

OmniFed: Towards Configurable Cross-Silo Federated Learning

Sahil Tyagi

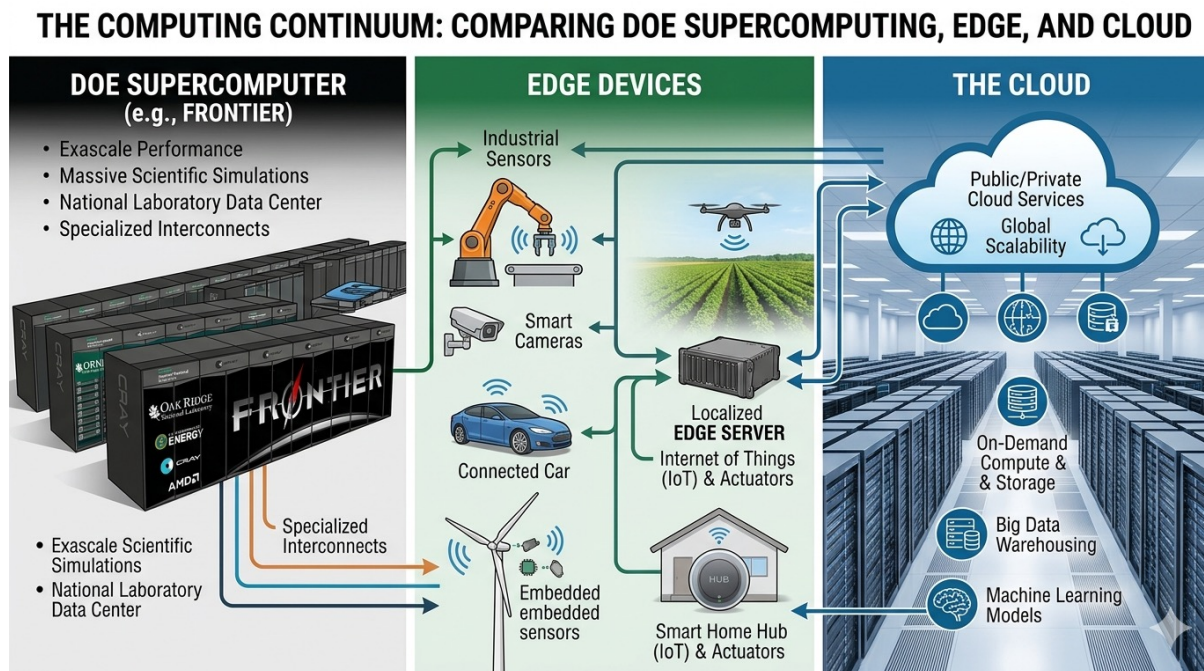
National Center for Computational Sciences

Oak Ridge National Laboratory



Trajectory of Federated Learning

- SoTA Models already pre-trained on the internet-scale data; next frontier lies behind silos
- Federated learning (FL) is more than mobile/edge devices training with a single server
- Expansion of the compute continuum; current research environment itself is more heterogeneous than ever!



Challenges

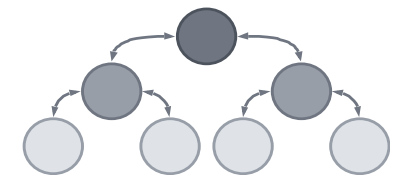
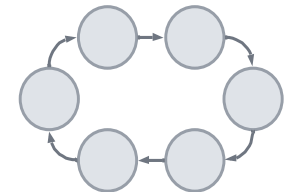
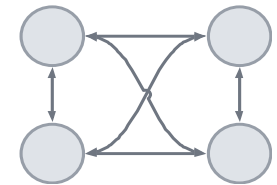
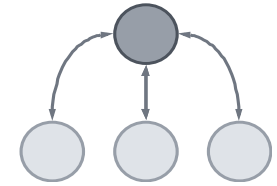


- Opportunities come with challenges
 - Data heterogeneity
 - Compute heterogeneity
 - Network heterogeneity
 - Privacy, security and authentication
 - Scalability

- To meet these challenges, tools/platforms need to make explicit design trade-offs:
 - production readiness vs. flexibility
 - Ease-of-use vs. configurability
 - Simulation vs. deployment

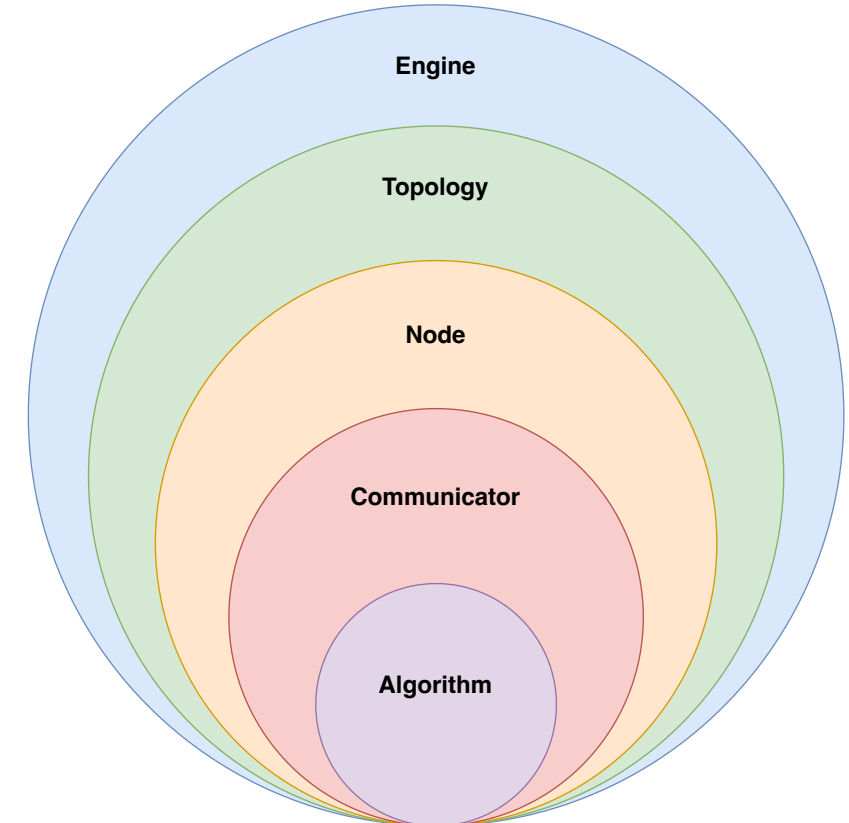
Motivations

- Rapid FL prototyping/testing with minimal boilerplate
- Flexibility for complex topologies
 - Centralized, decentralized, hierarchical
- Configuration-driven, schema-based deployment for reproducibility
 - Easily add/remove features like privacy, compression
- Swap-out design with clear separation of concerns
 - Separate compute resource, communication backend, algorithm strategy



Abstractions for configurable and flexible FL

- Core abstractions:
 - Algorithm
 - Communicator
 - Node
 - Topology
 - Engine
- Clear separation of concerns with minimal interfaces
- Override-what-you-need paradigm



Design Motivation Example

- Deploying FL algorithms with single file algorithm plugins

FedAvg

```
defaults:
- override topology: centralized
- override algorithm: fedavg
- override model: resnet152
- override datamodule: cifar100

topology:
_target_: src.omnifed.topology.CentralizedTopology
num_clients: 8
inner_comm:
_target_: src.omnifed.communicator.GrpcCommunicator
master_port: 50051
master_addr: 127.0.0.1

algorithm:
_target_: src.omnifed.algorithm.FedAvg
lr: 0.01

global_rounds: 2
```

Scaffold

```
defaults:
- override topology: centralized
- override algorithm: fedavg
- override model: resnet152
- override datamodule: cifar100

topology:
_target_: src.omnifed.topology.CentralizedTopology
num_clients: 8
inner_comm:
_target_: src.omnifed.communicator.GrpcCommunicator
master_port: 50051
master_addr: 127.0.0.1

algorithm:
_target_: src.omnifed.algorithm.Scaffold
lr: 0.01

global_rounds: 2
```

Design Motivation Example cont'd...

- Testing compression w.r.t. communication savings vs. convergence quality

```
inner_comm:
  _target_: src.omnifed.communicator.TorchDistCommunicator
  port: 28670
  compression:
    _target_: src.omnifed.communicator.compression.TopK
    k: 0.001
```

Sparsification

Quantization

```
inner_comm:
  _target_: src.omnifed.communicator.GrpcCommunicator
  port: 28670
  compression:
    _target_: src.omnifed.communicator.compression.QSGD
    bit_width: 8
```

```
inner_comm:
  _target_: src.omnifed.communicator.TorchDistCommunicator
  port: 28670
  compression:
    _target_: src.omnifed.communicator.compression.PowerSGD
    rank: 16
```

Low-rank approximation

Design Motivation Example

- Evaluating privacy-preserving mechanisms

```
inner_comm:  
  _target_: src.omnifed.communicator.TorchDistCommunicator  
  port: 28670  
  privacy:  
    _target_: src.omnifed.privacy.DifferentialPrivacy  
    epsilon: 1.0  
    delta: 0.00001
```

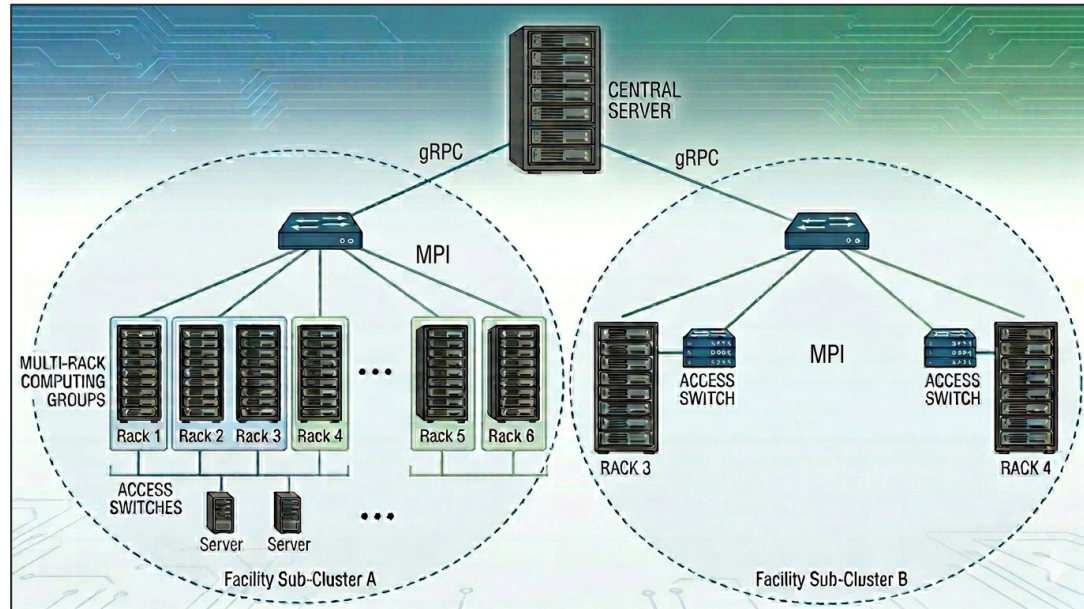
Noise injection-based

```
inner_comm:  
  _target_: src.omnifed.communicator.TorchDistCommunicator  
  port: 28670  
  privacy:  
    _target_: src.omnifed.privacy.HomomorphicEncryption  
    polymod_degree: 16384
```

Cryptography-based

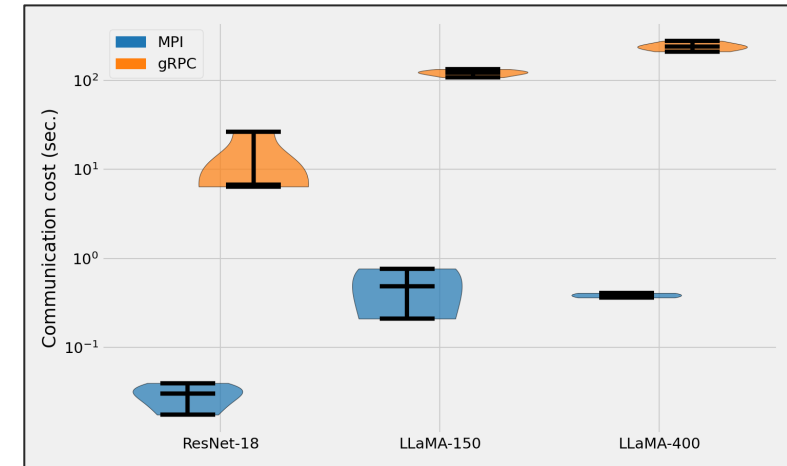
Use-case: Simulating Cross-silo FL

Hierarchical cross-silo training

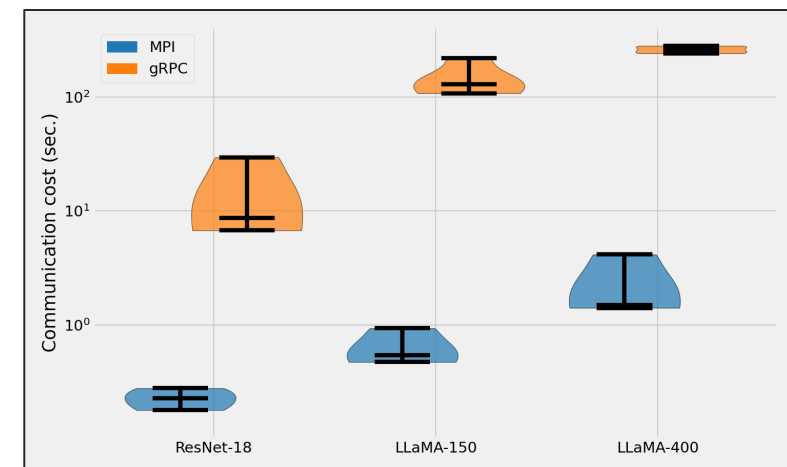


- Model cross-silo training as hierarchical FL
- Mix multiple communication protocol
- Tested on 16/128 nodes on Frontier @OLCF

16 Frontier nodes

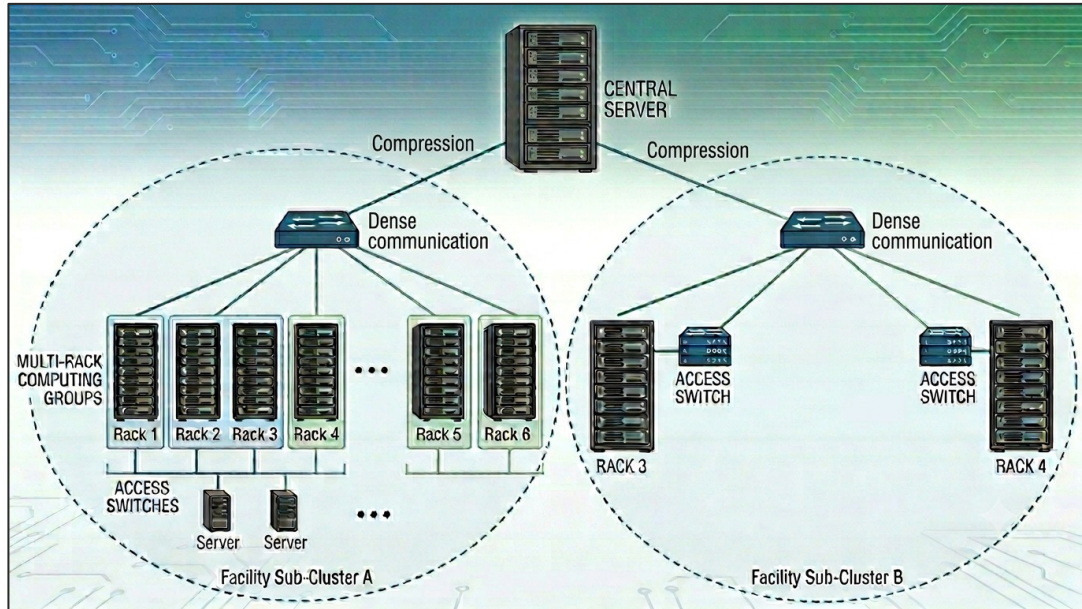


128 Frontier nodes

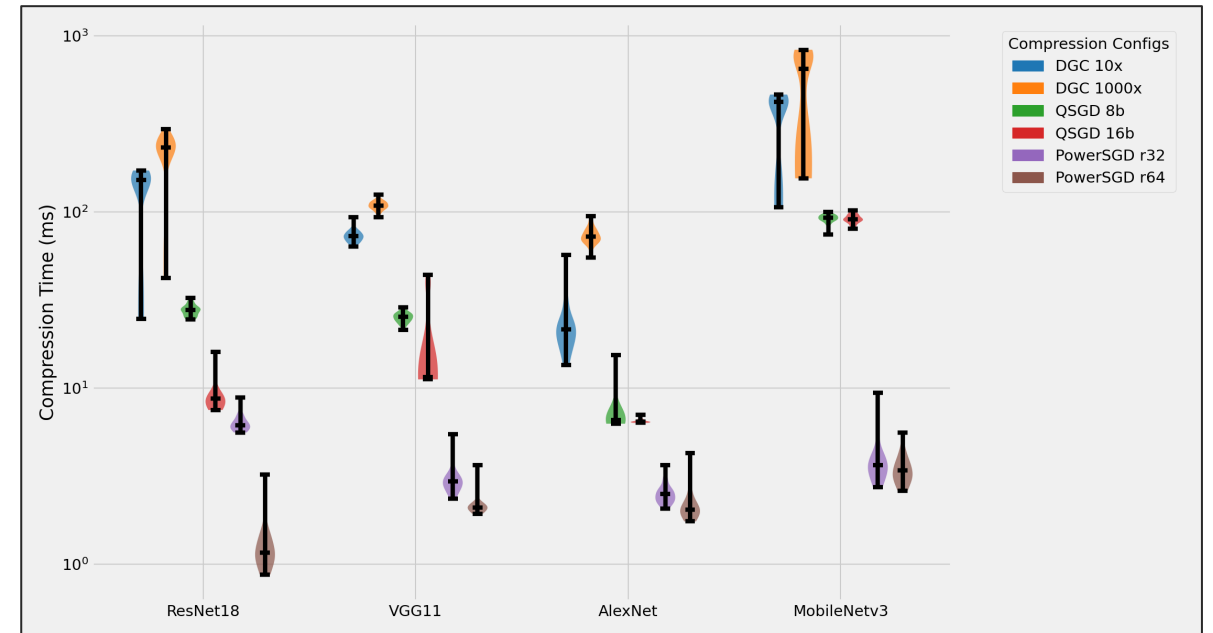


Use-case: Simulating Cross-silo FL cont'd...

Varying compression over sub-aggregators



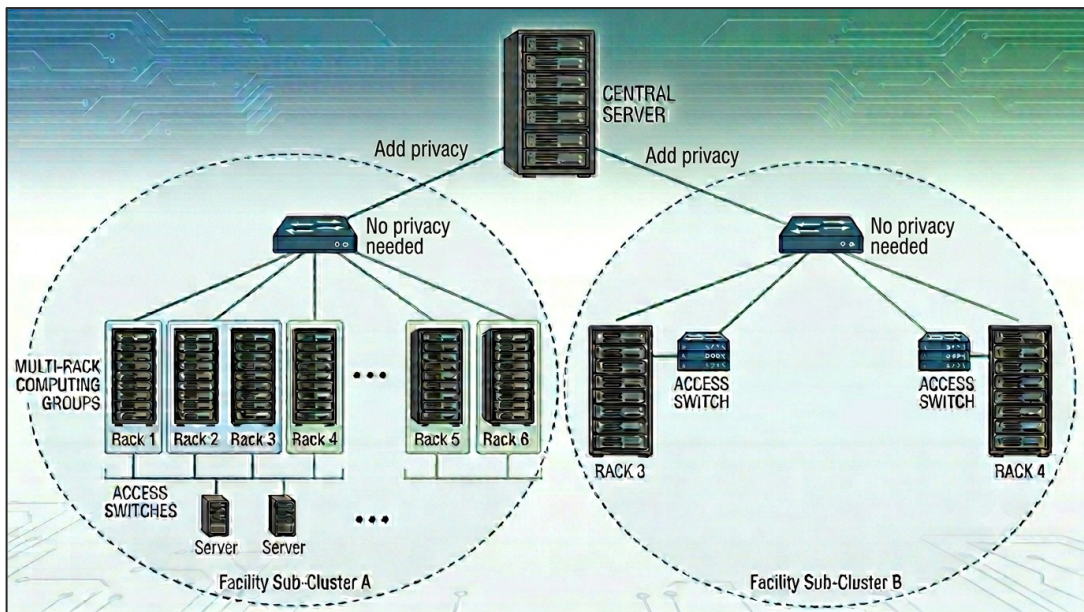
Comparing compression overhead



- Adaptive communication across sub-aggregators
- Tested on NVIDIA DGX server

Use-case: Simulating Cross-silo FL cont'd...

Adding privacy ONLY at infrastructure boundaries



Comparing overhead of privacy mechanisms

Model	Privacy overhead (seconds)	
	Differential Privacy	Homomorphic Encryption
ResNet-18	1.45	68.72
VGG-11	14.4	786
AlexNet	6.9	458.7
MobileNetv3	1.2	29.8

Varying Different Privacy parameters to compare model utility

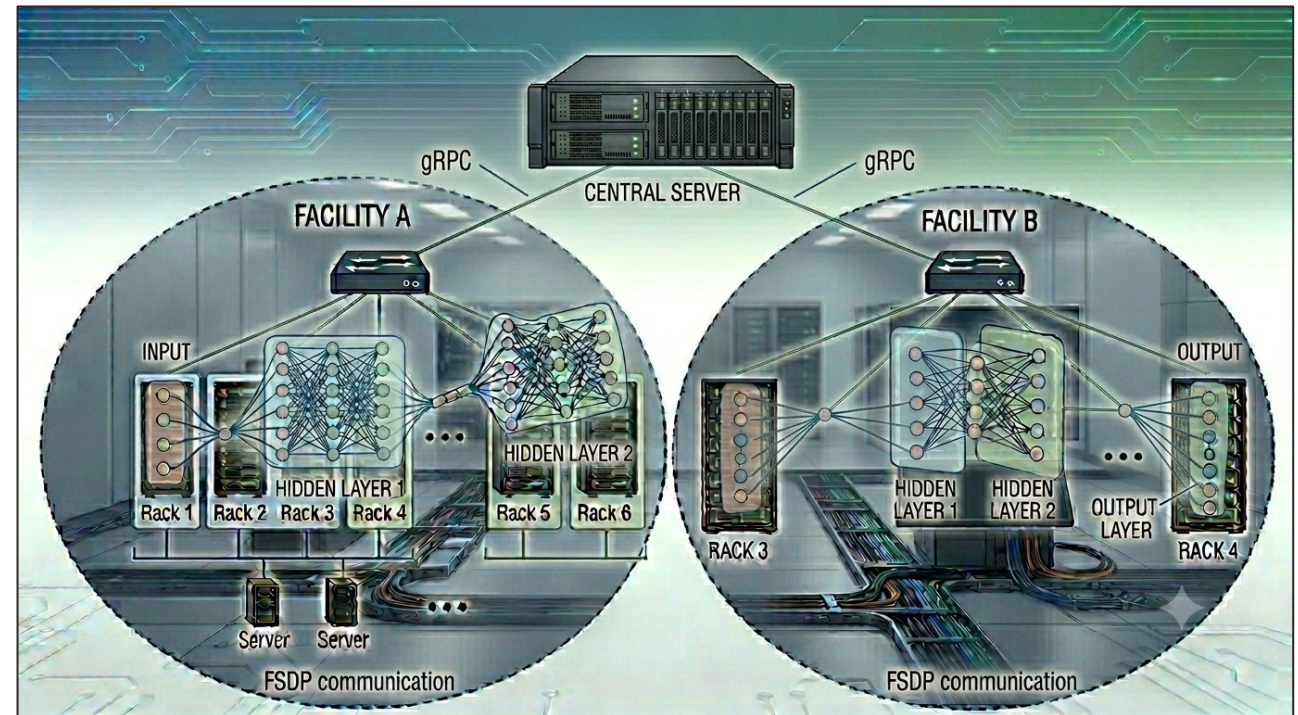
Model	DP Test accuracy	
	$\epsilon=1$	$\epsilon=10$
ResNet-18	97.98%	98.06%
VGG-11	24.1%	28.6%
AlexNet	16.31%	39.54%
MobileNetv3	23.72%	58.8%

- Adaptive privacy mechanisms across silos
- Data secure within a closed system; vulnerable outside

Use-case: Federated AI for LLMs

- Frontier/foundation models fine-tuned on scientific data stored behind secure silos
- Model-size not just a communication, but a computation bottleneck as well!
- Already optimized in HPC AI training!

Training/fine-tuning models too large for individual compute endpoints



Takeaways



- FL is underexplored on SoTA AI over scientific data residing behind closed silos
- Although always relevant, FL is scaling beyond clients training locally, and a central server aggregating globally
- Training configurations may comprise of *many* clients with *fewer* resources, and more recently, *fewer* clients with *larger* resource pools
- Research environment itself is heterogeneous, spanning edge, cloud and leadership-class HPC, and thus, should be *modular, flexible* and *configurable*
- Users/engineers/researchers should be able to *easily* reconfigure the entire FL stack (topology, communication, engine, algorithm, privacy)



Tyagi, S., Cozma, A., Kotevska, O., & Wang, F. (2025). *OmniFed: A Modular Framework for Configurable Federated Learning from Edge to HPC*. SC25-W: Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis, 516-523.