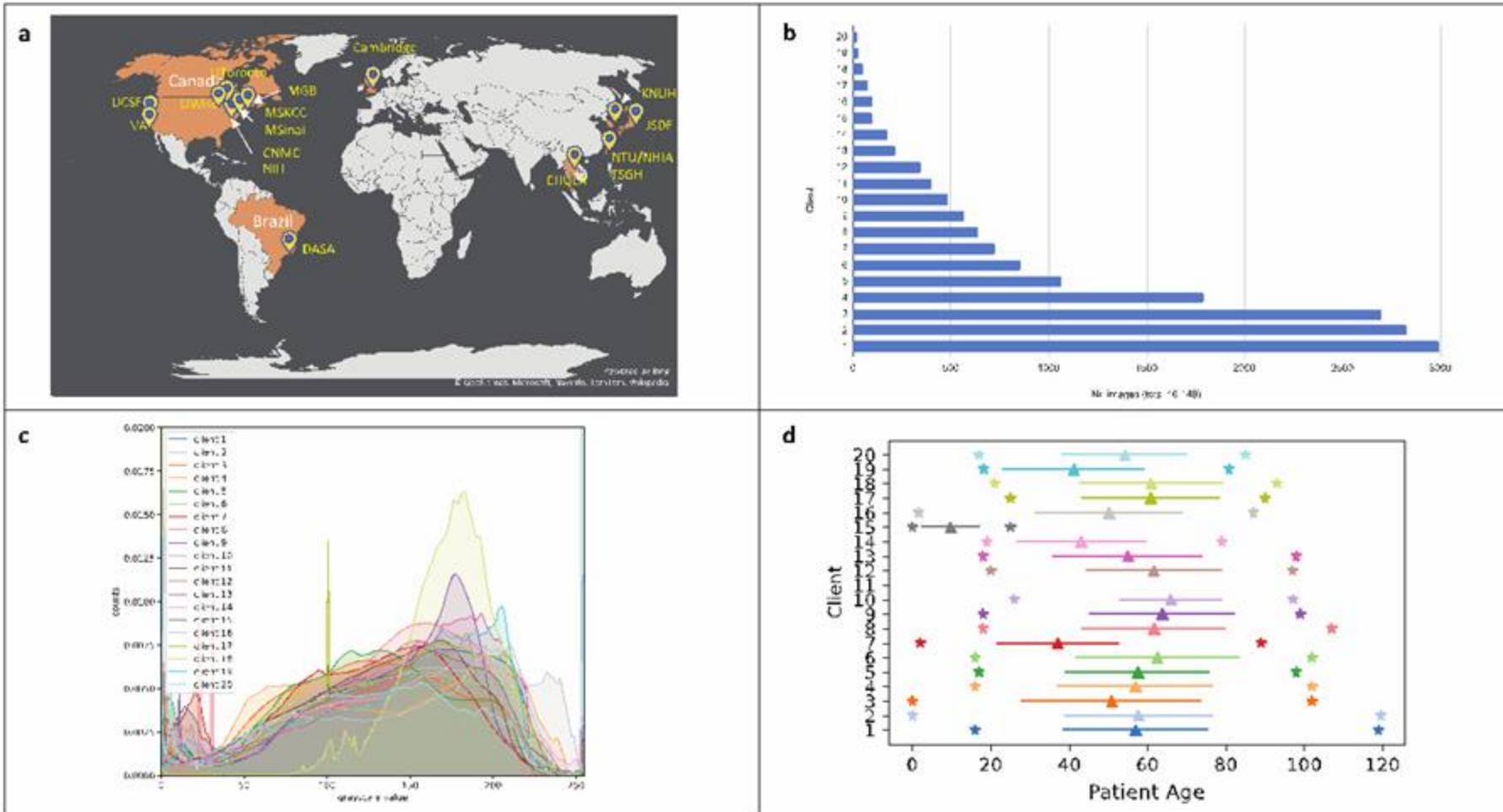
Ittai Dayan^{1,56}, Holger R. Roth^{1,2,56}, Aoxiao Zhong^{1,3,4,56}, Ahmed Harouni², Amilcare Gentili⁵, Anas Z. Abidin², Andrew Liu², Anthony Beardsworth Costa^{1,6}, Bradford J. Wood^{1,7,8}, Chien-Sung Tsai⁹, Chih-Hung Wang^{1,10,11}, Chun-Nan Hsu^{1,12}, C. K. Lee², Peiying Ruan², Daguang Xu², Dufan Wu³, Eddie Huang², Felipe Campos Kitamura^{1,13}, Griffin Lacey², Gustavo César de Antônio Corradi^{1,13}, Gustavo Nino¹⁴, Hao-Hsin Shin^{1,15}, Hirofumi Obinata¹⁶, Hui Ren³, Jason C. Crane¹⁷, Jesse Tetreault², Jiahui Guan², John W. Garrett^{1,18}, Joshua D. Kaggie¹⁹, Jung Gil Park^{1,20}, Keith Dreyer^{1,21}, Krishna Juluru¹⁵, Kristopher Kersten², Marcio Aloisio Bezerra Cavalcanti Rockenbach^{1,21}, Marius George Linguraru^{22,23}, Masoom A. Haider^{24,25}, Meena AbdelMaseeh²⁵, Nicola Rieke^{1,2}, Pablo F. Damasceno^{1,17}, Pedro Mario Cruz e Silva², Pochuan Wang^{1,26,27}, Sheng Xu^{7,8}, Shuichi Kawano¹⁶, Sira Sriswasdi^{1,28,29}, Soo Young Park³⁰, Thomas M. Grist³¹, Varun Buch²¹, Watsamon Jantarabenjakul^{32,33}, Weichung Wang^{26,27}, Won Young Tak³⁰, Xiang Li^{1,3}, Xihong Lin^{1,34}, Young Joon Kwon⁶, Abood Quraini², Andrew Feng², Andrew N. Priest^{1,35}, Baris Turkbey^{1,8,36}, Benjamin Glicksberg^{1,37}, Bernardo Bizzo^{1,21}, Byung Seok Kim³⁸, Carlos Tor-Díez²², Chia-Cheng Lee³⁹, Chia-Jung Hsu³⁹, Chin Lin^{40,41,42}, Chiu-Ling Lai⁴³, Christopher P. Hess¹⁷, Colin Compas², Deepeksha Bhatia², Eric K. Oermann⁴⁴, Evan Leibovitz²¹, Hisashi Sasaki¹⁶, Hitoshi Mori¹⁶, Isaac Yang², Jae Ho Sohn¹⁷, Krishna Nand Keshava Murthy^{1,15}, Li-Chen Fu⁴⁵, Matheus Ribeiro Furtado de Mendonça^{1,13}, Mike Fralick⁴⁶, Min Kyu Kang^{1,20}, Mohammad Adil², Natalie Gangai¹⁵, Peerapon Vateekul^{1,47}, Pierre Elnajjar¹⁵, Sarah Hickman¹⁹, Sharmila Majumdar¹⁷, Shelley L. McLeod^{48,49}, Sheridan Reed^{7,8}, Stefan Gräf^{1,50}, Stephanie Harmon^{1,8,51}, Tatsuya Kodama¹⁶, Thanyawee Puthanakit^{32,33}, Tony Mazzulli^{52,53,54}, Vitor Lima de Lavor¹³, Yothin Rakvongthai⁵⁵, Yu Rim Lee³⁰, Yuhong Wen², Fiona J. Gilbert^{1,19,56}, Mona G. Flores^{1,2,56}✉ and Quanzheng Li^{1,3,56}



NVIDIA FEDERATED LEARNING FOR COVID-19 PATIENT CARE



USING FEDERATED LEARNING TO LEVERAGE A GLOBAL DATASET

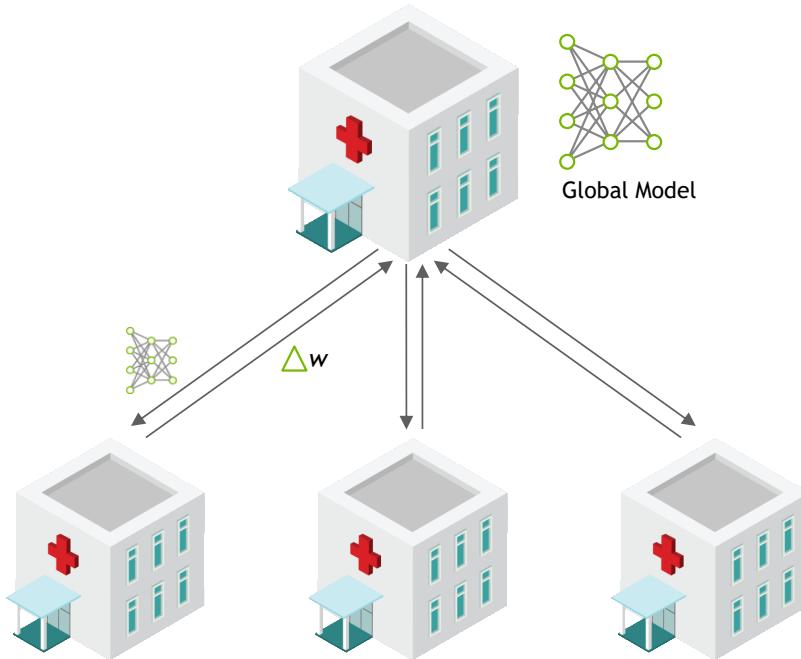


a, World map indicating the 20 different client sites contributing to the EXAM study. **b**, Number of cases contributed by each institution or site (client 1 represents the site contributing the largest number of cases). **c**, Chest X-ray intensity distribution at each client site. **d**, Age of patients at each client site, showing minimum and maximum ages (asterisks), mean age (triangles) and standard deviation (horizontal bars).

CONTROLLER FOR MODEL TRAINING

Typical workflow for FedAvg, FedOpt, FedProx, etc.

FedAvg

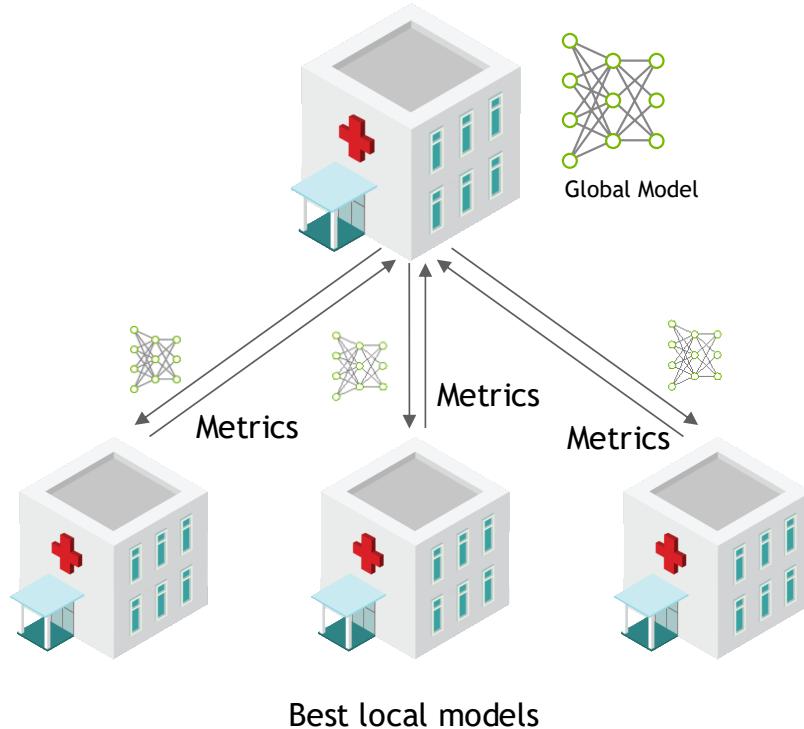


1. Server initializes model
2. For number of rounds:
 1. Server broadcasts global model to workers
 2. Workers validate global model and train on their data
 3. Workers keep track on their locally best model (Personalization)
 4. Workers send back updated model or updates
 5. Server Gathers (Aggregates) updates and updates the global model

Source: https://github.com/NVIDIA/NVFlare/blob/main/nvflare/app_common/workflows/fedavg.py

CONTROLLER FOR MODEL EVALUATION

Global model evaluation, Cross-site model evaluation



FedEval (Global Model Validation/Cross-Site Validation)

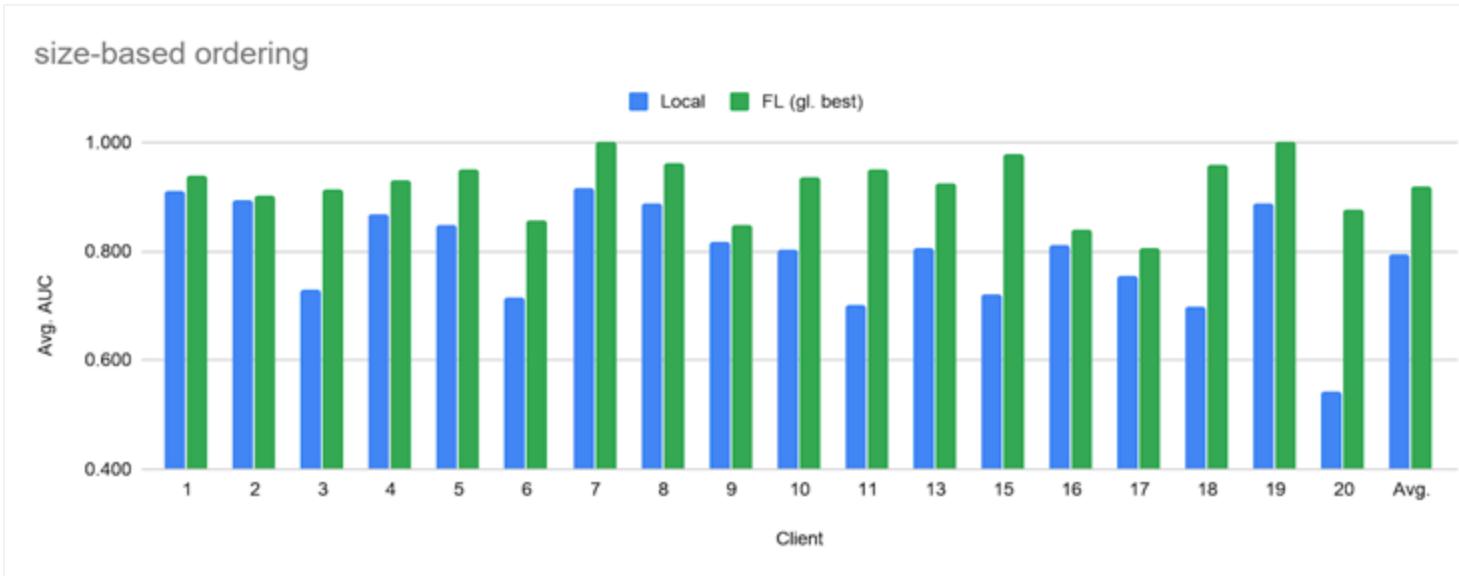
1. Server sends models (e.g., global model and registered best local models) to each worker for evaluation
2. Server gathers the resulting metrics

Models	Metrics	Evaluation sites			
		Site-1	Site-2	...	Site-N
	Global (Final)
	Global (Best)
	Site-1
	Site-2

	Site-N

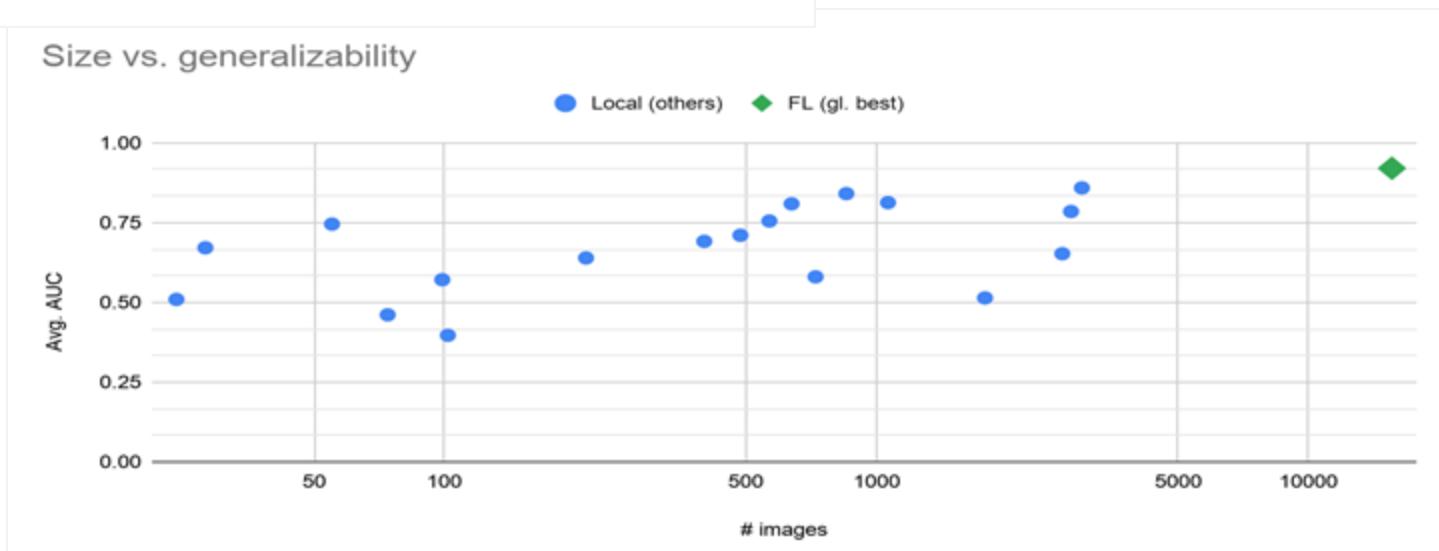
Source: https://github.com/NVIDIA/NVFlare/blob/main/nvflare/app_common/workflows/global_model_eval.py
https://github.com/NVIDIA/NVFlare/blob/main/nvflare/app_common/workflows/cross_site_model_eval.py

NVIDIA FL FOR COVID-19



FL resulted on average in:

- 16% performance improvement
- 38% generalizability improvement



NVIDIA FL FOR COVID-19: EXTERNAL VALIDATION



Predicting need for mechanical ventilation

EXAM Model Released on NGC:

24h avg. AUC: 0.94

72h avg. AUC: 0.91

https://ngc.nvidia.com/catalog/models/nvidia:med:clara_train_covid19_exam_ehr_xray

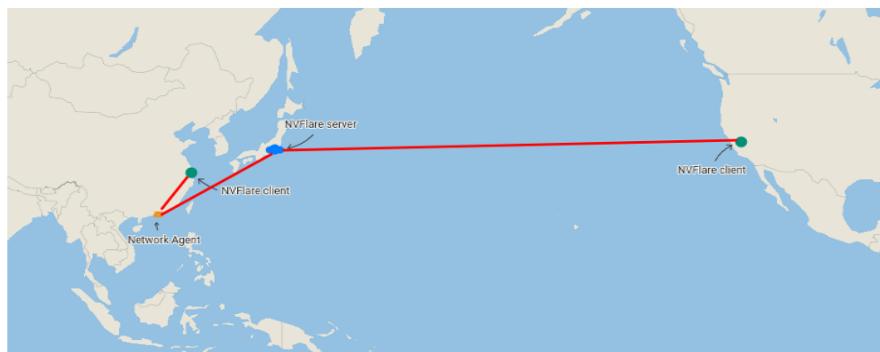
YouTube: <https://youtu.be/cOXVrtkv6FE>

NVIDIA AV Federated Learning

AV team leveraging FLARE to training object detection & tracking models

NVIDIA Autonomous Vehicle model training

- **Scenario:** small number of clients, and each client host a big dataset.
- **Workflow:** cyclic weight transfer
- **Goals:** Build a unified global model with the same or better accuracy. Reduce the effort for model approval process
- **Total Users:** 29 active users in the platform
- **Models trained:**
- **MLMCF** - A model to detect vehicles, person, bicycle, free space, parking area
- **DoNET** - A model to detect the status of vehicles, for example, lamp status (is light on/off), door status (is door open or not), etc.
- **Waitnet** - A model to detect static objects such as traffic lights, traffic signs, road marks, stop lines, and crosswalks.
- **PathNet/RoadNet** — Detecting the path the autonomous vehicle takes.
- **RadarNet** — Radar sensor data as input to predict the obstacle around the vehicle.
- **PredetionNet** — Object tracking and trajectory prediction
- **EGM** — Multi-camera input-based model for detecting the barrier in the parking lot such as height limit pole, pillar.



Cross-border federated learning setup: Clients in China & USA
Servers in Japan, Hong Kong



US Model: only capture the whole sign

FL Global Model: capture the individual objects and whole sign



Agenda

- What is Federated Learning?
 - NVIDIA Key Technologies for Federated Learning
 - Real-world Use cases of Federated Learning
 - **Getting Started with NVIDIA FLARE**
 - Research: Addressing Key challenges in Federated Learning
 - Summary & Announcements
-
-
-
-
-
-

Server Code: Controller

```
18 class FedAvg(BaseFedAvg):
19
20     def run(self) -> None:
21         self.info("Start FedAvg.")
22
23         model = self.load_model()
24         model.start_round = self.start_round
25         model.total_rounds = self.num_rounds
26
27         for self.current_round in range(self.start_round, self.start_round + self.num_rounds):
28             self.info(f"Round {self.current_round} started.")
29             model.current_round = self.current_round
30
31             clients = self.sample_clients(self.min_clients)
32
33             results = self.send_model_and_wait(targets=clients, data=model)
34
35             aggregate_results = self.aggregate(
36                 results, aggregate_fn=None
37             ) # if no `aggregate_fn` provided, default `WeightedAggregationHelper` is used
38
39             model = self.update_model(model, aggregate_results)
40
41             self.save_model(model)
42
43         self.info("Finished FedAvg.")
```

Client Code: Convert PyTorch to NVFlare

PyTorch CIFAR-10 Tutorial

```
6 from net import Net  
7  
8  
9 def main():  
10    transform = transforms.Compose([...])  
11    trainset = torchvision.datasets.CIFAR10(...)  
12    trainloader = torch.utils.data.DataLoader(...)  
13    testset = torchvision.datasets.CIFAR10(...)  
14    testloader = torch.utils.data.DataLoader(...)  
15    net = Net()  
16    criterion = nn.CrossEntropyLoss()  
17    optimizer = optim.SGD(...)  
18  
19    # Train loop  
20    for epoch in range(epochs):  
21        ...  
22  
23    print("Finished Training")
```



```
6 from net import Net  
7  
8 import nvflare.client as flare  
9  
10  
11 def main():  
12     transform = transforms.Compose([...])  
13  
14     trainset = torchvision.datasets.CIFAR10(...)  
15     trainloader = torch.utils.data.DataLoader(...)  
16  
17     testset = torchvision.datasets.CIFAR10(...)  
18     testloader = torch.utils.data.DataLoader(...)  
19  
20     net = Net()  
21  
22     criterion = nn.CrossEntropyLoss()  
23     optimizer = optim.SGD(...)  
24  
25     flare.init()                                2. Initialize  
26  
27     while flare.is_running():  
28         input_model = flare.receive()  
29         print(f"current_round={input_model.current_round}")      3. Receive global model  
30  
31         net.load_state_dict(input_model.params)           4. Load global model  
32  
33         for epoch in range(epochs): # loop over the dataset multiple times  
34             ...  
35  
36         print("Finished Training")                      5. Send back the updated model  
37  
38         output_model = flare.FLModel(  
39             params=net.cpu().state_dict(),  
40             metrics={"accuracy": accuracy},  
41             meta={"NUM_STEPS_CURRENT_ROUND": epochs * len(trainloader)},  
42         )  
43         flare.send(output_model)
```

Create a FedJob and Run Simulation

```
7 ► if __name__ == "__main__":
8     n_clients = 2
9     num_rounds = 2
10    train_script = "src/cifar10_fl.py"
11
12    # Create basic fed Job with initial model
13    job = BaseFedJob(
14        name="cifar10_pt_fedavg",
15        initial_model=Net(),
16    )
17
18    # Define the controller and send to server
19    controller = FedAvg(
20        num_clients=n_clients,
21        num_rounds=num_rounds,
22    )
23    job.to_server(controller)
24
25    # Add clients
26    for i in range(n_clients):
27        runner = ScriptRunner(script=train_script) # script_args=f"--batch_size 32 --data_path /data/site-{i}"
28        job.to(runner, target: f"site-{i}")
29
30    # job.export_job("/tmp/nvflare/jobs/job_config") # Exported jobs can be used in real deployment!
31    job.simulator_run( workspace: "/tmp/nvflare/jobs/workdir", gpu="0")
```

Client code: Lightning client API

Transform your script to FL with a few lines of code changes:

1. Import NVFlare lightning API
2. Patch your lightning trainer
3. (Optionally) validate the current global model
4. Train as usually

```
from nemo.core.config import hydra_runner
from nemo.utils import AppState, logging
from nemo.utils.exp_manager import exp_manager
from nemo.utils.model_utils import inject_model_parallel_rank

# (0): import nvflare lightning api
import nvflare.client.lightning as flare

mp.set_start_method("spawn", force=True)

...
# (1): flare patch
flare.patch(trainer)

while flare.is_running():

    # (2) evaluate the current global model to allow server-side model selection
    print("--- validate global model ---")
    trainer.validate(model)

    # (3) Perform local training starting with the received global model
    print("--- train new model ---")
    trainer.fit(model)
```

Demo

NVIDIA FLARE Website
<https://nvidia.github.io/NVFlare>



Agenda

- What is Federated Learning?
 - NVIDIA Key Technologies for Federated Learning
 - Real-world Use cases of Federated Learning
 - Getting Started with NVIDIA FLARE
 - **Research: Addressing Key challenges in Federated Learning**
 - Summary & Announcements
-
-
-
-
-
-
-

Medical Imaging

ADDRESSING DATA HETEROGENEITY

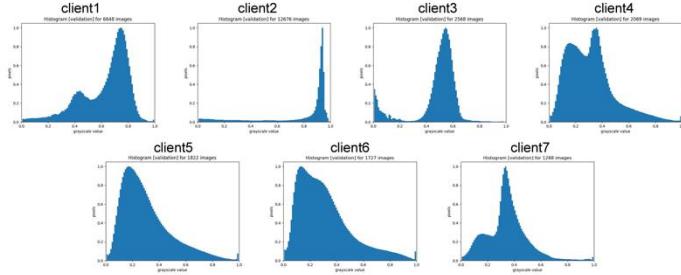


Fig. 4: Intensity distribution at different sites.

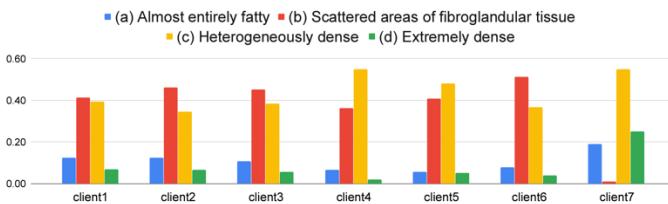
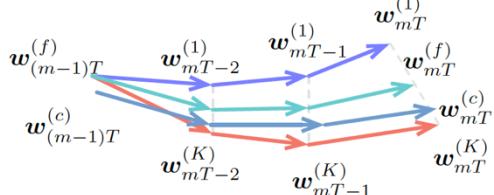


Fig. 3: Class distribution at different client sites as a fraction of their total data.

IID Settings:



Non-IID Settings:

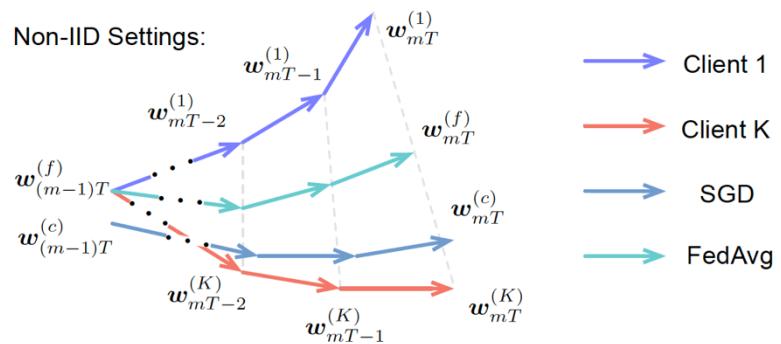


Figure sources: Intensity, label & quantity distribution of Mammography images across sites, [Roth et al. DCL 2020](#); Zhao, et. al, Federated Learning with Non-IID Data, 2018

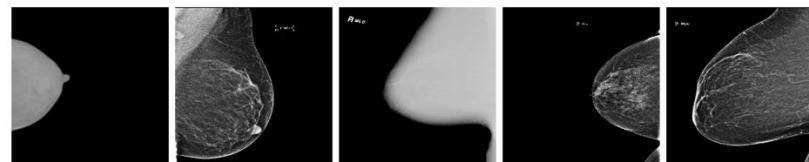


Fig. 1: Mammography data examples from different sites after resizing the original images to a resolution of 224×224 . No special normalization was applied in order to keep the scanners' original intensity distribution that can be observed in 4.

Institution	Train	# Val.	# Test
client1	22933	3366	6534
client2	8365	1216	2568
client3	44115	6336	12676
client4	7219	1030	2069
client5	6023	983	1822
client6	6874	853	1727
client7	4021	664	1288

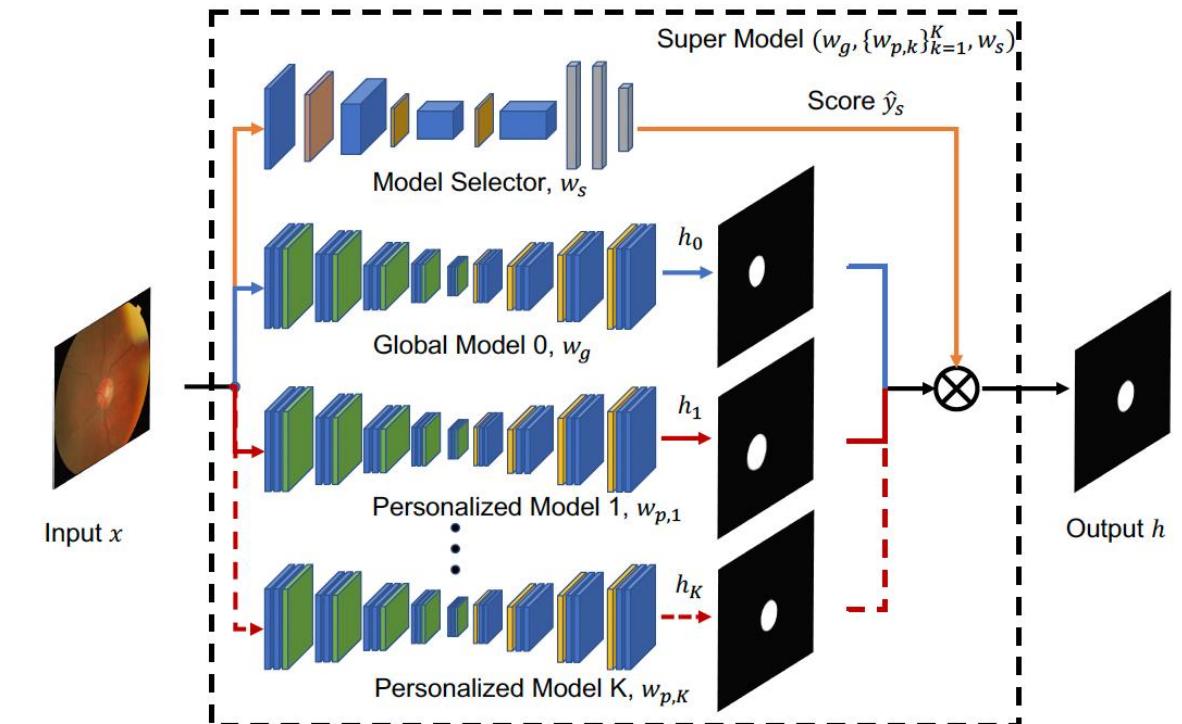
Standard FedAvg is sub-optimal if the clients have heterogeneous (non-i.i.d.) distributions

- **Personalization** can help close this gap.
- Carefully tuning hyperparameter plays a critical role to achieve optimal performance (**Auto-ML/FL**).

FEDERATED SUPER MODEL (FedSM)

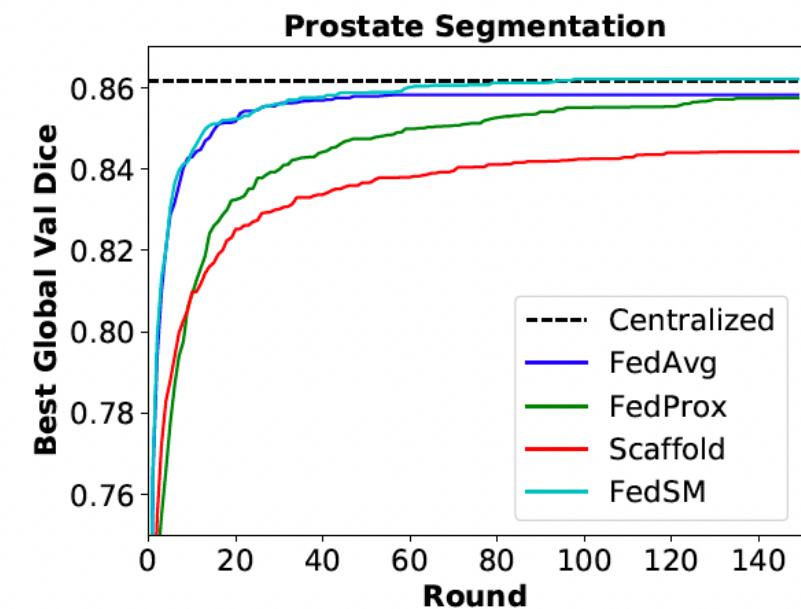
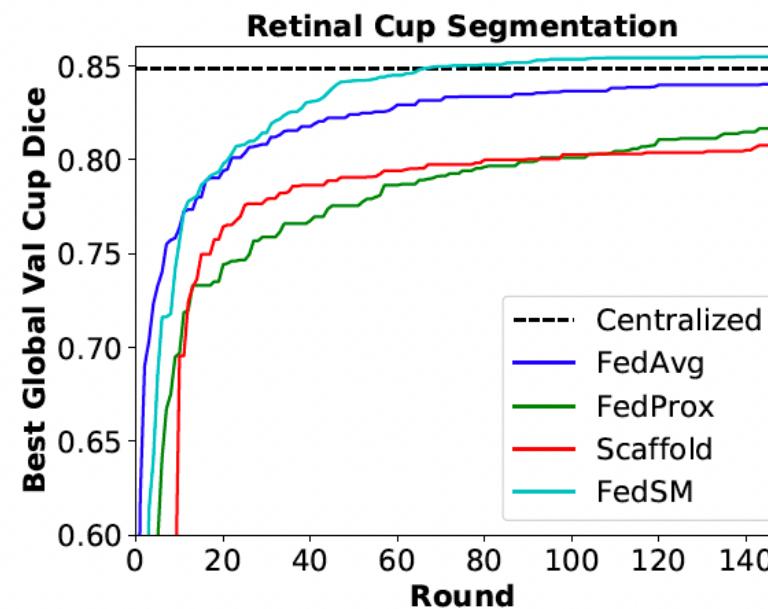
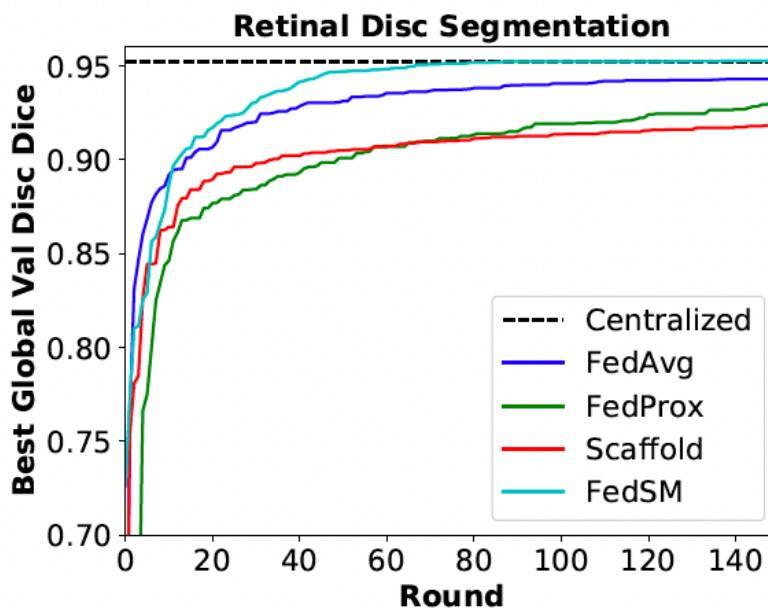
Closing the gap between FL and centralized training CVPR 2022

- The goal of FL is to collaboratively train **one global model** that generalizes well on all clients' joint data distribution – which is hard.
- Instead of finding one global model that tries to fits all clients' data distribution, we propose:
 - **Personalized models** to fit different data distributions well
 - A **soft-pull mechanism** to harmonize the global and local information
 - A **model selector** to decide the closest model/data distribution for any unseen test data from **one super model**.
- If input is similar to a site N, use model N otherwise, use global model.



FEDSM: RESULTS

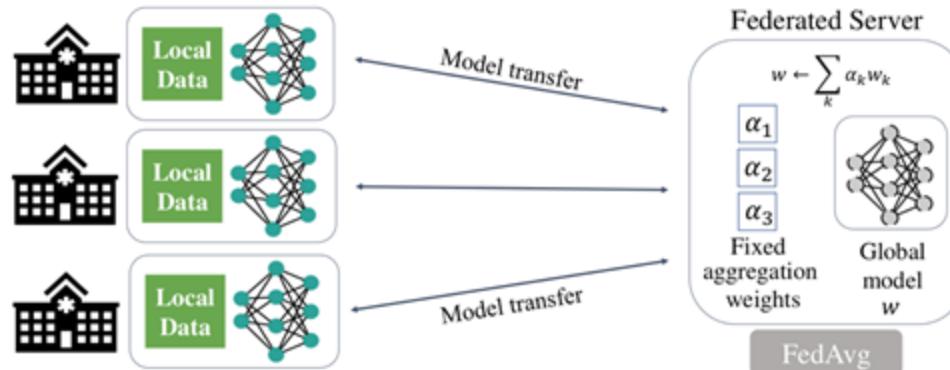
Closing the validation gap to centralized training



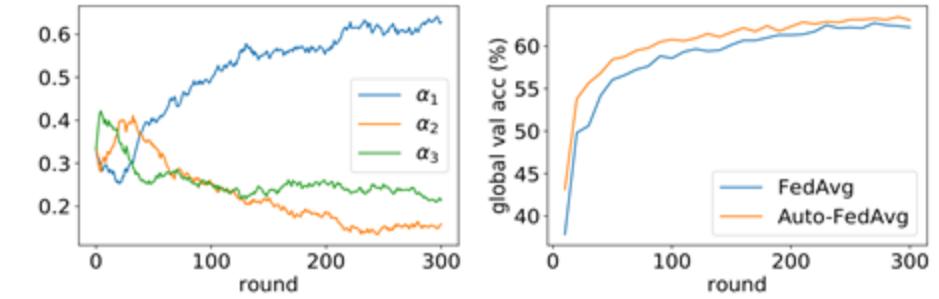
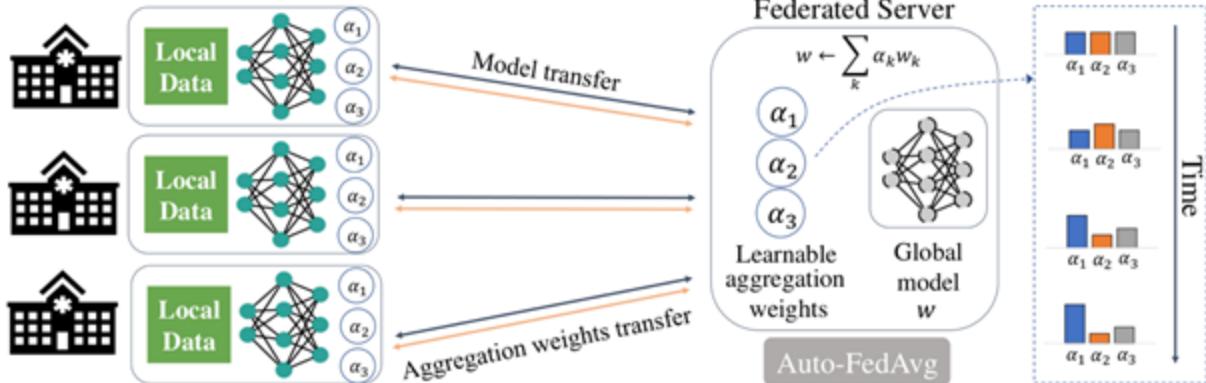
AUTO-FEDAVG: LEARNABLE FEDERATED AVERAGING FOR MULTI-INSTITUTIONAL MEDICAL IMAGE SEGMENTATION

Yingda Xia et al.

FedAvg:

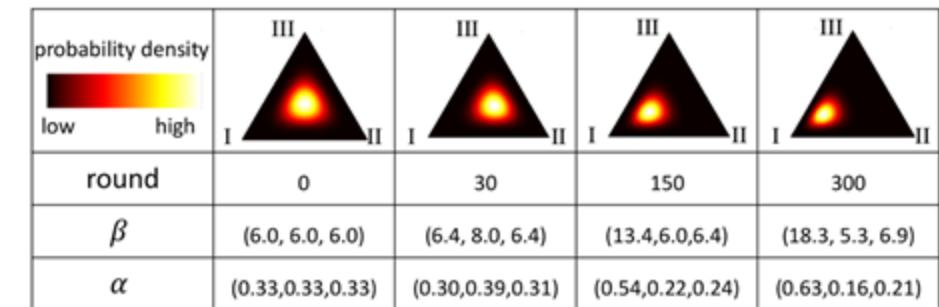


Auto-FedAvg:



(a) The learning curve of α .

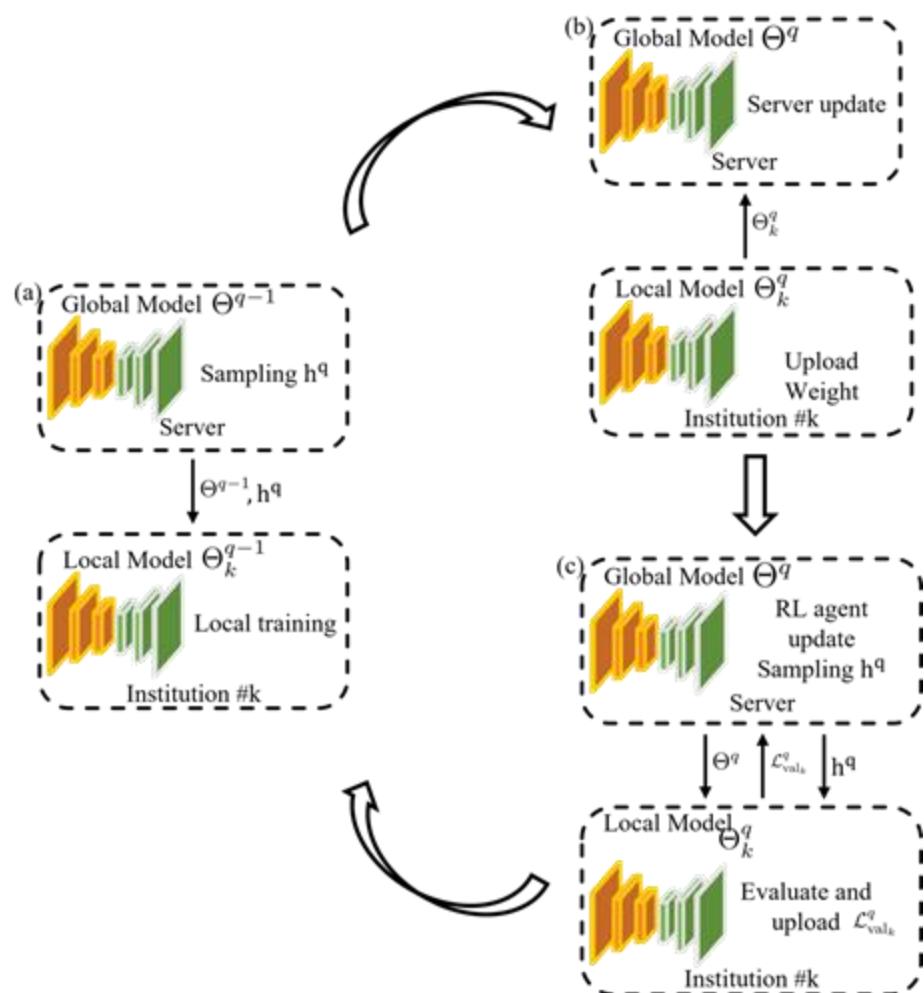
(b) Validation accuracy growth.



(c) Visualizations of Dirichlet distribution.

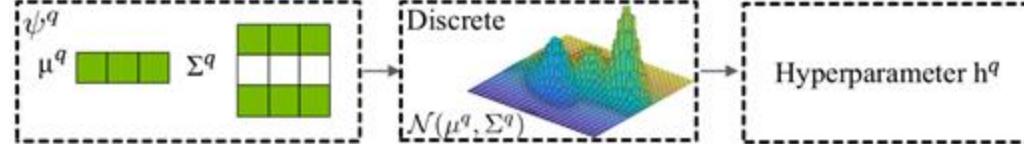
We show benefits for COVID lesion and pancreas segmentation in FL.

AUTO-FEDRL: OPTIMIZE ALL FL HYPERPARAMETERS

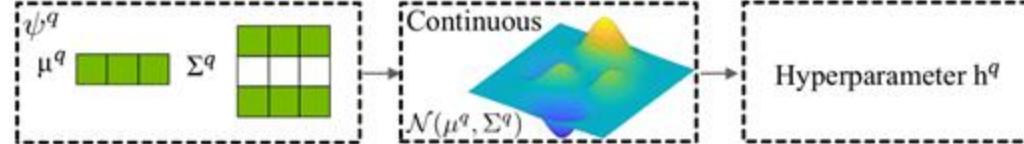


Search strategies:

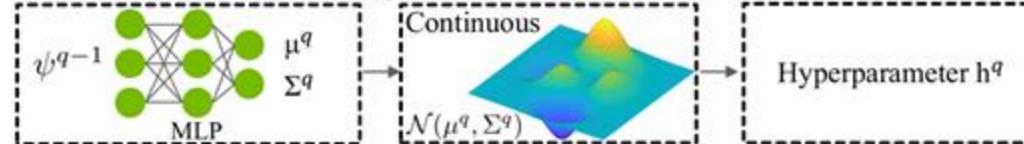
(a) Discrete search



(b) Continuous search



(c) Continuous search with Deep RL



The proposed FL method consists of the following three steps:

- a) The server broadcasts the current **global model** and **hyperparameters** to all **clients**. Clients perform **local training**
- b) Clients upload the trained local models, and the server updates the **global model**
- c) Clients evaluate the received the **new global model** and upload their **validation loss**. The server updates the **RL agent** and distributes the a **new set of hyperparameter** to all clients

AUTO-FEDRL: EXPERIMENTS

Results

CIFAR-10

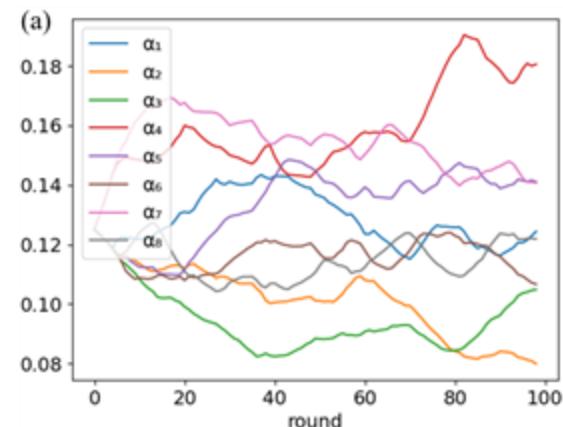
Method	Accuracy (%)
FedAvg [29]	88.43
FedProx [22]	89.45
Mostafa <i>et al.</i> [32]	89.86
Auto-FedAvg [48]	89.16
Auto-FedRL(dss)	90.70
Auto-FedRL(css)	90.85
Auto-FedRL(css MLP)	91.27
Centralized	92.56

Auto-FedRL search space:

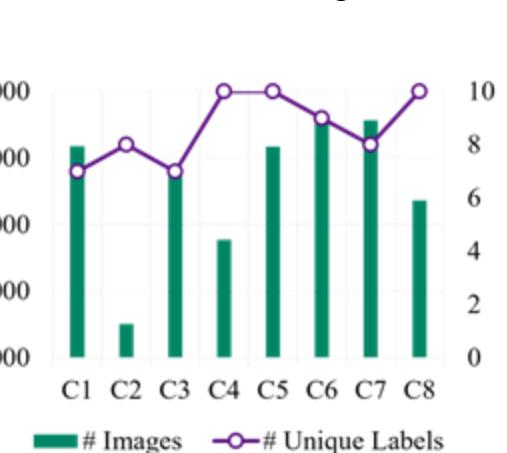
- Learning Rate (LR)
- # Local iterations (LI)
- Aggregation weights (AW)
- Server learning rate (SLR)

Search Space Type	Memory Usage	Running Time for Search
Discrete	42.8 GB	8.246 s
Continuous	3.00 GB	0.012 s
Continuous MLP	<u>3.13 GB</u>	<u>0.019 s</u>

Learning process of Auto-FedRL(css) in CIFAR-10



Non-IID setting



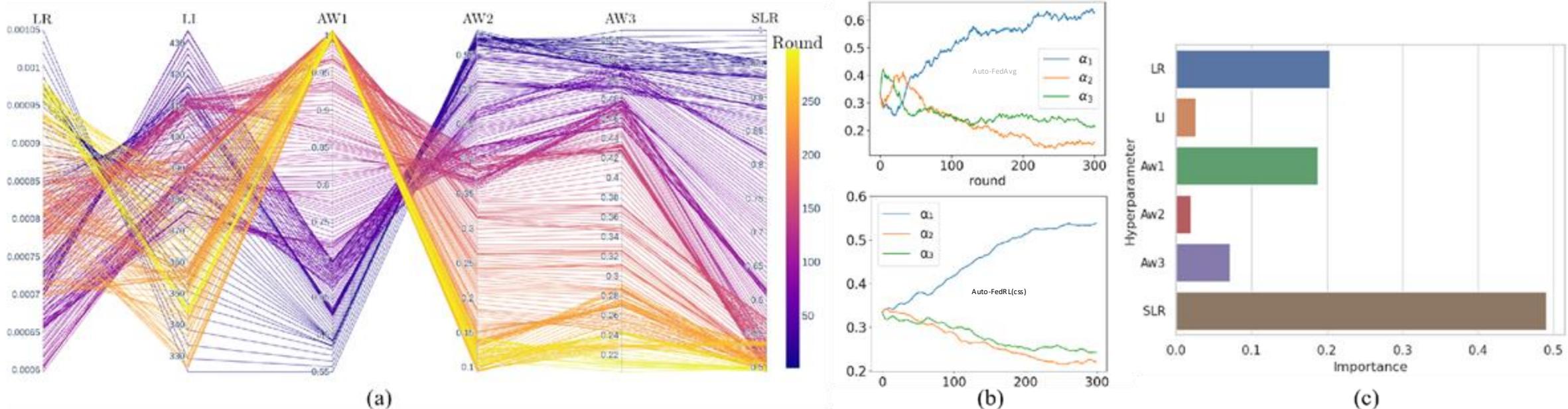
Analysis of the learning process of Auto-FedRL(css) in CIFAR-10. (a) the evolution of the distribution mean of aggregation weights during the training. (b) the statistic of different clients.

Mostafa et al.

AUTO-FEDRL: VISUALIZATIONS

Learned Hyperparameters

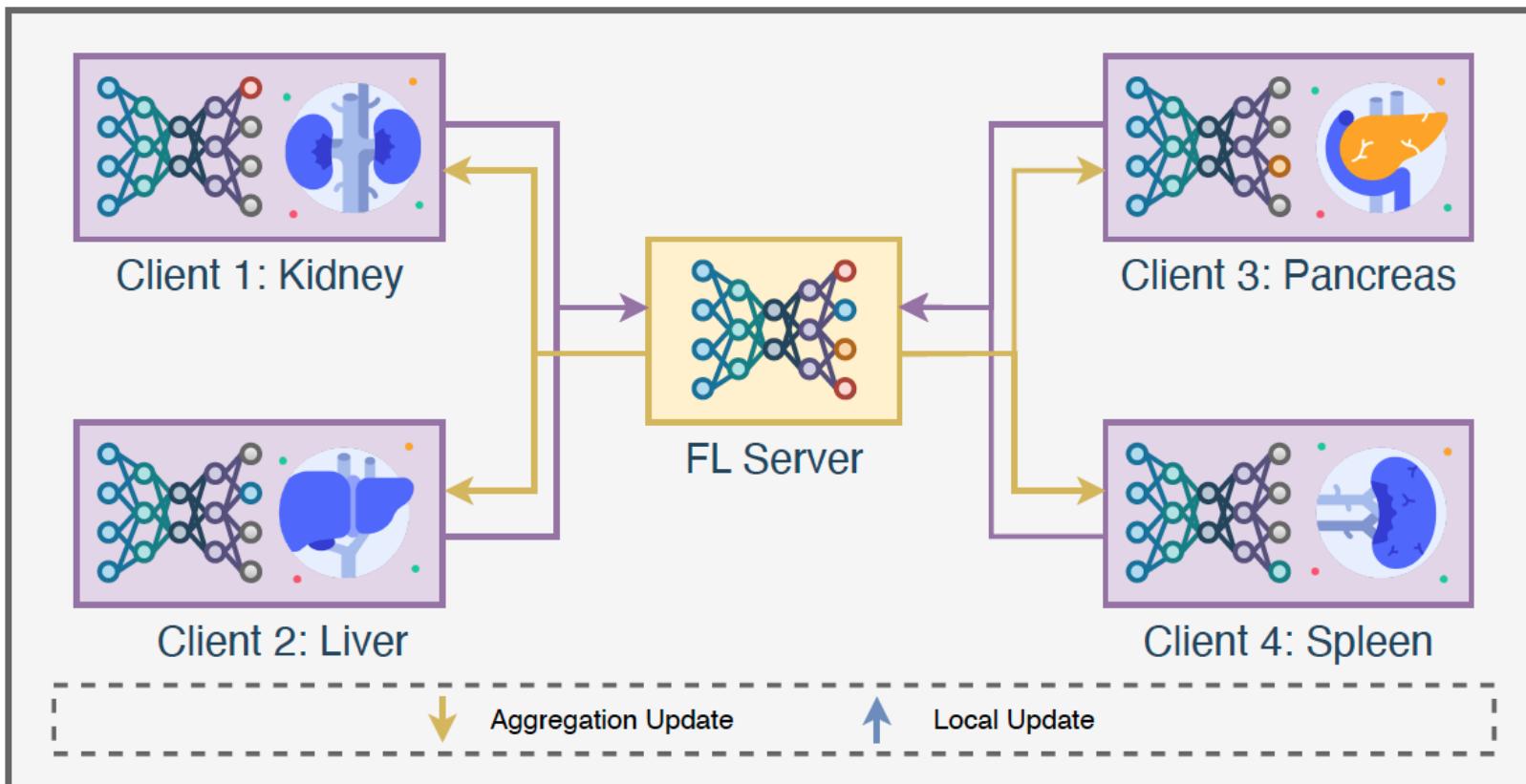
COVID-19 Lesion Segmentation



Analysis of the learning process of Auto-FedRL(css) in COVID-19 lesion segmentation. (a) The parallel plot of hyperparameter change during the training. LR, LI, AW, and SLR demotes learning rate, local iterations, aggregation weights, and server learning rate, respectively. (b) The aggregation weights evolution of Auto-FedAvg in top row and Auto-FedRL(css) in bottom row. (c) The importance analysis of different hyperparameter.

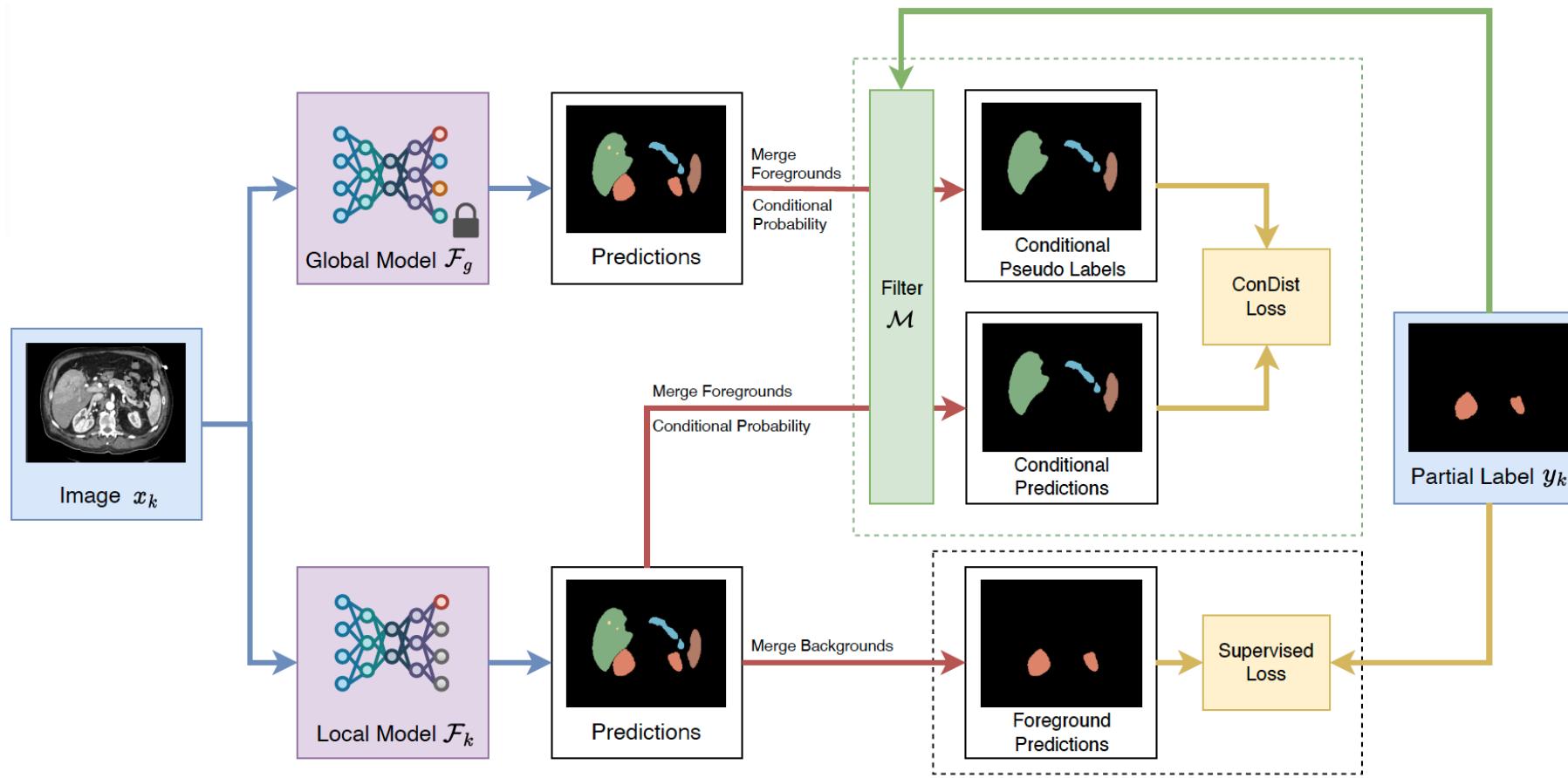
FEDERATED LEARNING WITH PARTIALLY LABELED DATA

Utilize diverse annotated data



FEDERATED LEARNING WITH PARTIALLY LABELED DATA

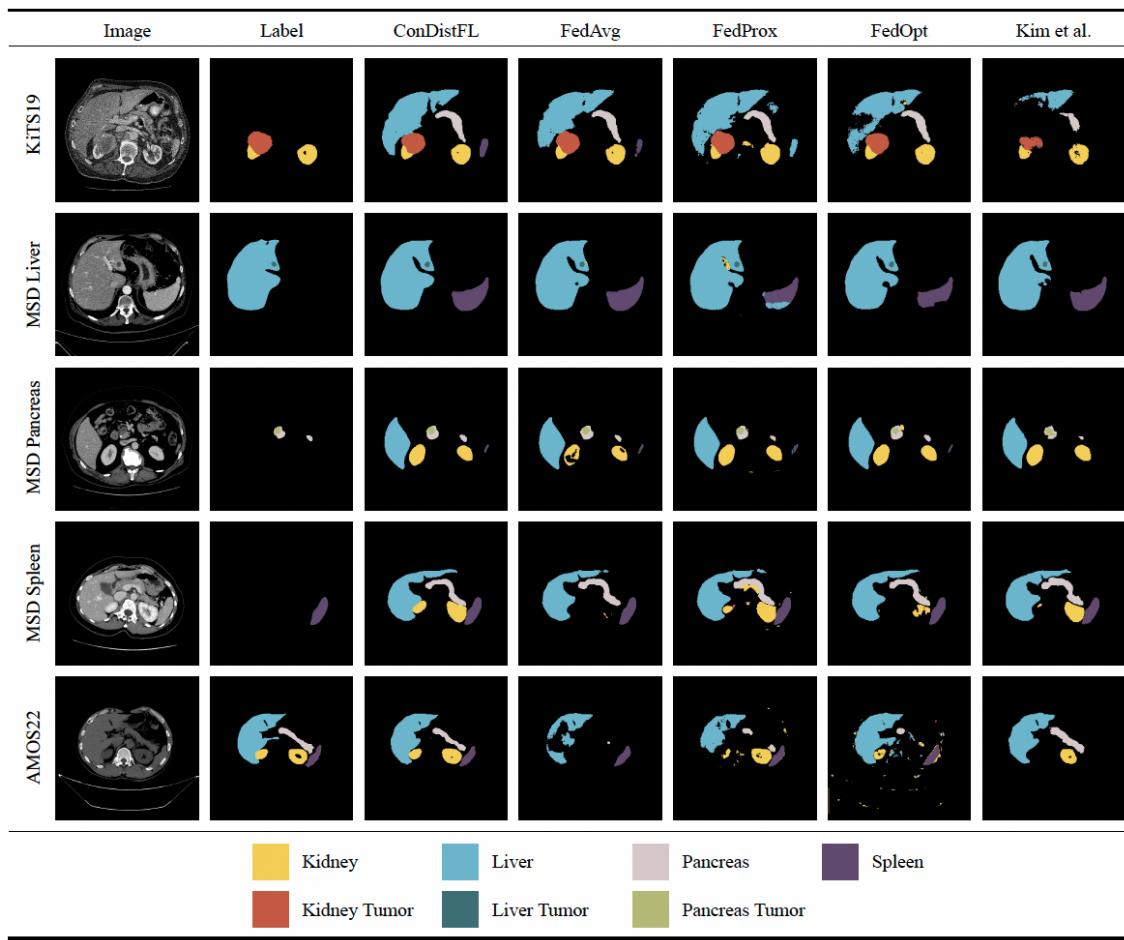
A Conditional Distillation Approach



FEDERATED LEARNING WITH PARTIALLY LABELED DATA

Results

Method	Average Dice ↑	Kidney		Liver		Pancreas		Spleen
		Organ	Tumor	Organ	Tumor	Organ	Tumor	Organ
Standalone	0.8167	0.9563	0.8116	0.9520	0.7265	0.7846	0.5126	0.9631
FedAvg	0.8091	0.9536	0.7766	0.9608	0.7241	0.7824	0.5041	0.9618
FedProx	0.7432	0.7736	0.7481	0.8973	0.5968	0.7893	0.5049	0.8923
FedOpt	0.7598	0.6767	0.7751	0.8915	0.6908	0.7841	0.5435	0.9568
Kim et al.	0.7330	0.9430	0.6734	0.9386	0.6439	0.7207	0.3778	0.8334
ConDistFL (ours)	0.8235	0.9547	0.8247	0.9613	0.7322	0.7975	0.5278	0.9664



Large Language Models

Parameter-efficient Fine-tuning (PEFT)

Adapt Foundational LLMs in FL

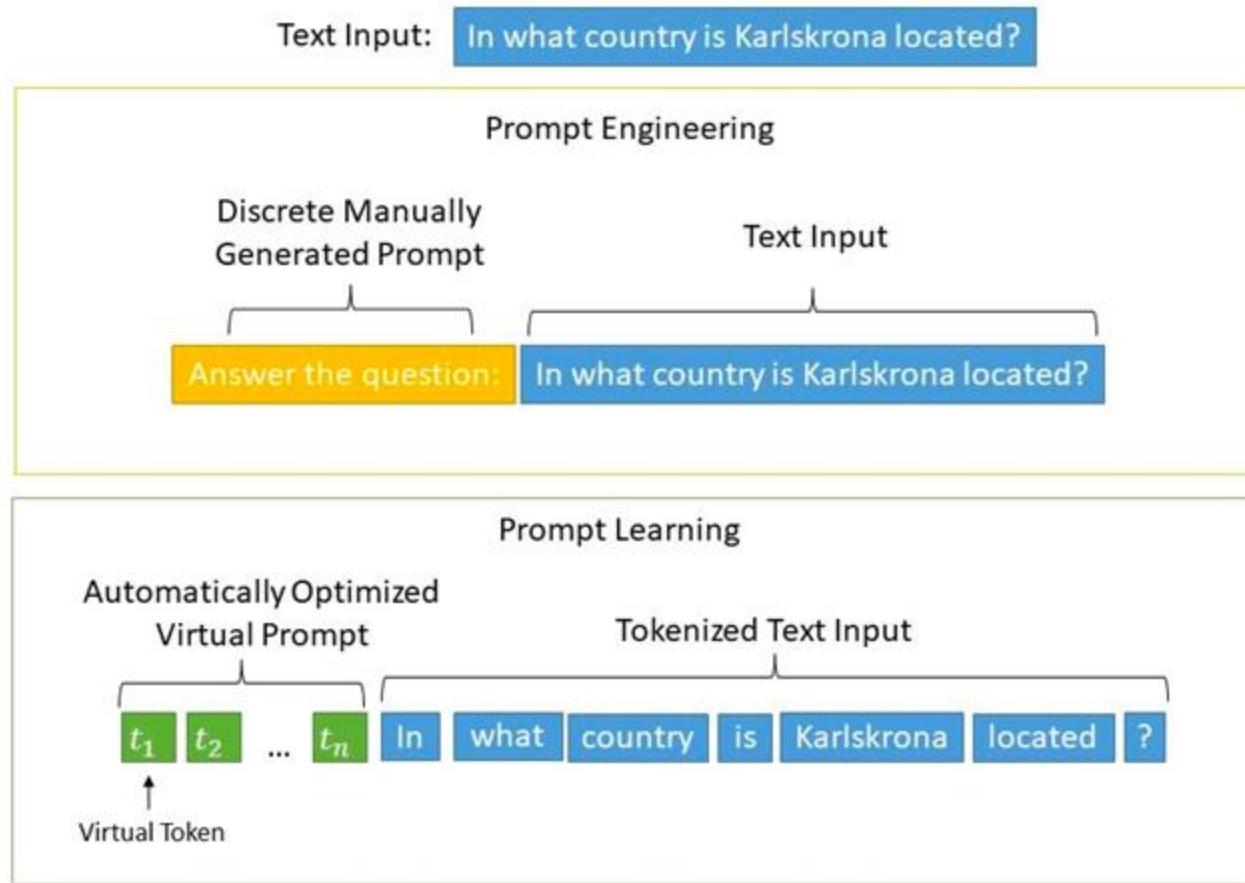
Parameter-efficient fine-tuning

Fine-tuning with a task-specific module

- **Most LLM layers fixed;** Only few dozen million params are being exchanged
- Tech: prompt-tuning/p-tuning/adapter/LoRA/others
- NVFlare example: sentiment analysis example with NeMo GPT model (**345M/5B/20B**)

Prompt Learning

Parameter-efficient adaptation of LLMs to downstream tasks



Tasks: brainstorming, classification, closed QA, generation, information extraction, open QA, summarization, etc.

P-Tuning for Sentiment Analysis

Downstream task example:

- Financial PhraseBank dataset ([Malo et al.](#)) for sentiment analysis.
- The Financial PhraseBank dataset contains the sentiments for financial news headlines from a retail investor's perspective.

Example prompts and predictions:

The products have a low salt and fat content . ***sentiment: neutral***

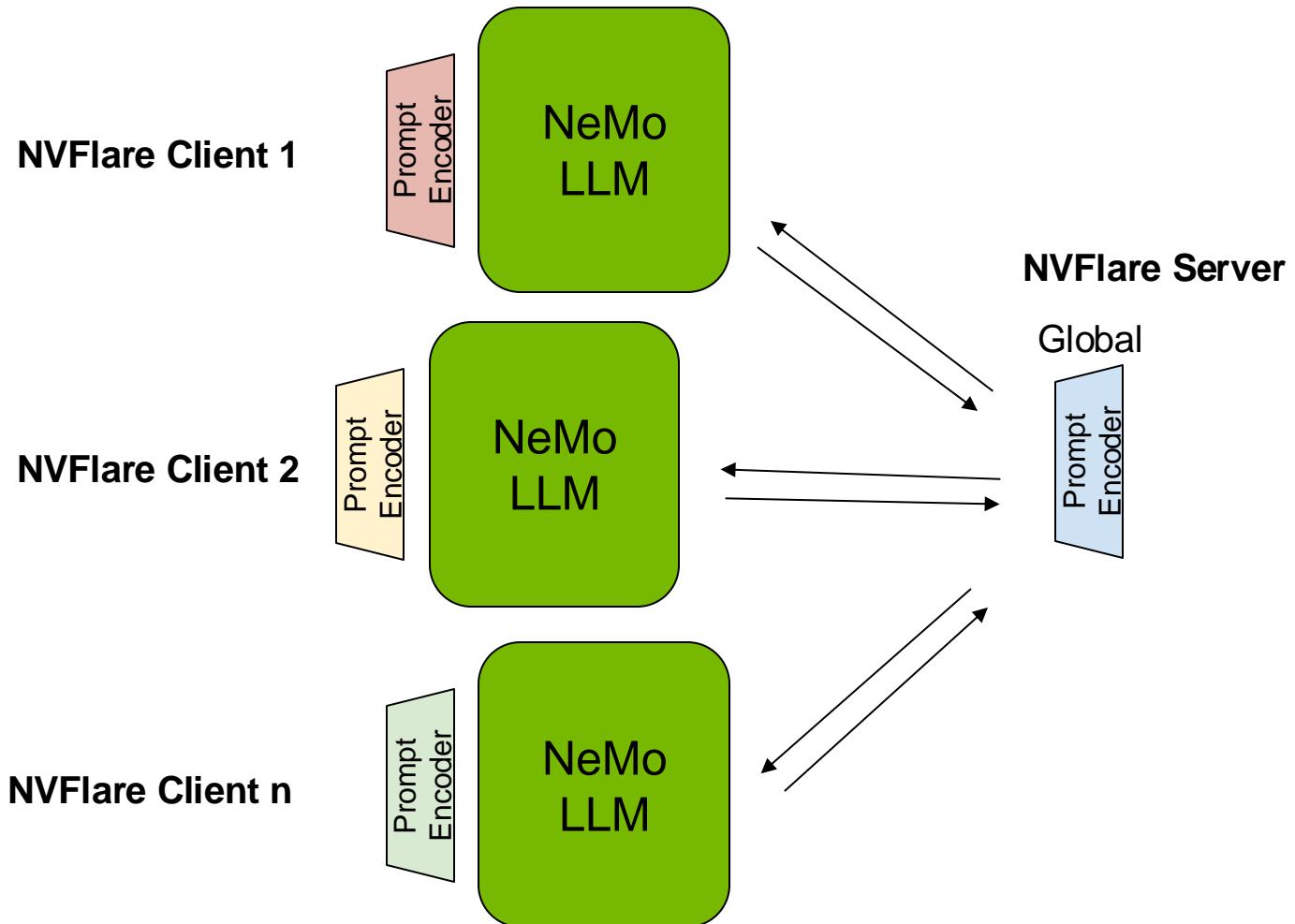
The agreement is valid for four years . ***sentiment: neutral***

Diluted EPS rose to EUR3 .68 from EUR0 .50 . ***sentiment: positive***

The company is well positioned in Brazil and Uruguay . ***sentiment: positive***

Profit before taxes decreased by 9 % to EUR 187.8 mn in the first nine months of 2008 , compared to EUR 207.1 mn a year earlier . ***sentiment: negative***

NVFlare for P-Tuning With NeMo



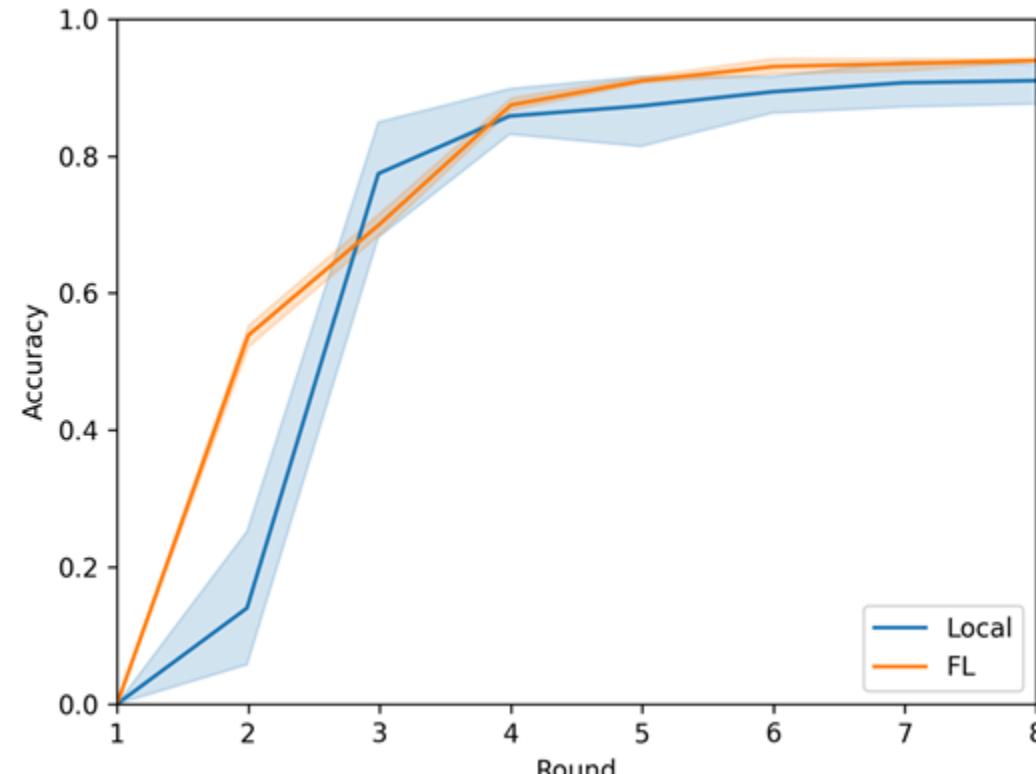
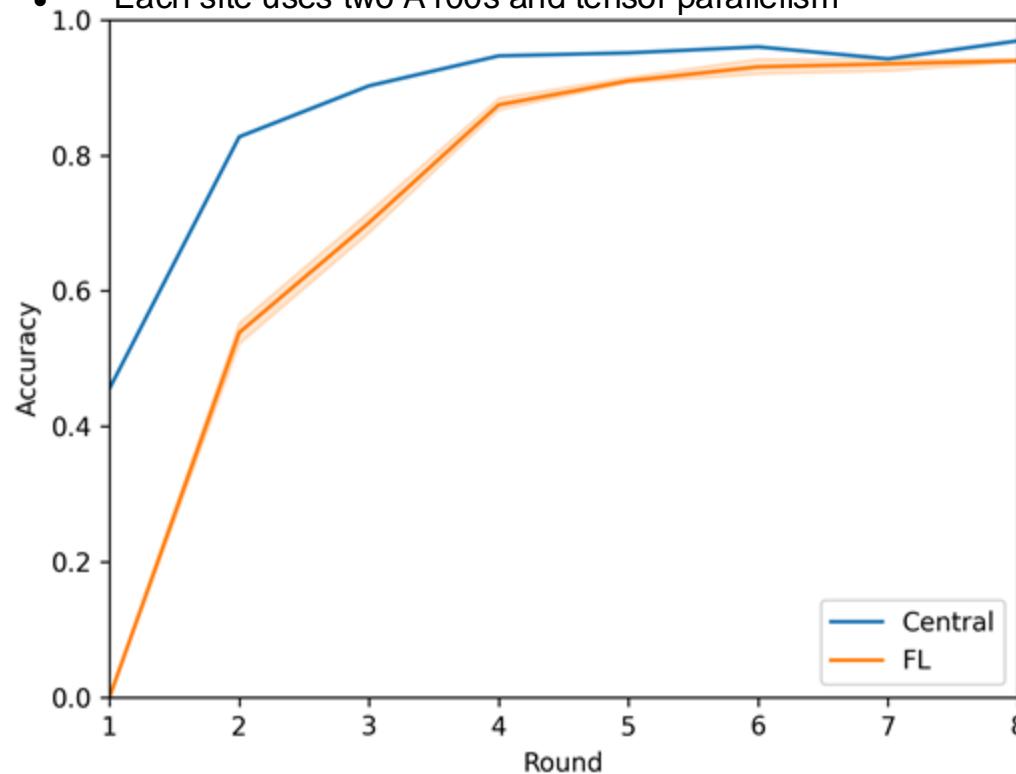
LLM parameters stay fixed; Prompt encoder parameters are trained/updated

P-Tuning for Sentiment Analysis

FL can achieve performance comparable to centralized training

Federated p-tuning experiment:

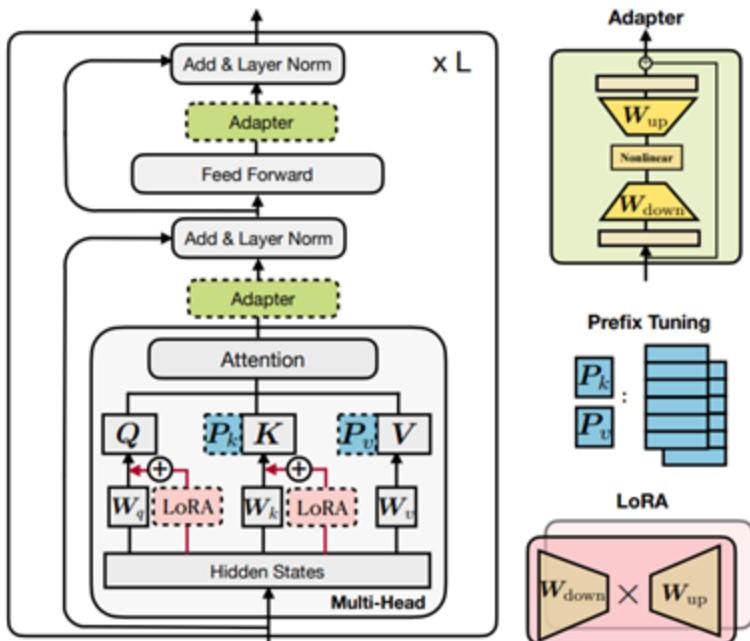
- Using **20B NeMo Megatron-GPT** model hosted on HuggingFace
- **50M** parameters are updated (0.25%)
- 1800 pairs of statement and sentiment
- 600 for each site; shared validation set for direct comparison
- Each site uses two A100s and tensor parallelism



Compare PEFT Methods With NeMo

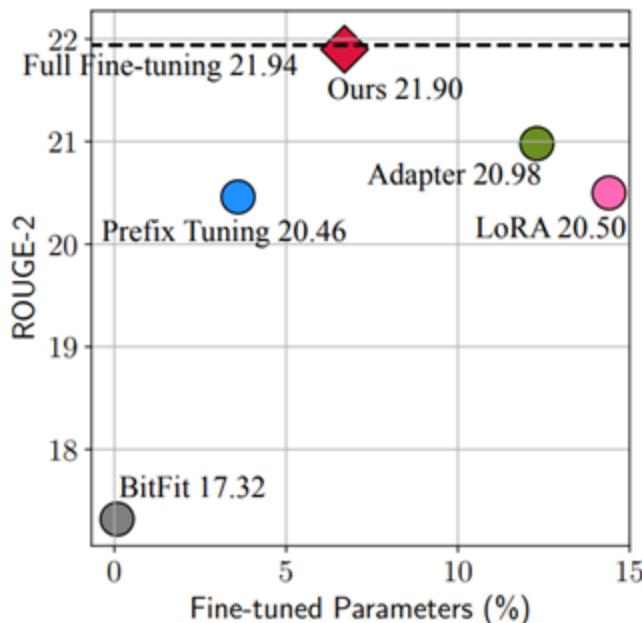
Lightning Client API + 1 line configuration change

Transformer and PEFT methods:



Source: <https://arxiv.org/abs/2110.04366>

Different PEFT methods on the XSum summarization task:

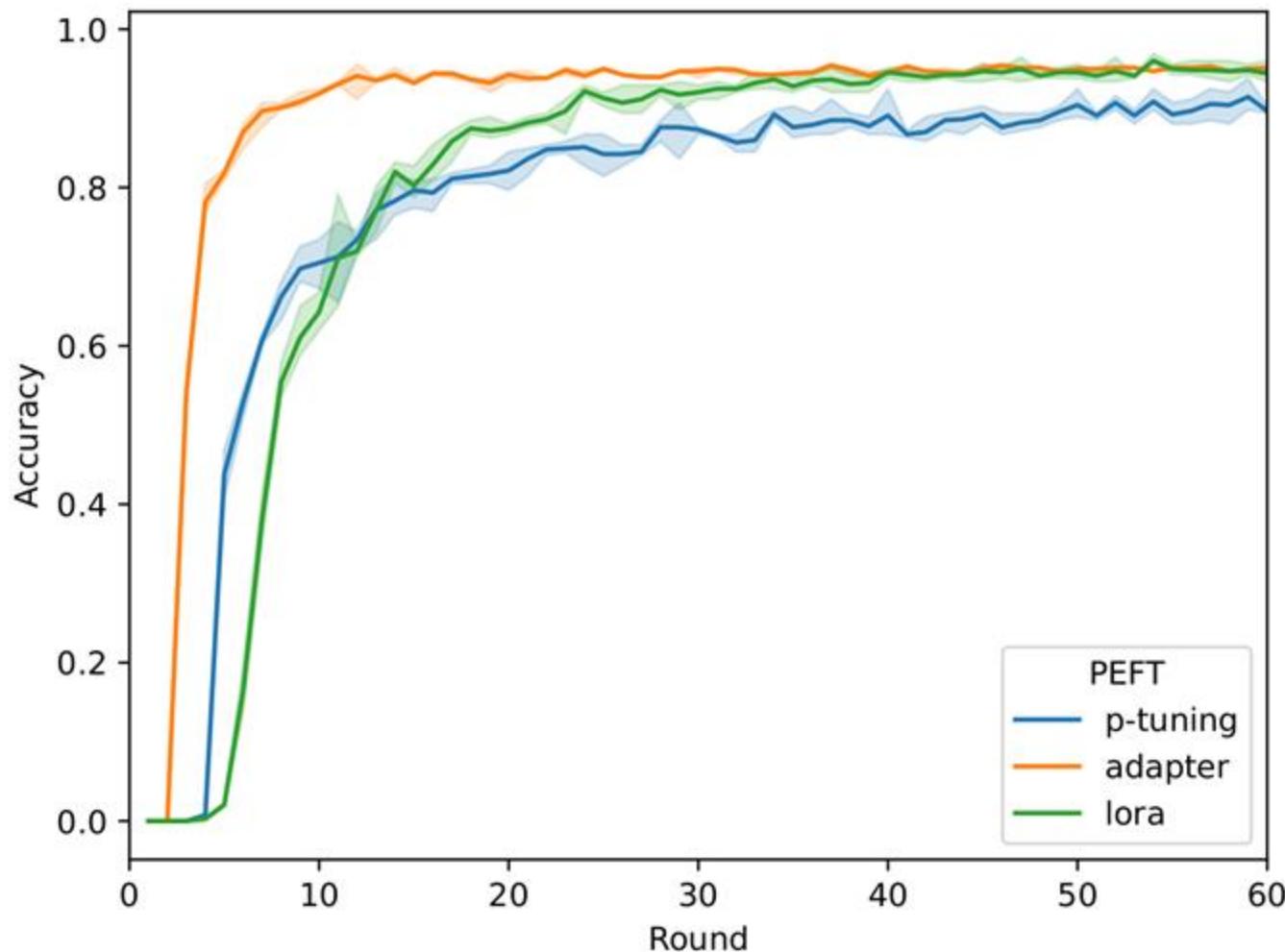


```
peft:  
  peft_scheme: "adapter" # can be either 'adapter', 'ia3', or 'ptuning'  
  restore_from_path: null  
  
  # Used for adapter peft training  
  adapter_tuning:  
    type: 'parallel_adapter' # this should be either 'parallel_adapter'  
    adapter_dim: 32  
    adapter_dropout: 0.0  
    norm_position: 'pre' # This can be set to 'pre', 'post' or null  
    column_init_method: 'xavier' # IGNORED if linear_adapter is used  
    row_init_method: 'zero' # IGNORED if linear_adapter is used, otherwise  
    norm_type: 'mixedfusedlayernorm' # IGNORED if layer_adapter is used  
    layer_selection: null # selects in which layers to add adapter  
    weight_tying: False  
    position_embedding_strategy: null # used only when weight_tying  
  
  lora_tuning:  
    adapter_dim: 32  
    adapter_dropout: 0.0  
    column_init_method: 'xavier' # IGNORED if linear_adapter is used  
    row_init_method: 'zero' # IGNORED if linear_adapter is used, otherwise  
    layer_selection: null # selects in which layers to add lora adapter  
    weight_tying: False  
    position_embedding_strategy: null # used only when weight_tying  
  
  # Used for p-tuning peft training  
  ptuning:  
    virtual_tokens: 10 # The number of virtual tokens the prompt encoder has  
    bottleneck_dim: 1024 # the size of the prompt encoder mlp bottleneck  
    embedding_dim: 1024 # the size of the prompt encoder embedding  
    init_std: 0.023  
  
  ia3_tuning:  
    layer_selection: null # selects in which layers to add ia3 adapter
```

[NeMo YAML configuration](#)

Compare PEFT Methods With NeMo

P-tuning vs. Adapter vs. LoRa



Tensor parallel with 2 GPUs per client

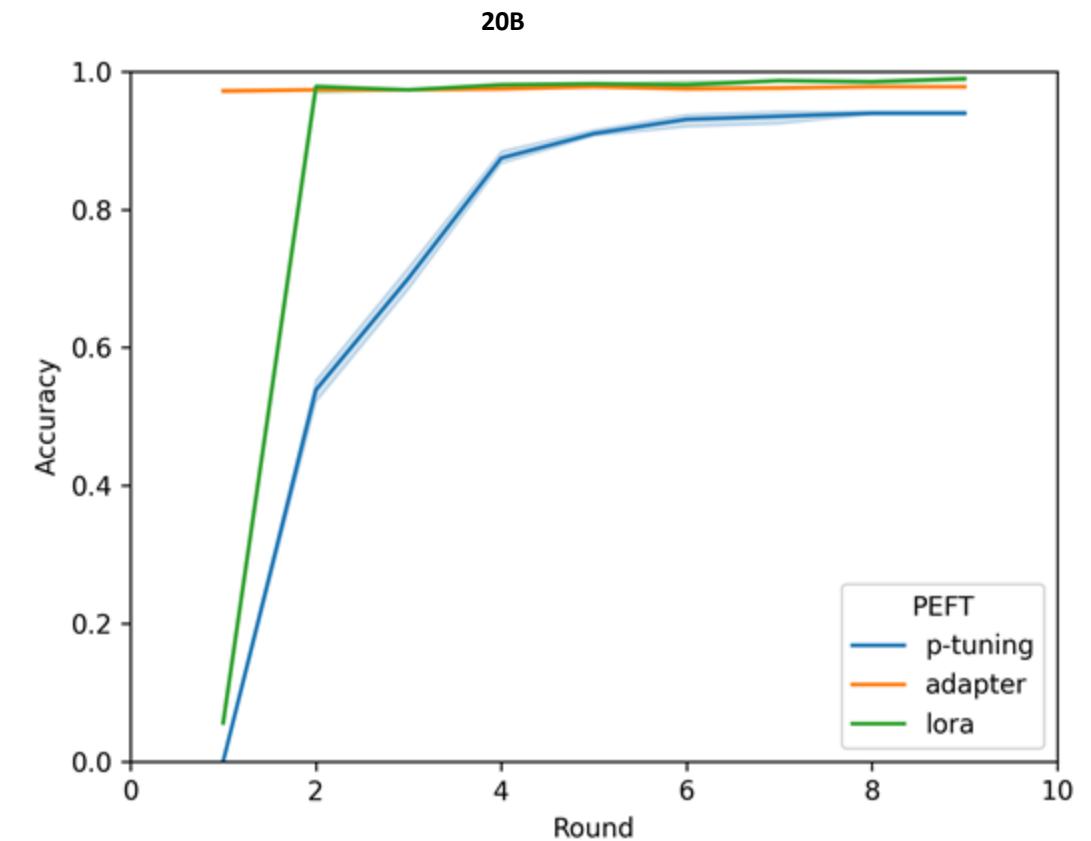
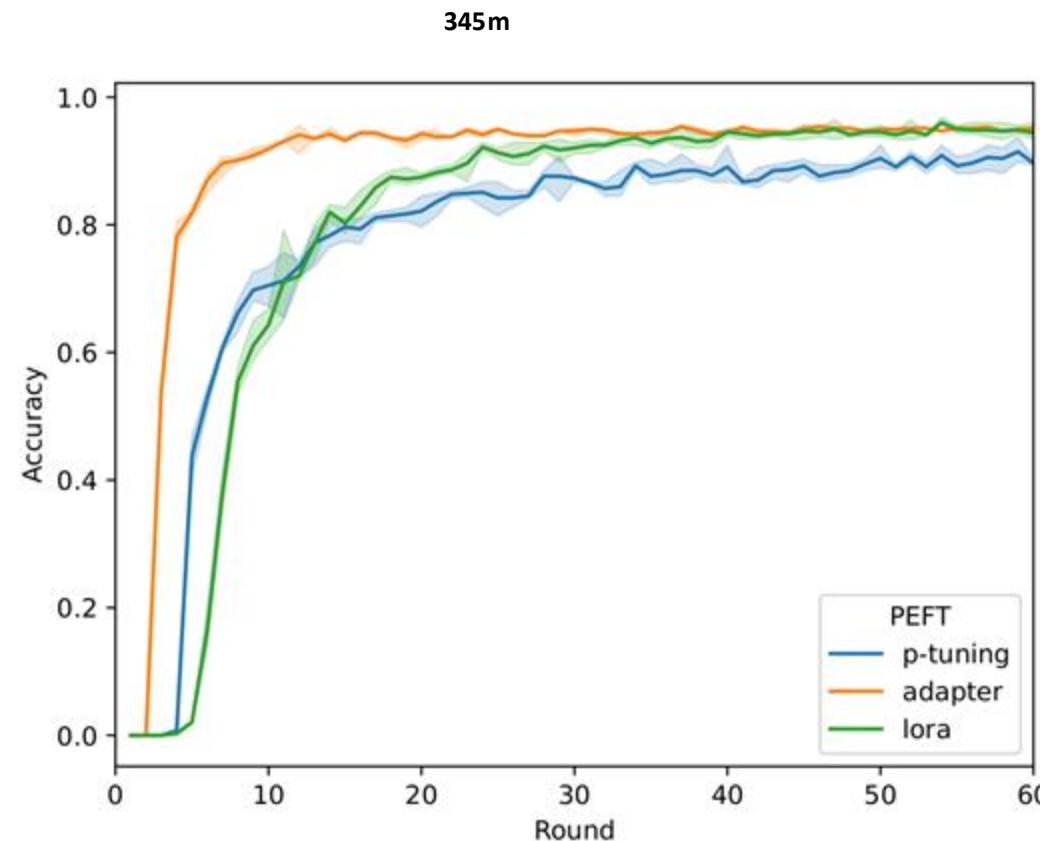
345M Param NeMo GPT Megatron model

PEFT Method	Execution time
P-tuning	4h 59m
Adapter	11h 25m
LoRA	7h 27m

Example [notebook](#)

Compare PEFT Methods With NeMo

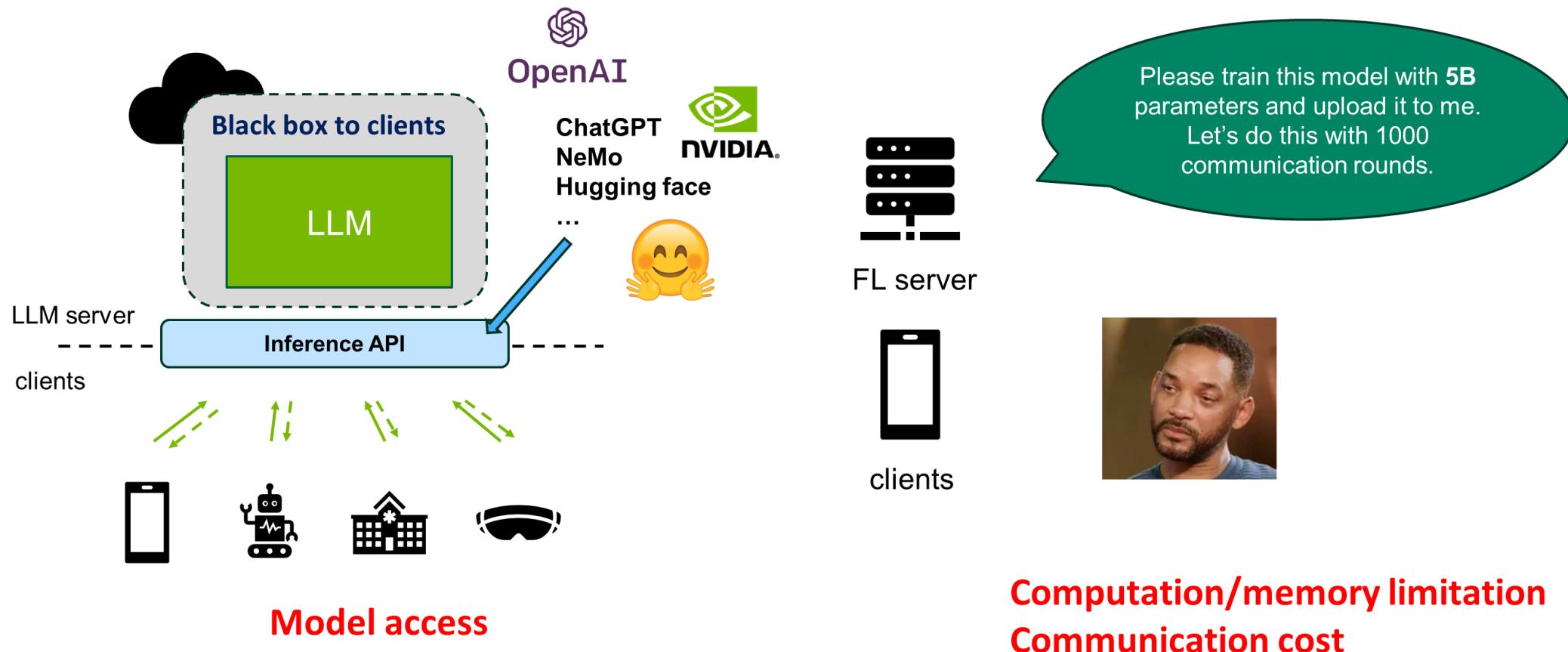
P-tuning vs. Adapter vs. LoRa



Example [notebook](#)

FedBPT: Efficient Federated Black-box Prompt Tuning for Large Language Models

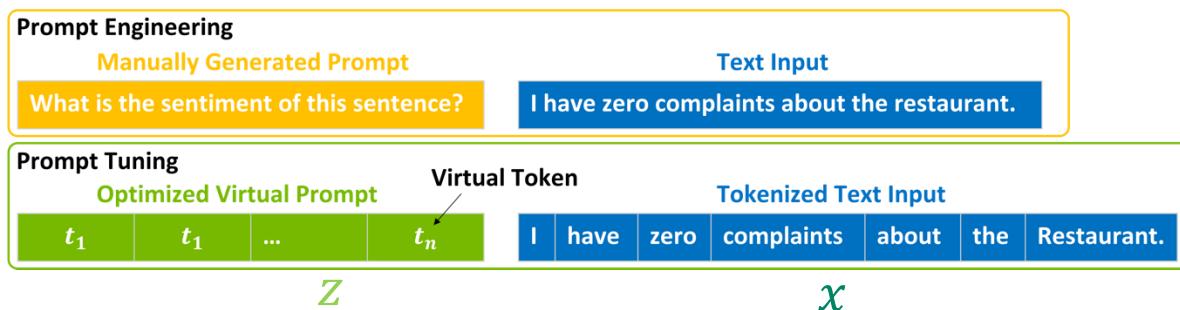
ICML 2024



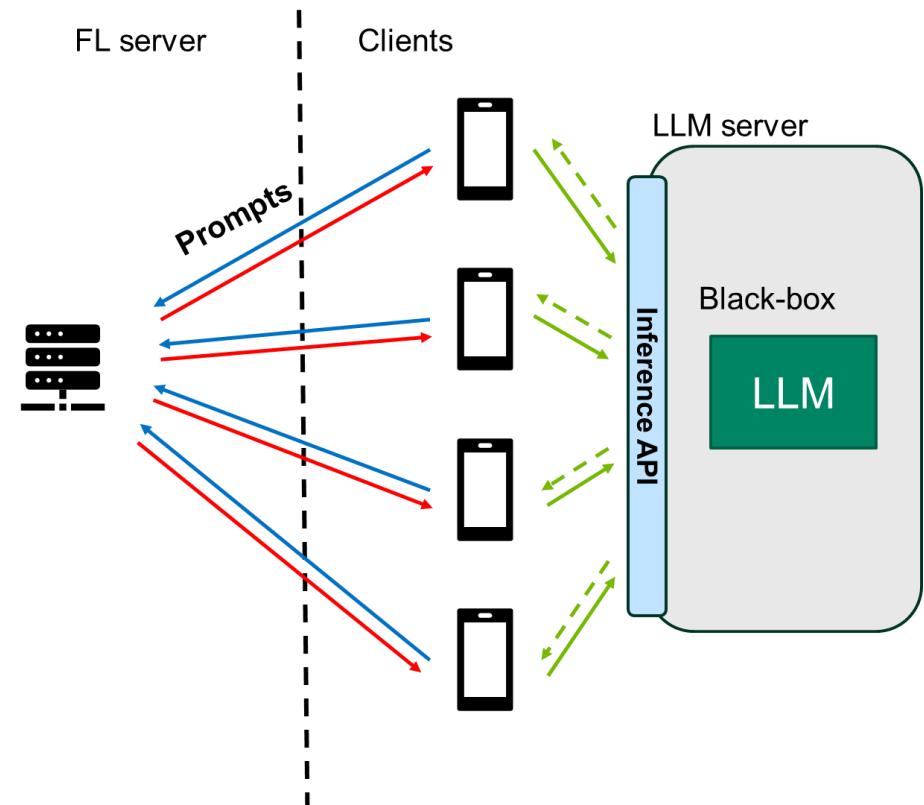
Unique Challenge in the Age of LLM

Efficient Federated Black-box Prompt Tuning

- The clients train prompts while treating the LLM as a black-box model.
 - The clients only conduct inference (with LLM server) without back-propagation.
 - The clients upload and download prompts (to and from FL server) rather than the whole model.
-
- Learning an optimized prompt encoding scheme in a federated fashion – prompt tuning
 - Utilizing gradient-free optimization methods



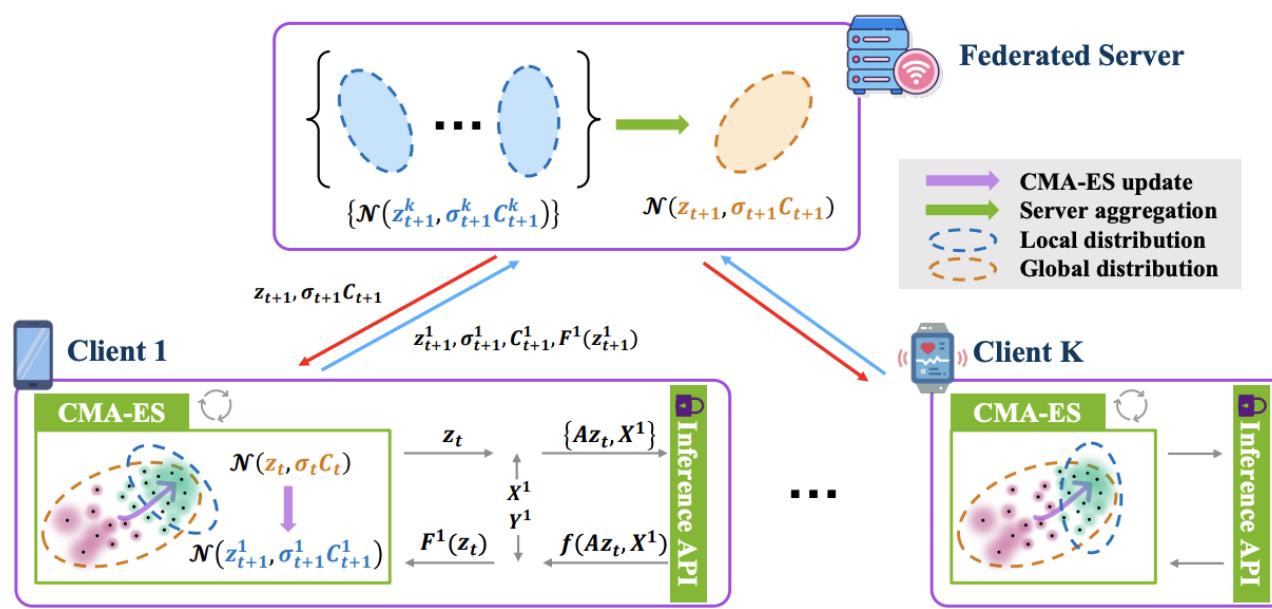
Dimension of z is much lower than Θ



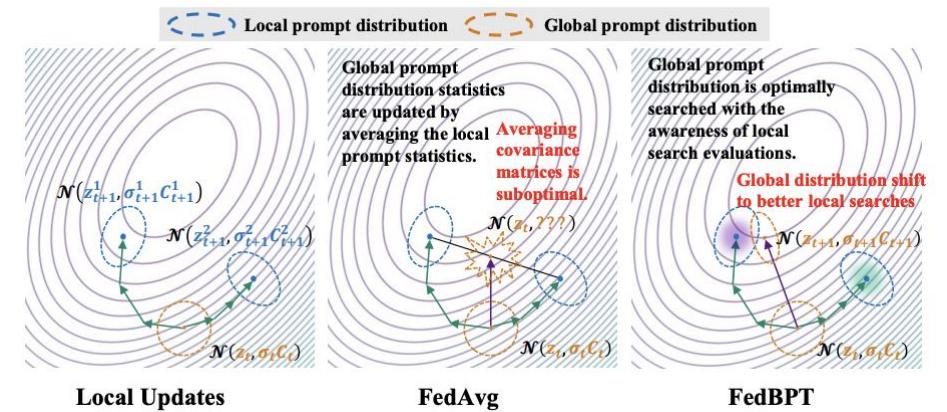
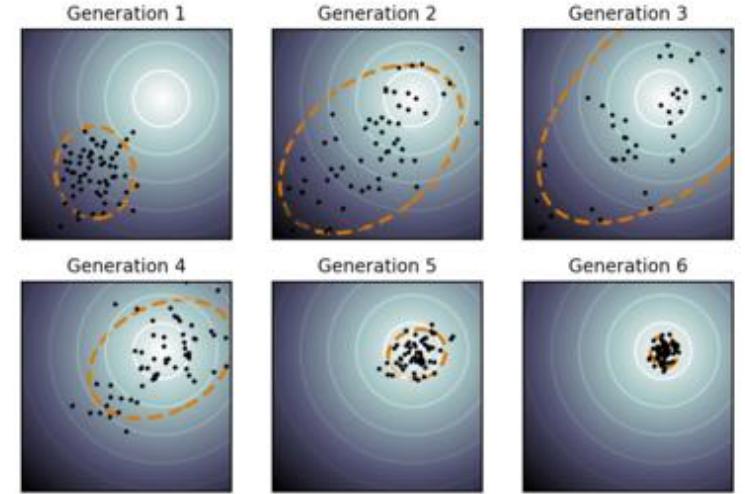
FedBPT: Gradient-free Optimization

Federated global estimation

- Use evolution strategy for optimization over the optimal prompt distribution
- Local optimization with inference only (“gradient-free”)
- Rather than simple averaging, global prompt distribution is searched with awareness of local prompt distributions
- Gradually converge to optimal solution



Inference is all you need!!!



Results

RoBERTa-350M

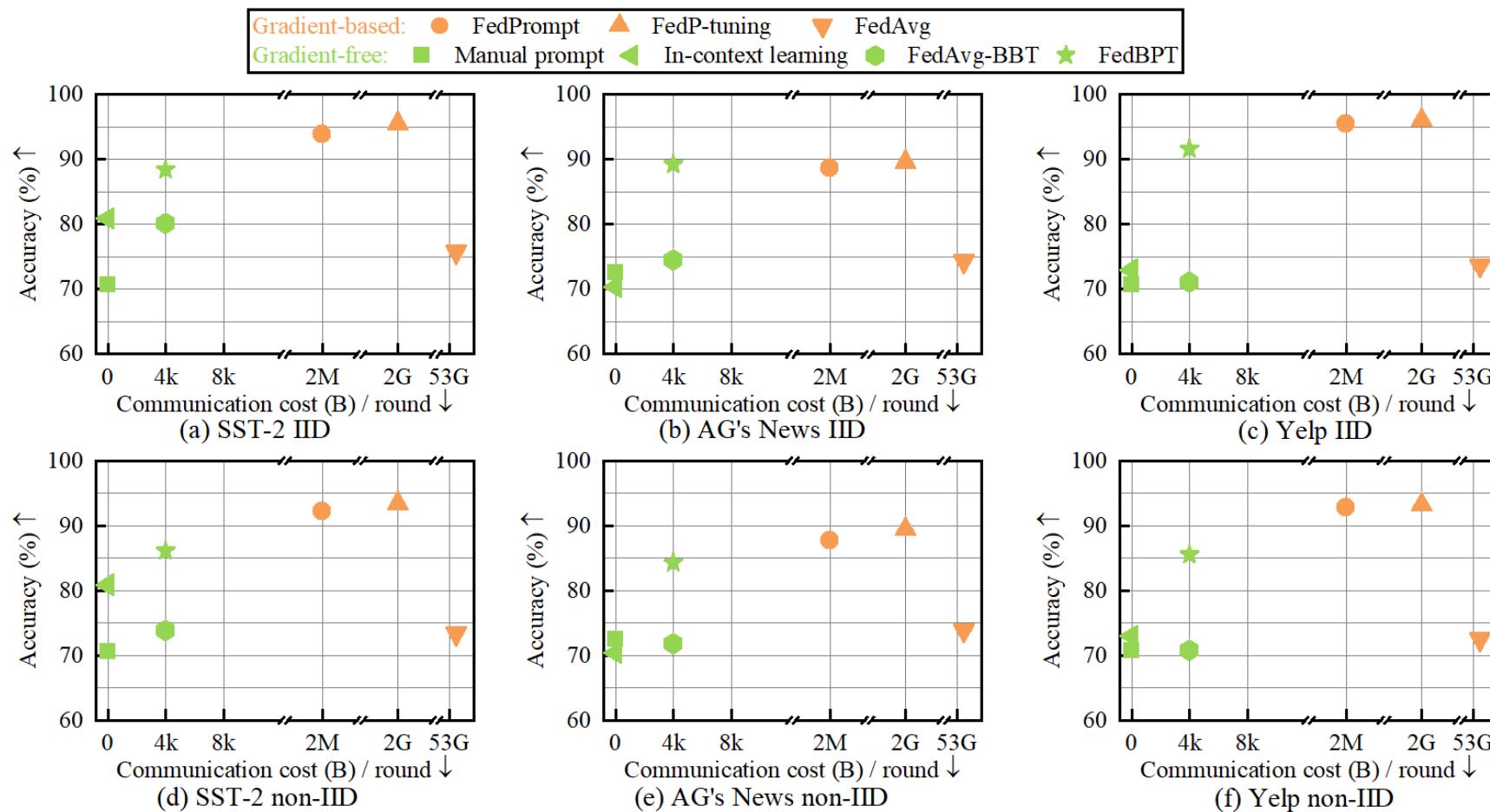
- Compare both gradient-based and gradient-free methods
- Achieves higher accuracies than gradient-free baselines under IID and non-IID settings for all the datasets.
- Decent performance among gradient-based methods
- Simply combining FedAvg and BBT cannot achieve decent performance

Method	Trainable Params.	SST-2		AG's NEWS		Yelp	
		Acc. (%) IID	Acc. (%) non-IID	Acc. (%) IID	Acc. (%) non-IID	Acc. (%) IID	Acc. (%) non-IID
<i>Gradient-based methods</i>							
FedPrompt	51K	90.25	85.55	87.72	85.62	91.44	91.47
FedP-tuning	15M	90.6	87.16	88.17	86.11	93.61	91.63
FedAvg	355M	84.7	82.4	77.43	76.54	88.25	88.03
FedLoRA	786K	84.6	84.53	77.85	75.9	88.52	88.2
<i>Gradient-free methods</i>							
Manual prompt	0	83.6		75.75		88.37	
In-Context Learning	0	79.7		76.96		89.65	
FedAvg-BBT	500	84.45	84.17	76.54	76.46	89.64	89.72
<i>FedBPT</i>	500	87.16	86.47	82.36	81.03	91.12	90.8

Results

Llama2-7B

- Improve the accuracy significantly compared with the gradient-free baselines and achieve comparable accuracies with the gradient-based methods in most settings
- Reduce the communication cost of one device in one round from nearly 2GB to 4KB (against P-tuning).

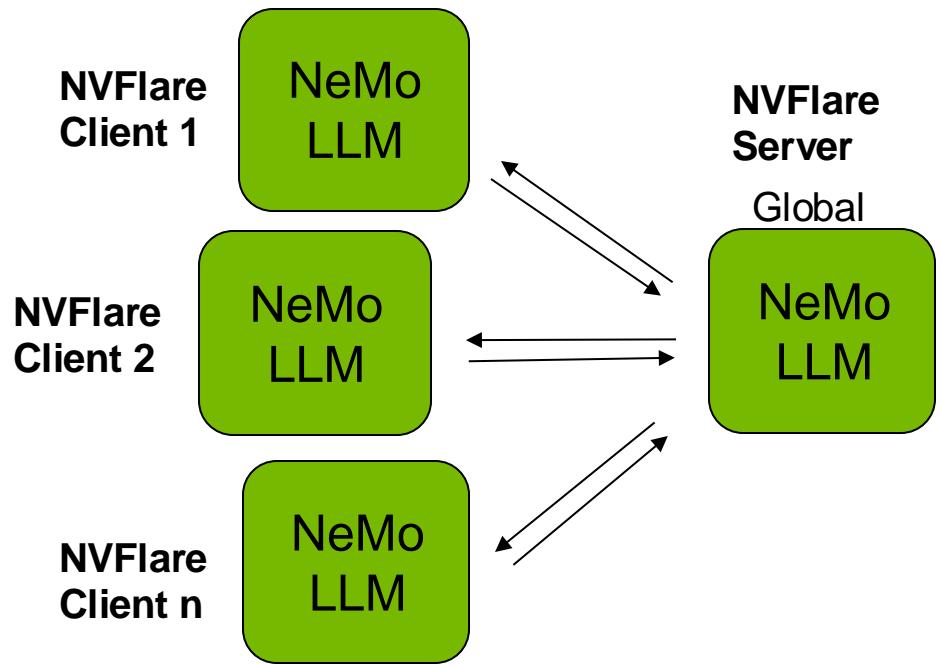
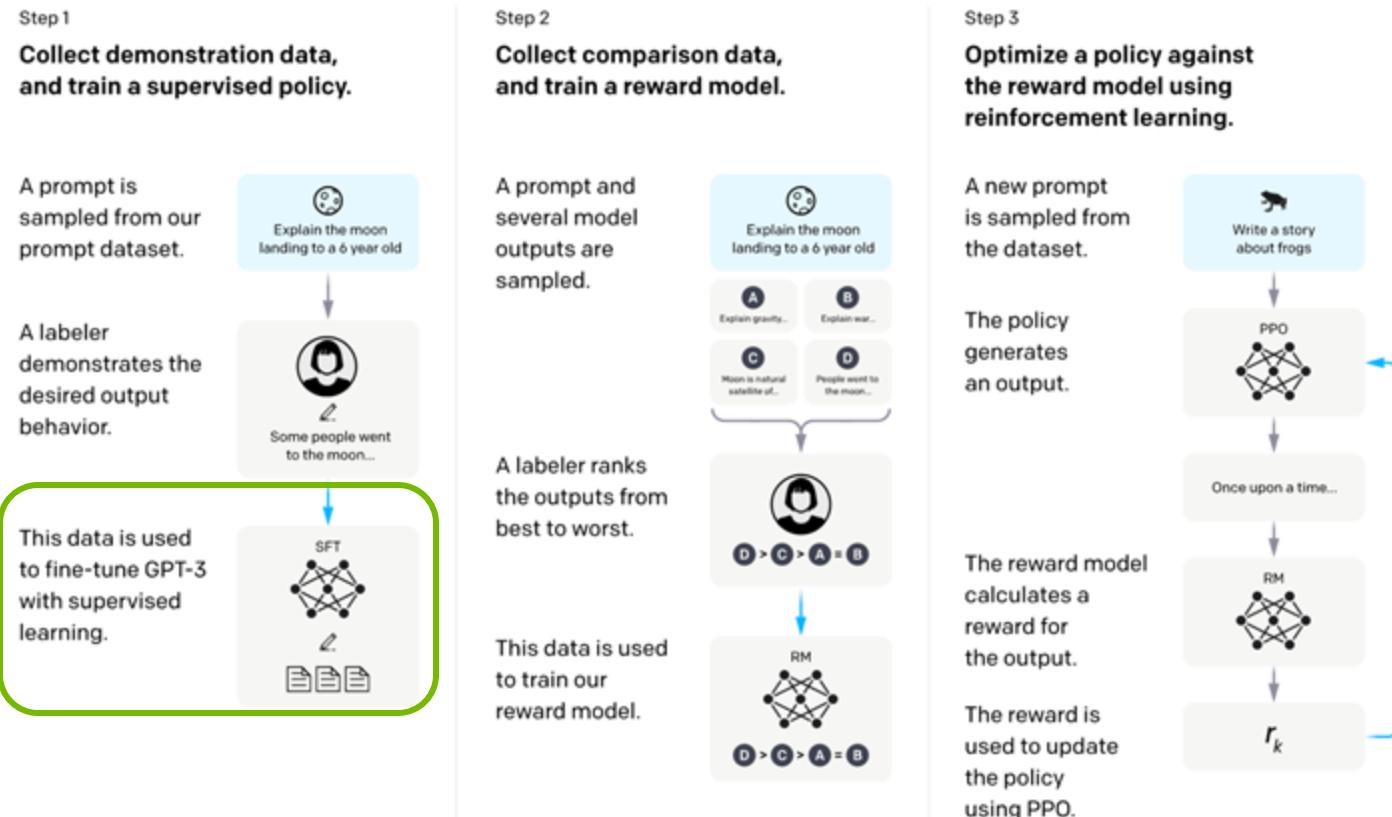


Supervised Fine-tuning

Supervised Fine-tuning (SFT)

Towards “instruction-following” LLM

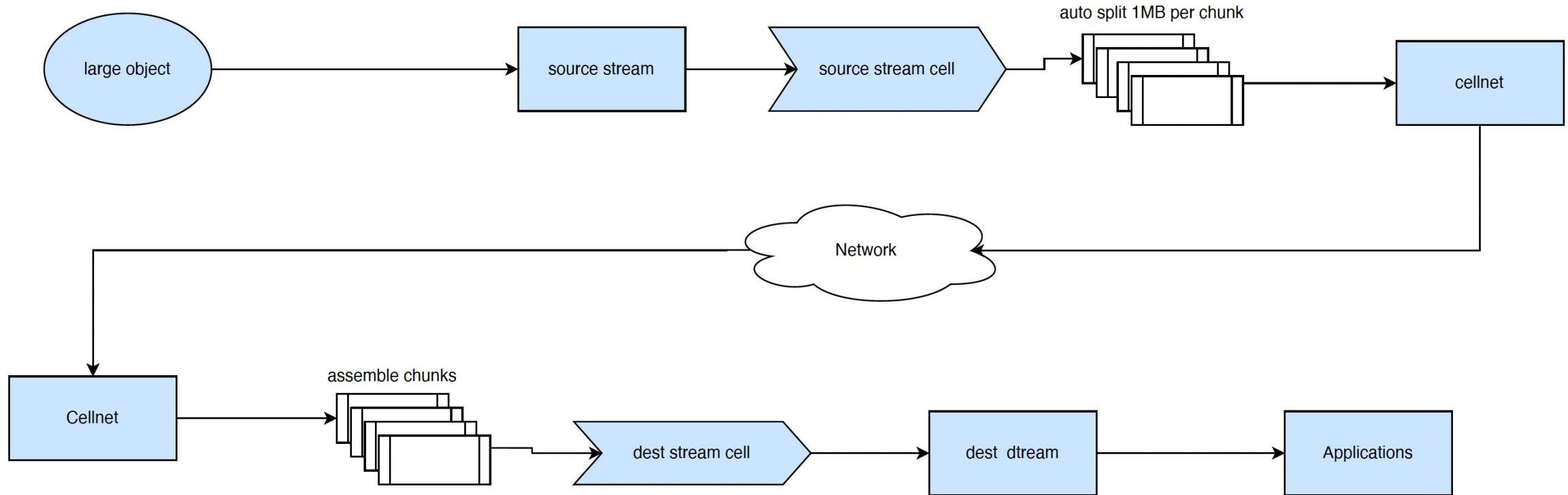
The first step of “Chat-GPT training scheme”. Unlike PEFT, SFT finetunes the entire network.



NVFlare Streaming

Support Large Model Transmission

- Model size of mainstream LLM can be huge: 7B -> 26 GB (beyond the 2 GB GRPC limit)
- In order to transmit LLMs in SFT, NVFlare supports **large object streaming**



SFT for Instruction Following

3 open datasets

We use three datasets:

- Alpaca
- databricks-dolly-15k
- OpenAssistant

Containing instruction-following data in different formats under different settings:

- oasst1 features a tree structure for full conversations
- other two are instruction(w/ or w/o context)-response pairs

Examples with NeMo and HuggingFace

https://github.com/NVIDIA/NVFlare/tree/main/integration/nemo/examples/supervised_fine_tuning
https://github.com/NVIDIA/NVFlare/tree/main/examples/advanced/llm_hf

SFT With FL

Achieving better performance

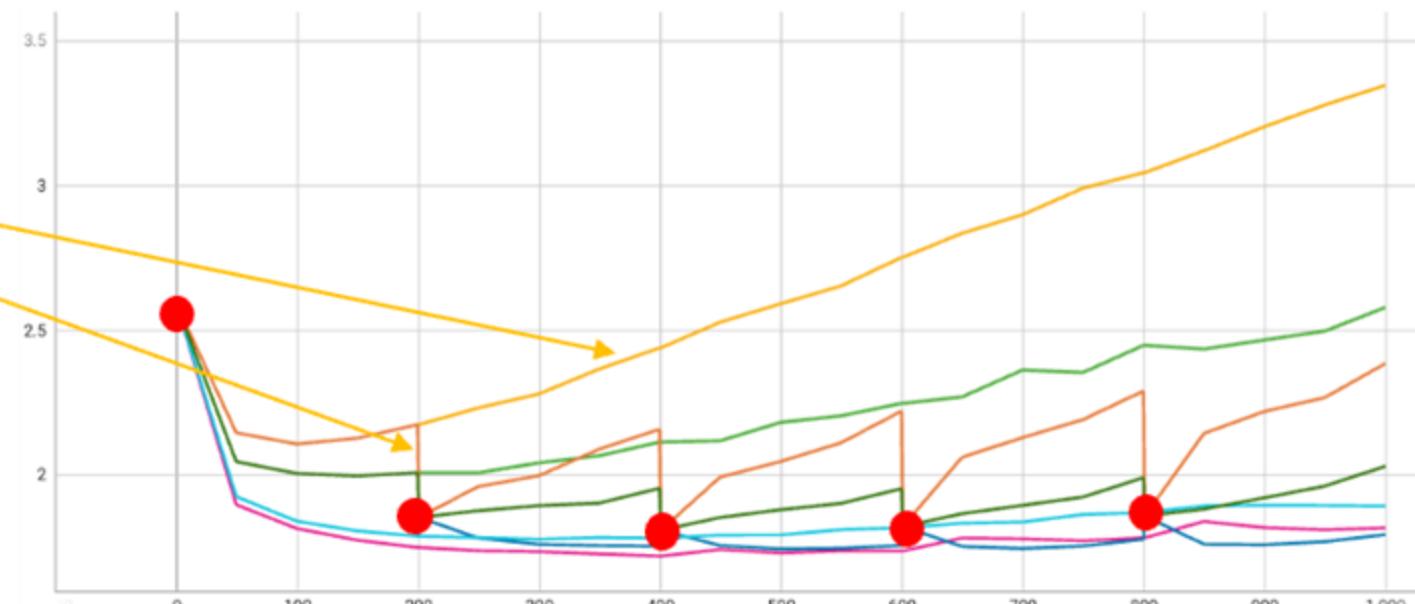
NeMo 1.3B model, SFT for 5 rounds

5 experiments in total: training on each client's own dataset, on combined dataset, and all three clients training together using the FedAvg algorithm implemented in NVFlare.

Validation loss curves:

- yellow: oasst1
- green: dolly
- blue: alpaca
- magenta: three datasets combined

Smooth curves for local training,
“Step curves” for FL - because of
global model sync and update



SFT Model Evaluation

LLM Performance

Non-trivial task compared with “fixed downstream tasks” where we usually have metrics like accuracy, F-1 scores, etc.

Common practice is to test the resulting LLMs on **benchmark tasks**, and/or human evaluations

We perform standard language modeling tasks under zero-shot setting, including HellaSwag(H), PIQA(P), and WinoGrande(W)

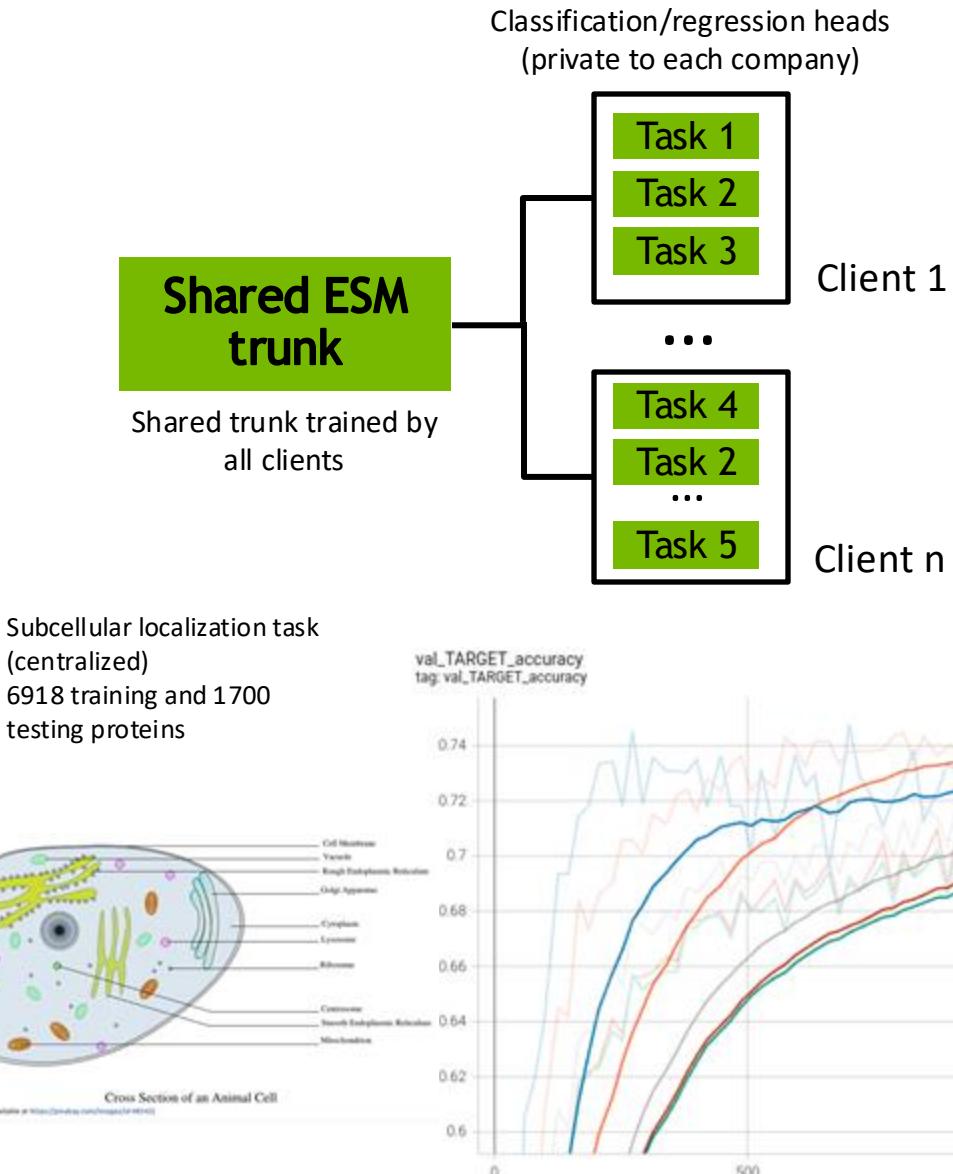
BaseModel - Before SFT

	H_acc	H_acc_norm	P_acc	P_acc_norm	W_acc	Mean
BaseModel	0.357	0.439	0.683	0.689	0.537	0.541
Alpaca	0.372	0.451	0.675	0.687	0.550	0.547
Dolly	0.376	0.474	0.671	0.667	0.529	0.543
Oasst1	0.370	0.452	0.657	0.655	0.506	0.528
Combined	0.370	0.453	0.685	0.690	0.548	0.549
FedAvg	0.377	0.469	0.688	0.687	0.560	0.556

Table 1. Model performance on three benchmark tasks: HellaSwag (H), PIQA (P), and WinoGrande (W)

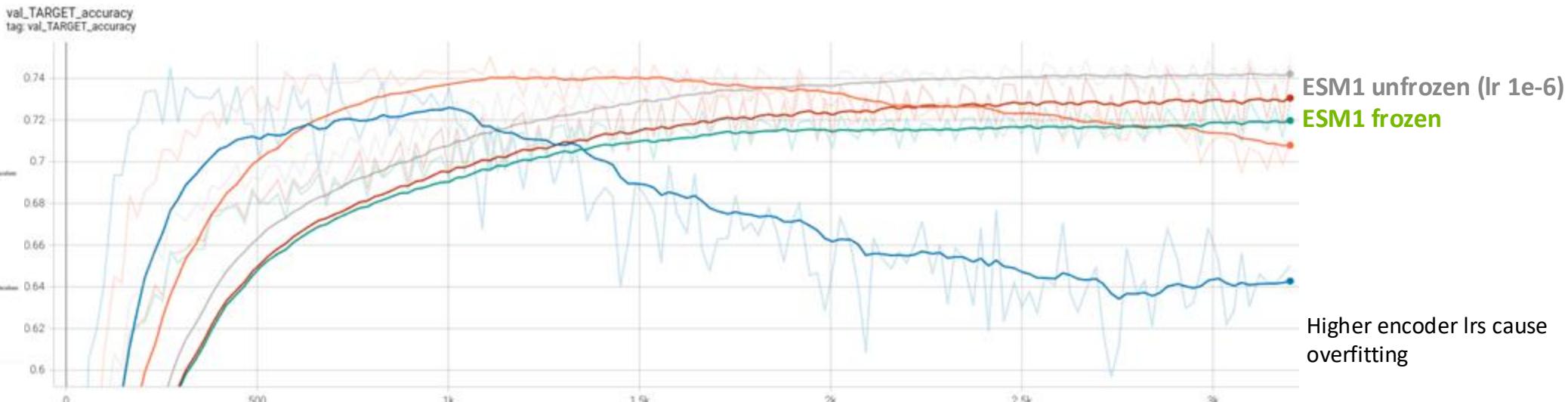
Biopharma Applications

BioNeMo ESM Fine-tuning



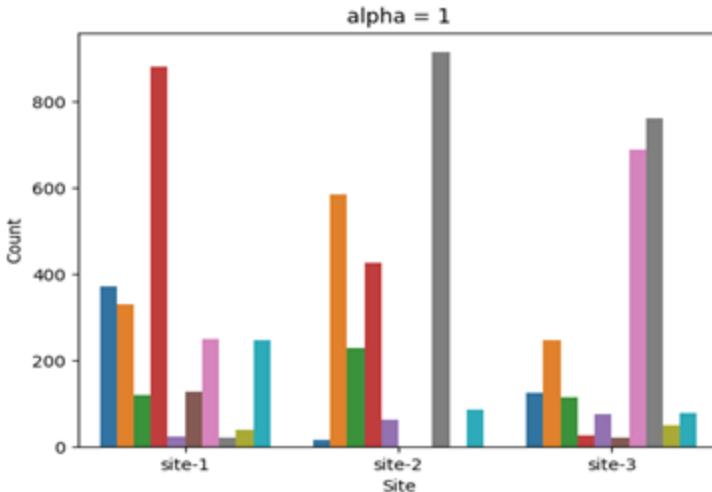
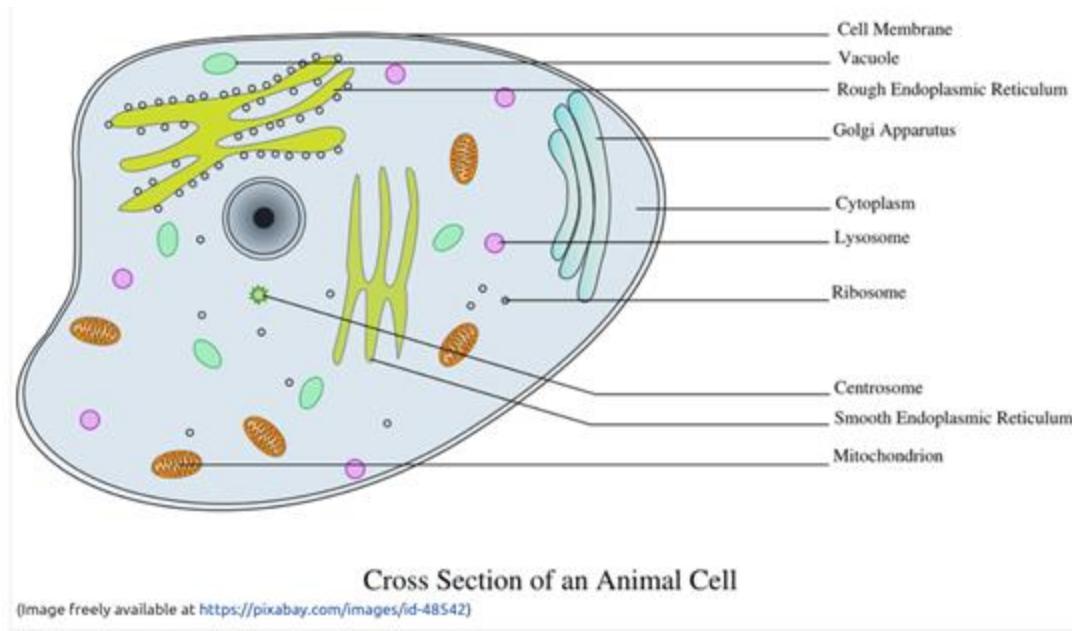
Model	Params	Frozen Encoder	Validation Accuracy
ESM1nv	43.6 M	Yes	0.7344
ESM1nv	43.6 M	No	0.7488
ESM2nv	650M	Yes	0.7883
ESM2nv	650M	No	0.7919

ESM-2 is a pre-trained, bi-directional encoder (BERT-style model) over amino acid sequences.

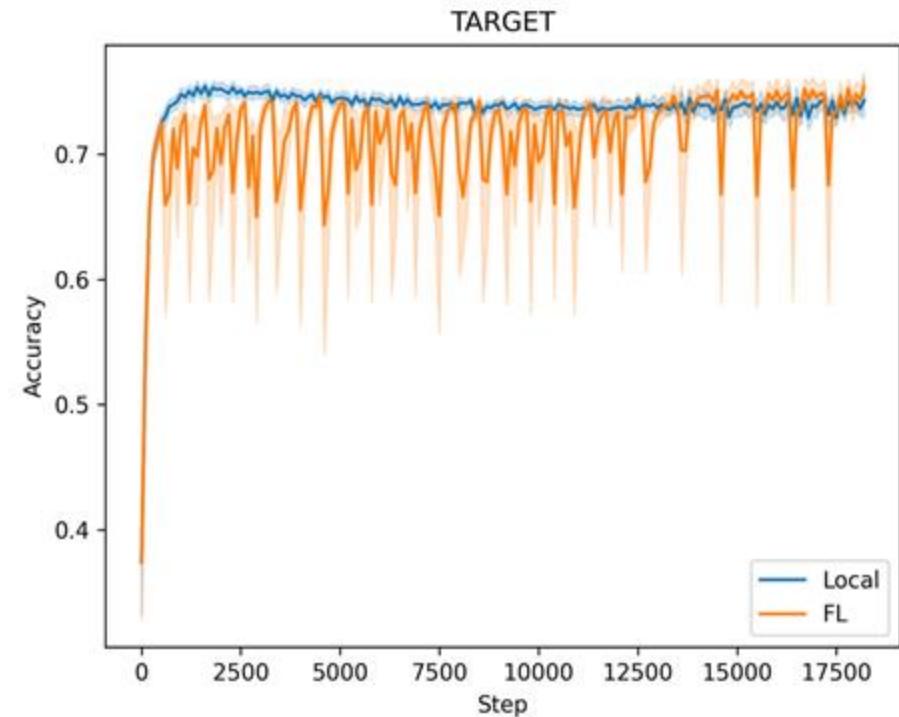


FedAvg With ESM2nv 650M

Subcellular Location Prediction



- 3 clients
- alpha 1.0
- 20 rounds/epochs
- ~2300 sequences per client
- Frozen ESM encoder



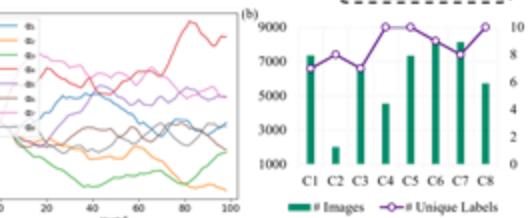
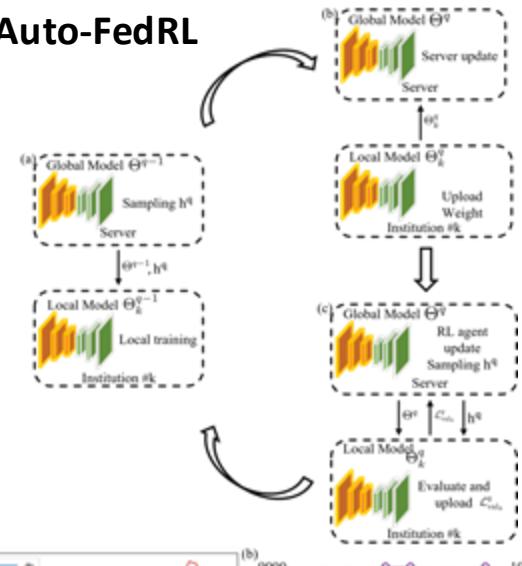
Setting	Accuracy
Local	0.773
FL	0.776

More examples: <https://github.com/NVIDIA/NVFlare/tree/main/examples/advanced/bionemo>

Research With NVIDIA FLARE

<https://github.com/NVIDIA/NVFlare/tree/dev/research>

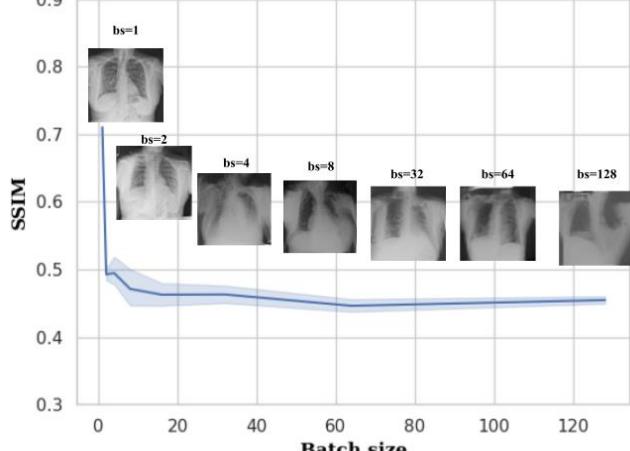
Auto-FedRL



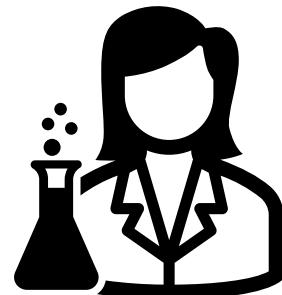
[P. Guo et al. ECCV 2022](https://arxiv.org/pdf/2204.07074.pdf)

Baseline Implementations

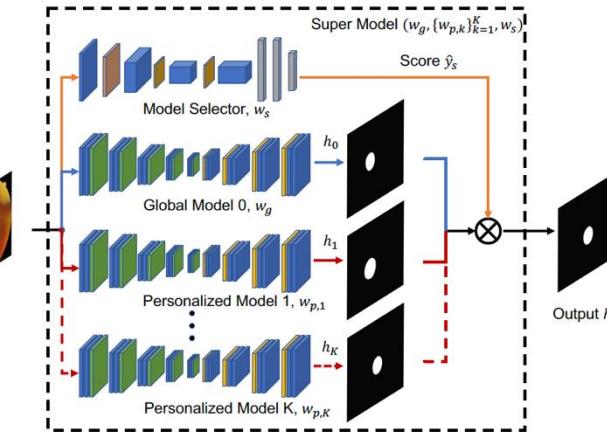
- FedAvg
- FedProx
- FedOpt
- SCAFFOLD
- Ditto
- Cyclic Weight Transfer
- Swarm Learning



[A. Hatamizadeh et al. TMI 2022](https://arxiv.org/pdf/2204.07074.pdf)

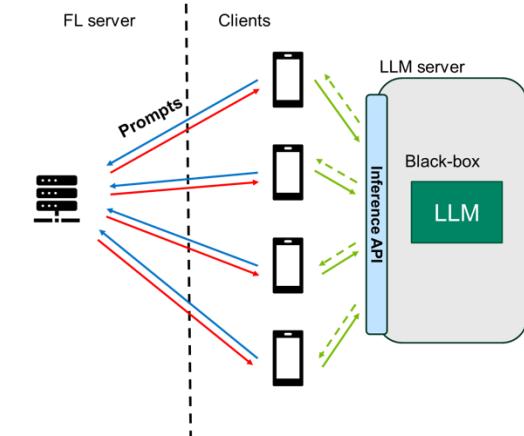


FedSM: Personalized FL

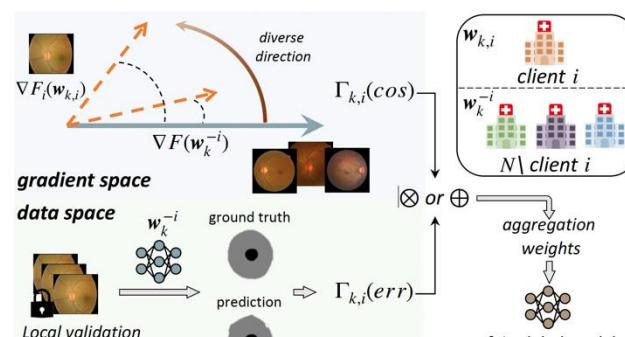


[A. Xu et al. CVPR 2022](https://arxiv.org/pdf/2204.07074.pdf)

Federated Black-Box Prompt Tuning



FedCE: Contribution Estimation



[M. Jiang et al. CVPR 2023](https://arxiv.org/pdf/2304.07074.pdf)



Agenda

- What is Federated Learning?
 - NVIDIA Key Technologies for Federated Learning
 - Real-world Use cases of Federated Learning
 - Getting Started with NVIDIA FLARE
 - Research: Addressing Key challenges in Federated Learning
 - **Summary & Announcements**
-
-
-
-
-
-
-

NVIDIA FLARE: Summary

A domain-agnostic, open-source, extensible FL framework

- **Federated Computing** -- a federated computing framework at core
- **End-to-End Confidential Federated AI** – combination of confidential computing and additional technologies to secure end-to-end process
- **Built for productivity** -- designed for maximum productivity, providing a range of tools to enhance user experience
- **Built for security & privacy** -- prioritizes robust security and privacy preservation
- **Built for concurrency & scalability** -- designed for concurrency, supporting resource-based multi-job execution
- **Built for customization** -- structured in layers, with each layer composed of customizable components
- **Built for integration** -- multiple integration options with third-party system
- **Built for production** -- robust, production-scale deployment in real-world federated learning and computing scenarios
- **Rich examples repository** -- wealth of built-in implementations, tutorials and examples
- **Growing application categories** -- medical imaging, medical devices, edge device application, financial services, HPC and autonomous driving vehicles

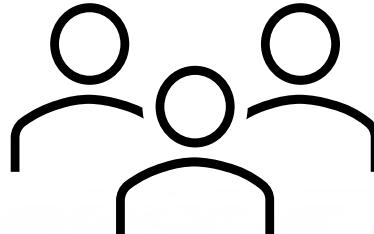
GitHub : <https://github.com/NVIDIA/NVFlare>

Web: <https://nvidia.github.io/NVFlare/>



Get Involved!

NVIDIA FLARE is open-source



NVIDIA FLARE DAY 2025

Exploring Real-world Examples of Federated Learning

September 2025 (US + EMEA)

September 2025 (APAC)

2024 Recordings available:

<https://nvidia.github.io/NVFlare/flareDay>

NVIDIA Academic Grant Program for Researchers:

Edge AI and Federated Learning

<https://www.nvidia.com/en-us/industries/higher-education-research/academic-grant-program>

NVIDIA FLARE GitHub : <https://github.com/NVIDIA/NVFlare>

Web Page : <https://nvidia.github.io/NVFlare/>



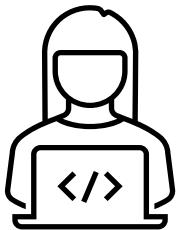
NVIDIA®

NVIDIA FLARE Summer of Code 2025

Select contributions will feature at NVIDIA FLARE DAY in September

Step 1: Register a Proposal

[Register Now](#)



Step 2: Develop your Contribution

- Work on your solution following the project's coding standards and best practices.
- Regularly engage with maintainers and the community for feedback.
- Test your code thoroughly and document it properly.

Step 3: Submit a Pull Request

- Fork the [NVFlare repository](#) and create a branch for your contribution.
- Submit a clear and well-documented PR with:
 - A description of your solution
 - Relevant benchmarks or evaluations
 - Any necessary instructions for running/testing the contribution
- Label your PR with "NVFlare Summer of Code" to ensure visibility.
- For more info, see the [CONTRIBUTION GUIDE](#).

Step 4: Review and Iterate

- Address feedback from the maintainers and community.
- Refine your work based on suggestions.

Step 5: Final Submission and Selection

- By September 1, 2025, ensure your PR is merged or under final review.
- Select contributions will be invited to present at NVIDIA FLARE DAY in September 2025 (5-min Lightning Talk).

Focus Areas

We welcome contributions in the following areas:

- **Energy-Efficient Communication & Learning** (e.g., compression, pruning, parameter-efficient finetuning, asynchronous FL, carbon footprint tracking)
- **Privacy-Preserving Techniques in the Age of LLMs** (e.g., differential privacy, memorization prevention, privacy attacks & defenses)
- **Training on the Edge (Real-Time & Multi-Sensor Data)** (e.g., mobile, automotive, robotics, surgery, dynamic sensor networks)
- **Collaborative Inference & Information Aggregation** (e.g., federated RAG, multi-agent learning, decentralized inference)

Contribution Types

We are looking for the following contributions:

- **Algorithms** – Novel methods to improve federated learning efficiency, scalability, and security
- **Core Framework & Utilities** – Enhancements to NVFlare's core functionalities
- **Examples & Use Cases** – Demonstrations of FL applications in real-world scenarios

NVIDIA FLARE Webinars

2025

- **Q1: Introduction to NVIDIA Federated Learning: Concepts, Technology, and Use Cases**

Kick off the series with a comprehensive dive into federated learning fundamentals, the technology that powers it, and real-world use cases.

- **Q2: Federated Learning Applications: Healthcare and Finance**

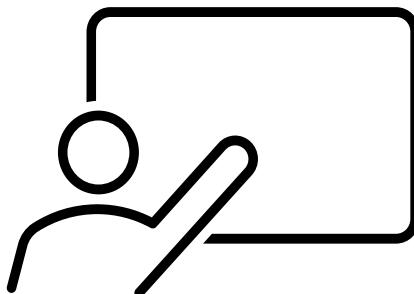
Explore how federated learning is transforming two of the most critical industries, driving innovation while preserving data privacy.

- **Q3: Provision, Run, and Monitor Federated Learning Applications**

How to setting up, executing, and managing federated learning system

- **Q4: Real-world Federated Learning Use Cases from Industries**

Delve deeper into how federated learning is shaping advanced use cases across diverse industries.





Thank you!

Questions?

Holger Roth <hroth@nvidia.com>