

Advances in APPFL: A Comprehensive and Extensible Federated Learning Framework

Zilinghan Li*, Shilan He[†], Ze Yang[†], Minseok Ryu[‡], Kibaek Kim*, Ravi Madduri*

*Argonne National Laboratory [†]University of Illinois at Urbana-Champaign [‡]Arizona State University
{zilinghan.li, kimk, madduri}@anl.gov, {shilanh2, zeyang2}@illinois.edu, minseok.ryu@asu.edu

Abstract—Federated learning (FL) is a distributed machine learning paradigm enabling collaborative model training while preserving data privacy. In today’s landscape, where most data is proprietary, confidential, and distributed, FL has become a promising approach to leverage such data effectively, particularly in sensitive domains such as medicine and the electric grid. Heterogeneity and security are the key challenges in FL, however; most existing FL frameworks either fail to address these challenges adequately or lack the flexibility to incorporate new solutions. To this end, we present the recent advances in developing APPFL, an extensible framework and benchmarking suite for federated learning, which offers comprehensive solutions for heterogeneity and security concerns, as well as user-friendly interfaces for integrating new algorithms or adapting to new applications. We demonstrate the capabilities of APPFL through extensive experiments evaluating various aspects of FL, including communication efficiency, privacy preservation, computational performance, and resource utilization. We further highlight the extensibility of APPFL through case studies in vertical, hierarchical, and decentralized FL. APPFL is open-sourced at <https://github.com/APPFL/APPFL>.

Index Terms—Federated Learning, Distributed Computing, Benchmarking, Privacy Preservation, Scheduling Algorithms

I. INTRODUCTION

Availability of extensive training data is becoming increasingly crucial for developing more capable machine learning (ML) models, especially as these models continue to grow in size and complexity. Nonetheless, most of the data in today’s landscape is confidential and distributed across various data silos [1]. This distribution makes it difficult to collect the data for centralized model training, posing significant challenges in fully leveraging the existing data to train more powerful ML models. In this context, federated learning (FL), a distributed ML paradigm, offers a promising solution to utilize data from multiple data owners without direct data sharing [2]–[5].

In FL, multiple data owners, referred to as clients, collaborate under a central server to train a shared ML model by iterating the following two steps: (1) each client trains an ML model using its local dataset and submits the updated model to the server, and (2) the server aggregates these local models to update the global model and then sends it back to the clients for further local training. In this way, FL leverages data from multiple data sources to build a more powerful and robust model without requiring data centralization, thereby protecting data privacy. Consequently, FL has been widely adopted in domains such as medicine [6]–[8], finance [9], [10], and electric grid [11], where data privacy is paramount.

Depending on the amount, capability, and availability of client devices, FL is broadly categorized into two types, cross-device FL and cross-silo FL [5]. In cross-device FL, numerous mobile or IoT devices with limited computing power and intermittent availability collaboratively train relatively small models such as keyboard suggestion and hot word detection models [12]–[14]. In contrast, cross-silo FL involves fewer but more reliable and powerful clients, typically represented by large data silos and institutions, to develop more complex ML models with extensive parameters.

While FL can be conceptually simplified to traditional machine learning with an additional global aggregation operation, its distributed nature introduces significant challenges in terms of heterogeneity and security when adopted for real-world applications. Data heterogeneity, stemming from the unbalanced, and non-independent and identically distributed (non-IID) nature of client local datasets, can lead to varied local training objectives across clients and potentially degrade the performance of the global model [15]. Additionally, the heterogeneity in computation and communication, caused by diverse computing capabilities and network connectivity of client devices, can severely impact the efficiency of FL training [16], [17]. This is particularly problematic in synchronous FL algorithms, where the server has to wait for all clients to submit their local models before global aggregation. With regard to security, FL is vulnerable to various attacks. Untrusted clients might maliciously attack FL experiments by submitting corrupted local models, and there is also a risk that training data could be reconstructed from the model updates sent by clients, thereby compromising data privacy [5], [18].

Most existing FL frameworks, such as FLOWER [19], FEDML [20], and FEDSCALE [21], do not adequately address the full spectrum of challenges inherent in FL. For example, some frameworks do not support asynchronous aggregation that could improve training efficiency, lack implementations of robust authentication and privacy preservation algorithms, or fail to offer user-friendly interfaces for developers to easily integrate new solutions or algorithms. To bridge these gaps, we developed the Advanced Privacy-Preserving Federated Learning (APPFL) framework, a comprehensive and extensible FL framework that builds on and improves the work presented in [22]. APPFL supports both single-node and multi-node simulations, as well as distributed deployment of FL experiments. It features advanced aggregation strategies to address data heterogeneity [15] and various asynchronous aggregation

strategies to boost training efficiency in environments with computation heterogeneity [23]. Additionally, APPFL incorporates versatile communication protocols [24], data transfer methods [25], and compression strategies to meet different communication requirements and enhance communication efficiency [26]. It also includes robust authentication via Globus [27]–[29], along with plugins for adding new authentication methods, and implements privacy preservation strategies to prevent the reconstruction of training data [30]. Moreover, APPFL is extensible; it follows a modular design that enables users and developers to seamlessly adapt the framework for different use cases and integrate custom algorithmic solutions to tackle various FL challenges.

The contributions of this work are outlined as follows:

- Advance APPFL, an open-source FL framework for both FL users and developers that provides established solutions to common FL challenges for FL users and offers flexible and modular interfaces that facilitate easy integration of new algorithmic solutions for FL developers
- Conduct comprehensive evaluations of various aspects of FL using APPFL, including the efficiency of the versatile communication protocols, data transfer methods, and compression strategies, as well as the performance of privacy preservation strategies and training effectiveness of different FL aggregation algorithms
- Provide case studies in vertical, hierarchical, decentralized FL to highlight the extensibility and adaptability of the APPFL framework in diverse FL scenarios

II. BACKGROUND AND RELATED WORK

A. Heterogeneity in Federated Learning

Heterogeneity is one of the significant challenges in FL due to its distributed nature across multiple clients [5], [31], [32]. This heterogeneity can be categorized into three primary types: data heterogeneity, computation heterogeneity, and communication heterogeneity.

Data heterogeneity arises from the fact that client datasets are often unbalanced and non-IID, meaning they may not be representative of the overall population distribution. This discrepancy leads to varying local training objectives among clients, causing their locally trained models to diverge from one another, a phenomenon known as client drift [33]. As a result, simple weighted averaging of local models, as in the FedAvg strategy [3], may degrade the performance of the global model as data heterogeneity increases [15]. Several solutions have been proposed to address this issue on both the server and client sides. For instance, server-side optimizations such as FedAvgM [15], FedAdam, FedAdagrad, and FedYogi [34] have been introduced to enhance FL performance on non-IID data. Client-side approaches such as FedProx [16] and SCAFFOLD [33] incorporate correction terms into the client’s local objective function to reduce excessive drift between local and global models.

Computation heterogeneity occurs when the computing devices of FL clients have varying computing power, resulting

in large variants in the local training times [35]. This variance poses challenges for synchronous FL strategies, where the server must wait for all clients to submit their local models before proceeding with global aggregation. Such delays from slower clients can reduce training efficiency and lead to underutilization of computing resources. In order to address this issue, various asynchronous aggregation strategies have been proposed. These strategies, including FedAsync [36], FedASO [37], FedBuff [38], AREA [39], and FedCompass [23], update the global model immediately upon receiving models from one or a few clients. These methods are particularly beneficial in environments with heterogeneous computing capabilities because they minimize client idle time and enhance resource utilization. Other approaches include disregarding contributions from straggler clients [40] or explicitly selecting clients for local training based on their computing capabilities [41]–[43]. Nonetheless, these methods are best suited for cross-device FL, where only a subset of clients participates in each training round. They do not align well with cross-silo FL where there are only a few FL clients and ensuring the participation of every client is usually vital for maintaining the robustness of the global model.

Communication heterogeneity is originally rooted in the intermittent availability of client devices due to the limited power and bandwidth in cross-device FL settings, which is less of an issue in cross-silo FL. However, as foundation models increasingly dominate various domains and applications, the interest in using FL to train or fine-tune these models has surged. This surge has led to a substantial increase in communication costs, which become a critical factor affecting the efficiency of FL training [44]–[46]. Consequently, improving the efficiency and robustness of transferring large model parameters has become critically important as well [47]. To address this situation, some client selection methods have been proposed to mitigate communication issues in cross-device FL by strategically selecting clients based on their availability, data quality, and performance [48]–[51]. Other approaches focus on generic FL settings by applying compression techniques to large model parameters, thereby reducing the overall communication workload [26], [52], [53].

B. Attacks and Security Concerns in Federated Learning

The distributed and uninspectable nature of FL exposes it to various adversarial attacks and security risks. These attacks generally fall into two broad categories: (1) inferring clients’ confidential training data from the model gradients and (2) degrading the performance of the trained global model [5].

Gradient inversion algorithms, for example, can reveal information about the private training data by iteratively updating a randomly initialized sample to match its gradient update with the actual model gradient. These algorithms are particularly effective in the early stages of training where the gradients contain more information about the training data [60]–[63]. Countermeasures to these inversion attacks include increasing training batch sizes [62], implementing differential privacy techniques to add noise to model gradients [8], [18], compress-

TABLE I: Comparison of popular open-source federated learning frameworks.

Framework	Data Hetero.	Sync. FL	Async. FL	Compression	Versatile Comm.	Privacy	Auth.	Real Deployment	FL Variants
LEAF [54]	✗	✓	✗	✗	✗	✗	✗	✗	✗
TFF [40]	✓	✓	✗	✗	✗	✓	✗	✗	✗
APPFL-V0 [22]	✓	✓	✗	✗	✗	✓	✗	✓	✗
FEDERATEDSCOPE [55]	✓	✓	✗	✗	✗	✓	✗	✓	VFL
FLARE [56]	✓	✓	✗	✗	✗	✓	✓	✓	VFL
OPENFL [57]	✓	✓	✗	✓	✗	✓	✓	✓	VFL
FEDSCALE [21]	✓	✓	✓	✓	✗	✓	✗	✓	✗
FLGO [58]	✓	✓	✓	✗	✗	✗	✗	✓	VFL
FEDLAB [59]	✓	✓	✓	✓	✗	✗	✗	✓	✗
FLOWER [19]	✓	✓	✗	✗	✓	✓	✓	✓	VFL
FEDML [20]	✓	✓	✗	✗	✓	✓	✓	✓	VFL, HierFL, DFL
APPFL (this work)	✓	✓	✓	✓	✓	✓	✓	✓	VFL, HierFL, DFL

ing model gradients [64], [65], and employing cryptographic strategies such as homomorphic encryption to secure the gradients [66].

Since the FL server cannot access and inspect the client training data, FL is also vulnerable to attacks from Byzantine clients [67], which may either submit corrupted model parameters (model poisoning) or use tampered data for training (data poisoning) [68]–[73]. To counter these, several algorithms have been developed to exclude models whose parameters significantly deviate from the norm [70], [74], [75]. Other approaches use the client models to inversely derive training data and then exclude models if their derived data significantly diverges from other models. Alternatively, some solutions assume that the FL server holds a clean and secret validation dataset to evaluate and potentially exclude poorly performing client models from the aggregation process [76], [77].

Beyond algorithmic defenses, addressing malicious attacks in FL can also be achieved through system-level enhancements, particularly by integrating with identity and access management (IAM) services [78]. Such integration enables the creation of secure federations that permit only trusted and known clients to participate in FL experiments, thus fundamentally resolving security concerns at their root.

C. Existing Federated Learning Frameworks

We conduct a brief survey of twelve popular open-source federated learning frameworks, including the previous version of APPFL without advancements [22] (APPFL-V0), focusing on their approaches to heterogeneity and security challenges, usability, and extensibility for various application scenarios. The results are summarized in Table I.

In addressing data heterogeneity, most frameworks implement advanced client training and server aggregation strategies to mitigate the client drift issue, with the exception of LEAF [54], one of the earliest FL frameworks. Regarding computation heterogeneity, while all frameworks include synchronous FL algorithms, only FEDSCALE [21], FLGO [58], FEDLAB [59], and APPFL offer asynchronous communication stack and the corresponding asynchronous aggregation strategies. For communication heterogeneity concerns, we evaluate whether the frameworks feature lossless or lossy compression algorithms to reduce the communication loads and whether they

provide versatile communication stacks that support multiple protocols, enhancing efficiency and adaptability to different deployment requirements and scenarios.

In terms of the privacy and security challenges, our investigation focuses on whether the frameworks incorporate privacy preservation algorithms such as differential privacy, and integrate IAM services for user authentication and authorization. Most existing frameworks support privacy preservation to some extent, with FLARE [56], OPENFL [57], FLOWER [19], FEDML [20], and APPFL featuring IAM integration for verifying user identities and managing access to specific FL experiments.

As for the usability, most of the frameworks facilitate both the simulation of FL experiments within the same machine or cluster and the real deployment among distributed clients. However, LEAF and TFF [40] are limited to simulation environments only. To evaluate the extensibility and customization capabilities of the frameworks, we assess their support for different FL variants beyond the traditional federated learning, specifically vertical federated learning (VFL) [79], hierarchical federated learning (HierFL) [80], and decentralized federated learning (DFL) [81]. FEDERATEDSCOPE [55], FLARE, OPENFL, FLGO, and FLOWER provide use cases in VFL settings, and FEDML and APPFL extend support to all three variants.

III. ARCHITECTURE AND IMPLEMENTATION

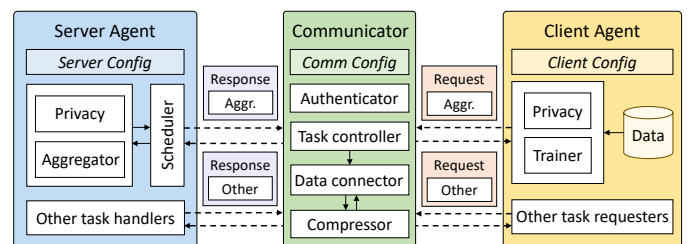


Fig. 1: Overview of the APPFL framework’s new software architecture design. *Server agent* and *client agent* act on behalf of the FL server and client, respectively, to fulfill various tasks for FL experiments. *Communicator* exchanges task control signals and model parameters between the server and client.

The APPFL framework is a Python package available on PyPI. Figure 1 provides an overview of its new software architecture. APPFL defines a *server agent* and a *client agent*, connected by the *communicator*, to represent the FL server and clients in performing the primary aggregation task and other necessary tasks for running FL experiments. The *server agent* is mainly composed of a scheduler module that orchestrates the aggregation of client local models under various synchronicity settings, an aggregator module that aggregates the local models passed from the scheduler to update global model, and a privacy module for additional privacy protection. The *client agent* consists of a trainer module responsible for training the ML model using the confidential local dataset and a privacy module for the privacy preservation algorithms. The *communicator* facilitates robust communication between the server and clients, supporting multiple communication protocols for exchanging task control signals and data, with an option to separate the transmission of control signals and data via a data connector. Additionally, the communicator incorporates several compressors for improved efficiency and authenticators for securing the FL experiments. Overall, APPFL incorporates solutions for various challenges in FL and is designed to be modular and extensible, facilitating easy integration of new algorithms and strategies to address FL challenges. The following subsections detail the key components of APPFL.

A. FL Experiment Configuration

APPFL provides a straightforward way to configure FL experiments: each experiment utilizes a configuration YAML file for the FL server and individual YAML files for each FL client. Listing 1 presents an example of the server configuration file, which includes server-specific settings, such as the aggregation algorithm and the number of global epochs, along with general configurations for the clients like the trainer and compressor types. These general configurations are distributed to all clients at the beginning of each FL experiment, simplifying the setup by ensuring that shared configuration fields do not need to be individually set by each client. In addition to the configurations shared by the server, each client possesses its own YAML configuration file that defines client-specific settings, as shown in Listing 2. The client-specific settings include a Python loader file, which defines a function for loading the client’s local datasets, and some training-related configurations such as the device to use and directories for logging and checkpoints.

This configuration also facilitates integration of new algorithms by allowing developers to directly add necessary settings to the relevant configuration files and use them in their respective module blocks. For instance, to create a trainer for a particular application, a developer simply needs to define a new trainer within the APPFL trainer module and include all necessary arguments in the *client_configs.train_configs* section of the configuration file.

```

1 # Server configurations
2 server_configs:
3   aggregator: FedAvgAggregator
4   num_global_epochs: 10
5   ...

```

```

6 # General client configurations for all clients
7 client_configs:
8   train_configs:
9     trainer: VanillaTrainer
10    lr: 0.001
11    ...
12   comm_configs:
13     compressor_configs:
14       lossy_compressor: SZ2Compressor
15     ...

```

Listing 1: An example server configuration YAML file, containing both server configurations and general client configurations to be shared among all clients.

```

1 # Information needed to load local data
2 data_configs:
3   dataset_path: ./dataset/covid_dataset.py
4   dataset_name: get_covid #function to load data
5   dataset_kwargs: #optional function arguments
6   ...
7 # Client-specific training settings
8 train_configs:
9   device: cpu
10  logging_dir: ./appfl_logging
11  checkpoint_dir: ./appfl_checkpoint
12  ...

```

Listing 2: An example client configuration YAML file, containing client-specific configurations such as the data loader file.

B. Communication Stack

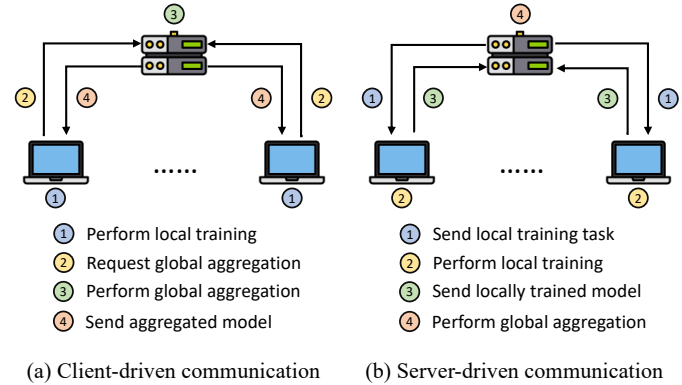


Fig. 2: Running one local training and global aggregation iteration using (a) client-driven and (b) server-driven communication protocols.

In FL, communication protocols can be broadly classified into two types based on the driven side of the FL process: (1) *client-driven*: the clients control the FL process and interact with the server for aggregation and other tasks by sending various requests and (2) *server-driven*: the server controls the FL process by dispatching various types of tasks to the clients. Figure 2 illustrates the differences between these two types of communication protocols during a local training and global aggregation FL iteration. Client-driven protocols offer clients greater autonomy over the FL process, whereas server-driven protocols simplify the coordination of FL experiments, with the central server itself managing the whole distributed training process. The APPFL communicator supports MPI and gRPC as client-driven communication protocols and Globus

Compute [24] as the server-driven protocol. Specifically, MPI is for simulation purposes only, while gRPC and Globus Compute can be used for real deployments. Notably, gRPC requires the server to open a specific port for inbound TCP connections, which is typically restricted in high-performance computing environments and institutional computing facilities. Conversely, Globus Compute only necessitates outbound connections to the Globus service, thus enabling a broader range of computing resources to serve as the FL server. The versatile communication protocols supported by APPFL make it capable of meeting diverse communication needs in FL deployments.

For client-driven communication protocols, the APPFL communicator provides a server communicator that defines handlers for various types of requests, such as sending general client configurations and performing global aggregation, by interacting with the server agent. Additionally, a client communicator assists the client agents in sending requests to the server. As for Globus Compute, the server-driven communication protocol, it is a distributed function-as-a-service platform that can dispatch Python functions to run on remote machines. The APPFL communicator provides a Globus Compute server communicator to send various tasks, such as local training, to run on the remote client machines and collect results back for conducting FL experiments. Overall, APPFL supports commonly used server request handlers and client task implementations and provides a user-friendly interface that enables developers to easily define new request handlers or tasks without in-depth knowledge of the underlying communication protocol.

While the communication protocols can transfer the task control signals (i.e., requests in client-driven and tasks in server-driven protocols) along with the associated data, APPFL provides an option to separate the transfer of task controls from the associated model parameters through the integration with ProxyStore [25]. ProxyStore can create a *proxy* for any target Python object, providing a lightweight reference that can remotely resolve the target object when used. When the *proxy* is resolved, the object is transferred via an underlying data connector. APPFL currently supports two connectors: an S3 connector, which uses AWS S3 buckets for data transfer, and a ProxyStore endpoint connector, which employs the ProxyStore-hosted relay server. The integration with ProxyStore offers two main benefits: (1) it prevents exceeding the maximum data size limits imposed by certain communication protocols (e.g., Globus Compute restricts task arguments and result sizes to 10 MB to reduce server load, thus making data transfer separation a must when exchanging large model parameters), and (2) it offers users a variety of data transmission options for different communication scenarios and facilitates easy integration of other efficient data transmission methods suitable for their specific use cases to accelerate the FL communication, regardless of the communication protocol in use.

Furthermore, APPFL incorporates a range of data compressors to enhance communication efficiency, crucial for transferring parameters of large models or operating in environments

with limited network bandwidth. It supports various lossless compressors including zstd [82], gzip [83], and blosc [84], as well as lossy data compressors including SZ2 [85], SZ3 [86], and ZFP [87]. These compressors can help reduce the communication load, enabling faster data transfer between the server and clients.

C. Server Scheduling and Aggregation

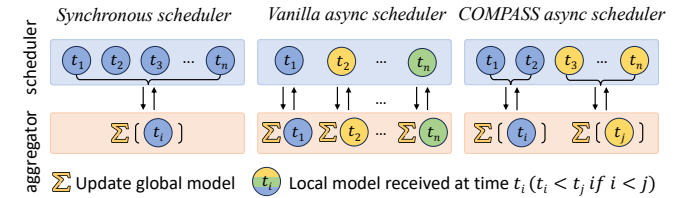


Fig. 3: Scheduling of the aggregation for client local models under three schedulers with different synchronicity settings.

In order to tackle the computation heterogeneity in FL where clients have varying computing capabilities, many asynchronous aggregation algorithms have been proposed to reduce client idle times and enhance resource utilization. To support aggregation with different synchronicity settings, APPFL introduces a server-side scheduler that acts as an interface between the communicator and the aggregator. Upon receiving a local model from a client, the communicator forwards it to the scheduler, which determines the appropriate time to pass the local model(s) to the aggregator for updating the global model. For synchronous aggregation strategies, such as FedAvg [3], a synchronous scheduler buffers each client’s local model until all models are received, at which point it forwards them to the aggregator to update the global model. Conversely, for asynchronous strategies like FedAsync [36], a vanilla asynchronous scheduler immediately sends the client model to the aggregator and returns the updated global model back to the communicator. Additionally, the scheduler module is designed to be extensible for the incorporation of more advanced scheduling algorithms. Specifically, APPFL supports the state-of-the-art Compass asynchronous scheduler [23], aiming to alleviate the drift of the global model toward faster clients. Such drift is prevalent in other asynchronous FL algorithms, where faster clients update the global model more frequently and the models from slower clients become stale. The Compass scheduler synchronizes the arrival of a group of client local models by assigning different amounts of local training tasks to different clients to enable a grouped global aggregation and avoid stale local models, mitigating the client drift issue. Figure 3 illustrates the scheduling processes under the three different schedulers.

As for the aggregator module, APPFL supports a broad range of aggregation strategies, going beyond the widely used FedAvg. These include FedAvgM, FedAdam, and FedYogi, which address data heterogeneity, PLFL [88] for personalized FL, as well as ICEADMM [89] and IIDADMM [22], which focus on efficient privacy preservation. Additionally, for asynchronous aggregation, APPFL includes strategies such as

FedAsync, FedBuff, AREA [39], and FedCompass. This diverse suite of options ensures that APPFL can accommodate a variety of needs and scenarios in FL, illustrating its adaptability and comprehensive approach to FL challenges.

D. Privacy Preservation and Authentication

To tackle security and privacy concerns in FL, APPFL offers solutions that span both algorithmic and system-level measures. Algorithmically, APPFL incorporates the differential privacy (DP) algorithms [30] into the FL process that perturbs the client model parameters with noises before sending to the server, protecting against the reconstruction of confidential training data. A study utilizing APPFL showcases that the usage of DP in FL can effectively mitigate the risk of data reconstruction [8].

At the system level, APPFL enhances security through the integration of identity and access management (IAM) services into its communication stack for user authentication and access control for FL experiments. Specifically, Globus Compute itself is already integrated with the Globus authentication service, ensuring that the server dispatches training functions only to clients within a specified Globus group. This setup helps create a secure federation of trusted collaborators, authenticated via institutional emails linked to Globus accounts.

As for gRPC, APPFL utilizes token-based authenticators to verify users. Clients have to attach an access token to each remote procedure call (RPC) request over an SSL-encrypted channel, allowing the server to confirm the user’s identity before processing the request. The token-based authenticator consists of two primary functions: one invoked by the client to generate the token prior to sending the RPC request and another invoked by the server to verify the validity of the token upon receipt. This straightforward interface allows developers to effortlessly integrate their own authentication methods tailored to specific use cases and applications. Currently, APPFL supports a Globus authenticator, with its login flow depicted in Figure 4. Users can employ APPFL’s command line interface (CLI), `appfl-auth`, to perform a one-time login. Depending on the selected role during login, either as an FL server or client, the appropriate Globus access token (Group Service or Identity Service) is requested. The access tokens, along with the corresponding refresh tokens, are securely stored in the client’s local token storage. Whenever an FL client makes an RPC request, it attaches its Globus Identity Service token. The FL server uses this token to retrieve the client’s Globus ID and, leveraging its Globus Group Service token, verifies whether the client belongs to the specified Globus group. This robust authentication process ensures a secure and controlled federation for FL experiments.

IV. PERFORMANCE EVALUATION

In this section, we employ APPFL to comprehensively benchmark a broad spectrum of aspects within FL to highlight the capabilities of its new software design. Specifically, we utilize APPFL to evaluate the communication efficiency of different communication protocols, data transfer methods,

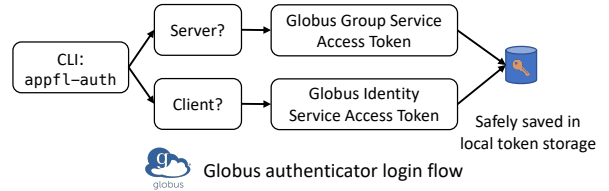


Fig. 4: Login flow for the Globus authenticator.

and compression algorithms. We also explore the impacts of privacy preservation algorithms on the performance of FL-trained models, as well as the training efficiency and resource utilization of various FL strategies under different synchronicity settings. These comprehensive assessments highlight the capabilities of the newly designed APPFL framework in benchmarking a broad spectrum of aspects within FL.

A. Communication Efficiency

We evaluate the communication efficiency for different communication and data transfer protocols across different numbers of clients and various model sizes, ranging from a few-byte 1×1 fully connected (FC) layer to a Vision Transformer (ViT) with hundreds of megabytes. Table II details the sizes of all models used in our experiments.

TABLE II: Sizes of the models used in the experiments.

Model	Params	Size
1×1 FC	2	8 B
CNN [90]	1.20M	4.58 MB
ResNet18 [91]	11.17M	42.66 MB
ResNet50 [91]	23.52M	89.93 MB
ResNet101 [91]	42.51M	162.58 MB
Vision Transformer [92]	88.22M	336.55 MB

In the experiments, each FL client runs on a single core of the CPU-only compute nodes of the Delta supercomputer at the National Center for Supercomputing Applications. Each CPU node contains two 64-core AMD EPYC 7763 “Milan” CPUs with PCIe Gen4 interfaces and 256 GB of RAM. Delta is connected to the NPCF core router and exit infrastructure via two 100 gigabits per second (Gbps) connections. The FL server is hosted on an AWS EC2 `x2iedn.2xlarge` instance, equipped with 8 virtual CPUs, 256 GB of RAM, and up to 25 Gbps connections. We exponentially increase the number of clients from 2 to 128 across all models, except for the ViT model, which scales only from 2 to 64 because of memory constraints on the client and server hardware. We evaluate gRPC and Globus Compute communication protocols as well as two data transfer methods, AWS S3 buckets and ProxyStore endpoints. Because of the 10 MB data transfer limit with Globus Compute, it is integrated with other data transfer protocols rather than being tested in isolation, resulting in five distinct communication pattern combinations.

Figure 5 shows the epoch-wise two-way communication time in seconds for various models using different communication and data transfer protocols. From the plots, we note

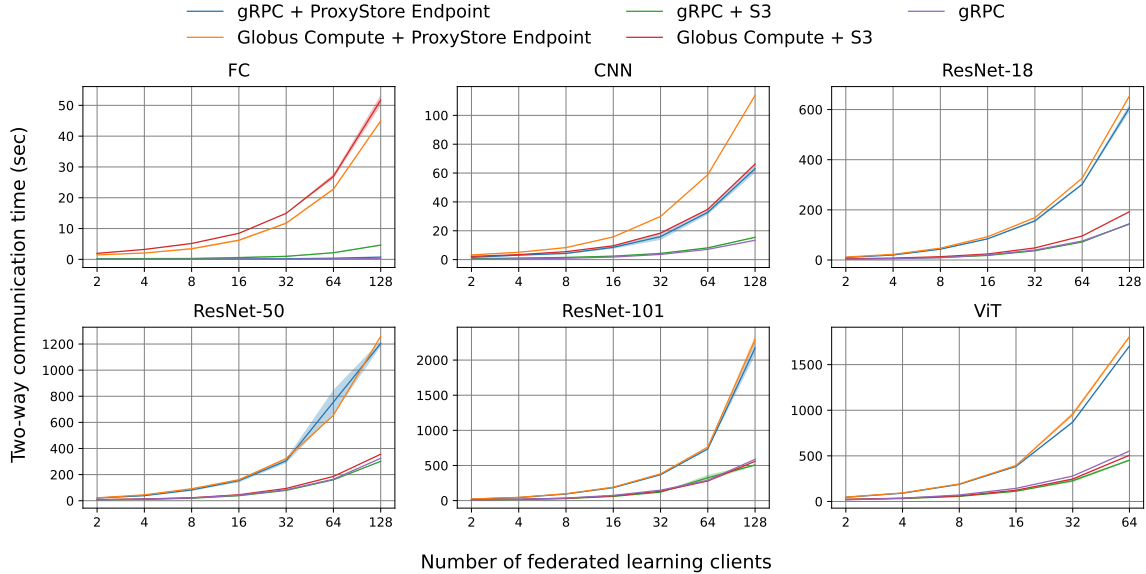


Fig. 5: Efficiency comparison of communication and data transfer protocols: Average two-way communication time per global epoch and the corresponding standard deviation as the number of clients increases exponentially across various models.

the following key points: (1) Separating the transmission of data (i.e. model parameters) from task control signals helps communication protocols exceed their maximum data size limitations. (2) Globus Compute consistently incurs longer overheads than does gRPC in transmitting control signals, which is a significant factor when the FL model size is small. (3) While data transfer via ProxyStore endpoints generally results in longer communication times, it offers a free and straightforward solution for protocols such as Globus Compute that have message size restrictions. (4) Data transfer through S3 features relatively low latency and also provides a secure, reliable means to store model checkpoints during training, although it incurs some costs.

B. Compression Efficiency

We assess the efficiency of various data compression algorithms integrated within the APPFL framework. Specifically, we utilize the lossless compressor blosc [84] for tensors with less than 1,024 parameters and lossy compressors SZ2 [85], SZ3 [86], and ZFP [87], each with a relative error bound of 0.01, for larger tensors. The experiments adhere to the same hardware configurations described in Subsection IV-A, with clients running on Delta and the server on an AWS EC2 instance. We conduct experiments on ResNet-50 and ResNet-101 models, scaling client numbers from 2 to 128. Figure 6a illustrates the reduction in model sizes by 3 to 5 times using different lossy compressors. Notably, previous studies have shown that such levels of lossy compression can preserve model accuracy within a 0.5% margin of uncompressed results [26]. Figure 6b presents the two-way communication times via gRPC for the two models using various compressors. Solid lines represent times with compression and decompression overheads, whereas dotted lines depict times without. The

comparison reveals significant overhead, particularly with SZ2 and SZ3. Despite this, the use of compressors notably reduces communication costs and overall two-way communication times, even under high-bandwidth conditions for both clients and the server.

C. Privacy Preservation

In this subsection we study the impact of differential privacy (DP) techniques on the performance of models trained via FL. We select four tasks in medical domains, where data privacy is paramount, from the FLamby benchmark containing naturally split medical datasets [93]. Table III provides an overview of these tasks. We assess model performance across varying values of privacy loss parameter ϵ , a measure of how much privacy is lost when using DP algorithms, with lower ϵ values signifying larger added noises and enhanced privacy. Figure 7 shows the change of model performance throughout the FL training process for these tasks at different ϵ values. The performance metrics represent the average outcomes of five independent trials with different random seeds. The results indicate that a decrease in ϵ values, corresponding to increased privacy preservation, leads to varying degrees of performance degradation across various models and training tasks.

TABLE III: Overview of selected tasks from FLamby.

	Fed-TCGA-BRCA	Fed-Heart-Disease	Fed-IXI	Fed-ISIC2019
Input	Patient info	Patient info	T1WI	Dermoscopy
Prediction	Risk of death	Heart disease	Brain mask	Melanoma class
Task type	Regression	Classification	3D Segmentation	Classification
Model	Cox model [94]	Logistic Reg.	3D U-Net [95]	EfficientNet [96]
Metric	C-index	Accuracy	DICE [97]	Balanced Acc.
# Clients	6	4	3	6

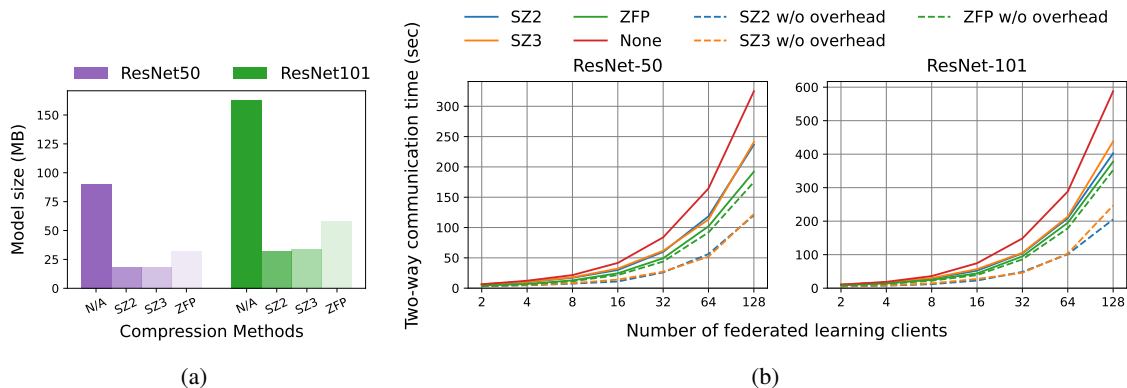


Fig. 6: (a) Model sizes for ResNet-50 and ResNet-101 using different lossy compression methods with a relative error bound of 0.01. (b) gRPC two-way communication time for ResNet-50 and ResNet-101 using different lossy compressors, with and without the compression and decompression overhead.

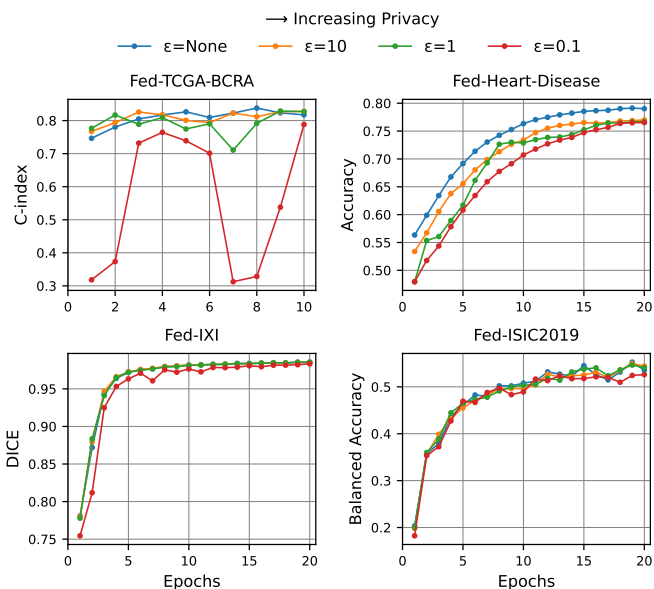


Fig. 7: Change of model performance throughout the FL training on the selected FLamby tasks at different ϵ values.

D. Addressing Heterogeneous Clients

In this subsection we evaluate the performance and efficiency of different FL algorithms under various synchronicity settings. Specifically, we benchmark five FL algorithms: (1) FedAvg, a widely used synchronous algorithm that updates the global model by averaging all client local models; (2) FedAvgM, another synchronous algorithm, which incorporates momentum on top of FedAvg; (3) FedAsync, which asynchronously updates the global model upon receipt of any local model; (4) FedBuff, which is similar to FedAsync but buffers multiple local models before updating the global model; and (5) FedCompass, which introduces a *COMputing Power Awareness Scheduler* (Compass) that dynamically adjusts the number of client local training steps based on real-time estimates of client computing power to synchronize the

training completion for groups of clients. As for the datasets, we partition the CIFAR-10 dataset [98], one of the most commonly used datasets in evaluating FL algorithms [99], into ten client splits in a non-IID manner, with each client holding data from five to seven classes out of ten classes. All clients use Nvidia A100 GPUs for training, and we simulate a group of heterogeneous clients by assigning different average batch processing times from an exponential distribution.

Figure 8a presents the average validation accuracy and the corresponding standard deviation across five independent runs for each FL algorithm during training. Key observations from the figure include the following. (1) Asynchronous FL algorithms like FedAsync and FedBuff, which use the vanilla asynchronous scheduler, converge to significantly lower global model accuracy compared with synchronous methods, primarily due to the drifting toward faster clients, as the global model gets more updates from faster clients and slower clients' models become stale. (2) Synchronous algorithms exhibit slower convergence as the server has to wait for the slow clients for aggregation. (3) FedCompass effectively addresses substantial client drift issues and attains high global model accuracy by ensuring nearly simultaneous model arrivals for grouped aggregation. It also achieves quicker convergence than synchronous methods without extensive waiting.

Figure 8b shows the average training time per batch for the ten clients involved in the FL training, as well as the resource utilization, calculated as the ratio of client compute time to total training time, for algorithms using the synchronous, vanilla asynchronous, and Compass asynchronous scheduler. The synchronous scheduler shows the lowest client resource utilization, correlating with training time per batch: the quicker the client, the lower the utilization. In contrast, the vanilla asynchronous scheduler, which immediately sends any received local model for aggregation and returns the updated global model, allows client resource utilization to approach 100%. Despite full utilization, however, this method results in poorly performing models due to client drift. The Compass scheduler, by estimating client speeds and adjusting train-

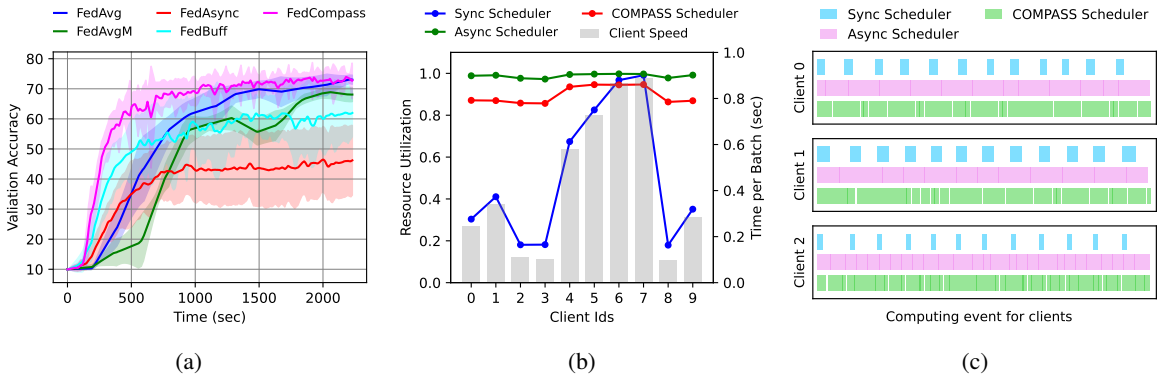


Fig. 8: (a) Average validation accuracy and the corresponding standard deviation on the partitioned CIFAR-10 dataset for different FL algorithms during the training process. (b) Client resource utilization for algorithms using different schedulers and the average training time per batch for different clients. (c) Visualization of computing resource utilization for three clients under different schedulers, where colored bar represents the computing period and the blank places means idle times.

ing steps accordingly, maintains approximately 90% resource utilization and reduces client drift through timely grouped aggregations. Figure 8c visualizes the resource utilization for three clients under different scheduling scenarios, highlighting the significant resource underutilization of the synchronous scheduler compared with the asynchronous alternatives when client computing resources vary widely.

V. CASE STUDY: EXTENSIBILITY DEMONSTRATION

To highlight the versatility and extensibility of the APPFL framework across various FL applications, we present case studies on three distinct FL variants: vertical FL, hierarchical FL, and decentralized FL, all built upon the APPFL framework, illustrating how it can be adapted to different FL paradigms.

A. Vertical Federated Learning

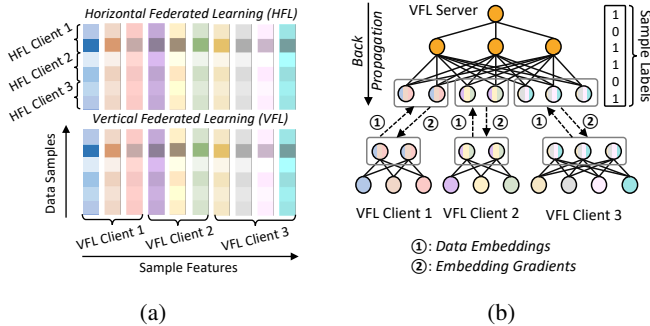
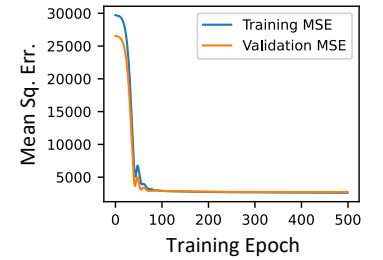


Fig. 9: (a) Comparison of client training data distribution in HFL and VFL. (b) Overview of the VFL process.

Vertical federated learning (VFL) is a specialized paradigm of FL where different clients hold distinct features from the same dataset [79]. Unlike traditional FL (i.e., horizontal FL) dealing with the same feature space across diverse data samples, VFL enables collaboration among clients that have partially overlapping or non-overlapping features but share the same sample IDs, as illustrated in Figure 9a. Figure 9b

Dimension	In	Hidden	Out
Client 1	3	24	8
Client 2	3	24	8
Client 3	4	24	8
Server	24	12	1

(a)



(b)

Fig. 10: (a) Input, hidden, and output dimensions of two-layer perceptrons for the VFL clients and server. (b) Training and validation MSE during the VFL training process.

depicts a typical VFL process. In VFL, rather than training the same model architecture and sharing model parameters, each client possesses its embedding model to process its local data sample features and then sends their embeddings to the server. The server, holding the labels of the client data samples, concatenates the received embeddings to train a central model. It then sends the gradients of the feature embeddings back to the corresponding clients, enabling them to update their local embedding models accordingly.

APPFL seamlessly supports VFL by providing the VFL trainer and aggregator in the corresponding modules. In this case study, we use the diabetes datasets [100] from the `scikit-learn` library, which contains ten features of 442 data samples. The labels, ranging from 25 to 346, are the responses of interest that quantitatively measure the disease progression. We split the dataset into 80% for training and 20% for validation and use three VFL clients, where clients 1 and 2 possess three patient features and client 3 possesses four. Each of the three clients as well as the server employs a two-layer perceptron with ReLU nonlinear activation [101] as their embedding models. Figure 10a presents the input, hidden, and output dimensions of these models. During the training,

the server model is updated based on the mean squared error (MSE) loss between the labels and predictions, using the Adam optimizer [102] with a learning rate of 0.01. Figure 10b shows the training and validation MSE throughout the training.

B. Hierarchical Federated Learning

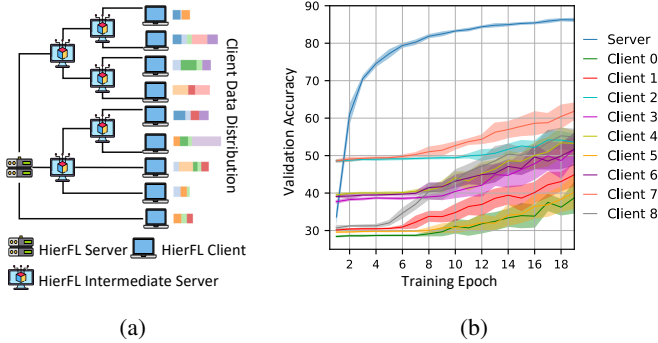


Fig. 11: (a) Topology of the multi-layer HierFL experiments with nine clients, five intermediate servers, and one root server. (b) HierFL validation accuracy for the server and client models, where the accuracy of client models is evaluated after each local training step.

Hierarchical federated learning (HierFL) is also a special type of FL that introduces an additional role, the intermediate server (edge server). This server first aggregates local model parameters from connected clients or child intermediate servers and then forwards the aggregated model to the parent server for further aggregation [80]. HierFL is particularly beneficial when FL clients are geographically clustered, since placing an intermediate server for these clusters can significantly improve overall communication efficiency. To support HierFL in APPFL, in addition to the general server agent for the root server and the client agent for the clients, we define an intermediate server agent, similar to the server agent, which handles FL-related requests from connected clients or child intermediate servers by interacting with its parent server.

In this case study we conduct four-tier HierFL experiments involving nine clients, five intermediate servers, and one root server. The MNIST training dataset [103] is partitioned into nine heterogeneous client splits, with each client containing training data for only three to five classes. Figure 11a illustrates the topology of the HierFL experiments and the client data distribution. Training is conducted over 20 global epochs, with each client performing 100 local steps per epoch using a batch size of 64 and the Adam optimizer with a learning rate of 0.001. The experiments are repeated five times with different random seeds. Figure 11b presents the average validation accuracy and standard deviation for both the server model and each client’s local model on the MNIST validation set. We note that the client models are evaluated after local training. Since each client has data for only three to five classes, their local models perform significantly worse than the global model, highlighting the advantages of federated learning in leveraging data from distributed clients to train a more robust ML model.

C. Decentralized Federated Learning

Decentralized federated learning (DFL) is yet another variant of FL that eliminates the need for a central server to aggregate models. Instead, each node trains its local model using its own data, requests model parameters from neighboring clients, and aggregates these with its local model [81], [104]. APPFL supports DFL by implementing a DFL node agent that inherits most functionalities of both an FL client and server, enabling it to train local models and handle requests from neighboring clients. In this case study we set up DFL experiments with six nodes, where each node has three neighbors, as shown in Figure 12a. Each node holds a heterogeneously partitioned MNIST dataset with six to eight classes and trains the model for 20 epochs. During each epoch, the node updates its model for 100 steps with a batch size of 64 using the Adam optimizer with a learning rate of 0.001, then aggregates its local model with those of its three neighbors. The experiment is repeated five times with different random seeds. Figure 12b presents the average validation accuracy and its standard deviation on the MNIST validation set across the training process for the six DFL nodes.

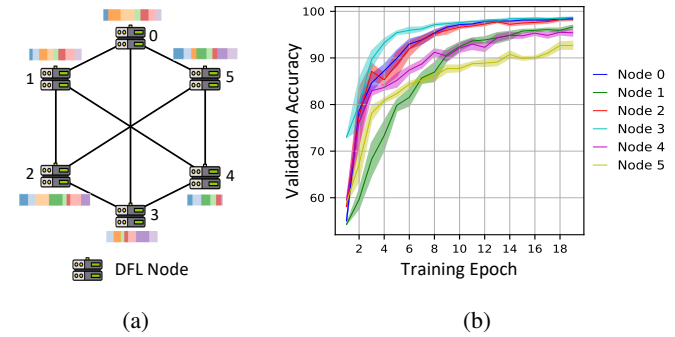


Fig. 12: (a) Topology of the DFL experiments. (b) DFL validation accuracy for the DFL nodes, where the accuracy is evaluated after aggregating the local models of the neighbor DFL nodes.

VI. CONCLUSION

In this paper we present the recent advancements in APPFL, a redesigned federated learning framework to simplify FL usage by offering comprehensive solutions to various challenges and to advance FL research through an easy-to-use, modular interface that facilitates the seamless integration of new algorithms. We demonstrate the capability and extensibility of APPFL by employing it to benchmark various FL aspects and provide case studies across different FL variants. APPFL is open-sourced under the MIT License, and we actively encourage contributions from the community.

In our future work we plan to incorporate more advanced privacy-enhancing technologies into the framework, such as secure multi-party computation, homomorphic encryption, and trusted execution environments, to further ensure the security of FL experiments. We will also focus on secure storage and deployment of the FL-trained models, enabling involved clients to efficiently and safely access the trained models for

inference. Furthermore, APPFL will be used to backend a web platform that aims to provide federated learning as a service, streamlining the end-to-end process of AI model development for domain scientists, from data preparation through model training to model deployment.

ACKNOWLEDGMENT

This research utilizes computing resources provided by the National Artificial Intelligence Research Resource (NAIRR) Pilot, supported by award NAIRR240008. We also gratefully acknowledge Amazon Web Services (AWS) for providing cloud computing credits that were used to assist with benchmarking efforts for this paper.

REFERENCES

- [1] K. Crawford, *The atlas of AI: Power, politics, and the planetary costs of artificial intelligence*. Yale University Press, 2021.
- [2] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint arXiv:1610.05492*, 2016.
- [3] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [4] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 10, no. 2, pp. 1–19, 2019.
- [5] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.
- [6] S. Pati, U. Baid, B. Edwards, M. Sheller, S.-H. Wang, G. A. Reina, P. Foley, A. Gruzdev, D. Karkada, C. Davatzikos *et al.*, "Federated learning enables big data for rare cancer boundary detection," *Nature Communications*, vol. 13, no. 1, p. 7346, 2022.
- [7] G. Kaissis, A. Ziller, J. Passerat-Palmbach, T. Ryffel, D. Usynin, A. Trask, I. Lima Jr, J. Mancuso, F. Jungmann, M.-M. Steinborn *et al.*, "End-to-end privacy preserving deep learning on multi-institutional medical imaging," *Nature Machine Intelligence*, vol. 3, no. 6, pp. 473–484, 2021.
- [8] T.-H. Hoang, J. Fuhrman, R. Madduri, M. Li, P. Chaturvedi, Z. Li, K. Kim, M. Ryu, R. Chard, E. Huerta *et al.*, "Enabling end-to-end secure federated learning in biomedical research on heterogeneous computing environments with APPFLx," *arXiv preprint arXiv:2312.08701*, 2023.
- [9] G. Wang, C. X. Dang, and Z. Zhou, "Measure contribution of participants in federated learning," in *2019 IEEE international conference on Big Data (Big Data)*. IEEE, 2019, pp. 2597–2604.
- [10] G. Wang, "Interpret federated learning with shapley values," *arXiv preprint arXiv:1905.04519*, 2019.
- [11] S. Bose and K. Kim, "Federated short-term load forecasting with personalization layers for heterogeneous clients," *arXiv preprint arXiv:2309.13194*, 2023.
- [12] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [13] S. Ramaswamy, R. Mathews, K. Rao, and F. Beaufays, "Federated learning for emoji prediction in a mobile keyboard," *arXiv preprint arXiv:1906.04329*, 2019.
- [14] D. Leroy, A. Coucke, T. Lavril, T. Gisselbrecht, and J. Dureau, "Federated learning for keyword spotting," in *ICASSP 2019-2019 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2019, pp. 6341–6345.
- [15] T.-M. H. Hsu, H. Qi, and M. Brown, "Measuring the effects of non-identical data distribution for federated visual classification," *arXiv preprint arXiv:1909.06335*, 2019.
- [16] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.
- [17] C. Xu, Y. Qu, Y. Xiang, and L. Gao, "Asynchronous federated learning on heterogeneous devices: A survey," *arXiv preprint arXiv:2109.04269*, 2021.
- [18] G. Kaissis, A. Ziller, J. Passerat-Palmbach, T. Ryffel, D. Usynin, A. Trask, I. Lima Jr, J. Mancuso, F. Jungmann, M.-M. Steinborn *et al.*, "End-to-end privacy preserving deep learning on multi-institutional medical imaging," *Nature Machine Intelligence*, vol. 3, no. 6, pp. 473–484, 2021.
- [19] D. J. Beutel, T. Topal, A. Mathur, X. Qiu, J. Fernandez-Marques, Y. Gao, L. Sani, K. H. Li, T. Parcollet, P. P. B. de Gusmão *et al.*, "Flower: A friendly federated learning research framework," *arXiv preprint arXiv:2007.14390*, 2020.
- [20] C. He, S. Li, J. So, X. Zeng, M. Zhang, H. Wang, X. Wang, P. Vepakomma, A. Singh, H. Qiu *et al.*, "FedML: A research library and benchmark for federated machine learning," *arXiv preprint arXiv:2007.13518*, 2020.
- [21] F. Lai, Y. Dai, S. Singapuram, J. Liu, X. Zhu, H. Madhyastha, and M. Chowdhury, "Fedscale: Benchmarking model and system performance of federated learning at scale," in *International conference on machine learning*. PMLR, 2022, pp. 11 814–11 827.
- [22] M. Ryu, Y. Kim, K. Kim, and R. K. Madduri, "APPFL: open-source software framework for privacy-preserving federated learning," in *2022 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. IEEE, 2022, pp. 1074–1083.
- [23] Z. Li, P. Chaturvedi, S. He, H. Chen, G. Singh, V. Kindratenko, E. A. Huerta, K. Kim, and R. Madduri, "FedCompass: efficient cross-silo federated learning on heterogeneous client devices using a computing power aware scheduler," *arXiv preprint arXiv:2309.14675*, 2023.
- [24] R. Chard, Y. Babuji, Z. Li, T. Skluzacek, A. Woodard, B. Blaiszik, I. Foster, and K. Chard, "Funcx: A federated function serving fabric for science," in *Proceedings of the 29th International symposium on high-performance parallel and distributed computing*, 2020, pp. 65–76.
- [25] J. G. Pauloski, V. Hayot-Sasson, L. Ward, N. Hudson, C. Sabino, M. Baughman, K. Chard, and I. Foster, "Accelerating communications in federated applications with transparent object proxies," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2023, pp. 1–15.
- [26] G. Wilkins, S. Di, J. C. Calhoun, Z. Li, K. Kim, R. Underwood, R. Mortier, and F. Cappello, "FedSZ: Leveraging error-bounded lossy compression for federated learning communications," in *2024 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. IEEE, 2024, pp. 1187–1188.
- [27] I. Foster, "Globus Online: Accelerating and democratizing science through cloud-based services," *IEEE Internet Computing*, vol. 15, no. 3, pp. 70–73, 2011.
- [28] B. Allen, J. Bresnahan, L. Childers, I. Foster, G. Kandaswamy, R. Kettimuthu, J. Kordas, M. Link, S. Martin, K. Pickett *et al.*, "Software as a service for data scientists," *Communications of the ACM*, vol. 55, no. 2, pp. 81–88, 2012.
- [29] S. Tuecke, R. Ananthkrishnan, K. Chard, M. Lidman, B. McCollam, S. Rosen, and I. Foster, "Globus Auth: A research identity and access management platform," in *2016 IEEE 12th International Conference on e-Science (e-Science)*. IEEE, 2016, pp. 203–212.
- [30] C. Dwork, "Differential privacy," in *International colloquium on automata, languages, and programming*. Springer, 2006, pp. 1–12.
- [31] A. Ghosh, J. Hong, D. Yin, and K. Ramchandran, "Robust federated learning in a heterogeneous environment," *arXiv preprint arXiv:1906.06629*, 2019.
- [32] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, M. Al-Shedivat, G. Andrew, S. Avestimehr, K. Daly, D. Data *et al.*, "A field guide to federated optimization," *arXiv preprint arXiv:2107.06917*, 2021.
- [33] S. P. Karimireddy, S. Kale, M. Mohri, S. Reddi, S. Stich, and A. T. Suresh, "Scaffold: Stochastic controlled averaging for federated learning," in *International conference on machine learning*. PMLR, 2020, pp. 5132–5143.
- [34] S. Reddi, Z. Charles, M. Zaheer, Z. Garrett, K. Rush, J. Konečný, S. Kumar, and H. B. McMahan, "Adaptive federated optimization," *arXiv preprint arXiv:2003.00295*, 2020.
- [35] C. Xu, Y. Qu, Y. Xiang, and L. Gao, "Asynchronous federated learning on heterogeneous devices: A survey," *Computer Science Review*, vol. 50, p. 100595, 2023.
- [36] C. Xie, S. Koyejo, and I. Gupta, "Asynchronous federated optimization," *arXiv preprint arXiv:1903.03934*, 2019.

- [37] Y. Chen, Y. Ning, M. Slawski, and H. Rangwala, "Asynchronous online federated learning for edge devices with non-iid data," in *2020 IEEE International Conference on Big Data (Big Data)*. IEEE, 2020, pp. 15–24.
- [38] J. Nguyen, K. Malik, H. Zhan, A. Yousefpour, M. Rabbat, M. Malek, and D. Huba, "Federated learning with buffered asynchronous aggregation," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 3581–3607.
- [39] C. Iakovidou and K. Kim, "Asynchronous federated stochastic optimization with exact averaging for heterogeneous local objectives," *arXiv preprint arXiv:2405.10123*, 2024.
- [40] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konečný, S. Mazzocchi, B. McMahan *et al.*, "Towards federated learning at scale: System design," *Proceedings of Machine Learning and Systems*, vol. 1, pp. 374–388, 2019.
- [41] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *ICC 2019-2019 IEEE international conference on communications (ICC)*. IEEE, 2019, pp. 1–7.
- [42] J. Hao, Y. Zhao, and J. Zhang, "Time efficient federated learning with semi-asynchronous communication," in *2020 IEEE 26th International Conference on Parallel and Distributed Systems (ICPADS)*. IEEE, 2020, pp. 156–163.
- [43] Z. Chen, W. Liao, K. Hua, C. Lu, and W. Yu, "Towards asynchronous federated learning for heterogeneous edge-powered Internet of Things," *Digital Communications and Networks*, vol. 7, no. 3, pp. 317–326, 2021.
- [44] W. Zhuang, C. Chen, and L. Lyu, "When foundation model meets federated learning: Motivations, challenges, and future directions," *arXiv preprint arXiv:2306.15546*, 2023.
- [45] W. Kuang, B. Qian, Z. Li, D. Chen, D. Gao, X. Pan, Y. Xie, Y. Li, B. Ding, and J. Zhou, "FederatedScope-LLM: a comprehensive package for fine-tuning large language models in federated learning," *arXiv preprint arXiv:2309.00363*, 2023.
- [46] Z. Li, S. He, P. Chaturvedi, V. Kindratenko, E. A. Huerta, K. Kim, and R. Madduri, "Secure federated learning across heterogeneous cloud and high-performance computing resources – a case study on federated fine-tuning of LLaMA 2," *Computing in Science & Engineering*, 2024.
- [47] C. Chen, X. Feng, J. Zhou, J. Yin, and X. Zheng, "Federated large language model: A position paper," *arXiv preprint arXiv:2307.08925*, 2023.
- [48] Y. J. Cho, J. Wang, and G. Joshi, "Client selection in federated learning: Convergence analysis and power-of-choice selection strategies," *arXiv preprint arXiv:2010.01243*, 2020.
- [49] Y. J. Cho, D. Jhunjunwala, T. Li, V. Smith, and G. Joshi, "To federate or not to federate: Incentivizing client participation in federated learning," *arXiv preprint arXiv:2205.14840*, 2022.
- [50] D. Jhunjunwala, P. Sharma, A. Nagarkatti, and G. Joshi, "FedVARP: Tackling the variance due to partial client participation in federated learning," in *Uncertainty in Artificial Intelligence*. PMLR, 2022, pp. 906–916.
- [51] Y. J. Cho, J. Wang, and G. Joshi, "Towards understanding biased client selection in federated learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2022, pp. 10351–10375.
- [52] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek, "Robust and communication-efficient federated learning from non-iid data," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 9, pp. 3400–3413, 2019.
- [53] F. Haddadpour, M. M. Kamani, A. Mokhtari, and M. Mahdavi, "Federated learning with compression: Unified analysis and sharp guarantees," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 2350–2358.
- [54] S. Caldas, S. M. K. Duddu, P. Wu, T. Li, J. Konečný, H. B. McMahan, V. Smith, and A. Talwalkar, "Leaf: A benchmark for federated settings," *arXiv preprint arXiv:1812.01097*, 2018.
- [55] Y. Xie, Z. Wang, D. Gao, D. Chen, L. Yao, W. Kuang, Y. Li, B. Ding, and J. Zhou, "FederatedScope: A flexible federated learning platform for heterogeneity," *arXiv preprint arXiv:2204.05011*, 2022.
- [56] H. R. Roth, Y. Cheng, Y. Wen, I. Yang, Z. Xu, Y.-T. Hsieh, K. Kersten, A. Harouni, C. Zhao, K. Lu *et al.*, "NVIDIA FLARE: Federated learning from simulation to real-world," *arXiv preprint arXiv:2210.13291*, 2022.
- [57] P. Foley, M. J. Sheller, B. Edwards, S. Pati, W. Riviera, M. Sharma, P. N. Moorthy, S.-h. Wang, J. Martin, P. Mirhaji *et al.*, "OpenFL: the open federated learning library," *Physics in Medicine & Biology*, vol. 67, no. 21, p. 214001, 2022.
- [58] Z. Wang, X. Fan, Z. Peng, X. Li, Z. Yang, M. Feng, Z. Yang, X. Liu, and C. Wang, "LGo: A fully customizable federated learning platform," *arXiv preprint arXiv:2306.12079*, 2023.
- [59] D. Zeng, S. Liang, X. Hu, H. Wang, and Z. Xu, "FedLab: a flexible federated learning framework," *Journal of Machine Learning Research*, vol. 24, no. 100, pp. 1–7, 2023.
- [60] B. Zhao, K. R. Mopuri, and H. Bilen, "idlg: Improved deep leakage from gradients," *arXiv preprint arXiv:2001.02610*, 2020.
- [61] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients-how easy is it to break privacy in federated learning?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 16937–16947, 2020.
- [62] H. Yin, A. Mallya, A. Vahdat, J. M. Alvarez, J. Kautz, and P. Molchanov, "See through gradients: Image batch recovery via gradient inversion," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 16337–16346.
- [63] A. Hatamizadeh, H. Yin, P. Molchanov, A. Myronenko, W. Li, P. Dogra, A. Feng, M. G. Flores, J. Kautz, D. Xu *et al.*, "Do gradient inversion attacks make federated learning unsafe?" *IEEE Transactions on Medical Imaging*, vol. 42, no. 7, pp. 2044–2056, 2023.
- [64] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," *arXiv preprint arXiv:1712.01887*, 2017.
- [65] Y. Tsuzuku, H. Imachi, and T. Akiba, "Variance-based gradient compression for efficient distributed deep learning," *arXiv preprint arXiv:1802.06058*, 2018.
- [66] Y. Aono, T. Hayashi, L. Wang, S. Moriai *et al.*, "Privacy-preserving deep learning via additively homomorphic encryption," *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 5, pp. 1333–1345, 2017.
- [67] L. Lamport, R. Shostak, and M. Pease, "The Byzantine generals problem," in *Concurrency: the works of Leslie Lamport*. Machinery, 2019, pp. 203–226.
- [68] V. Shejwalkar and A. Houmansadr, "Manipulating the byzantine: Optimizing model poisoning attacks and defenses for federated learning," in *NDSS*, 2021.
- [69] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *International conference on artificial intelligence and statistics*. PMLR, 2020, pp. 2938–2948.
- [70] X. Cao, M. Fang, J. Liu, and N. Z. Gong, "FLTrust: Byzantine-robust federated learning via trust bootstrapping," *arXiv preprint arXiv:2012.13995*, 2020.
- [71] X. Zhou, M. Xu, Y. Wu, and N. Zheng, "Deep model poisoning attack on federated learning," *Future Internet*, vol. 13, no. 3, p. 73, 2021.
- [72] V. Tolpegin, S. Truex, M. E. Gursoy, and L. Liu, "Data poisoning attacks against federated learning systems," in *Computer security—ESORICS 2020: 25th European symposium on research in computer security, ESORICS 2020, guildford, UK, September 14–18, 2020, proceedings, part i 25*. Springer, 2020, pp. 480–501.
- [73] M. Fang, X. Cao, J. Jia, and N. Gong, "Local model poisoning attacks to {Byzantine-Robust} federated learning," in *29th USENIX security symposium (USENIX Security 20)*, 2020, pp. 1605–1622.
- [74] P. Blanchard, E. M. El Mhamdi, R. Guerraoui, and J. Stainer, "Machine learning with adversaries: Byzantine tolerant gradient descent," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [75] D. Yin, Y. Chen, R. Kannan, and P. Bartlett, "Byzantine-robust distributed learning: Towards optimal statistical rates," in *International conference on machine learning*. Pmlr, 2018, pp. 5650–5659.
- [76] C. Xie, S. Koyejo, and I. Gupta, "Zeno: Distributed stochastic gradient descent with suspicion-based fault-tolerance," in *International Conference on Machine Learning*. PMLR, 2019, pp. 6893–6901.
- [77] S. Prakash and A. S. Avestimehr, "Mitigating Byzantine attacks in federated learning," *arXiv preprint arXiv:2010.07541*, 2020.
- [78] Z. Li, S. He, P. Chaturvedi, T.-H. Hoang, M. Ryu, E. Huerta, V. Kindratenko, J. Fuhrman, M. Giger, R. Chard *et al.*, "APPFLx: providing privacy-preserving cross-silo federated learning as a service," in *2023 IEEE 19th International Conference on e-Science (e-Science)*. IEEE, 2023, pp. 1–4.
- [79] Y. Liu, Y. Kang, T. Zou, Y. Pu, Y. He, X. Ye, Y. Ouyang, Y.-Q. Zhang, and Q. Yang, "Vertical federated learning: Concepts, advances, and challenges," *IEEE Transactions on Knowledge and Data Engineering*, 2024.

- [80] M. S. H. Abad, E. Ozfatura, D. Gunduz, and O. Ercetin, "Hierarchical federated learning across heterogeneous cellular networks," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 8866–8870.
- [81] A. Lalitha, S. Shekhar, T. Javidi, and F. Koushanfar, "Fully decentralized federated learning," in *Third workshop on bayesian deep learning (NeurIPS)*, vol. 2, 2018.
- [82] Y. Collet and M. Kucherawy, "Zstandard compression and the application/zstd media type," Tech. Rep., 2018.
- [83] P. Deutsch, "GZIP file format specification version 4.3," Tech. Rep., 1996.
- [84] Blosc Development Team. (2009-2023) A fast, compressed and persistent data store library. <https://blosc.org>.
- [85] X. Liang, S. Di, D. Tao, S. Li, S. Li, H. Guo, Z. Chen, and F. Cappello, "Error-controlled lossy compression optimized for high compression ratios of scientific datasets," in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018, pp. 438–447.
- [86] "SZ—3: A modular framework for composing prediction-based error-bounded lossy compressors, author=Liang, Xin and Zhao, Kai and Di, Sheng and Li, Sihuan and Underwood, Robert and Gok, Ali M and Tian, Jiannan and Deng, Junjing and Calhoun, Jon C and Tao, Dingwen and others, journal=IEEE Transactions on Big Data, volume=9, number=2, pages=485–498, year=2022, publisher=IEEE."
- [87] P. Lindstrom, "Fixed-rate compressed floating-point arrays," *IEEE Transactions on Visualization and Computer Graphics*, vol. 20, no. 12, pp. 2674–2683, 2014.
- [88] S. Bose, Y. Zhang, and K. Kim, "Addressing heterogeneity in federated load forecasting with personalization layers," *arXiv preprint arXiv:2404.01517*, 2024.
- [89] S. Zhou and G. Y. Li, "Communication-efficient ADMM-based federated learning," *arXiv preprint arXiv:2110.15318*, 2021.
- [90] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [91] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [92] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [93] J. Ogier du Terrail, S.-S. Ayed, E. Cyffers, F. Grimberg, C. He, R. Loeb, P. Mangold, T. Marchand, O. Marfoq, E. Mushtaq *et al.*, "FLamby: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings," *Advances in Neural Information Processing Systems*, vol. 35, pp. 5315–5334, 2022.
- [94] D. R. Cox, "Regression models and life-tables," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 34, no. 2, pp. 187–202, 1972.
- [95] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, "3D U-Net: learning dense volumetric segmentation from sparse annotation," in *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2016: 19th International Conference, Athens, Greece, October 17-21, 2016, Proceedings, Part II 19*. Springer, 2016, pp. 424–432.
- [96] M. Tan and Q. Le, "EfficientNet: Rethinking model scaling for convolutional neural networks," in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.
- [97] L. R. Dice, "Measures of the amount of ecologic association between species," *Ecology*, vol. 26, no. 3, pp. 297–302, 1945.
- [98] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [99] X. Ma, J. Zhu, Z. Lin, S. Chen, and Y. Qin, "A state-of-the-art survey on solving non-iid data in federated learning," *Future Generation Computer Systems*, vol. 135, pp. 244–258, 2022.
- [100] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani, "Least angle regression," *Ann. Statist.*, vol. 32, no. 2, pp. 407–499, 2004.
- [101] V. Nair and G. E. Hinton, "Rectified linear units improve restricted Boltzmann machines," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 807–814.
- [102] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [103] Y. LeCun, "The MNIST database of handwritten digits," <http://yann.lecun.com/exdb/mnist/>, 1998.
- [104] T. Sun, D. Li, and B. Wang, "Decentralized federated averaging," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 4, pp. 4289–4301, 2022.