



NVIDIA FEDERATED LEARNING

BUILDING AI FOR REAL-WORLD CLINICAL PERFORMANCE

Taking Algorithms Beyond Proof-of-Concept

REAL-WORLD AI DESIGN

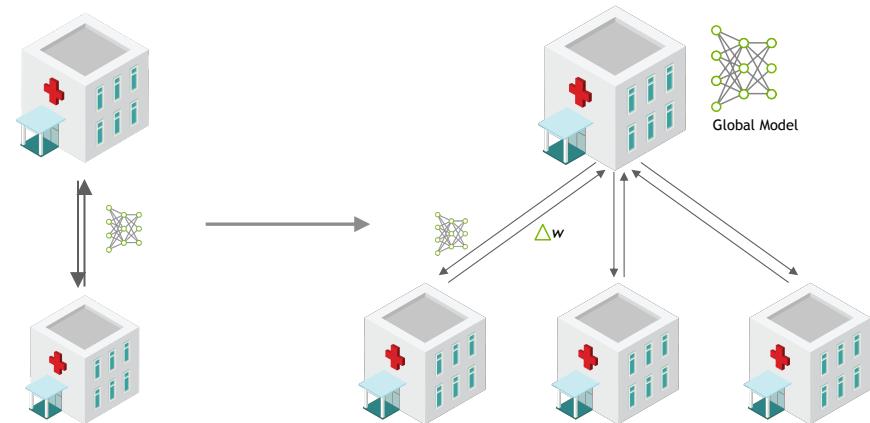
External Validation, Multiple Institutions, Prospective Data

Design Characteristic	All Articles (n = 516)	Articles Published in Medical Journals (n = 437)
External validation		
Used	31 (6.0)	27 (6.2)
Not used	485 (94.0)	410 (93.8)
In studies that used external validation		
Diagnostic cohort design	5 (1.0)	5 (1.1)
Data from multiple institutions	15 (2.9)	12 (2.7)
Prospective data collection	4 (0.8)	4 (0.9)
Fulfillment of all of above three criteria	0 (0)	0 (0)
Fulfillment of at least two criteria	3 (0.6)	3 (0.7)
Fulfillment of at least one criterion	21 (4.1)	18 (4.1)

Only 6% of published AI studies have external validation
Few included multiple institutions

FEDERATED LEARNING PARADIGM

Model to Data | Generalize Model

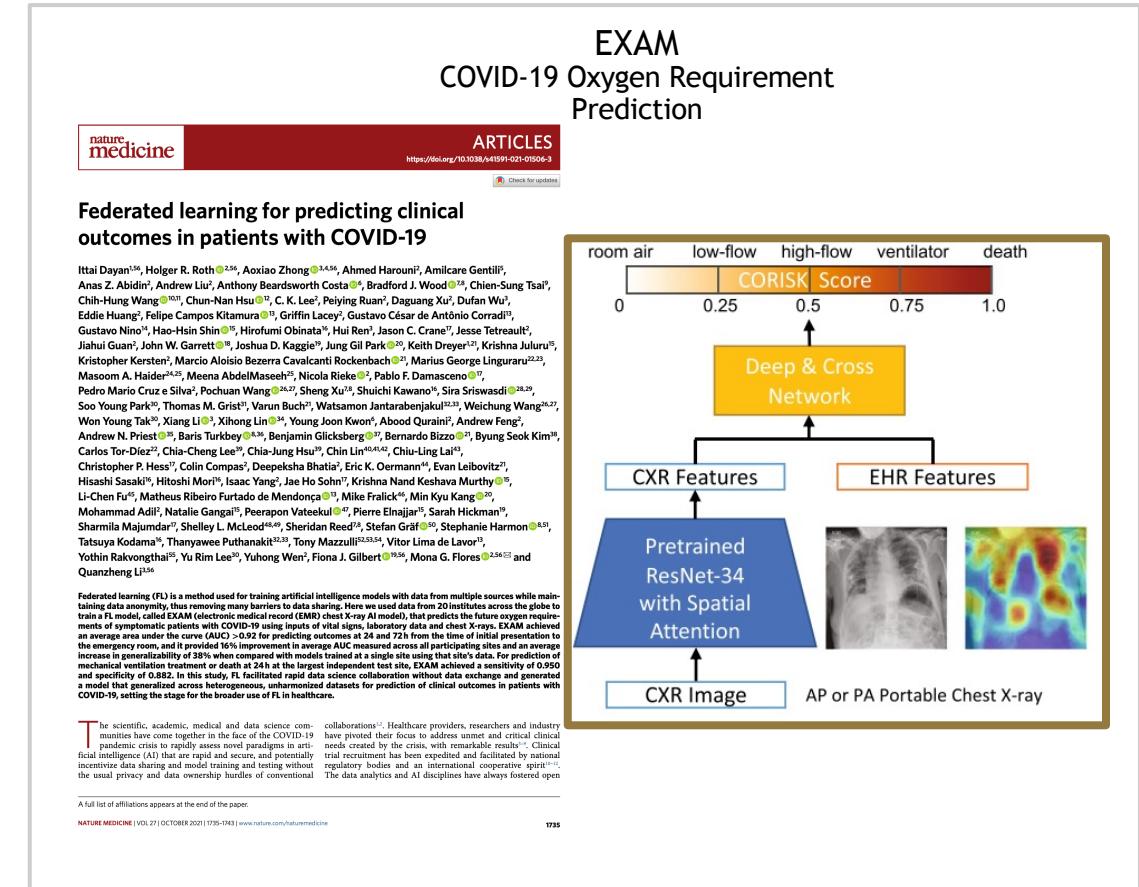
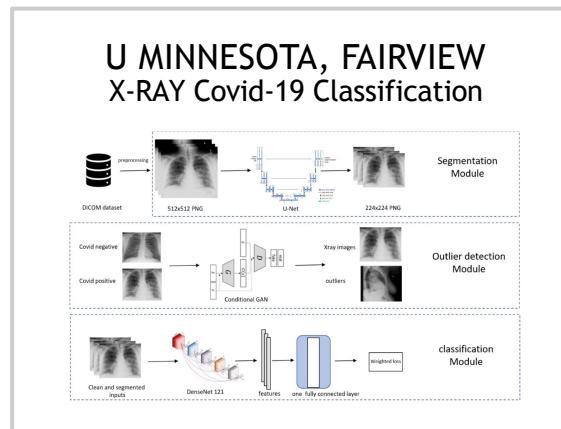
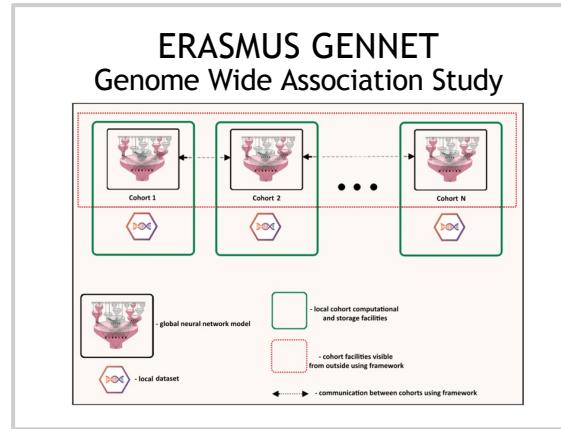
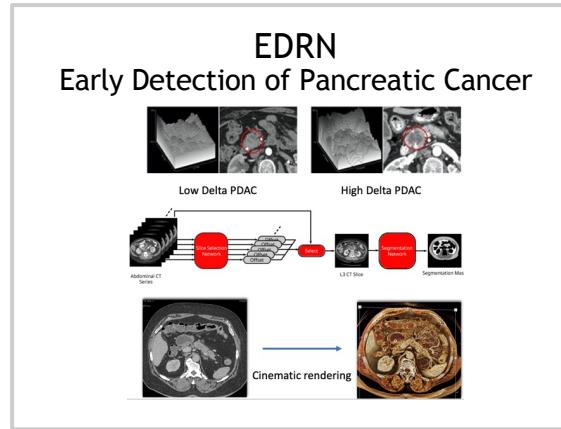


Transfer Learning
“Adapt”

Federated Learning
“Generalize”

FEDERATED LEARNING MOMENTUM

Breaking Healthcare Data Siloes



NVIDIA FEDERATED LEARNING

Applications across industries

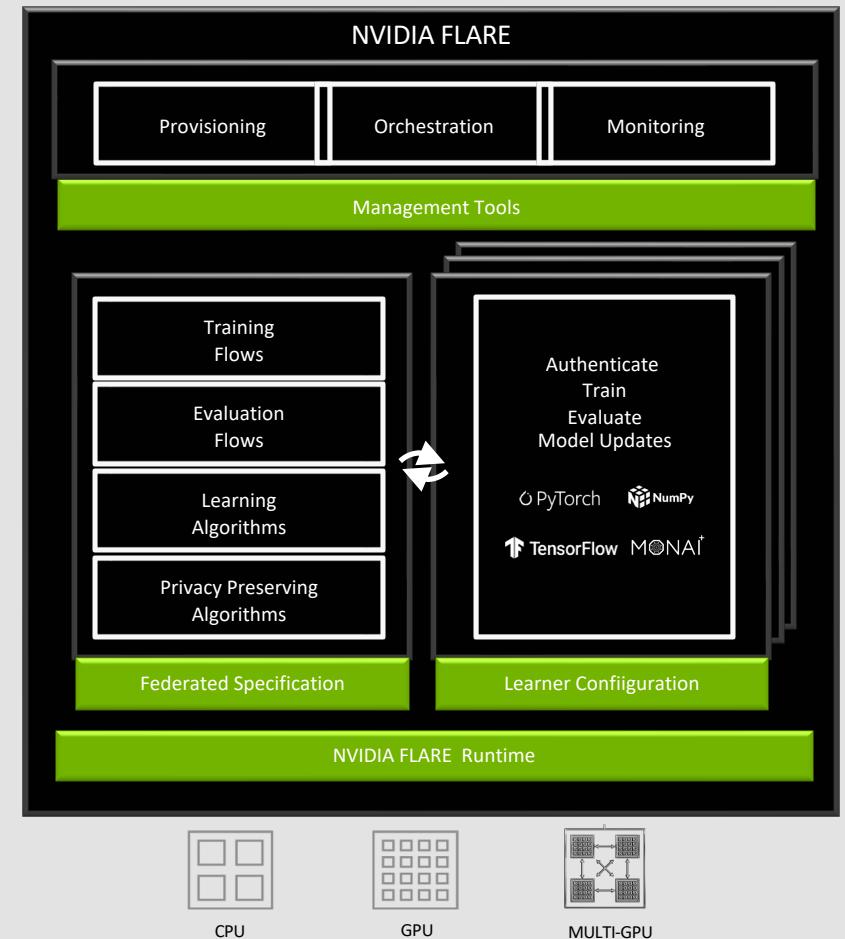


NVIDIA FLARE

Open-Source SDK for Federated Learning

- Apache License 2.0 to catalyze FL research & development
- Enables Distributed, Multi-Party Collaborative Learning
- Production Scalability with high availability and multi-task execution
- Adapt existing ML/DL workflows to a Federated paradigm
- Privacy Preserving Algorithms
 - Homomorphic Encryption & Differential Privacy
- Secure Provisioning, Orchestration & Monitoring
- Programmable APIs for Extensibility

Available on Github: <https://github.com/nvidia/nvFlare>



PERSONAS (WHO & VALUE PROP FOR EACH)

FL RESEARCHERS



Enables ease of getting started with FL experiments execution & evaluation in real world.

Extensible APIs for ease of creating custom implementations for new federated workflows, learning & privacy preserving algorithms.

DATA SCIENTISTS



Extend existing DL/ML workflows with a Federated paradigm and explore potential of Federated learning.

Ready to use FL specification and management tools enabling seamless execution.

PLATFORM DEVELOPERS

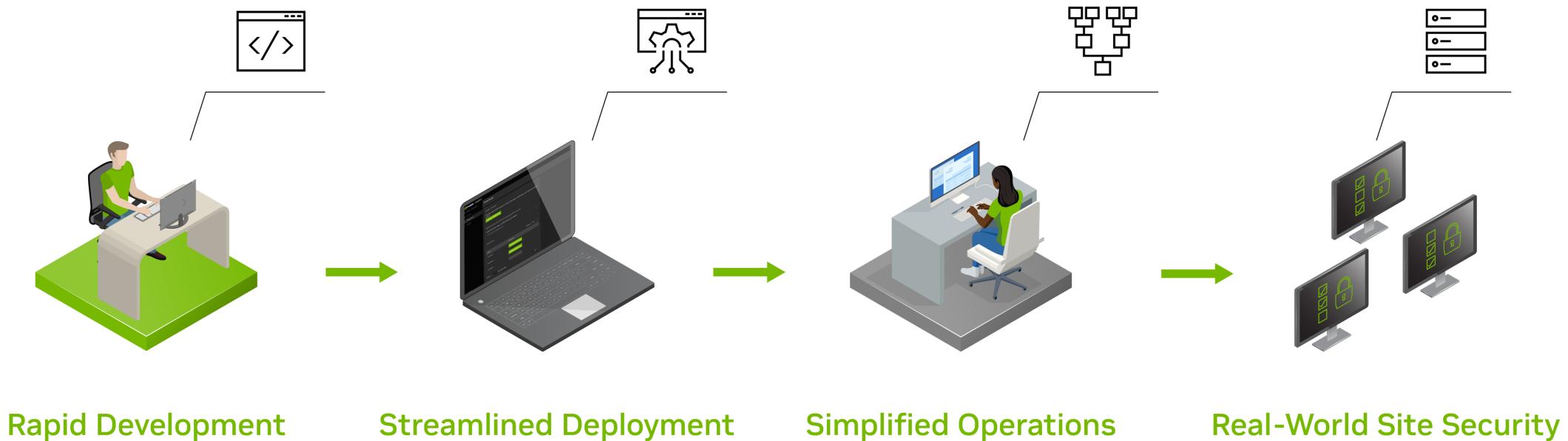


A robust, extensible foundation to customize a platform offering for end users.

Built-in implementations of Federated learning spec & Aux APIs to build custom offerings.

NVIDIA FLARE Workflow

From rapid research prototyping to streamlined real world deployment



FEDERATED 2.2 NEW FEATURES

From Research Simulation to Real World Deployment

FL SIMULATOR

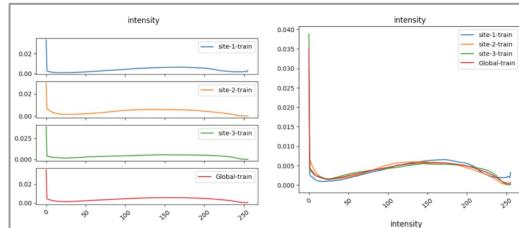
Rapid Development and Debugging

```
def run_simulator(simulator_args):
    simulator = SimulatorRunner(
        job_folder=simulator_args.job_folder,
        workspace=simulator_args.workspace,
        clients=simulator_args.clients,
        n_clients=simulator_args.n_clients,
        threads=simulator_args.threads,
        gpu=simulator_args.gpu,
        max_clients=simulator_args.max_clients,
    )
    run_status = simulator.run()

    return run_status
```

FEDERATED STATS

Analyze data distributions



FRAMEWORKS INTEGRATION

MONAI & XGBOOST



FLARE DASHBOARD

Streamlined operation & deployment



UNIFIED CLI

Multi-task Learning Chemical Assays



CLIENT CONTROLLED PP

Genome Wide Association Study



NVIDIA FLARE KEY CAPABILITIES

Runtime-ready and extensible suite of features

Privacy-Preserving Algorithms

NVIDIA FLARE provides privacy-preserving algorithms that ensure each change to the global model stays hidden and prevent the server from reverse-engineering the submitted weights and discovering any training data.

Training and Evaluation Workflows

Built-in workflow paradigms use local and decentralized data to keep models relevant at the edge, including learning algorithms for FedAvg, FedOpt, and FedProx.

Extensible Management Tools

Management tools help secure provisioning using SSL certifications, orchestration through an admin console, and monitoring of federated learning experiments using TensorBoard for visualization.

Supports Popular ML/DL Frameworks

Flexible in design, the SDK can be used with PyTorch, Tensorflow, and even Numpy, which allows for integrating federated learning into your current workflow.

Extensive API

Its extensive and open-source API enables researchers to develop new federated workflow strategies, innovative learning, and privacy-preserving algorithms.

Reusable Building Blocks

NVIDIA FLARE provides an easy way to perform federated learning experiments by utilizing the reusable building blocks and example walkthroughs.

<https://developer.nvidia.com/flare>



SECURITY & PRIVACY

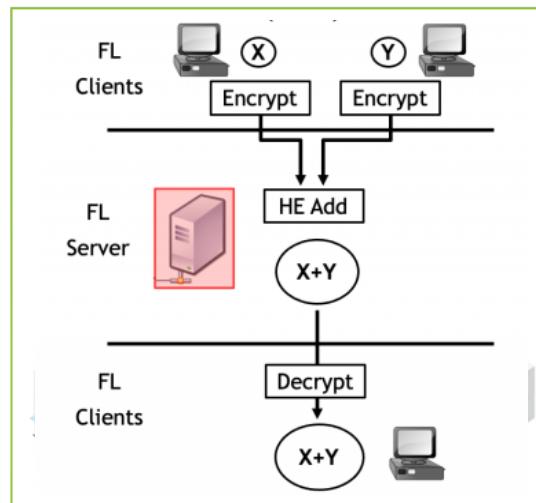
Homomorphic Encryption & Differential Privacy

Federated Learning with Homomorphic Encryption

What if I don't trust the server?

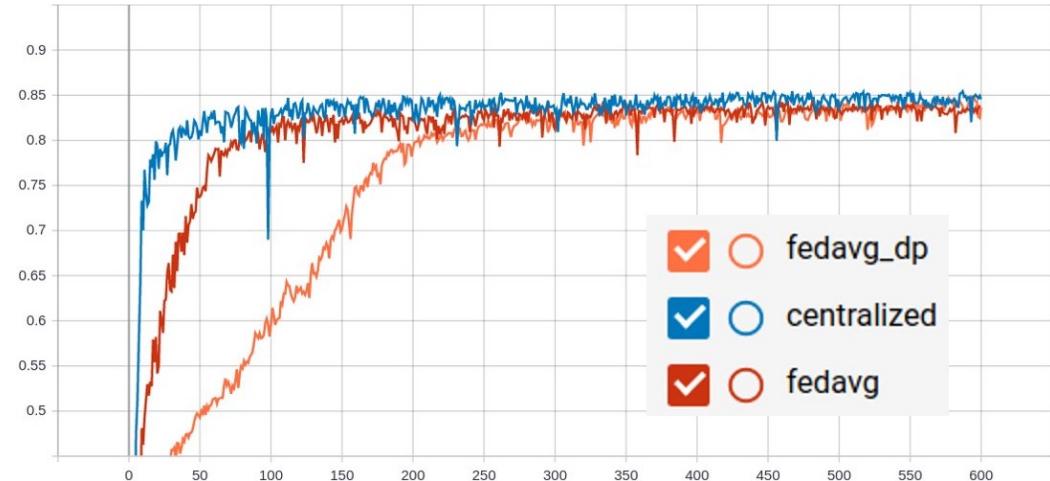
Homomorphic encryption (HE)

A form of encryption that permits users to perform computations on encrypted data



Differential Privacy for BraTS18 Segmentation

validation Dice scores of the global model for 600 training epochs:



Blog: <https://developer.nvidia.com/blog/federated-learning-with-homomorphic-encryption/>

Example: <https://github.com/NVIDIA/NVFlare/tree/main/examples/cifar10>

Example: <https://github.com/NVIDIA/NVFlare/tree/main/examples/brats18>

END-TO-END EXAMPLES (CIFAR10, BRATS18, PROSTATE)

▪ Comprehensive example for researchers to compare algorithms

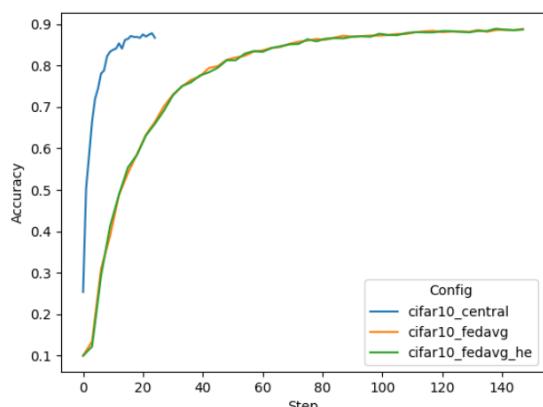
1. Set up a virtual environment
2. Create your FL workspace
3. Run automated experiments
 1. Varying data heterogeneity of data splits
 2. Centralized training
 3. FedAvg on different data splits
 4. Advanced FL algorithms (FedProx and FedOpt)
 5. Secure aggregation using homomorphic encryption
 6. Differential privacy

4. Results

4.1 Central vs. FedAvg

With a data split using `alpha=1.0`, i.e. a non-heterogeneous split, we achieve the following performance. We can see that FedAvg achieves a similar performance to central training and that HE does not impact the performance.

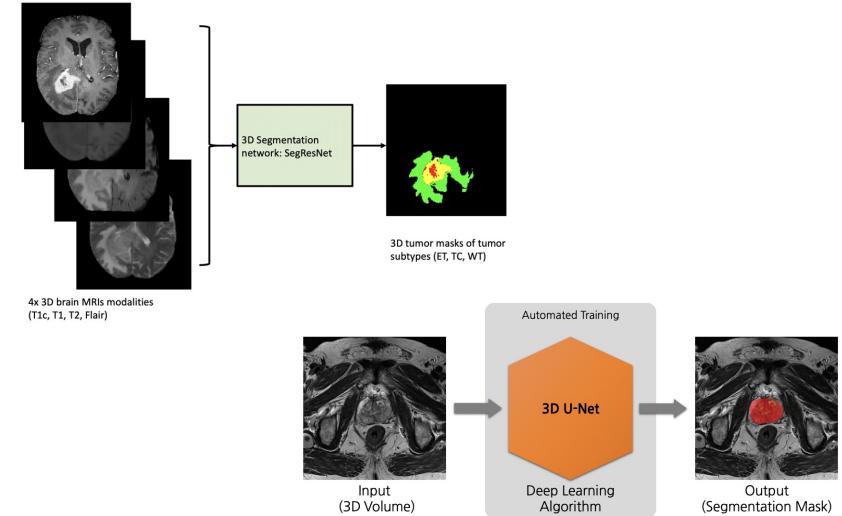
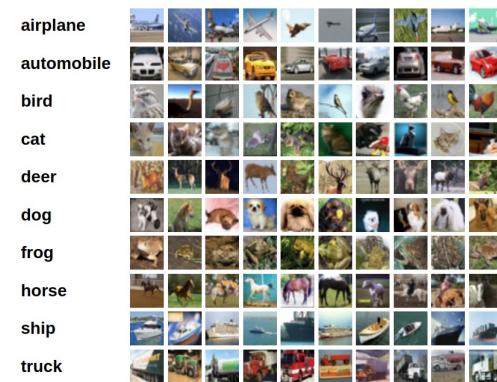
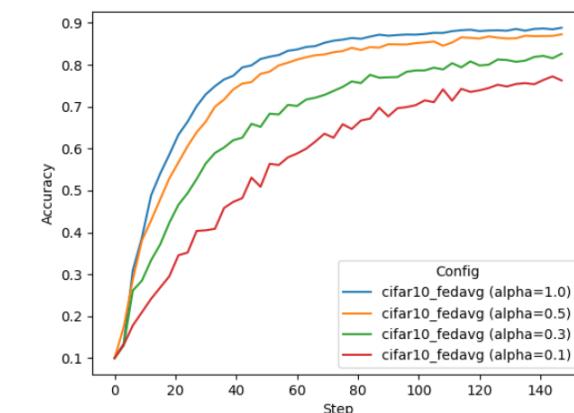
Config	Alpha	Val score
cifar10_central	1.0	0.8668
cifar10_fedavg	1.0	0.8840
cifar10_fedavg_he	1.0	0.8868



4.2 Impact of client data heterogeneity

We also tried different `alpha` values, where lower values cause higher heterogeneity. This figure shows the convergence rate of the FedAvg algorithm.

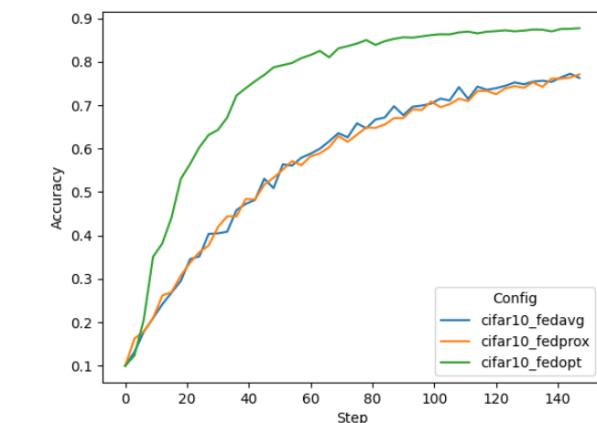
Config	Alpha	Val score
cifar10_fedavg	1.0	0.8840
cifar10_fedavg	0.5	0.8727
cifar10_fedavg	0.3	0.8264
cifar10_fedavg	0.1	0.7626



4.3 FedProx vs. FedOpt

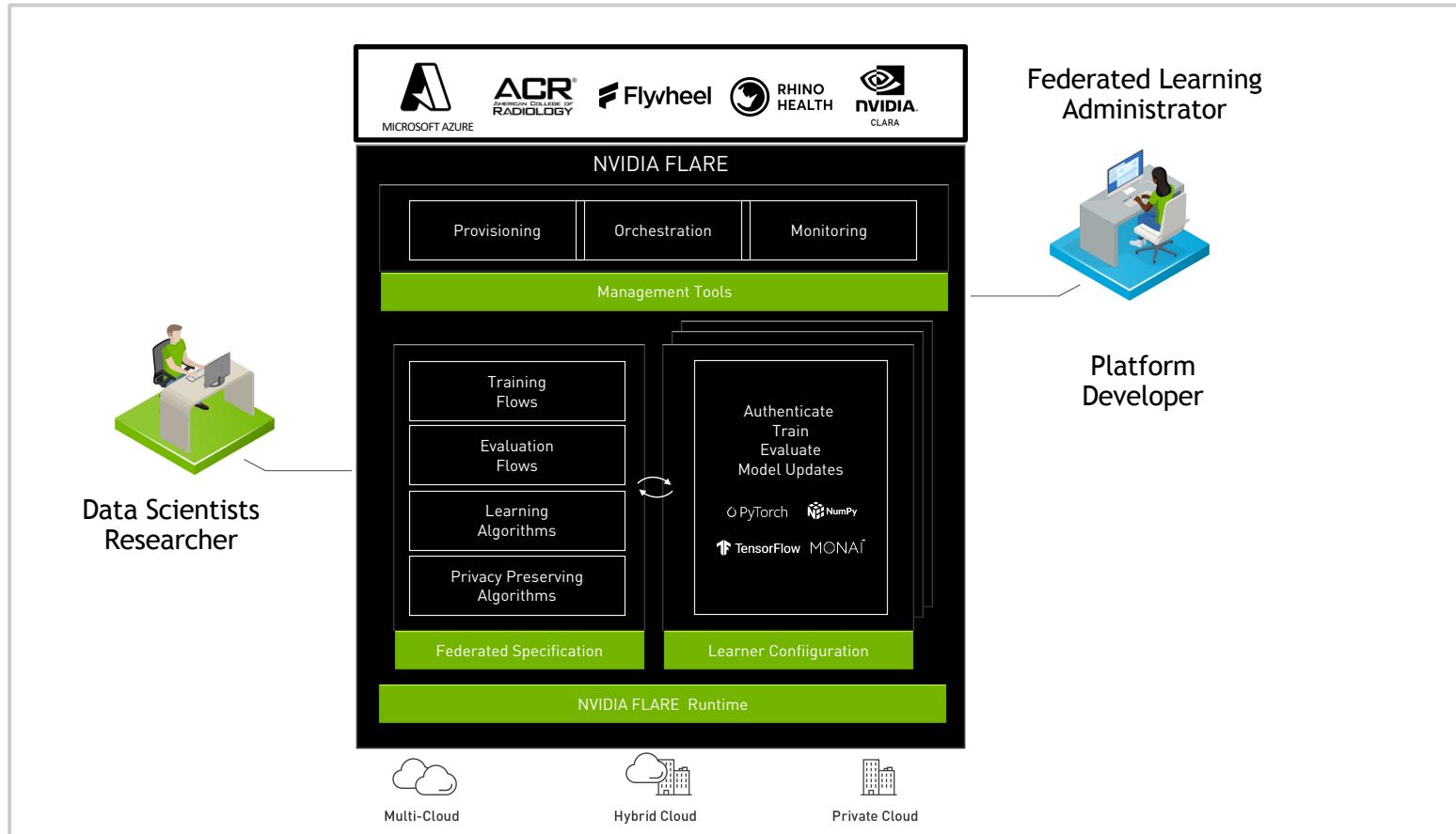
Finally, we are comparing an `alpha` setting of 0.1, causing a high client data heterogeneity, namely FedProx and FedOpt. Both achieve a better performance compared to FedAvg in terms of convergence rate by utilizing SGD with momentum to update the global model on the amount of training steps.

Config	Alpha	Val score
cifar10_fedavg	0.1	0.7626
cifar10_fedprox	0.1	0.7709
cifar10_fedopt	0.1	0.7963



NVIDIA FLARE 2.0

Accelerating the Federated Learning Paradigm



The background of the image is a dark gray or black color. It features a dense, abstract pattern of green glowing particles. These particles are arranged in a grid-like structure, with some lines being more prominent than others, creating a sense of depth and motion. The particles appear to be small, glowing spheres that are scattered across the entire frame.

NVIDIA FLARE

FALL 2021

LEARNING ALGORITHMS

Federated Averaging (FedAvg)

- Weighted average to update global model

[McMahan et al.](#)

Federated Proxy (FedProx)

- Clients add a loss to stay close the global model
- Avoids models drifting away from global model in heterogenous datasets

[Li et al.](#)

Adaptive Federated Optimization (FedOpt)

- Global model is updated using an optimizer (SGD w. momentum, Adam, Yogi, Adagrad, etc.)

[Reddi et al.](#)

Cyclic Weight Transfer

- Models are continuously fine-tuned and circulated around institutions

[Chang et al.](#)

SCAFFOLD

- Adds correction terms during training to deal with non-IID

[Karimireddy et al.](#)

Ditto

- Fairness through personalization

[Li et al.](#)

Algorithms can be extended

- Differential privacy
- Homomorphic Encryption

More to come...

FLARE 2.1: BUILT FOR SCALABILITY

High-Availability & Multi-Task Execution

- High availability (HA) supports multiple FL servers and automatically activates a backup server when the currently active server becomes unavailable.
- This is managed by a new entity in the federation, the overseer, that's responsible for monitoring the state of all participants and orchestrating the cutover to a backup server when needed.
- Multi-job execution supports resource-based multi-job execution by allowing for concurrent runs, provided that the resources required by the jobs are satisfied

<https://developer.nvidia.com/blog/experimenting-with-novel-distributed-applications-using-nvidia-flare-2-1/>

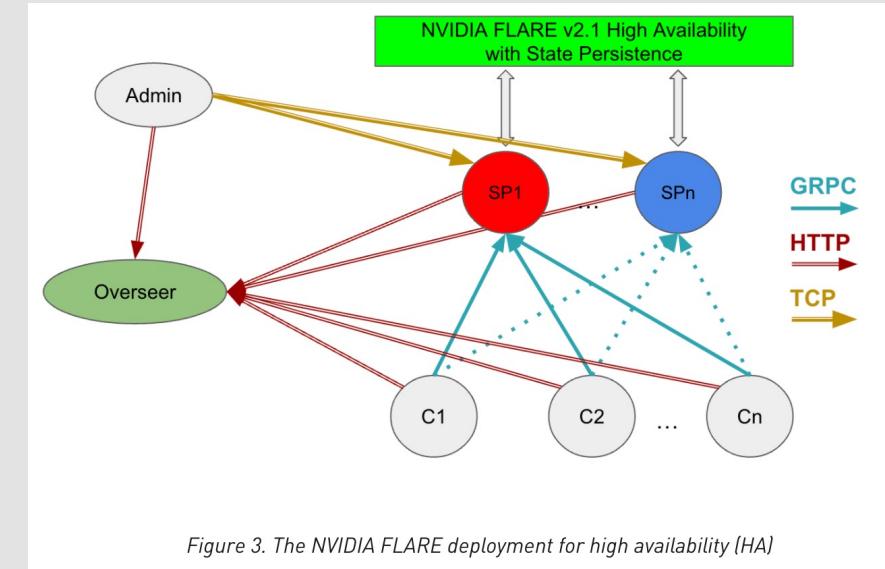


Figure 3. The NVIDIA FLARE deployment for high availability (HA)

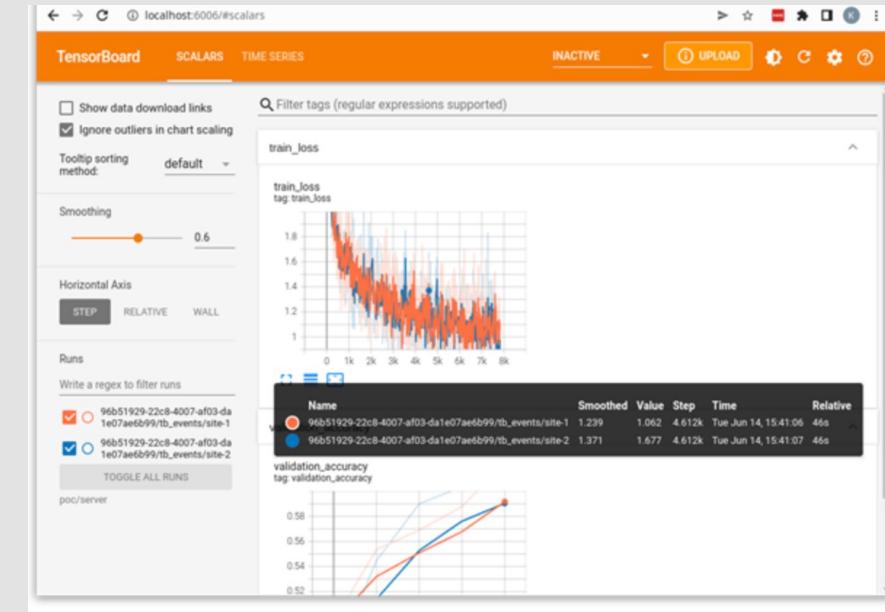


Figure 2. Example TensorBoard output from the hello-pt-tb application