

LATTE: Layer Algorithm-aware Training Time Estimation for Heterogeneous Federated Learning

Kun Wang¹, Zimu Zhou², and Zhenjiang Li¹

¹Department of Computer Science, City University of Hong Kong, Hong Kong SAR, China

²School of Data Science, City University of Hong Kong, Hong Kong SAR, China

Abstract

Accurate estimation of on-device model training time is increasingly required for emerging learning paradigms on mobile edge devices, such as heterogeneous federated learning (HFL). HFL usually customizes the model architecture according to the different capabilities of mobile edge devices to ensure efficient use of local data from all devices for training. However, due to oversimplification of latency modeling, existing methods rely on a single coefficient to represent computational heterogeneity, resulting in sub-optimal HFL efficiency. We find that existing methods ignore the important impact of runtime optimization of deep learning frameworks, which we call development-chain diversity. Specifically, layers of a model may have different algorithm implementations, and deep learning frameworks often have different strategies for selecting the algorithm they believe is the best based on a range of runtime factors, resulting in different training latencies and invalid predictions from existing methods. In this paper, in addition to considering this diversity to ensure synchronized completion time of model training, we also study how to select the best algorithm each time to reduce the latency of the per-round training, thereby further improving the overall efficiency of federated training. To this end, we propose LATTE, which consists of a novel selector that identifies the best algorithm at runtime based on relative runtime factors. By further integrating it into our training latency model, LATTE provides accurate training time estimation. We develop LATTE as middleware, compatible with different deep learning frameworks. Extensive results show significantly improved training convergence speed and model accuracy compared to state-of-the-art methods.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
ACM MobiCom '24, November 18–22, 2024, Washington D.C., DC, USA
© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0489-5/24/11

<https://doi.org/10.1145/3636534.3690705>

CCS Concepts

• **Human-centered computing** → **Mobile computing**; • **Computing methodologies** → **Machine learning**.

Keywords

Heterogeneous Federated Learning, Training Time Estimation

ACM Reference Format:

Kun Wang¹, Zimu Zhou², and Zhenjiang Li¹. 2024. LATTE: Layer Algorithm-aware Training Time Estimation for Heterogeneous Federated Learning. In *The 30th Annual International Conference on Mobile Computing and Networking (ACM MobiCom '24)*, November 18–22, 2024, Washington D.C., DC, USA. ACM, New York, NY, USA, 15 pages. <https://doi.org/10.1145/3636534.3690705>

1 Introduction

As mobile edge devices, such as NVIDIA Jetson series, become increasingly powerful, in addition to efficiently performing *inference* tasks [18, 63–65, 70], they can now also perform certain model *training* tasks [24]. However, since the computing capability of mobile edge devices is far less than that of personal computers or servers, on-device training at the mobile edge is not positioned to replace traditional model training. It mainly focuses on updating the model according to the requirements of many emerging learning paradigms, such as federated learning (FL) [22, 35, 37, 51, 52, 59, 62].

Training time is a key factor in determining the efficiency of federated learning on mobile edge devices as different devices may have vastly different computing and communication capabilities, known as *system heterogeneity* [73]. Specifically, multiple edge devices (as clients) collaborate to train a model under the coordination of the server [49]. In each round of training, the server assigns the current model to the devices. Each device then trains the model using its own data and returns updates to the server, which are aggregated into the next round of models. Primarily due to the heterogeneity in *computing capability* [53], devices usually cannot complete training of the same model at the same time, and the server must wait for updates of the slowest device in each round [49], leading to significant overall training latency measured in wall-clock time. This latency can seriously harm the availability of federated learning systems, especially in time-sensitive scenarios such as personal assistants, drone fleets, and robot swarms [21, 48, 68].

When the number of federated devices is large (*e.g.*, hundreds or thousands), recent work has found that training efficiency can be improved through client selection [34, 37, 40, 42] or asynchronous FL [59, 69], mitigating the impact of weak devices on training. However, a more general solution, known as heterogeneous federated learning (HFL), is to allocate appropriate sub-models according to the computing capability of the device [14, 22, 23, 43, 51, 56]. The HFL literature shows that having each device complete each round of model training simultaneously as much as possible can lead to better performance [39, 71]. This not only utilizes the training data on each participating device, but also significantly improves the efficiency of the overall federated training. It is particularly useful in most mobile edge applications in practice with medium-sized device networks, where each device may contain a unique set of data.

Since training time is critical for federated learning, the training time for each device needs to be accurately estimated to allocate appropriate sub-models. To this end, existing work typically predicts the latency by performing linear regression on the computational load of the model [22, 23], such as the number of floating point operations (FLOPs). They characterize computational heterogeneity across devices by computing capability related coefficients (*e.g.*, float point per second (FLOPS)) determined offline. Some recent studies have further proposed finer-grained layer-wise modeling [14, 43], acknowledging that different layer types may exhibit different execution times on the same platform but follow the same estimation principle. In this paper, we find that existing approaches ignore an important impact of different deep learning frameworks (*e.g.*, PyTorch and TensorFlow) exhibit distinct behaviors in invoking acceleration libraries (*e.g.*, cuDNN) during runtime optimization [41, 44, 72], which we refer to as *development-chain diversity*. It can cause significant deviations in training latency for the same model, even on the same device. This can make the characterization based on hardware parameters learned in advance inaccurate.

By delving into the runtime behaviors of deep learning frameworks, we discover that the core reason is the diversity of layer algorithms. For example, the convolutional layers, essential in many neural networks, can be implemented through various algorithms, such as direct convolution, matrix multiplication, and fast Fourier transform. However, different deep learning frameworks have different strategies to select what they think is the best layer algorithm at runtime based on hardware specifications, layer configuration, and available runtime resources. This selection needs to be made for both *forward* and *backward* propagation of layers before training, so the choice of layer algorithm can significantly affect the overall training latency. As shown in § 2.2, training the same model can have a 30.1% time difference and be 2.68× slower than the optimal for this reason.

If we deliberately measure the latency differences caused by the diversity of algorithm choices at each layer and take this into account when assigning sub-models, different devices can achieve synchronized completion times of model training. However, in addition to uncovering this key question, we aim to go one step further: Can every device also choose the best algorithm each time? If possible, training can not only be completed synchronously across devices, but the latency in each round can also be minimized (so that we can allocate the largest sub-model), which will further improve the overall training efficiency and performance [23].¹

To this end, we propose LATTE, a Layer Algorithm-aware Training Time Estimator for HFL, which adapts to the diversity of layer algorithms in the development chain to achieve accurate training time estimation. A key challenge in designing LATTE is that layer algorithm selection is affected by certain runtime parameters, such as resource availability and model configuration, which results in the inability to enumerate the consequences of choices in advance. To solve this problem, LATTE provides a lightweight and accurate layer algorithm selector. A main novelty in its design is the generation of training data and ground truth for the selector. Specifically, we design an acquisition tool that can derive a set of training data and leverage the functional modules from existing deep learning frameworks to generate corresponding ground truth for training the selector in a device-independent manner so that the generation can be executed just once offline before deployment. Trained on comprehensive, high-quality data, the selector can accurately select the best algorithm to match the model configurations and runtime resources of the development chain. By further integrating the selected algorithm into our fine-grained training latency model, LATTE can accurately estimate the latency for each single-pass (forward and backward) propagation.

With accurate training latency estimation, LATTE facilitates configuring sub-models across devices to complete local training simultaneously. Due to our local training time estimation capability, we can move from traditional *server-driven assignment* to *client-side sub-model allocation*, which enhances adaptability to resource dynamics within and across training rounds. We develop LATTE as middleware using Python 3.6, which is compatible with PyTorch and TensorFlow. We further integrate LATTE into the advanced and mainstream federated learning framework Flower [16], and conduct experiments on edge devices with different computing capabilities and deep learning frameworks. We evaluate LATTE on popular deep learning tasks by considering both

¹For ease of description, we do not discuss communication latency here, which is a relatively independent problem from training latency, but we do incorporate communication latency into our sub-model allocation design and experimental evaluations.

IID and non-IID data distributions, including image classification on CIFAR-10/100 and OpenImage datasets, speech recognition using the Google Speech dataset, and human activity recognition using the HARBox dataset. We compare LATTE with seven classical or state-of-the-art methods, such as FedRolex [14] and TailorFL [22]. Extensive results show that LATTE not only significantly speeds up convergence across tasks by 1.17–2.76 \times , but also improves the accuracy of the central model by 2.44–9.78%. Moreover, LATTE exhibits unique orthogonal capabilities to personalized FL methods such as Hermes [35], demonstrating versatility in non-IID scenarios. Our main contributions are summarized below.

- We reveal the problem of development-chain diversity in federated learning systems and identify diverse layer algorithms as the key to explain the variability in training time. Based on this, accurate estimation of model training time can be achieved without complex operator or kernel-level modeling.
- We devise LATTE, with a novel layer algorithm selector and training time estimator, to accurately estimate the single-pass (forward/backward) propagation latency of a model given its architecture, expected hardware and runtime memory. We further showcase its usability in a client-side sub-model selection for HFL.
- We conduct extensive experiments to evaluate LATTE in five typical HFL scenarios. The results show significant improvements in performance compared to seven classical or state-of-the-art methods.

2 Background & Motivation

We first review how system heterogeneity affects FL (§ 2.1) and the conventional training time modeling in FL (§ 2.2).

2.1 Primer of HFL

System Heterogeneity. Real-world FL deployments faces system heterogeneity [53]. One main reason is the diverse *computing capabilities* of participating devices. Such variations can lead to notable discrepancies in *training latency* of the same model, which affects the overall efficiency of FL.

In a typical FL process, such as the widely-used FedAvg [49], training occurs in *rounds*. Each device first downloads the current model from the server, trains the model on its local data for multiple epochs, and then uploads the updated model back to the server. Then the server aggregates these updates as the model for the next round. However, not all devices can complete the same amount of computation in each round. Slower devices can significantly delay the entire training process, leading to long wall-clock training time.

HFL Strategies. HFL is proposed to handle the system heterogeneity, and an effective approach studied in recent work [14, 22, 23, 43, 51, 56] is to tailor the model architecture

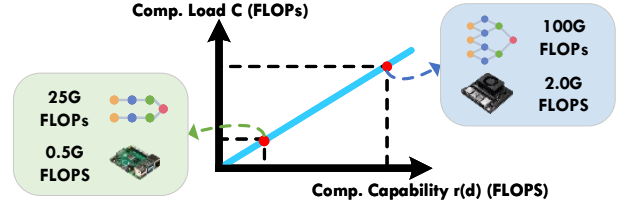


Figure 1: Allocating a sub-model with computational load proportional to the capability $r(d)$ of device d .

to fit the capabilities of each device, *i.e.*, larger models are allocated to devices with higher capabilities, while smaller models are designated for less powerful devices. The goal is to synchronize the training completion time across all devices. These model-level operations lead to a general solution to cope with the system heterogeneity in FL systems, while we find that a key issue in their training time modeling that still limits the effectiveness of federation.²

2.2 Existing Training Time Modeling in FL

Most HFL studies [14, 22, 23, 43, 51, 56] focus on the allocation and aggregation of heterogeneous models for effective training. They match the models with devices by assuming a linear mapping between the network architecture and its training latency, expecting that the computational loads of the allocated models are within the devices' capacities and that all devices finish their training simultaneously. However, this method can lead to inconsistent training time estimates due to the *development-chain diversity*, as explained below.

Conventional Modeling. Assume local training uses standard mini-batch gradient descent. The training time T_{train} of a model architecture M on a device d can be estimated as:

$$T_{train} = E \times (b \times n) \times T, \quad (1)$$

where E is the number of epochs, b is the batch size, n is the number of batches, and T represents the latency of single-pass backpropagation (forward and backward) during training. Although E , b , and n are deterministic, latency T can vary from device to device. Specifically, the latency of a single pass of backpropagation is typically assumed to be proportional to the computational load C_M of the model M , *i.e.*,

$$T = \frac{C_M}{r(d)}, \quad (2)$$

where $r(d)$ is a computing capability related coefficient, indicating the average processing speed of device d . For a fixed model architecture M , existing work assumes that C_M remains *constant*, suggesting that measuring the $r(d)$ suffices to accurately estimate the training latency [22, 23], which

²When the number of edge devices is quite large in FL systems, *e.g.*, hundreds or even thousands, some other orthogonal methods are also studied in the literature such as client selection and asynchronous/semi-asynchronous FL, which are detailed in related work (§5).

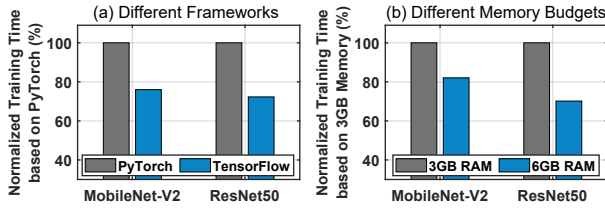


Figure 2: Wall-clock training times of (a) using different frameworks on same models, and (b) using same framework on the same model but with different memory resources. Since absolute training times differ for each setting, we normalize them for clarity.

can be performed offline as a one-time effort, given known hardware specifications. We can then allocate sub-models to devices based on Equation (2), as shown in Figure 1.

Specifically, to ensure roughly the same training latency $T(M_1) = T(M_2)$, the server allocates sub-models M_1 and M_2 proportional to the capability of device d_1 and d_2 , as follows:

$$\frac{C_{M_1}}{C_{M_2}} = \frac{r(d_1)}{r(d_2)}, \quad (3)$$

because the computational load C_{M_i} in Equation (3) of a model architecture M_i is constant given its topology.

Limitations. However, this modeling of training time is oversimplified because computing capability is not the only source of system heterogeneity. Diversity in the development-chain of neural networks also matters. To reveal this, we conduct preliminary experiments on NVIDIA Jetson devices.

- We first run the same neural network on the same device, but using different deep learning frameworks. Figure 2(a) shows that across two popular models, the wall-clock training time of a single round with TensorFlow is only 72.3–76.0% compared to using PyTorch.
- Even if developed using the same deep learning framework, models show different training times on the same hardware due to varied runtime resources (e.g., memory). Figure 2(b) shows that the wall-clock training time of single epoch round with 6GB memory budget is only 69.9–82.0% compared to a 3GB budget.

Implications. These experiments imply that training time is also related to the *runtime optimizations* of the deep learning framework. Furthermore, it links to *runtime resources* such as *memory*. These observations challenge the conventional training latency estimation and requires a new, more fine-grained modeling *beyond hardware specifications* to achieve more accurate training time estimation.

3 System Design

This section presents LATTE to solve the problems above. Figure 3 illustrates the functional modules of LATTE.

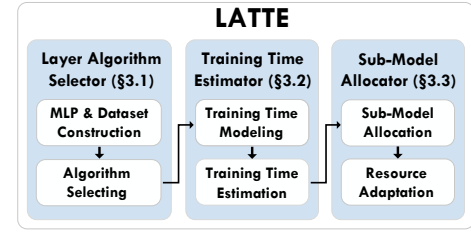


Figure 3: Architecture of the LATTE design.

- **Layer algorithm selector** (§ 3.1). Based on the key design insight to be described, this module first ranks candidate algorithms by considering the layer configuration, runtime memory, and hardware specifications. The selector then selects the fastest algorithm that satisfies the resource constraint as the output.
- **Training time estimator** (§ 3.2). This module further estimates the training time by explicitly incorporating the algorithm selection. We integrate the selected layer algorithm into our latency modeling and extend it to the entire model training. To achieve this estimation, we also propose effective ways to profile the relationship between computation and latency.
- **Sub-model allocator** (§ 3.3). This module finally employs the accurate training time estimation to allocate suitable sub-models for devices. It allows client-side sub-model adaptation to dynamic resource with a mechanism to ensure that all model parameters can be equally trained despite the varying sub-model architectures across different training rounds.

3.1 Layer Algorithm Selector

The same layer can have different implementations at runtime. For example, NVIDIA’s cuDNN library include multiple algorithms (see Table 1) to implement forward convolution (`cudaConvolutionBwdFilterAlgo_t`), e.g., direct method (`CUDNN_CONVOLUTION_FWD_ALGO_DIRECT`), matrix product (`CUDNN_CONVOLUTION_FWD_ALGO_GEMM`), to name a few. Different algorithms are also available for backward filter convolution (`cudaConvolutionBwdFilterAlgo_t`), backward data convolution (`cudaConvolutionBwdDataAlgo_t`), etc.³ Devices that primarily use CPUs (e.g., Raspberry Pi) also have diverse layer algorithms provided by their vendors [1, 7].

3.1.1 Diversity of Layer Algorithms. The same layer implemented by different algorithms results in different trade-offs. Some algorithms are general, which are compatible with most types of layer configurations, but their latency is always moderate, and their memory consumption is not small

³We use convolutions as an example to explain our insights and methods because they are among the most of optimized operations by deep learning frameworks and GPU vendors, which makes their latency estimation more challenging than less optimized operations.

Layer Algorithms	Generality	Memory Efficiency
GEMM	++	+
FFT	+	++
FFT_TILING	+	+++
IMPLICIT_GEMM	+++	++
IMPLICIT_PRECOMP_GEMM	++	++
DIRECT	++	+
WINOGRAD	+	+++
WINOGRAD_NONFUDED	++	++

Table 1: Candidate forward convolution layer algorithms and their characteristics. More '+' signs indicate better values for this characteristic category.

(e.g., DIRECT in Table 1). In contrast, some algorithms are tailored for certain layer configurations, with optimized performance (e.g., WINOGRAD in Table 1 is optimized for small kernels), but suffer substantial latency when applied to relatively incompatible configurations. Due to this diversity, all frameworks require algorithm selection before model training. Although strategies in different frameworks can vary significantly, they generally follow two steps [10, 11]:

- **1) Ranking:** rank candidate algorithms for a given layer by the predicted latency based on their layer estimation strategies, and
- **2) Selecting:** estimate the corresponding memory footprint and select the top-ranked (fastest) algorithm that satisfies the predefined memory constraint.

Thus, the layer algorithm selection could affect the training time estimation in two aspects. First, given enough memory for the same model, layer algorithm selection strategies in different DL frameworks may produce different results, primarily because these strategies obtain rankings of the candidate algorithms in different ways (e.g., heuristic strategy in PyTorch may not return the best ranking). Second, given limited memory budget, the same strategy may also return different algorithm ranking results even in the same device, since the selection must adhere to the memory limit. These two reasons lead to the observations of Figure 2(a) and (b) in §2.

Design Insight. The diversity of layer algorithms invalidates the conventional practice for sub-model allocation, i.e., Equation (3), because different layer algorithms also lead to different latencies. Therefore, given a model architecture M , its computational load is a *function* of the layer algorithms $\{algo\}$ rather than a constant, i.e., $C_M = C_M(algo)$. Hence, to ensure $T(M_1) = T(M_2)$ on two devices d_1 and d_2 , the sub-models M_1 and M_2 should satisfy:

$$\frac{C_{M_1}(algo_1)}{C_{M_2}(algo_2)} = \frac{r(d_1)}{r(d_2)}, \quad (4)$$

where $C_{M_1}(algo_1)$ and $C_{M_2}(algo_2)$ are the actual workloads for M_1 and M_2 given layer algorithms $algo_1$ and $algo_2$. That is, the sub-model architecture allocated to the device should be

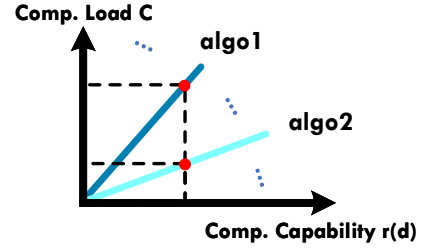


Figure 4: The relationship between the model’s computational load and the device’s computing capability is essentially described by a series of lines due to the diversity of layer algorithms. Accurate training time estimation requires selecting the correct line first.

calibrated by layer algorithm, as illustrated in Figure 4. Otherwise, training latency on devices with the same computing capabilities may differ significantly.

If we profile all layer algorithms in advance and consider their latency difference at each layer when allocating sub-models, devices can still complete per-round training simultaneously even when using different deep learning frameworks. However, beyond revealing this key design insight, we aim to go one step further in this paper: Can we pick the best (fastest) algorithm every time? If possible, we can use this strategy in different frameworks to ensure that each round of federated training can not only be completed synchronously across devices, but also the latency in each round can be minimized (to allocate the largest sub-model), which will further improve overall training efficiency and performance.

Design Challenge. However, it is difficult to achieve such best selection strategy for HFL systems. To understand this, we first analyze how layer algorithm selection is performed in existing deep learning frameworks, which usually form two categories: exhaustive testing (e.g., LaunchConv2Dop in TensorFlow) and black-box heuristics (e.g., cudnnGetConvolutionForwardAlgorithm in PyTorch).

The former, which exhaustively tests all layer algorithms on the target device and selects the fastest implementation, can achieve the best selection, but is very costly. Since sub-models could vary for each device and each training round in HFL systems, such high testing overhead may overwhelm the training latency. Figure 5 shows that even when all device memory is used, the latency of exhaustive testing (e.g., TensorFlow) in first round training still accounts for 18.8–38.4% of the overall training latency. The latter uses heuristic rules to predict the best algorithm without actual on-device testing, which is fast but inaccurate, e.g., its prediction accuracy is only around 73% of the exhaustive testing (see § 4.3.2), which can easily lead to inappropriate sub-model allocation.

Therefore, the main problem that challenges the design of new layer algorithm selection strategy is how to avoid the

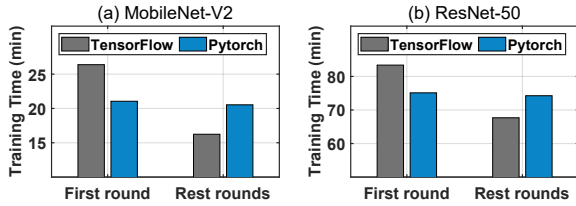


Figure 5: In the first round of training, exhaustive testing strategy (e.g., in TensorFlow) lead to much higher latencies than heuristic-based strategy (e.g., in Pytorch) on two models (a) MobileNet and (b) ResNet.

limitations of existing framework strategies and ensure the accuracy and efficiency of layer algorithm selection.

3.1.2 Design of Layer Algorithm Selector. The above challenge motivates a new design of the layer algorithm selection without runtime assessments, which can override and replace the corresponding module in commercial deep learning frameworks to achieve a *lightweight* and *accurate* selection.

Idea and Design. Since all candidate algorithms have been given in the optimization library of each layer [10, 11], our main idea to achieve an accurate and practical design is to treat this problem as a classification problem and propose an efficient classifier for selection. One of the advantages of this design is that we can shift unnecessary overhead from runtime to offline before deployment. To further improve the practicality of our approach, we also minimize the offline overhead and make it as a one-time and device-independent effort. Below, we first introduce our selector design.

1) Architecture of the selector. For lightweight selection, we design a simple multi-layer perceptron (MLP) classifier. It takes a layer configuration as input and outputs a score for each layer algorithm, indicating its latency ranking. We reuse the inputs of existing heuristic strategy’s API (e.g., I/O dimension, convolution kernel size, padding, stride, dilation *etc.* in `cudaGetConvolutionForwardAlgorithm` as inputs in our MLP, see Figure 6). The output of the softmax layer represents the probability of the fastest layer algorithm (and thus their rankings). We empirically design a 6-layer MLP ($13 \times 256 \times 128 \times 64 \times 32 \times 16 \times 8$) with CrossEntropy loss.

2) Feature enhancement. Inspired by the recent trend of explainable AI [28, 29], we further integrate SHAP [46], an explainable feature selection method to pick important features, into our design to improve performance. For example, `stride_height` and `stride_width` tend to be more important than `input_h`, `input_w`, and `output_c` for forward convolution. Thus, we add an attention layer after the input layer to assign greater weights to these important features.

Training Dataset Construction. With the above classifier design, we further propose a tool, responsible to acquire and label diverse layer configurations to train the classifier.

```

1: int perf_count;
2: std::unique_ptr<perf_t[]> perf_results(new perf_t[num_algos]);
3: if (!benchmark) {
4:   AT_CUDA_CHECK_WITH_SHAPES(cudaGetConvolution
5:     ForwardAlgorithm_v7(
6:       args.handle,
7:       args.idesc.desc(),
8:       args.wdesc.desc(),
9:       args.cdesc.desc(),
10:      args.odesc.desc(),
11:      num_algos,
12:      &perf_count,
13:      perf_results.get(), args);
14: }

```

Figure 6: Features of cuDNN’s official heuristic API.

1) Layer configuration generation. The acquisition tool first samples and then derives a wide spectrum of layer configurations, including I/O size, filter size, stride, padding, *etc.*, from a set of models widely used in HFL scenarios. In our current development, the tool samples from 28 widely used models (see Table 2), resulting in 30,000 configurations.

2) Ground-truth acquisition. After collecting sufficient configuration samples, the main issue is to obtain the ground-truth, *i.e.*, which algorithm is the best choice for each configuration that is difficult to analyze manually. To solve this issue, we propose to reuse the functional APIs of the exhaustive testing-based deep learning frameworks in our design. Although these APIs in deep learning frameworks are implemented as runtime processes, we re-engineering them so that they can be invoked independently before training. As discussed previously (§3.1), this type of strategy is very slow during on-device tests. However, since the latency ranking is based on the computational load of each algorithm rather than the absolute latency value, the ground-truth data is device-independent, and the overhead of collecting it is acceptable when it occurs during a one-time offline phase on a powerful server. Therefore, for ground-truth acquisition, the acquisition tool inputs configuration features into an exhaustive testing API, such as `LaunchConv2D0p` or `cudaFindConvolutionForwardAlgorithm`, to generate ground truth for training the classifier.

After training with layer configurations and their corresponding best choices collected by the tool, we can obtain an accurate selector that is also lightweight and does not require any runtime testing when selecting layer algorithms.

Model Family	Model Instances
VGG	VGG-11/13/16/19
MobileNet	MobileNet-V1/V2
ShuffleNet	ShuffleNet-V1/V2
ResNet	ResNet-18/34/50/50-V2/101/101-V2/152/152-V2
EfficientNet	EfficientNet-B0/B1/B2/B3/B4/B5/B6/B7
Others	GoogLeNet/SqueezeNet/Xception/DenseNet-121

Table 2: Models used to collect layer configurations.

Incorporating Memory Constraint. Given the algorithm rankings obtained by the classifier, the selector ultimately selects the fastest algorithm within memory budget *mem* as output (to be used by the training time estimator). We employ an

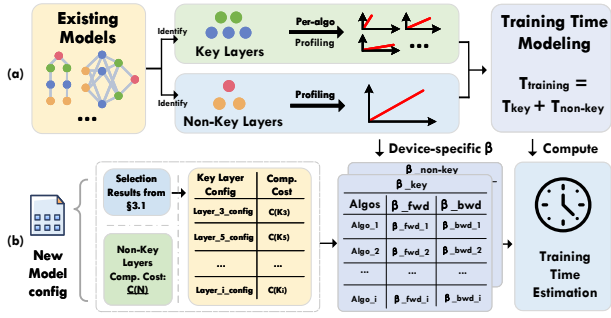


Figure 7: (a) Construction of training time estimator. (b) Training time estimation for a new model.

accurate estimator (`cudaGetConvolutionWorkspaceSize`) to retrieve the memory usage of each layer algorithm.

3.2 Training Time Estimator

The selector (§3.1) is designed to accommodate the diversity of layer algorithms to achieve accurate training time estimation. In this subsection, we describe how to integrate the selection results into our training time estimator.

In general, given a model M , a memory budget mem , and a target device d , the training time estimator (see Figure 7(a)) predicts its single-pass propagation latency T . We decompose the model into *key* layers and *non-key* layers (§ 3.2.1), integrate layer algorithms into the latency modeling of key layers (§ 3.2.2), and profile the coefficients for both key and non-key layers (§ 3.2.3). This approach can provide accurate latency estimates with low profiling overhead.

3.2.1 Key Layer Identification. Given a model, its key layers are the layers with more than one layer algorithms defined in the deep learning framework, *e.g.*, for CNNs, these mainly refer to the convolutional layers, while those with a single and same algorithm implementation are defined as non-key layers, *e.g.*, FC and activation layers. Therefore, we can easily classify each layer as key or non-key based on layer type.

3.2.2 Layer Algorithm Aware Training Time Modeling. We incorporate the layer algorithms by extending traditional latency modeling *i.e.*, Equation (2), as follows:

$$T = \left(\sum_i \frac{C_{fwd}^{key}(algo_i)}{r_{fwd}^{key}(d)} + \sum_i \frac{C_{bwd}^{key}(algo_i)}{r_{bwd}^{key}(d)} \right) + \left(\frac{C_{fwd}^{non}}{r_{fwd}^{non}(d)} + \frac{C_{bwd}^{non}}{r_{bwd}^{non}(d)} \right), \quad (5)$$

where $C_{fwd}^{key}(algo_i)$ and $C_{bwd}^{key}(algo_i)$ are the computational loads of key layer i of model M in the forward and backward pass, implemented by layer algorithm $algo_i$. Similar to existing deep learning frameworks, LATTE also selects the same algorithm for the same key layers of forward and

backward propagation. Similarly, C_{fwd}^{non} and C_{bwd}^{non} are the *total* computation of all the non-key layers of model M in the forward and backward pass. The latency is linear to the computational load when the layer algorithm can be determined, given coefficients $r_{fwd}^{key}(d)$, $r_{bwd}^{key}(d)$, $r_{fwd}^{non}(d)$, and $r_{bwd}^{non}(d)$. The rationale of our training latency modeling is as follows:

- We decompose the model M into key and non-key layers to reduce the profiling overhead while maintaining high latency estimation accuracy.
- As non-key layers have a single and same implementation, they are modeled the same as the traditional approach, *i.e.*, Equation (2), and can be profiled in advanced and applied in forward and backward passes.
- The modeling of key layers explicitly takes into account the diversity of layer algorithms.

For practical usage, we further rephrase Equation (5) as:

$$T = \sum_i (\beta_{fwd}^{key}(algo_i, d) + \beta_{bwd}^{key}(algo_i, d)) \times C(K_i) + (\beta_{fwd}^{non}(d) + \beta_{bwd}^{non}(d)) \times C(N), \quad (6)$$

where $C(K_i)$ is the computational load of key layer K_i with direct implementation (*e.g.*, the direct method for convolution), and $C(N)$ is the total computations loads of all non-key layers. Both $C(K_i)$ and $C(N)$ are deterministic given the model architecture M . In contrast, $\beta(algo, d) = C(algo)/r(d)$ absorbs all algorithm- and device-dependent variations. This formula makes it easy to estimate the single-pass training latency for each layer from a table (see Figure 7(b)).

3.2.3 Profiling Coefficient β . From Equation (6), we should profile $\beta_{fwd}^{key}(algo, d)$, $\beta_{bwd}^{key}(algo, d)$, $\beta_{fwd}^{non}(d)$ and $\beta_{bwd}^{non}(d)$ in advance for runtime training latency estimation. The profiling is performed for all device types in the federation. Furthermore, we need to profile each layer algorithm for each key layer $\beta_{fwd}^{key}(algo, d)$ and $\beta_{bwd}^{key}(algo, d)$. Recall that $\beta(algo, d) = C(algo)/r(d)$ without coupling between $algo$ and d . So, the profiling overhead scales linearly with $\max\{algo, d\}$ rather than their product, which is more manageable. We reuse the model family in Table 2 and extract 1,000 layer configurations to profile the coefficients in the current design of LATTE. We manage β by periodically monitoring processor utilization and updating β by multiplying the new utilization to avoid the overhead of re-estimating it.

3.3 Sub-Model Allocator

With accurate training latency estimation, we describe how LATTE allocates sub-models to devices to ensure synchronized local model training in each round for HFL systems.

Unlike the previous HFL schemes [14, 23], where clients *passively* train server-assigned sub-models, clients in LATTE can *actively* determine sub-models to better suit their own

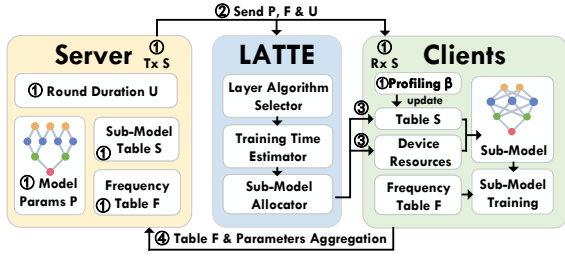


Figure 8: Overall HFL workflow with LATTE.

resource dynamics due to LATTE’s local training latency estimation capabilities. In the following, we first describe the entire client-server interaction in LATTE (§ 3.3.1) and then detail our client-side sub-model allocation design (§ 3.3.2).

3.3.1 Overall Client-Server Interactions. Figure 8 shows the overall HFL workflow with LATTE.

① On initialization, the server determines the round duration U , the global model architecture M , its complete model parameters P , and the frequency table F (introduced and used in §3.3.2). The server then generates candidate sub-model configurations (200 configurations in our case, which is a trade-off between the granularity of sub-model assignment and the estimation overhead, *e.g.*, it takes about 40 seconds to estimate 200 sub-models on a Jetson Nano device, and we will automate the selection of this number in the future.) via model scaling, storing them into the sub-model table S . Note that other methods of generating candidate architectures are also applicable. Afterwards, the server sends the sub-model configuration table S for all clients. Each client profiles its layer algorithm- and device-dependent coefficients β , then updates the sub-model table S , estimates its training latency and measures its memory usage. These estimation can be performed quickly and are a one-time effort.

② In each round, the server sends the current complete model parameters P , round duration U , and frequency table F to all clients for local training.

③ Each client polls the current memory budget mem and bandwidth B via the bandwidth monitoring tool *e.g.*, `bmon` [2]. It then calculates the overall training time T_{train} using the single-pass training latency T stored in the sub-model S and the communication latency T_{comm} based on the bandwidth B and the model size [22]. Each client then selects the largest sub-model configuration in S such that $T_{train} + T_{comm} \leq U$, instantiates the sub-model from the full model parameters P and updates the frequency table F , following the scheme in § 3.3.2, and then starts training on its local dataset.

④ The local parameters and frequency table F are sent back to the server for aggregation. We apply the standard aggregation scheme as in FedAvg [49]. Since the client submits a sub-model to the server, only the updated parameters are averaged [14, 22, 23, 74]. Due to accurate training time

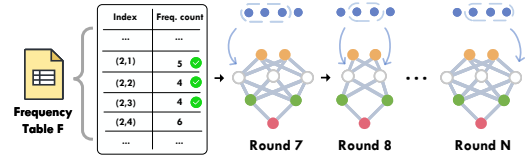


Figure 9: Parameter prioritizing via frequency table F .

estimates, all devices are expected to return updated model parameters within the deadline U . Steps ②-④ are iterated until the global model converges.

3.3.2 Sub-Model Allocation at Client. As mentioned above, LATTE allows dynamic client-side sub-model allocation to adapt to resource dynamics. However, such adaptation may affect the effectiveness of federated training. This is because if sub-model allocation takes resource dynamics into account, parameters in the global model may not be equally trained. In LATTE, we design a *priority-aware sliding window* scheme for sub-model allocation, which facilitates efficient training while maintaining adaptability to resource dynamics.

Our main idea is to train a subset of model parameters in rounds on a rolling basis, so that all parameters are trained evenly in the long run. Previous work [14] implemented this idea, and it was also shown to converge [76]. This approach uses a sliding window for each layer of fixed size and stride length, assuming that the sub-model size is the same across rounds, but this is invalid in our problem as the sub-model size on the client may change due to available resources.

Method. Since LATTE selects a different sub-model architecture in each round based on current resource availability to meet local training deadlines, the sub-model architecture naturally translates into a set of windows per layer, the size of which can vary from round to round. Therefore, the challenge is to ensure even training of all parameters under different window sizes over rounds. Our solution is to *prioritize* rarely trained model parameters into a sliding window for training. This is achieved by tracking how often the model parameters are trained. Specifically, we maintain a parameter trained frequency table F on each device, recording the training frequency of each parameter. The server also maintains a general table F_{server} . Only when the sub-model architecture changes due to resource dynamics, the client will upload its parameter trained frequency table to the server for aggregation and then distribute to each client. The size of each table is small, *e.g.*, only 2.3–24.5 MB, and the communication overhead of maintaining these tables is negligible. Figure 9 shows our parameter prioritizing process across rounds.

4 Evaluation

We implement LATTE using Python 3.6, C++ 14, CUDA 10.2 [4], PyTorch 1.6.0 [9] and TensorFlow 1.14.0. [13]. To train

the selector, we generate large amount of key layer configurations and then conduct benchmarking to obtain the ground truth with our acquisition tool. The development of acquisition tool is supported by cuDNN [5] and is compiled using NVCC [8]. Then LATTE overrides the layer algorithm selection module [10, 11] of deep learning frameworks with its layer algorithm selector output. For practical HFL scenarios and evaluation, LATTE is integrated into the advanced and mainstream FL framework Flower [16].

4.1 Experimental Setups

4.1.1 Tasks, Datasets, and ML Models. We test on IID datasets and real-world non-IID datasets. Each dataset corresponds to a task, and we test different models for each task.

i) Tasks with IID Datasets.

- **Image Classification.** We use CIFAR-10 and CIFAR-100 [3] and train MobileNet-V2 [57] and ResNet-50 [25] respectively for these tasks.
- **Human Activity Recognition.** We utilize the classical HAR [15] dataset, and train a simple customized CNN with 3 convolutional layers and 2 dense layers.

ii) Tasks with non-IID Datasets.

- **Image Classification.** We utilize the OpenImage [33] dataset, where we select 50 image classes, and train MobileNet-V2 for this task.
- **Speech Recognition.** We choose the Google Speech dataset [66], and train ResNet-50 for this task.
- **Human Activity Recognition.** We choose the HAR-Box [52] dataset and adopt the same pre-processing method as [37]. We train the same customized CNN as in i) the IID scenario for this task.

4.1.2 Baselines. We compare LATTE with the following:

- **HeteroFL [23]:** a classic HFL method that assigns different sub-models to devices according to their computing capabilities.
- **FedRolex [14]:** a state-of-the-art HFL scheme with a sliding window sub-model allocation mechanism to improve the training accuracy.
- **FedAvg [49]:** the classic generic FL algorithm without accounting for system heterogeneity.
- **TailorFL [22]:** the state-of-the-art FL method which considers both data and system heterogeneity.
- **Hermes [35]:** a state-of-the-art personalized FL method to handle data heterogeneity.
- **FedSEA [59]:** the state-of-the-art semi-asynchronous FL method for edge devices.
- **FedAsync [69]:** a classical asynchronous FL method.

4.1.3 Test-bed. We evaluate on a client pool of 10 devices, consisting of three tiers as shown in Table 3. We include 2

Category	Type	RAM	CPU	GPU
Higher-end	Dell Inspiron 5577	8GB	i5-7300HQ	GTX 1050
Mid-tier	Jetson Xavier NX	8GB	ARMv8.2	Pascal
	Jetson TX2	8GB	ARM A57	Volta
Lower-end	Jetson Nano	4GB	ARM A57	Maxwell
	Raspberry Pi 4B	8GB	ARM A72	–

Table 3: Device configuration in the client pool.

devices per type, running different frameworks, *i.e.*, TensorFlow for one and PyTorch for the other. The central server is a computer equipped with Intel i7-9700K CPU and NVIDIA RTX 2080Ti GPU. Default bandwidth is 100Mbps and we use Wondershaper [12] to control the bandwidth of each client. Device processor frequency is fixed during evaluation.

4.1.4 Evaluation Metrics. We compare different methods using the following metrics.

- **Time-to-Convergence.** It is the actual wall-clock time for training a model till convergence. Not all methods will converge within a specified time limit.
- **Model Test Accuracy.** It is the accuracy on the test datasets obtained by the model trained through FL.
- **Estimator Precision.** It is the precision between the estimated latency and the actual training time, *i.e.*,

$$\text{Estimate Precision} = 1 - \left| \frac{T_{\text{estimate}} - T_{\text{real}}}{T_{\text{real}}} \right|, \quad (7)$$

where T_{estimate} is from the training time estimator (§ 3.2) of LATTE, and T_{real} is the real wall-clock training time. A high precision means an accurate estimator.

4.2 Overall Performance

4.2.1 Performance on IID Datasets. This experiment compares LATTE with the-state-of-the-art HFL schemes to handle system heterogeneity alone.

Setups. We compare LATTE with HeteroFL [23] and FedRolex [14], two representative HFL schemes that allocate diverse sub-models to devices according to their capabilities. To isolate the impact of system heterogeneity, we conduct experiments on three IID datasets (CIFAR-10, CIFAR-100, and HAR). We also include FedAvg [49] as the baseline that does not explicitly address system heterogeneity.

Results. Figure 10 shows the results on the three IID datasets.

Time-to-Convergence. LATTE consistently achieves the fastest convergence. When training MobileNet on CIFAR-10, LATTE converges in 1.19 hours, 2.36× and 2.76× faster than FedRolex and HeteroFL, respectively. FedAvg fails to converge within the 20-hour limit, highlighting the necessity to explicitly deal with system heterogeneity. When training ResNet-50 on CIFAR-100, LATTE converges in 3.03 hours, shortening the convergence by 2.03× and 2.2× than FedRolex and HeteroFL, respectively. LATTE’s time-to-convergence is

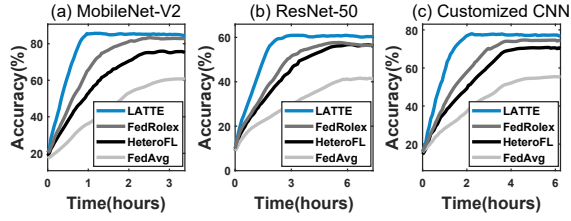


Figure 10: Overall performance on IID datasets.

2.27 hours when training Customized CNN on HAR, 2.52× and 2.64× faster than FedRolex and HeteroFL, respectively.

Model Test Accuracy. LATTE also achieves the highest test accuracy at convergence, reaching an accuracy of 85.72%, 60.94%, and 78.00% on MobileNet, ResNet, and Customized CNN respectively, which is 2.44% to 3.33% higher than FedRolex, and 4.19% to 9.78% higher than HeteroFL.

Figure 10 shows that LATTE notably improves both the efficiency and effectiveness of federated learning under system heterogeneity than the state-of-the-art HFL methods.

4.2.2 Performance on non-IID Datasets. This experiment aims to further evaluate the effectiveness of LATTE under both data and system heterogeneity. As a transparent design, LATTE can be easily integrated with FL solutions that handle data heterogeneity to improve their efficiency.

Setups. We integrate LATTE into Hermes [35] and TailorFL [22], two representative FL schemes that combat data heterogeneity by training personalized models. We test their performance on three non-IID datasets (OpenImage, Google Speech, and HARBox) with heterogeneous edge devices with and w/o the enhancement of LATTE.

Results. Figure 11 shows performance gains with LATTE.

Time-to-Convergence. The personalized FL methods converges notably faster after adding LATTE. When training MobileNet on OpenImage, the LATTE-enhanced TailorFL and Hermes converge in 4.50 hours and 7.58 hours respectively, which are 2.36× and 1.25× faster than the standalone versions. When training ResNet-50 on Google Speech, LATTE+TailorFL converges 1.81× faster than TailorFL, and LATTE+Hermes is 1.19× faster than Hermes. When training Customized CNN on HARBox, TailorFL and Hermes are accelerated by 2.1× and 1.17×, respectively.

Model Test Accuracy. Integration with LATTE also improves the test accuracy. LATTE improves the accuracy of TailorFL by up to 2.83%, reaching 48.34%, 52.78%, and 62.16% on MobileNet, ResNet, and Customized CNN, respectively. Similarly, LATTE improves the accuracy of Hermes by up to 6.38%, reaching 41.96%, 46.68%, and 56.17% on MobileNet, ResNet, and Customized CNN, respectively.

These results affirm that LATTE can complement existing FL methods for data heterogeneity, achieving both accurate and fast training under both data and system heterogeneity.

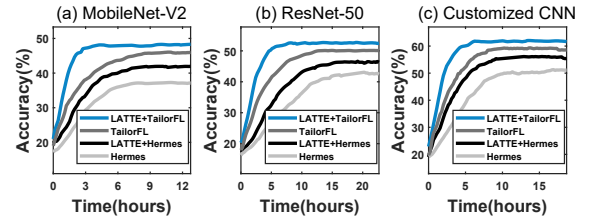


Figure 11: Overall performance on non-IID Datasets.

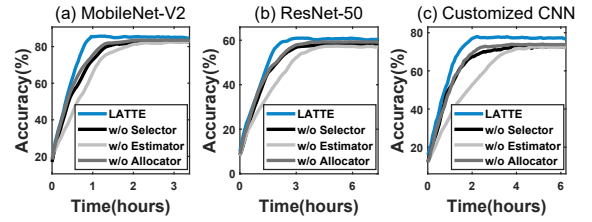


Figure 12: Contributions of individual modules.

4.3 Ablation Study

4.3.1 Performance Breakdown. This experiment evaluates the contributions of each module to the overall performance.

Setups. We implement three variants of LATTE:

- **LATTE w/o Selector:** it replaces the layer algorithm selector (§ 3.1) with the default APIs [10, 11] provided by deep learning frameworks.
- **LATTE w/o Estimator:** it replaces the training time estimator (§ 3.2) with the naive modeling in Equation (2).
- **LATTE w/o Allocator:** it replaces the sub-model allocator (§ 3.3) with the sliding window scheme in [14].

Results. Figure 12 shows the results of different variants.

Time-to-Convergence. The variant w/o predictor suffers from notable slowdown (0.5× to 0.57× that of LATTE). This is because the default APIs may not recommend the fastest layer algorithm, leading to inaccurate time estimates. The variant w/o estimator also experiences significant slowdown (0.36× to 0.57× that of LATTE). This is because the naive latency modeling ignores the impact of layer algorithms, which is the key to inconsistent training time estimates. The variant w/o selector shows a mild slowdown (0.58× to 0.67× that of LATTE), mainly due to the lack of adaptation to resource dynamics. This can lead to sub-model mismatch when resources on the device change frequently.

Model Test Accuracy. Overall, the degradation in test accuracy is less severe than the degradation in convergence speed. On the three datasets, for the variants w/o selector, estimator, and allocator, the accuracy drops by 2.67% to 4.66%, 2.97% to 5.32%, and 1.64% to 4.00%, respectively.

In summary, removing individual modules from LATTE mainly affects the convergence speed. The convergence slowdown is less drastic for the variants w/o selector and w/o allocator than that w/o estimator. This is because the default

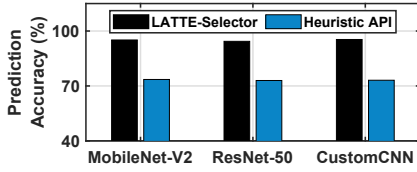


Figure 13: Effectiveness of the layer algorithm selector.

APIs still provide about 73% layer algorithm prediction accuracy, and system resources do not vary all the time. In contrast, the variant w/o estimator falls back to the HFL scheme ignoring the diversity of layer algorithms, which is the primary motivation of our work.

4.3.2 Effectiveness of Layer Algorithm Selector. A key to the convergence speedup of LATTE is the layer algorithm selector (§ 3.1). This experiment zooms into its performance against the default APIs in deep learning frameworks.

Setups. We compare the accuracy of our layer algorithm selector with the heuristic API [10] provided by PyTorch for the 200 sub-models listed in the sub-model table on the three IID datasets, and use the results from the on-device measurement API [11] of TensorFlow as ground truth (§ 3.1).

Results. In Figure 13, our layer algorithm selector achieves 95.11%, 94.37%, and 95.38% accuracy on MobileNet, ResNet and Customized CNN, respectively. Heuristic method only provides 73.54%, 72.98%, and 73.23% accuracy respectively, which is not sufficient for subsequent training time estimates.

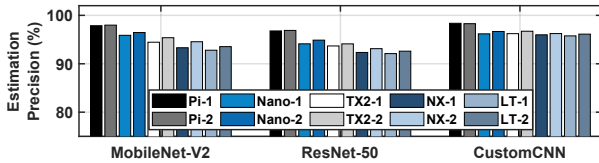


Figure 14: Effectiveness of the training time estimator.

4.3.3 Effectiveness of Training Time Estimator. The training time estimator (§ 3.2) is also essential for the convergence speedup of LATTE. This experiment aims to show the applicability of our training time estimator across diverse devices.

Setups. We measure the estimator precision of the training latencies on ten different devices using the three IID datasets. It covers three device types, with two devices per type.

Results. In Figure 14, our training latency estimator leads to a precision of 92.80–97.98%, 92.10–96.88%, and 95.76–98.36% on MobileNet, ResNet and Customized CNN, respectively. The estimator precision of ResNet-50 on CIFAR-100 is slightly lower due to the more complex network architectures. In contrast, the precision of Customized CNN on HAR is higher, as there are fewer key layers. However, as shown in § 4.2, the precision of our training time estimator already significantly improves the overall training convergence.

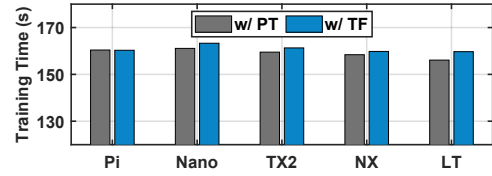


Figure 15: LATTE alleviates the variations in training time across deep learning frameworks.

4.4 Micro-benchmarks

4.4.1 Impact of Deep Learning Frameworks. This experiment validates that the layer algorithm is the key reason for variations in observed training times due to development-chain diversity, which can be resolved by LATTE. Specifically, we integrate LATTE with either TensorFlow or PyTorch (*i.e.*, replacing their default APIs) and measure the training latency of the same model on two identical devices, one developed with TensorFlow and the other with PyTorch. As shown in Figure 15, the wall-clock training times on the two devices using different deep learning frameworks are similar to each other. Recall that the wall-clock training time of a single epoch round with TensorFlow is only 72.3–76% of that of PyTorch without LATTE calibration (see § 2.2).

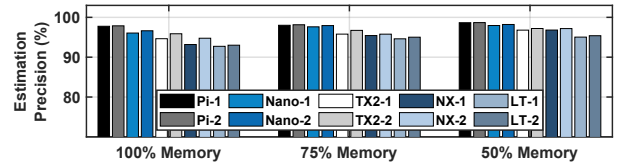


Figure 16: Impact of the memory resource budget.

4.4.2 Impact of Dynamic Memory. This experiment tests LATTE under various runtime memory budgets. We employ memhog [6] to periodically change the available memory resources on each device, cycling from 100% available memory to 75% to 50%, and measure the precision of our training time estimator. Figure 16 shows that when the memory budget decreases, the estimation precision increases. This is because a tighter memory budget means a simpler sub-model should be selected. In Figure 18(a), the convergence time under dynamic memory is 2.18× that of the normal state, and the model test accuracy decreases by 3.95%. This is because the drastic change in available memory force clients to conservatively select smaller sub-models, leading to a decrease in convergence speed and model test accuracy.

4.4.3 Impact of Dynamic Bandwidth. This experiment tests LATTE under different bandwidths. We use Wondershaper [12] to periodically change the communication bandwidth on each device, cycling from 100 Mbps to 50 Mbps to 10 Mbps. Figure 17 shows that the precision of training time estimator increases when the bandwidth decreases. This is because,

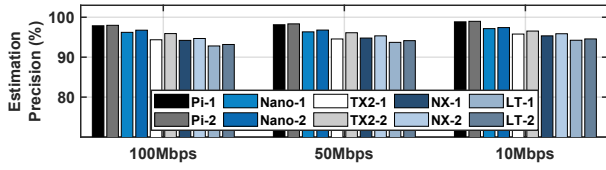


Figure 17: Impact of the communication bandwidth.

given a fixed deadline, as the communication latency increases, the budget for model training is shorter, causing LATTE to choose simpler sub-models and thus estimate the training time more accurately. Figure 18(a) shows that under dynamic bandwidth conditions, both the convergence speed and model test accuracy are similar to the normal state. The convergence time increases slightly to 1.09 \times of the stable state, and the model test accuracy decreases by 1.32%.

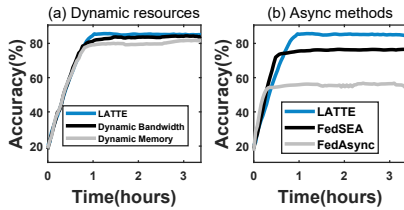


Figure 18: Time-to-accuracy (a) under dynamic resources and (b) compared between different schemes.

4.4.4 Comparison with Semi- and Fully-Asynchronous FL. This experiment compares the performance of LATTE with FedAsync [69], a fully asynchronous FL method, and FedSEA [59], a state-of-the-art semi-asynchronous FL method. Figure 18(b) shows that both FedSEA and FedAsync converge faster than LATTE, since they do not wait for slow devices. However, their test accuracy is lower for aggregating slate model updates. For example, the test accuracy of FedSEA decreases by 9.19%, while the test accuracy of FedAsync suffers from a drastic drop by 28.96%.

4.4.5 Impact of Client Heterogeneity Distribution. To understand this impact, we conduct a simulation-based evaluation. We simulate 100 clients using a server with 24 Xeon-Gold-6142 CPUs and 4 Tesla V100 32GB GPUs. All clients are initialized by Flower’s VCE (Virtual Client Engine) function, half of which have TensorFlow installed and the other half have PyTorch installed. Clients are divided into three categories: low-end, mid-tier, and high-end, with a computational power ratio of 1:2:5 respectively. In this experiment, we train MobileNet on CIFAR-10 and perform the evaluation using three common device distributions, as shown in Figure 19.

Results. Figure 20(a) shows that under the uniform distribution, LATTE converges 2.07 \times and 2.35 \times faster than FedRolex and HetetoFL, respectively. Its test accuracy is also 2.3% and 9.36% higher than FedRolex and HetetoFL, respectively.

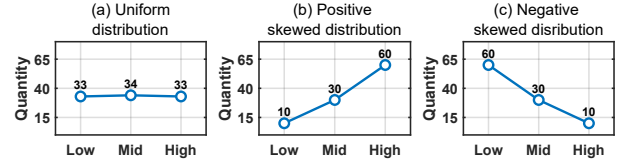


Figure 19: Three distributions used in the simulation.

When the devices follow a positively skewed distribution, Figure 20(b) shows that all three methods can accelerate convergence and improve test accuracy. However, LATTE still outperforms the two baselines, with 1.69–2.15 \times faster convergence and 2.65–6.35% higher test accuracy. This is because more clients have sufficient computational power, allowing them to select larger and better sub-models, thus contributing more to the overall system performance.

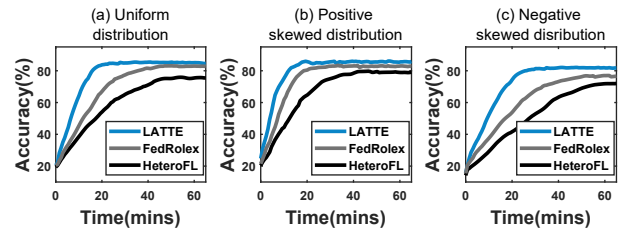


Figure 20: Performance comparison under different device distributions in the simulation.

When the devices follow a negatively skewed distribution, LATTE still maintains its advantages. Compared with the two baselines, LATTE converges 1.89–2.23 \times faster and improves test accuracy by 5.09–10.1% in Figure 20(c). Under this distribution, more clients have insufficient computational power, forcing them to select smaller sub-models. However, LATTE’s more accurate training time estimation enables it to select larger sub-models, leading to better performance. Overall, LATTE achieves faster convergence and higher model accuracy under different device heterogeneity distributions.

	ResNet-50	MobileNetV2	Customized CNN
Selector	191KB	191KB	191KB
Sub-model Table	509.6KB	197.6KB	52KB
Frequency Table	24.5MB	3.4MB	2.3MB

Table 4: LATTE memory overhead on different models.

4.5 System Overhead

Memory Overhead. LATTE introduces three types of memory overhead, including the parameters of selector, sub-model table, and frequency table, as shown in Table 4. In particular, each model used requires 191 KB memory of selector parameters. Sub-model tables take 52 KB – 509.6 KB, on ResNet, MobileNet and Customized CNN, respectively. Frequency tables are compressed with 8-bit, which occupy 2.3 MB – 24.5

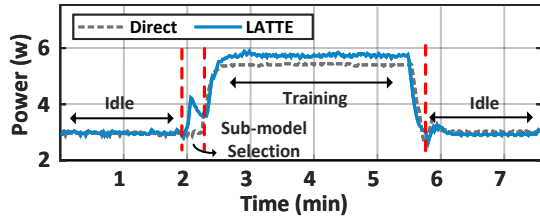


Figure 21: Power consumption of LATTE.

MB on the aforementioned models in the same order. Such memory overhead is also considered within the memory budget during sub-model selection.

Power Consumption. Low power consumption is important to mobile edge devices, and we measure the power consumption of LATTE on Jetson NX devices using the INA3221 power monitor [61]. As a baseline, we first measure the power consumption of direct training. We then measure LATTE for same duration of time. Figure 21 shows that the idle state of device consumes about 3 W. Before training, the power consumption of LATTE is slightly higher than the baseline, reaching 4.28 W, which is due to the fast sub-model allocation. During training, the power consumption is 5.84 W (by LATTE) and 5.6 W (by direct training), where LATTE itself consumes about 0.24 W only during resource polling.

5 Related Work

Heterogeneous Federated Learning. Federated learning with edge and IoT devices face the challenges of *data* and *system* heterogeneity [73]. Compared to centralized learning, data heterogeneity affects model accuracy [45, 77] and is often resolved by training personalized models [22, 36, 52, 60, 62, 75]. However, data heterogeneity is not the only factor affecting the model performance, and system heterogeneity also has a significant impact. System heterogeneity comes from the diverse computation and communication capabilities of devices, which can affect the training efficiency and is the focus of heterogeneous federated learning (HFL) [53].

HFL roughly falls into *model* and *system* level solutions. Model-level approaches [14, 22, 23, 27, 43, 51, 55, 58] assign sub-models tailored to the computation power of each device so that they can return the locally trained models almost simultaneously. The sub-models are extracted via masking [14, 22, 23, 51, 56], dropout [27, 55], knowledge distillation [38, 78], *etc.* Orthogonally, system-level strategies either perform client selection [34, 37, 40, 50] or adopt semi-asynchronous [47, 59, 67] and fully asynchronous [19, 30, 74] model aggregation schemes to exclude or mitigate the impact of slow devices. Our work targets at model-level HFL solutions because client selection might ignore valuable data on slow devices whereas asynchronous model aggregation might affect model convergence. We focus on masking-based sub-model

extraction for its widespread adoption and provide accurate local training time estimates so that the allocated sub-models match with the capabilities of devices. With the LATTE design, future methods to address data heterogeneity can be directly integrated to further improve FL performance.

Model Latency Estimation. Evaluating the latency of model execution on specific devices is crucial for the optimization of model inference [20] and training [31]. Due to various model architectures, device types, and deep learning development chains, there is an increasing interest to predict the execution latency rather than measuring it exhaustively. These predictors model latency at the *network*, *layer*, or *operator* level [41]. For example, the number of FLOPs or MAC of the entire network is a common proxy for its latency [26, 42]. Mainstream predictors [17, 32, 54] resort to the layer level, where different layer features (*e.g.*, FLOPs or layer types) and hardware features are leveraged to train a regressor to predict layer-wise latencies, which are then summed up as the overall model latency. Recent predictors [72] dive into the operator level to explicitly account for runtime optimizations such as operator fusion. However, all these studies are designed for latency prediction of model *inference*.

Our work is inspired by them yet aims at accurate latency prediction of model *training*. ElasticTrainer[29] proposed more fine-grained modeling of training time, but still overlooked the diversity of layer algorithm. By identifying layer algorithms as the previously overlooked feature, we devise a lightweight training latency estimator at the layer level. Our evaluations show such layer-level modeling matches with the current runtime optimization for model training on edge devices and delivers high accuracy despite its simplicity.

6 Conclusion

This paper presents LATTE, a new middleware design for accurate estimation of on-device training of deep learning models on mobile edge devices. Our core design insight is that even for the same model on the same device, training times can vary significantly due to runtime optimizations of deep learning frameworks. We solve this problem by designing a novel layer algorithm selector and incorporating it into LATTE for accurate delay estimation. We further showcase the usability of our design in heterogeneous federated learning. Extensive experiments demonstrate significant performance improvements compared to state-of-the-art methods.

Acknowledgments

We thank all reviewers. This work is supported by the GRF grant from Research Grants Council of Hong Kong (CityU 11202623 and CityU 11205624) and the APRC grant from City University of Hong Kong (Project No. 9610633). Zimu Zhou and Zhenjiang Li are the corresponding authors.

References

- [1] Arm compute library. <https://github.com/ARM-software/ComputeLibrary>.
- [2] bmon bandwidth monitor. <https://github.com/tgraf/bmon>.
- [3] Cifar datasets. <https://www.cs.toronto.edu/~kriz/cifar.html>.
- [4] CUDA Toolkit Documentation. <https://docs.nvidia.com/cuda/cuda-runtime-api/>.
- [5] cudnn. <https://developer.nvidia.com/cudnn>.
- [6] memhog. <https://github.com/afeinberg/memhog>.
- [7] Mkl-dnn. <https://oneapi-src.github.io/oneDNN/v0/index.html>.
- [8] NVIDIA CUDA Compiler Driver NVCC. <https://docs.nvidia.com/cuda/cuda-compiler-driver-nvcc/>.
- [9] PyTorch. <https://github.com/pytorch/pytorch>.
- [10] Pytorch runtime optimization. https://github.com/pytorch/pytorch/blob/main/aten/src/ATen/native/cudnn/Conv_v7.cpp.
- [11] Tensorflow runtime optimization. https://github.com/tensorflow/tensorflow/blob/v1.14.0/tensorflow/stream_executor/cuda/cuda_dnn.cc.
- [12] Wondershaper. <https://github.com/magnific0/wondershaper>.
- [13] Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: a system for Large-Scale machine learning. In *Proc. of USENIX OSDI*, 2016.
- [14] Samiul Alam, Luyang Liu, Ming Yan, and Mi Zhang. Fedrolex: Model-heterogeneous federated learning with rolling sub-model extraction. *Advances in Neural Information Processing Systems*, 2022.
- [15] Davide Anguita, Alessandro Ghio, Luca Oneto, Xavier Parra, and Jorge Luis Reyes-Ortiz. A public domain dataset for human activity recognition using smartphones. In *Proc. of ESANN*, 2013.
- [16] Daniel J Beutel, Taner Topal, Akhil Mathur, Xinchu Qiu, Javier Fernandez-Marques, Yan Gao, Lorenzo Sani, Hei Li Kwing, Titouan Parcollet, Pedro PB de Gusmão, and Nicholas D Lane. Flower: A friendly federated learning research framework. *arXiv preprint arXiv:2007.14390*, 2020.
- [17] Ermao Cai, Da-Cheng Juan, Dimitrios Stamoulis, and Diana Marculescu. Neuralpower: Predict and deploy energy-efficient convolutional neural networks. In *Proc. of PMLR ACML*, 2017.
- [18] Jiani Cao, Chengdong Lin, Yang Liu, and Zhenjiang Li. Gaze tracking on any surface with your phone. In *Proc. of ACM SenSys*, 2022.
- [19] Ming Chen, Bingcheng Mao, and Tianyi Ma. Efficient and robust asynchronous federated learning with stragglers. In *Proc. of ICLR*, 2019.
- [20] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Tvm: An automated end-to-end optimizing compiler for deep learning. In *Proc. of USENIX OSDI*, 2018.
- [21] Allan de Barcelos Silva, Marcio Miguel Gomes, Cristiano André da Costa, Rodrigo da Rosa Righi, Jorge Luis Victoria Barbosa, Gustavo Pessin, Geert De Doncker, and Gustavo Federizzi. Intelligent personal assistants: A systematic literature review. *Elsevier Expert Systems with Applications*, 2020.
- [22] Yongheng Deng, Weining Chen, Ju Ren, Feng Lyu, Yang Liu, Yunxin Liu, and Yaoxue Zhang. Tailorfl: Dual-personalized federated learning under system and data heterogeneity. In *Proc. of ACM SenSys*, 2022.
- [23] Enmao Diao, Jie Ding, and Vahid Tarokh. Heterofl: Computation and communication efficient federated learning for heterogeneous clients. In *Proc. of ICLR*, 2020.
- [24] In Gim and JeongGil Ko. Memory-efficient dnn training on mobile devices. In *Proc. of ACM MobiSys*, 2022.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proc. of IEEE CVPR*, 2016.
- [26] Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. Amc: Automl for model compression and acceleration on mobile devices. In *Prof. of Springer ECCV*, 2018.
- [27] Samuel Horvath, Stefanos Laskaridis, Mario Almeida, Ilias Leontiadis, Stylianos Venieris, and Nicholas Lane. Fjord: Fair and accurate federated learning under heterogeneous targets with ordered dropout. *Advances in Neural Information Processing Systems*, 2021.
- [28] Kai Huang and Wei Gao. Real-time neural network inference on extremely weak devices: agile offloading with explainable ai. In *Proc. of ACM MobiCom*, 2022.
- [29] Kai Huang, Boyuan Yang, and Wei Gao. Elastictrainer: Speeding up on-device training with runtime elastic tensor selection. In *Proc. of ACM MobiSys*, 2023.
- [30] Dzmitry Huba, John Nguyen, Kshitiz Malik, Ruiyu Zhu, Mike Rabbat, Ashkan Yousefpour, Carole-Jean Wu, Hongyuan Zhan, Pavel Ustinov, Harish Srinivas, Kaikai Wang, Anthony Shoumikhin, Jesik Min, and Mani Malek. Papaya: Practical, private, and scalable federated learning. In *Proc. of MLSys*, 2022.
- [31] Zhihao Jia, Matei Zaharia, and Alex Aiken. Beyond data and model parallelism for deep neural networks. In *Proc. of MLSys*, 2019.
- [32] Daniel Justus, John Brennan, Stephen Bonner, and Andrew Stephen McGough. Predicting the computational cost of deep learning models. In *Proc. of IEEE BigData*, 2018.
- [33] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Mallocci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *Springer International Journal of Computer Vision*, 2020.
- [34] Fan Lai, Xiangfeng Zhu, Harsha V Madhyastha, and Mosharaf Chowdhury. Oort: Efficient federated learning via guided participant selection. In *Proc. of USENIX OSDI*, 2021.
- [35] Ang Li, Jingwei Sun, Pengcheng Li, Yu Pu, Hai Li, and Yiran Chen. Hermes: an efficient federated learning framework for heterogeneous mobile clients. In *Proc. of ACM MobiCom*, 2021.
- [36] Ang Li, Jingwei Sun, Xiao Zeng, Mi Zhang, Hai Li, and Yiran Chen. Fedmask: Joint computation and communication-efficient personalized federated learning via heterogeneous masking. In *Proc. of ACM SenSys*, 2021.
- [37] Chenning Li, Xiao Zeng, Mi Zhang, and Zhichao Cao. Pyramidfl: A fine-grained client selection framework for efficient federated learning. In *Proc. of ACM MobiCom*, 2022.
- [38] Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, 2019.
- [39] Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE signal processing magazine*, 2020.
- [40] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In *Proc. of MLSys*, 2020.
- [41] Ying Li, Yifan Sun, and Adwait Jog. Path forward beyond simulators: Fast and accurate gpu execution time prediction for dnn workloads. In *Proc. of IEEE/ACM MICRO*, 2023.
- [42] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *Proc. of ICLR*, 2019.
- [43] Ruixuan Liu, Fangzhao Wu, Chuhan Wu, Yanlin Wang, Lingjuan Lyu, Hong Chen, and Xing Xie. No one left behind: Inclusive federated learning over heterogeneous devices. In *Proc. of KDD*, 2022.

- [44] Sicong Liu, Bin Guo, Cheng Fang, Ziqi Wang, Shiyao Luo, Zimu Zhou, and Zhiwen Yu. Enabling resource-efficient aiot system with cross-level optimization: A survey. *IEEE Communications Surveys & Tutorials*, 2023.
- [45] Zili Lu, Heng Pan, Yueyue Dai, Xueming Si, and Yan Zhang. Federated learning with non-iid data: A survey. *IEEE Internet of Things Journal*, 2024.
- [46] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. *Advances in Neural Information Processing Systems*, 2017.
- [47] Qianpiao Ma, Yang Xu, Hongli Xu, Zhida Jiang, Liusheng Huang, and He Huang. Fedsea: A semi-asynchronous federated learning mechanism in heterogeneous edge computing. *IEEE Journal on Selected Areas in Communications*, 2021.
- [48] Patrick McEnroe, Shen Wang, and Madhusanka Liyanage. A survey on the convergence of edge computing and ai for uavs: Opportunities and challenges. *IEEE Internet of Things Journal*, 2022.
- [49] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguerre y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proc. of PMLR AISTATS*, 2017.
- [50] Takayuki Nishio and Ryo Yonetani. Client selection for federated learning with heterogeneous resources in mobile edge. In *Proc. of IEEE ICC*, 2019.
- [51] Chaoyue Niu, Fan Wu, Shaojie Tang, Lifeng Hua, Rongfei Jia, Chengfei Lv, Zhihua Wu, and Guihai Chen. Billion-scale federated learning on mobile clients: A submodel design with tunable privacy. In *Proc. of ACM MobiCom*, 2020.
- [52] Xiaomin Ouyang, Zhiyuan Xie, Jiayu Zhou, Jianwei Huang, and Guoliang Xing. Clusterfl: a similarity-aware federated learning system for human activity recognition. In *Proc. of ACM MobiSys*, 2021.
- [53] Kilian Pfeiffer, Martin Rapp, Ramin Khalili, and Jörg Henkel. Federated learning for computationally-constrained heterogeneous devices: A survey. *ACM Computing Surveys*, 2023.
- [54] Hang Qi, Evan R Sparks, and Ameet Talwalkar. Paleo: A performance model for deep neural networks. In *Proc. of ICLR*, 2016.
- [55] Xinchu Qiu, Javier Fernandez-Marques, Pedro PB Gusmao, Yan Gao, Titouan Parcollet, and Nicholas Donald Lane. Zerofl: Efficient on-device training for federated learning with local sparsity. In *Proc. of ICLR*, 2021.
- [56] Martin Rapp, Ramin Khalili, Kilian Pfeiffer, and Jörg Henkel. Distreal: Distributed resource-aware learning in heterogeneous systems. In *Proc. of AAAI*, 2022.
- [57] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proc. of IEEE CVPR*, 2018.
- [58] Leming Shen, Qiang Yang, Kaiyan Cui, Yuanqing Zheng, Xiao-Yong Wei, Jianwei Liu, and Jinsong Han. Fedconv: A learning-on-model paradigm for heterogeneous federated clients. In *Proc. of ACM MobiSys*, 2024.
- [59] Jingwei Sun, Ang Li, Lin Duan, Samiul Alam, Xuliang Deng, Xin Guo, Haiming Wang, Maria Gorlatova, Mi Zhang, Hai Li, and Yiran Chen. Fedsea: A semi-asynchronous federated learning framework for extremely heterogeneous devices. In *Proc. of ACM SenSys*, 2022.
- [60] Alysa Ziyang Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [61] Texas Instruments Inc. INA3221 power monitor. <https://www.ti.com/product/INA3221>, 2016.
- [62] Linlin Tu, Xiaomin Ouyang, Jiayu Zhou, Yuze He, and Guoliang Xing. Feddl: Federated learning via dynamic layer sharing for human activity recognition. In *Proc. of ACM SenSys*, 2021.
- [63] Kun Wang, Jiani Cao, Zimu Zhou, and Zhenjiang Li. Swapnet: Efficient swapping for dnn inference on edge ai devices beyond the memory budget. *IEEE Transactions on Mobile Computing*, 2024.
- [64] Manni Wang, Shaohua Ding, Ting Cao, Yunxin Liu, and Fengyuan Xu. Asymo: scalable and efficient deep-learning inference on asymmetric mobile cpus. In *Proc. of ACM MobiCom*, 2021.
- [65] Qipeng Wang, Mengwei Xu, Chao Jin, Xinran Dong, Jinliang Yuan, Xin Jin, Gang Huang, Yunxin Liu, and Xuanzhe Liu. Melon: Breaking the memory wall for resource-efficient on-device machine learning. In *Proc. of ACM MobiSys*, 2022.
- [66] Pete Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv preprint arXiv:1804.03209*, 2018.
- [67] Wentai Wu, Ligang He, Weiwei Lin, Rui Mao, Carsten Maple, and Stephen Jarvis. Safa: A semi-asynchronous protocol for fast federated learning with low overhead. *IEEE Transactions on Computers*, 2020.
- [68] Yu Xianjia, Jorge Peña Queralta, Jukka Heikkonen, and Tomi Westerlund. Federated learning in robotic and autonomous systems. *Elsevier Procedia Computer Science*, 2021.
- [69] Cong Xie, Sanmi Koyejo, and Indranil Gupta. Asynchronous federated optimization. *arXiv preprint arXiv:1903.03934*, 2019.
- [70] Huatao Xu, Pengfei Zhou, Rui Tan, Mo Li, and Guobin Shen. Limu-bert: Unleashing the potential of unlabeled data for imu sensing applications. In *Proc. of ACM SenSys*, 2021.
- [71] Qiang Yang, Yang Liu, Tianjian Chen, and Yongxin Tong. Federated machine learning: Concept and applications. *ACM Transactions on Intelligent Systems and Technology*, 2019.
- [72] Li Lina Zhang, Shihao Han, Jianyu Wei, Ningxin Zheng, Ting Cao, Yuqing Yang, and Yunxin Liu. Nn-meter: Towards accurate latency prediction of deep-learning model inference on diverse edge devices. In *Proc. of ACM MobiSys*, 2021.
- [73] Tuo Zhang, Lei Gao, Chaoyang He, Mi Zhang, Bhaskar Krishnamachari, and Salman Avestimehr. Federated learning for the internet of things: Applications, challenges, and opportunities. *IEEE Internet of Things Magazine*, 2022.
- [74] Tuo Zhang, Lei Gao, Sunwoo Lee, Mi Zhang, and Salman Avestimehr. Timelyfl: Heterogeneity-aware asynchronous federated learning with adaptive partial training. In *Proc. of IEEE CVPR*, 2023.
- [75] Wenhao Zhang, Zimu Zhou, Yansheng Wang, and Yongxin Tong. Dm-pfl: Hitchhiking generic federated learning for efficient shift-robust personalization. In *Proc. of KDD*, 2023.
- [76] Hanhan Zhou, Tian Lan, Guru Prasad Venkataramani, and Wenbo Ding. Every parameter matters: Ensuring the convergence of federated learning with dynamic heterogeneous models reduction. *Advances in Neural Information Processing Systems*, 2024.
- [77] Hangyu Zhu, Jinjin Xu, Shiqing Liu, and Yaochu Jin. Federated learning on non-iid data: A survey. *Elsevier Neurocomputing*, 2021.
- [78] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *Proc. of ACM ICML*, 2021.