

---

# Towards Diverse Device Heterogeneous Federated Learning via Task Arithmetic Knowledge Integration

---

Mahdi Morafah<sup>\*1</sup>, Vyacheslav Kungurtsev<sup>2</sup>, Hojin Chang<sup>1</sup>, Chen Chen<sup>3</sup>, Bill Lin<sup>1</sup>

<sup>1</sup>University of California San Diego (UCSD), <sup>2</sup>Czech Technical University in Prague,

<sup>3</sup>University of Central Florida (UCF)

\*Correspondence: mmorafah@ucsd.edu

## Abstract

Federated Learning (FL) has emerged as a promising paradigm for collaborative machine learning, while preserving user data privacy. Despite its potential, standard FL algorithms lack support for diverse heterogeneous device prototypes, which vary significantly in model and dataset sizes—from small IoT devices to large workstations. This limitation is only partially addressed by existing knowledge distillation (KD) techniques, which often fail to transfer knowledge effectively across a broad spectrum of device prototypes with varied capabilities. This failure primarily stems from two issues: the dilution of informative logits from more capable devices by those from less capable ones, and the use of a single integrated logits as the distillation target across all devices, which neglects their individual learning capacities and the unique contributions of each device. To address these challenges, we introduce TAKFL, a novel KD-based framework that treats the knowledge transfer from each device prototype’s ensemble as a separate task, independently distilling each to preserve its unique contributions and avoid dilution. TAKFL also incorporates a KD-based self-regularization technique to mitigate the issues related to the noisy and unsupervised ensemble distillation process. To integrate the separately distilled knowledge, we introduce an adaptive *task arithmetic* knowledge integration process, allowing each student model to customize the knowledge integration for optimal performance. Additionally, we present theoretical results demonstrating the effectiveness of task arithmetic in transferring knowledge across heterogeneous device prototypes with varying capacities. Comprehensive evaluations of our method across both computer vision (CV) and natural language processing (NLP) tasks demonstrate that TAKFL achieves state-of-the-art results in a variety of datasets and settings, significantly outperforming existing KD-based methods. Our code is released at <https://github.com/MMorafah/TAKFL>.

## 1 Introduction

Federated Learning (FL) has rapidly gained traction as a promising approach to train machine learning models collaboratively across multiple devices, while preserving the privacy of user data. Standard federated learning methods, such as FedAvg [33], however, are primarily designed for unrealistic *device-homogeneous* scenarios, where all devices are assumed to have identical compute resource and can train the same neural network architecture [28, 33, 47, 21, 27, 46, 31]. Therefore, standard FL cannot support the participation of *heterogeneous devices*, all of which could significantly contribute to model training due to their unique and invaluable local datasets. To address this gap, knowledge distillation (KD) techniques have emerged as a promising approach to establish federation among heterogeneous device prototypes and facilitate knowledge transfer between them. In this approach, locally updated client models from different device prototypes, collectively termed as ensembles, serve as teachers to distill their knowledge into each device prototype’s server student model using an unlabeled public dataset.

---

This work was supported in part by a research grant from Cisco Systems, Inc.

Despite their success, however, existing KD-based methods for device heterogeneous FL are primarily designed for scenarios where device prototypes are in the same-size with similar capabilities, i.e. same model and dataset sizes. However, in practice, *device capabilities vary widely*, ranging from small devices like IoTs with small models and small datasets to large devices like workstations with large models and large datasets. This diversity, often overlooked in the existing literature, results in device prototypes with varying strengths and information qualities. Unfortunately, existing methods struggle to establish effective knowledge transfer in these challenging, real-world device heterogeneous settings, primarily due to two reasons: ① Existing methods often disregard the individual strengths and information quality of each device prototype’s ensembles and integrate their logits into a single distillation target. *This approach dilutes the richer, more informative logits from larger, more capable devices with less informative logits from smaller, less capable ones.* ② Additionally, these methods employ this single integrated distillation target to transfer knowledge across all different size student models. *This one-size-fits-all approach fails to provide customized knowledge integration based on the unique learning capacities of each student and the specific helpfulness of each device prototype’s ensembles.*

Moreover, the heterogeneous ensemble distillation process can inadvertently lead student models into erroneous learning directions, causing them to forget their self-knowledge acquired through averaged locally updated parameters. This issue arises primarily due to two reasons: ① The distillation process introduces noise, *as the ensembles’ logits are inferred on an unfamiliar public dataset*, distinct from their original training data. Additionally, the presence of data heterogeneity and the insufficient training of some ensembles, due to computational constraints, can further exacerbate this noise. ② The distillation process *lacks supervision from the actual private datasets, which are the ultimate learning objectives*. Consequently, these factors, combined with the limitations outlined earlier, result in *suboptimal knowledge transfer* in device heterogeneous settings. This underscores the urgent need for a more effective knowledge transfer framework.

In this paper, we introduce TAKFL, a novel “*Task Arithmetic Knowledge Transfer Integration for Federated Learning*” framework, designed to overcome the fundamental limitations in the existing methods and improve knowledge transfer in scenarios where device prototypes vary in size—both model and dataset—and consequently, in strength. TAKFL treats knowledge transfer from each device prototype’s ensembles as separate tasks, distilling them independently to ensure that each prototype’s unique contributions are accurately distilled without interference. To tackle the challenges associated with noisy and unsupervised ensemble distillation, we incorporate a KD-based self-regularization technique into this individual knowledge transfer process. Subsequently, to selectively integrate the separately distilled knowledge from heterogeneous prototypes’ ensembles, we introduce an adaptive *task arithmetic* knowledge integration method by extending the notion of task vectors from centralized learning to federated learning. Our approach enables the student model to strategically customize the knowledge integration process based on the quality of knowledge from each prototype’s ensembles and its intrinsic capacity, aiming to achieve optimal performance. We present theoretical results, grounded on the established theoretical learning properties of over-parametrized neural networks, that conceptualize knowledge distillation as the allocation of device prototypes’ capacities to accurately fit the chosen logits. These results demonstrate the advantages of employing task arithmetic for knowledge transfer in terms of overall accuracy, coverage, and efficiency, as well as the adaptive knowledge integration based on the capacity of the student prototype. Furthermore, we comprehensively evaluate our method across both computer vision (CV) and natural language processing (NLP) tasks, utilizing various datasets and architectures, and demonstrate that TAKFL consistently achieves state-of-the-art (SOTA) performance.

The contribution of our paper is as follows:

1. We formalize and review the important considerations of the problem statement of federated learning with heterogeneous device prototypes.
2. We introduce TAKFL, a novel KD-based method designed to overcome the fundamental limitations of existing approaches, effectively facilitating knowledge transfer across diverse heterogeneous device prototypes with varying capabilities.
3. We present a theoretical model for device heterogeneous KD, and demonstrate the effectiveness and efficiency of TAKFL compared to the standard alternatives that do not adapt to the student’s self-knowledge quality and available learning capacity.
4. Our comprehensive experimental evaluations on both CV and NLP tasks, spanning various datasets and architectures, reveal that TAKFL consistently achieves SOTA performance, outperforming existing KD-based methods.

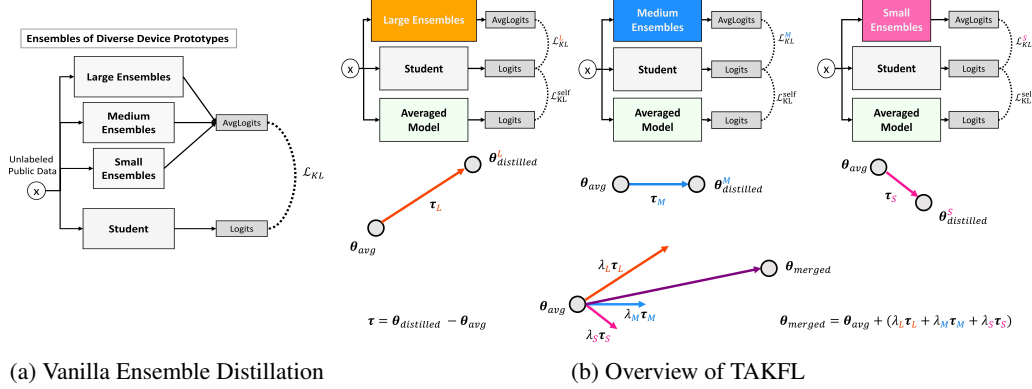


Figure 1: **Overview of our approach and its distinction from prior works.** (a) This figure illustrates the vanilla ensemble distillation process, where logits from ensembles of various sizes are averaged and used as the distillation target across all prototypes. This approach leads to the dilution of information and suboptimal knowledge transfer (refer to Sections 6 and 7 for details). (b) This figure depicts our approach, TAKFL, which treats knowledge transfer from each prototype’s ensemble as a separate task and distills them independently. Additionally, a KD-based self-regularization technique is introduced to mitigate issues related to the noisy and unsupervised ensemble distillation. Finally, the heterogeneously distilled knowledge is strategically integrated using an adaptive task arithmetic operation, allowing for customized knowledge integration based on each student prototype’s needs.

## 2 Related Works

**Device Heterogeneous FL.** Prior works on device heterogeneous FL have considered two distinct approaches with different objectives and settings. The first array of studies focuses on accommodating devices with varying compute resources, aiming to train a single global model. Techniques such as static and rolling-based partial model training allow devices to train a sub-model of the global model tailored to their compute resources [11, 18, 3, 1]. However, this approach does not fully reflect real-world scenarios. In practice, device prototypes such as IoTs and smartphones have unique neural network architectures designed for their specific configurations and underlying tasks, which may not support training varying neural architectures. This highlights a significant limitation in accommodating the full spectrum of device heterogeneity in this approach. The second array of studies addresses a more practical scenario where device prototypes with heterogeneous model architectures participate in FL to enhance their global model performance through mutual knowledge sharing [30, 41, 6]. In this context, KD techniques are used to transfer knowledge among prototypes, where locally updated client models, termed as ensembles, serve as teachers to distill their knowledge into each server’s student model using an unlabeled public dataset. For example, FedDF [30] uses vanilla logit averaging, while Fed-ET [6] applies an uncertainty-weighted logit averaging, enhanced by a diversity regularization technique. *However, existing works typically focus on settings where prototypes have similar capabilities—both model and dataset sizes—and thus neglecting the challenges in more diverse settings with varying capabilities. This oversight leaves their effectiveness in such settings largely unexplored. In this paper, we aim to study the underexplored diverse heterogeneous device settings.* See Appendix A for a more detailed discussion on the related works.

**Model Editing via Task Arithmetic.** Traditional methods for model editing often involve expensive joint fine-tuning across multiple tasks, which can limit scalability and democratization [60]. Recently, a promising technique called task arithmetic has emerged as a cost-effective and scalable method for updating pre-trained models with new information or refining undesired behavior [51, 37, 32]. The concept of “task vectors” introduced by Wortsman et al. [51] plays a pivotal role in these techniques. For any given task  $t$ , a task vector is derived by subtracting the model’s pre-trained weights  $\theta_{pre}$  from its fine-tuned weights  $\theta_{ft}^t$  on task  $t$ , i.e.  $\tau_t = \theta_{ft}^t - \theta_{pre}$ . These task vectors act as unique representations for specific tasks. Furthermore, researchers have demonstrated that by summing multiple task vectors  $\{\tau_t\}_{t=1}^T$ , and integrating them into a pre-trained model via  $\theta = \theta_{pre} + \lambda \sum_{t=1}^T \tau_t$ , one can effectively create a model capable of handling multiple tasks [51, 55]. *To the best of our knowledge, this work is the first to extend the notion of task vectors to the federated learning setting, introducing a task arithmetic for knowledge distillation across diverse heterogeneous device prototypes.*

### 3 Problem Statement: FL with Heterogeneous Device Prototypes

Consider a cross-device FL setup with a set of  $M$  distinct device prototypes  $\mathbb{M}$ , i.e.,  $M = |\mathbb{M}|$ . Each device prototype  $m_j \in \mathbb{M}$  has a distinct neural network architecture  $f^j(\cdot; \theta^j)$  parameterized by  $\theta^j \in \mathbb{R}^{n_j}$  and a set of clients  $\mathbb{C}^j$ , with  $N^j = |\mathbb{C}^j|$  clients in total. Each client  $c_k \in \mathbb{C}^j$  has a local private dataset  $\mathbb{D}_k^j = \{(x_i, y_i)\}_{i=1}^{n_{j,k}}$ , where  $n_{j,k} = |\mathbb{D}_k^j|$ , and locally trains the parameters  $\theta^j$  of the neural network architecture  $f^j$  on its local dataset. Furthermore, denote  $\mathbb{D}^j = \cup_{k \in \mathbb{C}^j} \mathbb{D}_k^j$  to be the union of the private datasets for device prototype  $j$ . We assume  $\mathbb{D}^j \sim \mathcal{D}^j$ , that is a subsample from the population distribution  $\mathcal{D}^j$  and similarly  $\mathbb{D}_k^j \sim \mathcal{D}_k^j$ . The union of the private datasets, i.e.  $\mathbb{D} = \cup_{j \in \mathcal{M}} \mathbb{D}^j$ , is sampled from the entire population  $\mathcal{D}$ , which is defined as an unknown mixture of the distributions each device prototype sampled its data from, i.e. generically non-i.i.d. We formalize this as a mixture of local clients data population, i.e.,  $\mathcal{D} = \sum_j \omega_j, : \mathcal{D}^j = \sum_j \sum_k \omega_{j,k} \mathcal{D}_k^j$ , where  $0 \leq \omega_{j,k} \leq 1$  and  $\sum_{j,k} \omega_{j,k} = 1$ , and  $\omega_{j,k}$  is unknown.

The ultimate objective is to minimize the test error and thus enable accurate inference for each device prototype  $j$ , aiming to obtain the optimal parameters for the population dataset:

$$\operatorname{argmin}_{\theta^j} \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell(f^j(x; \theta^j), y)] = \operatorname{argmin}_{\theta^j} \sum_{j=1}^M \sum_{k=1}^{N^j} \omega_{j,k} \mathbb{E}_{(x,y) \sim \mathcal{D}_k^j} [\ell(f^j(x; \theta^j), y)] \quad (1)$$

where  $\ell(\cdot, \cdot)$  is the sample-wise loss function (e.g. cross entropy for image classification) and we decompose by total population loss with the linearity of expectation in the mixture. See Fig 4b for a visual illustration of heterogeneous device prototype FL.

### 4 Background: Federated Ensemble Distillation

To address the limitations of standard FL in device heterogeneous settings, Lin et al. [30] proposed ensemble knowledge distillation to transfer knowledge between heterogeneous device prototypes in FL. This procedure consists of two stages: (1) local per-prototype FL, and (2) server-side vanilla ensemble distillation. The details of each stage discussed in the following paragraphs.

**Local Per-Prototype FL.** In this context, at each round  $r$  a subset of clients  $\mathbb{C}_r^j$  from each device prototype  $j \in \mathbb{M}$  is randomly selected by the server and download their corresponding model initialization  $\theta_r^j$ . Each client  $c_k^j \in \mathbb{C}_r^j$ , starting from this model initialization, locally train the model  $f^j$  on its local private data  $\mathbb{D}_k^j$  by taking multiple steps of stochastic gradient descent. Then, they send back their updated parameters  $\{\hat{\theta}_k^j\}_{k \in \mathbb{C}_r^j}$  to the server. The server aggregates the received clients parameters, and computes  $\theta_{avg}^j = \sum_{k \in \mathbb{C}_r^j} \frac{|\mathbb{D}_k^j|}{\sum_{k \in \mathbb{C}_r^j} |\mathbb{D}_k^j|} \hat{\theta}_k^j$ . In classic federated learning formalism, the parameters  $\theta_{avg}^j$  satisfy,

$$\theta_{avg}^j \in \operatorname{argmin}_{\theta^j} \sum_{k=1}^{N^j} \mathbb{E}_{(x,y) \sim \mathbb{D}_k^j} [\ell(f^j(x; \theta^j), y)] \quad (2)$$

**Vanilla Ensemble Distillation.** In this stage, each server model  $f^j$  gets initialized with  $\theta^j$ , and undergoes updates using ensemble knowledge distillation. Here, heterogeneous client models from heterogeneous device prototypes, collectively termed as ensembles, serve as teachers, i.e.  $\mathcal{T} := \{f^i(\cdot, \theta_k^i) | i \in \mathbb{M}, k \in \mathbb{C}^i\}$ , transferring their knowledge to each server student model, i.e.  $\mathcal{S}_i := f^i(\cdot, \theta^i)$ . For simplicity, we drop the index for each server student model, denoting it as  $\mathcal{S}$ . The ensemble distillation loss using a mini-batch of data from an unlabeled public dataset, i.e  $x \in \mathbb{D}^{public}$ , can be defined by the following equation:

$$\mathcal{L}_{ED} = \text{KL} \left[ \sigma \left( \frac{1}{|\mathcal{T}|} \sum_{\mathcal{F} \in \mathcal{T}} \mathcal{F}(x) \right), \sigma(\mathcal{S}(x)) \right], \quad (\text{AvgLogits}) \quad (3)$$

where  $\sigma(\cdot)$  is the softmax function. As illustrated in Eq. 3, vanilla ensemble distillation treats all heterogeneous device prototypes' ensembles equally by uniformly averaging their logits. This way of knowledge integration overlooks the individual strengths and informational value of each prototype's ensembles. As a result, the richer, more informative logits from stronger ensembles are diluted by less informative logits from weaker ensembles, leading to information loss. Furthermore,



this averaged logits is used as the distillation target across different-sized student models, irrespective of their intrinsic capacity and the helpfulness of each prototype’s ensembles. Consequently, this leads to suboptimal knowledge transfer in device heterogeneous FL. See Section 6 for theoretical analysis and Section 7 for experimental observations.

## 5 Task Arithmetic Knowledge Transfer and Integration

In this section, we introduce TAKFL, designed to overcome the fundamental limitations of previous approaches and enhance knowledge transfer across diverse heterogeneous device prototypes, which vary in size—in terms of both model and dataset size. TAKFL consists of two main components: (1) individually transferring knowledge from each device prototype’s ensembles, and (2) adaptively integrating knowledge via task arithmetic. Detailed descriptions of each component are provided in Section 5.1 and 5.2, respectively. An illustrative overview of TAKFL is presented in Figure 1b, and the full algorithm is detailed in Appendix B, Algorithm 1.

### 5.1 Knowledge Transfer from Individual Device Prototype

We begin by discussing our proposed knowledge transfer framework from each individual device prototype’s ensembles. This process consists of two main components: ensemble knowledge transfer and self-regularization, each detailed in the subsequent paragraphs.

**Ensemble Knowledge Transfer.** Vanilla ensemble distillation integrates the knowledge of varying strength ensembles by uniformly averaging their logits. This approach can potentially transform or even degrade the overall quality of the knowledge being transferred, leading to suboptimal knowledge transfer. To effectively distill the unique knowledge and contributions of each prototype’s ensembles, and to avoid dilution, information loss, and interference from other prototypes’ ensembles, we propose transferring the knowledge from each prototype’s ensembles separately and independently.

Specifically, let’s consider  $\mathcal{T}_i := \{f^i(\cdot, \hat{\theta}_k^i) | k \in \mathbb{C}^i\}$  denotes the ensembles of device prototype  $i$  as teacher and  $\mathcal{S}_j$  denotes the server student model of the device prototype  $j$ . Without loss of generality, we refer to each device prototype’s server student model as just student denoted as  $\mathcal{S}$ . Therefore, the knowledge distillation loss between the teacher ensembles  $\mathcal{T}_i$  and server student  $\mathcal{S}$  ( $\mathcal{T}_i \rightarrow \mathcal{S}$ ) is defined below:

$$\mathcal{L}_{KD}^{\mathcal{T}_i \rightarrow \mathcal{S}} = \text{KL} \left[ \sigma \left( \frac{1}{|\mathcal{T}_i|} \sum_{\mathcal{F} \in \mathcal{T}_i} \mathcal{F}(x) \right), \sigma(\mathcal{S}(x)) \right]. \quad (4)$$

**Scaffolding Student from Noisy Ensemble Distillation.** The ensemble distillation process may adversely impact the student, causing it to forget its own knowledge acquired through averaged locally updated parameters and be drifted into erroneous directions. This is primarily due to two key factors: (1) The ensemble distillation process introduces noise, mainly because the ensembles’ logits are inferred on an unfamiliar public dataset they have not been trained on. These ensembles are originally trained on local private datasets, which usually differ from the unlabeled public dataset used for distillation. Moreover, other factors such as the presence of data heterogeneity within FL and insufficient training of some ensembles due to limited computational resources can exacerbate this noise, particularly in the initial rounds of federation. (2) The ensemble distillation process lacks supervision from the actual private datasets, which is the ultimate learning objective.

To scaffold the student models from the noisy and unsupervised distillation process, which may cause them to drift into erroneous directions and forget their invaluable self-knowledge, we introduce a KD-based self-regularization technique. Our self-regularization technique mitigates these issues by enforcing similarity between the logits of the student and its initial logits (when the student is initialized with averaged parameters) using KL divergence loss defined below:

$$\mathcal{L}_S^{\text{self}} = \text{KL} \left[ \sigma(\mathcal{S}(x; \theta_{avg})), \sigma(\mathcal{S}(x)) \right]. \quad (5)$$

**Overall Knowledge Transfer Objective.** The overall knowledge transfer objective from teacher ensembles  $\mathcal{T}_i$  of device prototype  $i$  to the student  $\mathcal{S}$  is the combination of the ensemble knowledge distillation loss  $\mathcal{L}_{KD}^{\mathcal{T}_i \rightarrow \mathcal{S}}$  (Eq. 4) and the self-regularization loss  $\mathcal{L}_S^{\text{self}}$  (Eq. 5) defined in the following:

$$\mathcal{L}_S^{\mathcal{T}_i} = \mathcal{L}_{KD}^{\mathcal{T}_i \rightarrow \mathcal{S}} + \gamma \cdot \mathcal{L}_S^{\text{self}}. \quad (6)$$

Here,  $\gamma$  is a hyperparameter controlling the effect of self-regularization term. We associate the knowledge transfer from each device prototype  $i$  to a task  $T_i$  with the loss  $\mathcal{L}_S^{\mathcal{T}_i}$ .

## 5.2 Task Arithmetic Knowledge Integration

Herein, we delve into the details of our proposed method for customized integration of the separately distilled knowledge from heterogeneous ensembles. Drawing inspiration from recent advances in model editing via task arithmetic [51], where a pre-trained model’s knowledge can be edited via task-specific vectors using arithmetic operation, we propose a novel customizable knowledge integration method via task arithmetic. To do so we extend the notion of task vector from centralized learning to federated learning. We conceptualize the averaged locally updated parameters, i.e.  $\theta_{avg}$ , as a “pre-trained”, similar to those in centralized learning, and the parameters of the distilled model via knowledge transfer objective (Eq. 4), denoted as  $\theta_{distilled}$ , as a “fine-tuned” version of the model (see Fig. 2). Consequently, the task vector  $\tau_i$  associated with the knowledge transfer task  $\mathcal{L}_S^{\mathcal{T}_i}$  can be defined by subtracting the distilled parameters from the averaged locally updated parameters as follows:

$$\tau_i = \theta_{distilled}^{\mathcal{T}_i \rightarrow \mathcal{S}} - \theta_{avg}. \quad (7)$$

Essentially, task vectors serve as unique representations for the transferred knowledge from each prototype’s ensembles to the student and encapsulate the distinct contributions of each prototype’s ensembles to the student model. To selectively merge the knowledge of each prototype’s ensembles into the student, we employ an adaptive task arithmetic operation as follows:

$$\theta_{merged} = \theta_{avg} + \sum_{i \in \mathbb{M}} \lambda_i \tau_i, \quad (8)$$

where  $\lambda_i$  denotes the merging coefficient associated with task vector  $\tau_i$ , and they sum to one, i.e.  $\sum_{i \in \mathbb{M}} \lambda_i = 1$ . The merging coefficients determine the extent of knowledge integration from each prototype’s ensembles. Essentially, they enable the student to have customized knowledge integration to achieve maximum performance. The student can determine these merging coefficients based on its own learning capacity and the relative knowledge and helpfulness of other device prototypes’ ensembles. This approach provides an effective, low-cost, and scalable knowledge integration strategy in settings with diverse device heterogeneity. In our experiments, we considered this as a hyperparameter and tuned it manually or determined it using held-out validation sets which achieves similar results. More details can be found in Appendix F.3.

## 6 Theoretical Results

We present a theoretical understanding on the efficacy of knowledge distillation in device heterogeneous FL. We argue that vanilla ensemble distillation (VED) diffuses the information from logits, which presents a notable disadvantage for solving (1). This effect is particularly pronounced when the teacher ensembles are from a device prototype of small capacity, and the student model is from a device prototype of large capacity. By contrast, our proposed method of task arithmetic knowledge integration, mitigates the drawbacks of VED and is able to simultaneously incorporate information from differently sized heterogeneous ensembles, efficiently filling up the capacity of each student with the most informative knowledge, achieving optimal knowledge transfer.

**Assumptions and Preliminaries.** Standard practice, including the setting in consideration as well as the numerical experiments here, involves *overparametrized* neural networks, that is, the total number of weights far exceeds the training sample size. This implies that the set of weights that minimize the loss is non-unique, and moreover, it has been argued that they form a submanifold [8]. This submanifold structure of solution sets will provide the critical source of understanding the subsequent results. In particular, we shall consider knowledge distillation as filling up the capacity of device prototypes’ models with basis vectors corresponding to submanifolds that minimize as many device prototypes’ data distributions as possible.

Since we are interested in server-side distillation across heterogeneous device prototypes, we assume optimal conditions at the local per-prototype FL level, meaning that the perfect solution for local per-prototype FL is achieved. The formal details of the assumptions and statements are presented in Appendix C.

**Proposition 1. (information loss in VED, informal).** Consider the VED procedure in the form of solving (3). Consider two device prototypes with a device capacity and solution dimension of  $Q^1, Q^2$

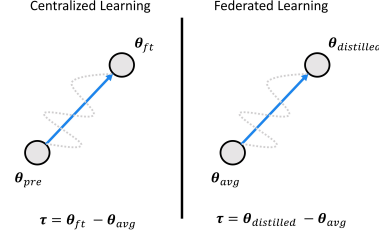


Figure 2: **Analogy between task vector in centralized learning and federated learning.**

and  $W^1, W^2$ , respectively, and with associated eigenbases  $Q^i, W^i$ . Denote  $W^{i,j}, i, j = 1, 2$  as the capacity allocated by student  $i$  in order to distill knowledge from teacher  $j$ 's logits.

1. **Case 1:** When the capacities are the same, that is  $Q^1 = Q^2$  and  $W^1 = W^2 = W^{1,2} = W^{2,1}$ , then with VED, there will be some capacity, in the sense of eigenspace, of student prototypes that will be allocated with parameters that do not minimize the student's its own data distribution.
2. **Case 2:** Assume that  $Q^1 > Q^2$  and  $W^1 = W^{1,2} > W^2$ . Then the phenomenon as for Case 1 holds. Moreover, there will be some capacity of student 1's model that will be allocated with parameters that do not minimize either of the teacher or student prototype's data distribution.

An interesting key mechanism of the proof is that when VED is applied in distilling logits from a small device prototype to a large one, the modeling capacity of  $W^{1,2}$  is structurally reduced to that of  $W^2 < W^{1,2}$ , i.e., it is an operation wasteful of the potential model capacity.

**Remark 1.** This proposition proves that in general, VED is prone to diffuse knowledge already present in students, and leads to inefficient and inaccurate use of model capacity. Furthermore, under the case that device prototypes have different capacities, VED ends up leading to more erroneous models entirely as the small information within the small teacher is transferred onto a larger capacity target.

**Proposition 2. (improve knowledge transfer with task arithmetic, informal).** Consider the TAKFL procedure as in the form of computing (8). Consider two device prototypes with a device capacity and solution dimension of  $Q^1, Q^2$  and  $W^1, W^2$ , respectively, and with associated eigenbases  $Q^i, W^i$ .

1. **Case 1:** In the case that that  $Q^1 \geq Q^2$  and  $W^1 \geq W^2$ , it holds that the TAKFL with prototype 1 as student preserves the eigenbasis associated to the parameters used to accurately fit the data  $\mathcal{D}^1$ .
2. **Case 2:** Assume that  $Q^1 = Q^2$  and  $W^1 = W^2$ . TAKFL yields a solution for the student that is at the intersection of the subspaces corresponding to minimizing the two data distributions.
3. **Case 3:** Assume that  $Q^1 > Q^2$  and  $W^1 > W^2$ . In the case of prototype 1 being the student, TAKFL yields a solution that:
  - (a) retains the approximation accuracy on device 1's data distribution,
  - (b) ensures approximation accuracy to the level of device 2's relative capacity
  - (c) fills the remaining local capacity device 1 has allocated for device 2's logits with no informative new knowledge, unless enforced otherwise.

**Remark 2.** This proposition proves that in general, TAKFL promotes the most efficient allocation of the devices' capacity in order to accurately fit a diverse set of data distributions. With TAKFL, the previously acquired knowledge is entirely preserved. Even under the case that device prototypes have different capacities, TAKFL smartly transfers the most informative knowledge to each prototype's student model based on its own intrinsic capacity. Still, the final statement indicates that in the case that there are many different teachers, while a small device prototype serving as teacher will not be necessarily compromise information, it would still be preferable to allocate that capacity to a more informative, larger, teacher model.

## 7 Experiments

### 7.1 Main Experimental Setup

**Dataset and Architecture.** We evaluate our method on computer vision (CV) and natural language processing (NLP) tasks. For CV, we train image classification using CIFAR10/100 [24], CINIC-10 [9], and TinyImagenet [25]. For NLP, we fine-tune pre-trained models for text classification on MNLI [50], SST-2 [43], MARC [22], and AG News [58]. Our architectures include ResNet [17], VGG [42], and ViT [12] for CV, and small BERT variants [45] (-Tiny, -Mini, -Small) for NLP. We simulate a federated non-i.i.d setting using a Dirichlet distribution  $Dir(\alpha)$ , where a lower  $\alpha$  indicates higher heterogeneity [27, 36]. Further details can be found in Appendix F.1 and F.2.

**Implementation Details.** We use the Adam optimizer for both CV and NLP tasks. For CV, local training involves 20 epochs with a learning rate of 0.001, weight decay of 5e-5, and a batch size of 64. NLP training is conducted over 1 epoch with a learning rate of 3e-5, no weight decay, and a batch size of 32. For distillation, Adam is used with a learning rate of 1e-5 and weight decay of 5e-4 for CV, and 3e-5 with no weight decay for NLP. Batch sizes for distillation are 128 for CV and 32 for NLP. The softmax temperature is set at 3 for both tasks, with a temperature of 20 for self-regularization. Further details are provided in Appendix F.1 and F.2.

Table 1: **Performance Results for CV task on CIFAR-10 and CIFAR-100.** Training data is distributed among S, M, and L device prototypes in a 1:3:6 ratio, subdivided among clients using Dirichlet distribution. Public datasets are CIFAR-100 [24] for CIFAR-10 [24] and ImageNet-100 [10] for CIFAR-100. Client configurations include 100, 20, and 4 clients for S, M, and L, with sampling rates of 0.1, 0.2, and 0.5. Architectures are ResNet-8, ResNet-14, and ResNet-18 [17] for S, M and L, respectively. All models are trained from scratch for 60 rounds. See Appendix D.1 for additional experiments using hetero-family architecture and more details.

Dataset	Baseline	Low Data Heterogeneity (Dir(0.3))				High Data Heterogeneity (Dir(0.1))			
		S	M	L	Average	S	M	L	Average
CIFAR-10	FedAvg	36.21 $\pm$ 2.24	46.41 $\pm$ 2.33	59.46 $\pm$ 6.17	47.36	22.01 $\pm$ 0.78	25.26 $\pm$ 3.89	51.51 $\pm$ 3.52	32.93
	FedDF	49.31 $\pm$ 0.15	50.63 $\pm$ 0.73	49.82 $\pm$ 0.98	49.92	34.71 $\pm$ 1.48	35.27 $\pm$ 4.74	51.08 $\pm$ 4.04	40.35
	FedET	49.21 $\pm$ 0.72	55.01 $\pm$ 1.81	53.60 $\pm$ 6.47	52.60	29.58 $\pm$ 3.00	30.96 $\pm$ 4.70	45.53 $\pm$ 6.46	35.36
	TAKFL	55.90 $\pm$ 1.70	57.93 $\pm$ 3.49	60.58 $\pm$ 2.35	58.14	37.40 $\pm$ 1.68	38.96 $\pm$ 0.17	51.49 $\pm$ 6.15	42.62
	TAKFL+Reg	<b>56.37<math>\pm</math>0.46</b>	<b>58.60<math>\pm</math>0.43</b>	<b>65.69<math>\pm</math>1.28</b>	<b>60.22</b>	<b>40.51<math>\pm</math>1.05</b>	<b>40.12<math>\pm</math>1.24</b>	<b>53.24<math>\pm</math>2.51</b>	<b>44.62</b>
CIFAR-100	FedAvg	13.22 $\pm$ 0.14	21.39 $\pm$ 1.11	29.47 $\pm$ 0.86	21.36	11.86 $\pm$ 0.08	14.63 $\pm$ 0.65	26.25 $\pm$ 1.64	17.58
	FedDF	19.54 $\pm$ 0.20	24.32 $\pm$ 0.45	29.29 $\pm$ 1.45	24.38	16.09 $\pm$ 0.32	19.80 $\pm$ 0.17	26.59 $\pm$ 0.25	20.83
	FedET	19.67 $\pm$ 0.35	25.27 $\pm$ 0.66	31.10 $\pm$ 1.53	25.35	11.18 $\pm$ 1.68	18.22 $\pm$ 0.35	26.40 $\pm$ 0.65	18.60
	TAKFL	24.48 $\pm$ 0.42	27.60 $\pm$ 0.25	29.84 $\pm$ 0.94	27.31	<b>22.90<math>\pm</math>0.18</b>	<b>23.63<math>\pm</math>0.72</b>	<b>26.98<math>\pm</math>0.13</b>	<b>24.50</b>
	TAKFL+Reg	<b>27.18<math>\pm</math>0.27</b>	<b>29.14<math>\pm</math>0.20</b>	<b>31.15<math>\pm</math>0.97</b>	<b>29.15</b>	<b>22.88<math>\pm</math>0.37</b>	<b>23.92<math>\pm</math>0.57</b>	<b>28.01<math>\pm</math>0.34</b>	<b>24.94</b>

**Baselines and Evaluation Metric.** We compare our method against standard FL, i.e. FedAvg [33] and SOTA KD-based methods designed for heterogeneous device prototypes FL, including FedDF [30] and FedET [6]. The evaluation metric is the final top-1 classification accuracy of each device prototype’s global model on the test dataset, as per the methodology described in [36]. We report the average results and the standard deviation over three independent runs, each with a different random seed. *A more detailed version of the experiments, alongside additional experiments and ablation studies, is presented in Appendix D and E.*

## 7.2 Main Experimental Results

In this section, we evaluate the performance of our method, TAKFL, in a federated learning environment that mirrors real-world scenarios with diverse, heterogeneous device prototypes, as illustrated in Fig. 4b. Our experimental setup includes three different device prototype sizes: Small (S) with a small model and small dataset, Medium (M) with a medium-sized model and medium-sized dataset, and large (L) with a large model and large dataset.

**Performance on CV Task.** Table 1 presents the performance of TAKFL in the homo-family architecture setting on the CIFAR-10 and CIFAR-100 [24] datasets (for hetero-family architecture results, see Appendix D.1, Table 4). TAKFL consistently enhances performance across all device prototypes in various scenarios, achieving SOTA results. Notably, in the Dir(0.3) setting on CIFAR-10, TAKFL improves average performance across all prototypes by 8%, and by 4% on CIFAR-100. From Table 1, inconsistent performance improvements are observed with prior KD-based methods, especially for the L prototype. While S and M prototypes achieve gains, the L prototype suffers up to a 10% degradation compared to vanilla FedAvg, highlighting the dilution issue where valuable information from larger, more capable device prototypes is diluted by less informative outputs from smaller devices. Moreover, the significant performance improvements TAKFL achieves for each device prototype, particularly for S and M prototypes, illustrate the ineffectiveness of the one-size-fits-all approach used in the existing KD methods. These observations confirm the shortcomings of vanilla ensemble distillation and corroborate our theoretical findings in Remark 1 and 2. The effectiveness of our self-regularization technique is further supported by these experimental results. For more detailed and insightful analysis see Appendix D.1.1.

**Performance on NLP Task.** Table 2 presents the results on MNLI [50] and SST-2 [43] datasets (see Appendix D.3 for further experiments). Similar to the CV task, TAKFL has consistently improved performance across all device prototypes of varying sizes, achieving SOTA results: a 3% average increase on MNLI and 2% on SST-2. The suboptimality of existing KD methods, is evident from the results presented here as well. Notably, FedET suffers from a significant performance degradation compared to vanilla FedAvg. This issue stems from FedET’s reliance on the confidence scores of neural networks for uncertainty estimates. However, neural networks, especially pretrained language models (PLMs), tend to be poorly calibrated and overconfident, undermining reliable uncertainty estimates [48, 15, 5, 53].

## 7.3 Scalability Evaluation

We evaluate the scalability of TAKFL across a spectrum of device prototypes, from extremely small (XXS) to extremely large (XXL), to see how well our method adapts from a uniform array of small-size prototypes to a diverse mix of sizes. Each prototype is equipped with appropriately scaled model and dataset sizes, simulating real-world variations in device capabilities.

Table 2: **Performance Results for NLP Task on MNLI and SST-2.** Training data distribution is similar to the CV task using only Dir(0.5) here. Public datasets are SNLI [2] for MNLI [50] and Sentiment140 [14] for SST-2 [43]. Client configurations are 8, 4, and 2 clients for S, M, and L, with sample rates of 0.3, 0.5, and 1.0, respectively. Architectures include Bert-Tiny, Bert-Mini, and Bert-Small [45] for S, M, and L, initialized from pre-trained parameters and fine-tuned for 20 communication rounds. See Appendix F.2 for more details.

Baseline	MNLI				SST-2			
	S	M	L	Average	S	M	L	Average
FedAvg	36.15 $\pm$ 0.46	54.47 $\pm$ 2.48	57.51 $\pm$ 2.79	49.37	54.98 $\pm$ 1.81	74.71 $\pm$ 8.22	86.69 $\pm$ 0.06	72.13
FedDF	54.21 $\pm$ 0.15	60.44 $\pm$ 1.91	66.71 $\pm$ 1.09	60.45	74.41 $\pm$ 2.62	80.71 $\pm$ 1.63	84.35 $\pm$ 1.66	79.82
FedET	48.03 $\pm$ 6.32	50.33 $\pm$ 7.87	53.80 $\pm$ 6.18	50.72	66.63 $\pm$ 9.14	65.89 $\pm$ 16.35	70.05 $\pm$ 15.83	67.52
TAKFL	57.43 $\pm$ 0.21	63.58 $\pm$ 0.31	68.74 $\pm$ 0.12	63.25	74.73 $\pm$ 0.55	82.17 $\pm$ 0.31	86.93 $\pm$ 0.42	81.28
TAKFL+Reg	<b>57.61<math>\pm</math>0.89</b>	<b>63.91<math>\pm</math>1.05</b>	<b>68.96<math>\pm</math>1.10</b>	<b>63.49</b>	<b>74.88<math>\pm</math>0.43</b>	<b>82.40<math>\pm</math>0.83</b>	<b>87.33<math>\pm</math>0.63</b>	<b>81.54</b>

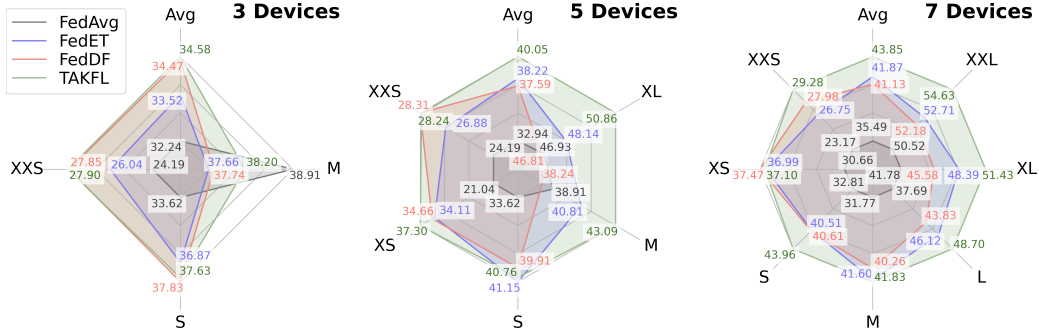


Figure 3: **Scalability Evaluation of TAKFL.** Image classification on CINIC-10 [9] dataset is used to evaluate TAKFL’s scalability across device prototypes ranging from XXS to XXL. Training data is distributed among prototypes in a 1:2:3:4:5:6:7 ratio, further subdivided using Dir(0.5). Client configurations range from 35 for XXS to 5 for XXL. Architectures span from ResNet10-XXS for XXS to ResNet50 for XXL prototype, all initialized from scratch and trained over 30 communication rounds. The public dataset is CIFAR-100 [24]. See Appendix D.4 for more details.

Figure 3 illustrates TAKFL’s ability to effectively scale from 3 to 7 device prototypes. In scenarios where all devices are similarly small, i.e. 3-device setup, TAKFL’s performance is slightly better than FedDF. This is because when devices are homogeneously small and similar in capability, they do not offer unique contributions that could benefit from more complex distillation strategies. However, as the scenario expands to include larger devices like XL and XXL in the 5- and 7-device configurations, TAKFL significantly outperforms existing KD-based methods. This improvement is driven by the larger devices’ ability to offer more significant and higher-quality knowledge, which TAKFL effectively distills across all prototypes, contrasting sharply with existing methods that fail to utilize this potential. These experimental observations, corroborated by our theoretical insights in Remark 2, demonstrate TAKFL’s superior scalability and effectiveness.

## 8 Conclusion and Discussion

In this work, we addressed a fundamental issue in standard federated learning: the lack of support for heterogeneous device prototypes. Existing KD-based methods often fall short in real-world scenarios, where device capabilities vary widely. To address this, we introduced TAKFL, a novel KD-based method that treats knowledge transfer from each prototype’s ensembles as separate tasks and distills them independently. TAKFL subsequently integrates the knowledge using an adaptive task arithmetic technique for optimized performance. We also introduced a KD-based self-regulation technique to mitigate issues arising from noisy and unsupervised ensemble distillation. The effectiveness of our method is substantiated by both theoretical results and extensive experimentation across CV and NLP tasks, using various datasets and models.

Limitations remain, notably in real-world applicability. While TAKFL’s effectiveness in an approximated real-world setup has been demonstrated, actual deployment on physical devices and in environments with extremely large models remains untested due to resource constraints. Experiencing TAKFL in genuine real-world settings could unveil additional challenges or limitations, providing further insights into its scalability and efficiency.



## References

- [1] 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*, 35:29677–29690, 2022. 3, 1
- [2] Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. A large annotated corpus for learning natural language inference, 2015. 9, 12, 18
- [3] Sebastian Caldas, Jakub Konečný, H Brendan McMahan, and Ameet Talwalkar. Expanding the reach of federated learning by reducing client resource requirements. *arXiv preprint arXiv:1812.07210*, 2018. 3, 1
- [4] Chuan Chen, Tianchi Liao, Xiaojun Deng, Zihou Wu, Sheng Huang, and Zibin Zheng. Advances in robust federated learning: Heterogeneity considerations. *arXiv preprint arXiv:2405.09839*, 2024. 2
- [5] Yangyi Chen, Lifan Yuan, Ganqu Cui, Zhiyuan Liu, and Heng Ji. A close look into the calibration of pre-trained language models. *arXiv preprint arXiv:2211.00151*, 2022. 8, 2, 12
- [6] Yae Jee Cho, Andre Manoel, Gauri Joshi, Robert Sim, and Dimitrios Dimitriadis. Heterogeneous ensemble knowledge transfer for training large models in federated learning. *arXiv preprint arXiv:2204.12703*, 2022. 3, 8, 1, 2, 17
- [7] Adam Coates, Andrew Ng, and Honglak Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pages 215–223, Fort Lauderdale, FL, USA, 2011. PMLR. 11, 16
- [8] Yaim Cooper. Global minima of overparameterized neural networks. *SIAM Journal on Mathematics of Data Science*, 3(2):676–691, 2021. 6
- [9] Luke N. Darlow, Elliot J. Crowley, Antreas Antoniou, and Amos J. Storkey. Cinic-10 is not imagenet or cifar-10, 2018. 7, 9, 11, 13, 16
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. 8, 9, 16
- [11] Enmao Diao, Jie Ding, and Vahid Tarokh. Heterofl: Computation and communication efficient federated learning for heterogeneous clients. *arXiv preprint arXiv:2010.01264*, 2020. 3, 1
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale, 2021. 7, 9, 17
- [13] Cong Fang, Hanze Dong, and Tong Zhang. Mathematical models of overparameterized neural networks. *Proceedings of the IEEE*, 109(5):683–703, 2021. 8
- [14] Alec Go, Richa Bhayani, and Lei Huang. Twitter sentiment classification using distant supervision. *CS224N project report, Stanford*, 1(12):2009, 2009. 9, 12, 18
- [15] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR, 2017. 8, 2, 12, 14
- [16] Farzin Haddadpour and Mehrdad Mahdavi. On the convergence of local descent methods in federated learning. *arXiv preprint arXiv:1910.14425*, 2019. 4
- [17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 7, 8, 9
- [18] 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*, 34:12876–12889, 2021. 3, 1
- [19] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications, 2017. 11, 17
- [20] Anil Kag, Durmus Alp Emre Acar, Aditya Gangrade, and Venkatesh Saligrama. Scaffolding a student to instill knowledge. In *The Eleventh International Conference on Learning Representations*, 2022. 17

- [21] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International conference on machine learning*, pages 5132–5143. PMLR, 2020. 1
- [22] Phillip Keung, Yichao Lu, György Szarvas, and Noah A. Smith. The multilingual amazon reviews corpus. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020. 7, 12, 18
- [23] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 11
- [24] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009. 7, 8, 9, 11, 14, 16
- [25] Ya Le and Xuan S. Yang. Tiny imagenet visual recognition challenge. 2015. 7, 11, 16
- [26] Lin Li, Jianping Gou, Baosheng Yu, Lan Du, and Zhang Yiand Dacheng Tao. Federated distillation: A survey. *arXiv preprint arXiv:2404.08564*, 2024. 2
- [27] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10713–10722, 2021. 1, 7
- [28] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine learning and systems*, 2:429–450, 2020. 1
- [29] Zhiyuan Li, Tianhao Wang, and Sanjeev Arora. What happens after sgd reaches zero loss?—a mathematical framework. In *International Conference on Learning Representations*, 2021. 8
- [30] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020. 3, 4, 8, 1, 2
- [31] Chang Liu, Chenfei Lou, Runzhong Wang, Alan Yuhan Xi, Li Shen, and Junchi Yan. Deep neural network fusion via graph matching with applications to model ensemble and federated learning. In *International Conference on Machine Learning*, pages 13857–13869. PMLR, 2022. 1
- [32] Tian Yu Liu and Stefano Soatto. Tangent model composition for ensembling and continual fine-tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 18676–18686, 2023. 3
- [33] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 1, 8
- [34] Sachin Mehta and Mohammad Rastegari. Mobilevit: Light-weight, general-purpose, and mobile-friendly vision transformer, 2022. 11, 17
- [35] Alessio Mora, Irene Tenison, Paolo Bellavista, and Irina Rish. Knowledge distillation for federated learning: a practical guide. *arXiv preprint arXiv:2211.04742*, 2022. 2
- [36] Mahdi Morafah, Weijia Wang, and Bill Lin. A practical recipe for federated learning under statistical heterogeneity experimental design. *IEEE Transactions on Artificial Intelligence*, 2023. 7, 8
- [37] Guillermo Ortiz-Jimenez, Alessandro Favero, and Pascal Frossard. Task arithmetic in the tangent space: Improved editing of pre-trained models. *Advances in Neural Information Processing Systems*, 36, 2024. 3, 4, 5
- [38] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimesheine, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library, 2019. 16
- [39] Kilian Pfeiffer, Martin Rapp, Ramin Khalili, and Jörg Henkel. Federated learning for computationally constrained heterogeneous devices: A survey. *ACM Computing Surveys*, 55(14s):1–27, 2023. 2
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 14

- [41] Felix Sattler, Tim Korjakow, Roman Rischke, and Wojciech Samek. Fedaux: Leveraging unlabeled auxiliary data in federated learning. *IEEE Transactions on Neural Networks and Learning Systems*, 2021. [3](#), [1](#)
- [42] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2015. [7](#), [9](#), [16](#)
- [43] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA, 2013. Association for Computational Linguistics. [7](#), [8](#), [9](#), [12](#), [18](#)
- [44] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15 (1):1929–1958, 2014. [1](#)
- [45] Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*, 2019. [7](#), [9](#), [12](#), [18](#)
- [46] Hongyi Wang, Mikhail Yurochkin, Yuekai Sun, Dimitris Papailiopoulos, and Yasaman Khazaeni. Federated learning with matched averaging. *arXiv preprint arXiv:2002.06440*, 2020. [1](#)
- [47] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33:7611–7623, 2020. [1](#)
- [48] Yuxia Wang, Daniel Beck, Timothy Baldwin, and Karin Verspoor. Uncertainty estimation and reduction of pre-trained models for text regression. *Transactions of the Association for Computational Linguistics*, 10:680–696, 2022. [8](#), [2](#), [12](#)
- [49] Ross Wightman, Hugo Touvron, and Herve Jegou. Resnet strikes back: An improved training procedure in timm. In *NeurIPS 2021 Workshop on ImageNet: Past, Present, and Future*. [17](#)
- [50] Adina Williams, Nikita Nangia, and Samuel R. Bowman. A broad-coverage challenge corpus for sentence understanding through inference, 2018. [7](#), [8](#), [9](#), [12](#), [18](#)
- [51] Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International Conference on Machine Learning*, pages 23965–23998. PMLR, 2022. [3](#), [6](#), [4](#)
- [52] Yebo Wu, Li Li, Chunlin Tian, and Chengzhong Xu. Breaking the memory wall for heterogeneous federated learning with progressive training. *arXiv preprint arXiv:2404.13349*, 2024. [1](#)
- [53] Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. Uncertainty quantification with pre-trained language models: A large-scale empirical analysis. *arXiv preprint arXiv:2210.04714*, 2022. [8](#), [2](#), [12](#), [14](#)
- [54] Ziyue Xu, Mingfeng Xu, Tianchi Liao, Zibin Zheng, and Chuan Chen. Fedbrb: An effective solution to the small-to-large scenario in device-heterogeneity federated learning. *arXiv preprint arXiv:2402.17202*, 2024. [1](#)
- [55] Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-merging: Resolving interference when merging models. *Advances in Neural Information Processing Systems*, 36, 2024. [3](#)
- [56] Tuo Zhang, Lei Gao, Sunwoo Lee, Mi Zhang, and Salman Avestimehr. Timelyfl: Heterogeneity-aware asynchronous federated learning with adaptive partial training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5064–5073, 2023. [1](#)
- [57] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2015. [12](#), [18](#)
- [58] Xiang Zhang, Junbo Jake Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In *NIPS*, 2015. [7](#), [12](#), [18](#)
- [59] Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification, 2016. [12](#), [18](#)
- [60] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*, 109(1):43–76, 2020. [3](#)

## Appendix

The supplementary materials are organized as follows:

- Appendix **A**: Provides more details on related works.
- Appendix **B**: Presents the full algorithm description of TAKFL.
- Appendix **C**: Presents formal theoretical statements, assumptions, and proofs supporting our method.
- Appendix **D**: Presents detailed experimental results including some additional experiments.
- Appendix **E**: Presents the ablation studies experiments.
- Appendix **F**: Presents hyper-parameters and implementation details.

### A More Detailed Related Works

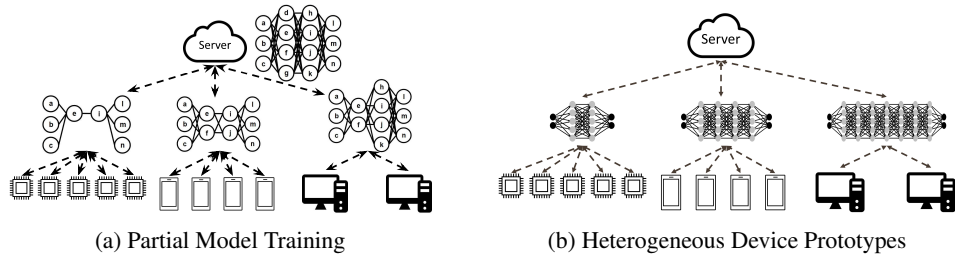


Figure 4: **Overview of Two Different Device Heterogeneous FL Settings.** (a) In the partial model training setting, the objective is to train a single global model where heterogeneous devices train a specific sub-model based on their computational resources. This approach necessitates device support for varying neural network architectures, which is impractical as devices typically have specialized architectures designed to match their hardware, software configurations, and underlying machine learning tasks. (b) In the heterogeneous device prototypes setting, device prototypes participate in FL to enhance the performance of their global model by transferring knowledge across prototypes. This setting is more feasible as it accommodates diverse device prototypes with their own specific configurations, including neural network architecture and dataset. However, establishing effective knowledge transfer between differently sized prototypes (like IoTs and workstations) and diverse configurations is challenging. In this paper, we address this issue.

Prior works on device heterogeneous FL have considered two distinct approaches with different objectives and settings. The first group of studies focuses on accommodating devices with varying compute resources, aiming to train a single global model [11, 3, 56, 52, 54]. Various partial model training techniques have been proposed for this setting, where devices are tasked with training a sub-model of a global model according to their compute resources. These include dropout-based [3], static [11, 18], and rolling-based sub-model extraction techniques [1]. Federated Dropout builds upon the concept of dropout [44] to extract smaller sub-models. Static sub-model extraction techniques like in HeteroFL [11] and FjORD [18] consistently extract designated portions of the global model, whereas FedRolex [1] introduces a more flexible rolling method for sub-model extraction. However, these approaches assume that devices can support various sub-model architectures for training, which does not fully reflect the real-world scenario. In practice, there exist a diverse spectrum of device prototypes such as IoT devices and smartphones each have unique and unhashable neural network architectures tailored to their specific hardware and software configurations and underlying machine learning tasks. Consequently, these device prototypes may not support training various neural network architectures, highlighting a significant limitation in accommodating the full spectrum of device heterogeneity in this setting.

The second array of studies tackles a more practical scenario where device prototypes with heterogeneous model architectures participate in FL to enhance their global model performance through mutual knowledge sharing. In this context, knowledge distillation techniques are employed to transfer knowledge among device prototypes [30, 6, 41]. Here, locally updated client models from various

device prototypes, collectively referred to as ensembles, serve as teachers to distill their knowledge into each server’s student model using an unlabeled public dataset. For instance, FedDF [30] utilizes vanilla averaging of all ensemble logits as the distillation target for all server student models. In contrast, FedET [6] employs an uncertainty-weighted average of ensembles’ logits as the distillation target for all server student models, complemented by a diversity regularization technique. However, methods like FedET rely on the neural networks’ confidence scores for uncertainty estimates, overlooking the fact that neural networks are often poorly calibrated and prone to overconfidence, which compromises their ability to provide reliable uncertainty estimates [48, 15, 5, 53]. *These existing works typically focus on settings where device prototypes have similar capabilities, i.e. similar model and dataset sizes, thus neglecting the challenges presented in more diverse settings where device prototypes vary significantly in terms of model and dataset size. This oversight limits the effectiveness of these methods in truly diverse and heterogeneous environments. In this paper, we introduce TAKFL, which is designed to address the limitations of existing methods in these underexplored diverse device heterogeneous settings.*

Figure 1 illustrates the distinctions between these two different settings studied in the literature. For more information, we refer the reader to recent surveys [35, 26, 39, 4].



## B Full Algorithm Description of TAKFL

The full algorithm description of TAKFL is presented in Algorithm 1.

---

### Algorithm 1 TAKFL Algorithm

---

**Require:** number of communication rounds ( $R$ ), public unlabeled dataset  $\mathbb{D}^{\text{public}}$ , server training iterations  $I$ , heterogeneous device prototypes ( $i \in \mathbb{M}$ ) with their associated clients ( $\mathbb{C}^i$ ) and local datasets ( $\{\mathbb{D}_k^i\}_{k \in \mathbb{C}^i}$ ), model architecture ( $f^i$ ), local training iterations ( $I_{\text{local}}$ ), local learning rate ( $\eta_{\text{local}}$ ), server distillation iterations ( $I_{\text{distill}}$ ), and server distillation learning rate ( $\eta_{\text{distill}}$ ).

```

1: Server Executes:
2: Randomly initialize all device prototype's server model  $\{\theta_0^i\}_{i \in \mathbb{M}}$ 
3: for each round  $r = 0, 1, \dots, R - 1$  do
4:    $\mathbb{C}_r^i \leftarrow$  (randomly select clients from each device prototype)  $\forall i \in \mathbb{M}$ 
5:   for each client  $k \in \mathbb{C}_r^i, \forall i \in \mathbb{M}$  in parallel do
6:      $\hat{\theta}_k^i \leftarrow \text{ClientUpdate}(k; \theta_r^i)$ 
7:   end for
8:    $\theta_{\text{avg}}^i = \sum_{k \in \mathbb{C}_r^i} \frac{|\mathbb{D}_k^i|}{\sum_{k \in \mathbb{C}_r^i} |\mathbb{D}_k^i|} \hat{\theta}_k^i$ 
9:   for each device prototype's server student  $i = 1, 2, \dots, M$  in parallel do
10:    for each device prototype's teacher ensembles  $j = 1, 2, \dots, M$  in parallel do
11:       $\theta \leftarrow \theta_{\text{avg}}^i$ 
12:      for each server distillation iteration  $t = 0, 1, 2, \dots, I_{\text{distill}}$  do
13:         $x \leftarrow$  sample a mini-batch of data from public dataset  $\mathbb{D}^{\text{public}}$ 
14:         $\theta^{t+1} \leftarrow \theta^t - \eta_{\text{distill}} \cdot \nabla \mathcal{L}_S^{\tau_i}$  defined in Eq. 6.
15:      end for
16:       $\tau_j \leftarrow \theta^{I_{\text{distill}}} - \theta_{\text{avg}}^i$ 
17:    end for
18:     $\theta_{r+1}^i \leftarrow \theta_{\text{avg}}^i + \sum_{j=1}^M \lambda_j \tau_j$ 
19:  end for
20:   $\theta_{r+1}^i \leftarrow \theta^i$ 
21: end for

22: function ClientUpdate( $k, \theta_r^i$ )
23:    $\theta \leftarrow \theta_r^i$ 
24:   for each local update iteration  $t = 0, 1, \dots, I_{\text{local}} - 1$  do
25:      $\{x, y\} \leftarrow$  sample a mini-batch of data from local dataset  $\mathbb{D}_k^i$ 
26:      $\theta^{t+1} \leftarrow \theta^t - \eta_{\text{local}} \cdot \nabla \ell(f^i(x; \theta^t), y)$ 
27:   end for
28:    $\hat{\theta}_k^i \leftarrow \theta^{I_{\text{local}}}$ 
29: end function

```

---

## C Theoretical Results

### C.1 Proofs of the Main Propositions

First we present the formal assumptions associated with our theoretical derivations.

**Assumption 1.** Local federated averaging is performed with perfect test accuracy, i.e.,

$$\operatorname{argmin}_{\boldsymbol{\theta}^j} \sum_{j=1}^M \sum_{k=1}^{N^j} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbb{D}_k^j} [\ell(f^j(\mathbf{x}; \boldsymbol{\theta}^j), y)] = \operatorname{argmin}_{\boldsymbol{\theta}^j} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}^j} [\ell(f^j(\mathbf{x}; \boldsymbol{\theta}), y)] \quad (9)$$

That is, the training error on the datasets  $\{\mathbb{D}_k^j\}$  for the computed  $\boldsymbol{\theta}_{avg}^j$  is the same as the test error on the population distribution  $\mathcal{D}^j$ . Moreover assume that we can write  $\mathcal{T}_i = \left\{ \sum_{k=1}^{N^i} f^i(\cdot, \hat{\boldsymbol{\theta}}_k^i) | k \in \mathbb{C}^i \right\} = \{f^i(\cdot, \boldsymbol{\theta}_{avg}^i)\}$ . Finally, we assume that the same population distribution  $\sum_j \omega_j \mathcal{D}^j$  is the same that the clients perform their testing on as the server performs distillation on.

These assumptions are made for mathematical practicality while at the same time not starkly unreasonable. The local FL the device prototypes perform is generically prone to imprecision, especially as the clients' data varies, but this discrepancy is bounded [16]. Similarly the difference in the average of logits and the logit of averages has a bounded difference norm [51]. Thus, violations of the Assumption add additional perturbations to quantities derived in the Theoretical analysis without having structural/qualitative effects, and thus would only present clutter in the presentation.

**Notations.** Now we present the notation defining the specific quantities we refer to in the derivations below. The set of important quantities is given in Table 3. Note that the formal definitions of the first two quantities are,

$$\boldsymbol{\Theta}^j := \operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}^j} [\ell(f^j(\mathbf{x}; \boldsymbol{\theta}), y)], \quad \boldsymbol{\Theta}^{j,k} := \operatorname{argmin}_{\boldsymbol{\theta}} \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}^i} [\ell(f^j(\mathbf{x}; \boldsymbol{\theta}), y)]$$

Table 3: Notation and Definitions

Notation	Definition
$\boldsymbol{\Theta}^j$	Parameters in $j$ 's device model that minimize the loss on its population distribution
$\boldsymbol{\Theta}^{j,k}$	Parameters in $j$ 's device model that minimize the loss on $i$ 'th population distribution
$Q^j = \dim(\boldsymbol{\Theta}^j)$	The total capacity of device prototype $j$
$\mathcal{Q}^j = \{e_k^j\}_{k=1, \dots, Q^j}$	Eigenbasis for the model of device prototype $j$
$W^j = \dim(\boldsymbol{\Theta}^j)$	Dimension of the solution submanifold $\boldsymbol{\Theta}^j$
$W^{j,k} = \dim(\boldsymbol{\Theta}^{j,k})$	Dimension of the solution submanifold $\boldsymbol{\Theta}^{j,k}$
$\mathcal{W}^j = \{e_k^j\}_{k=1, \dots, W^j}$	Eigenbasis the solution submanifold $\boldsymbol{\Theta}^j$
$\mathcal{W}^{j,k} = \{e_l^{j,k}\}_{l=1, \dots, W^{j,k}}$	Eigenbasis the solution submanifold $\boldsymbol{\Theta}^{j,k}$

We shall make use of the ‘‘Choose’’ combinatorial operator, defined to be  $Ch(n, p) = \frac{n!}{p!(n-p)!}$ . The standard  $O(\cdot)$  notation indicates  $a_k = O(b_k)$  to mean there exists  $K$  and  $C$  such that for  $k \geq K$ ,  $a_k \leq Cb_k$ .

A recent finding that inspired the methodology in this work is the discovery of the weight disentanglement phenomenon underlying task arithmetic [37]. Indeed the *task arithmetic property* provides the ideal circumstance for federated knowledge transfer as we shall see below. Formally, adapting their definition to our notation:

**(Task Arithmetic Property)** holds for a set of vectors  $\{\boldsymbol{\tau}_j\}$  if for all  $j$  it holds that,

$$f^j \left( \mathbf{x}; \boldsymbol{\theta}_{avg}^j + \sum_{i \neq j} \lambda_i \boldsymbol{\tau}_i \right) = \begin{cases} f^j(\mathbf{x}; \boldsymbol{\theta}_{avg}^j + \lambda_i \boldsymbol{\tau}_i) & \mathbf{x} \in \mathcal{D}^i \\ f^j(\mathbf{x}; \boldsymbol{\theta}_{avg}^j) & \mathbf{x} \in \mathcal{D}^j \setminus \cup_{i \neq j} \mathcal{D}^i \end{cases} \quad (10)$$

Let us define an important property of task arithmetic that we shall use in the sequel.

**(Weight disentanglement).**[37] A parametric function  $f : \mathcal{X} \times \Theta \rightarrow \mathcal{Y}$  is weight disentangled with respect to a set of task vectors  $T = \{\tau_j\}_{j \in \mathbf{T}}$  and the corresponding supports  $\mathcal{D}_T := \{\mathcal{D}_j\}_{j \in \mathbf{T}}$  if

$$f(\mathbf{x}; \theta_0 + \sum_{i \in \mathbf{T}} \alpha_i \tau_i) = \sum_{i \in \mathbf{T}} g_j(\mathbf{x}; \alpha_i \tau_i) + g_0(\mathbf{x}),$$

where  $g_i(\mathbf{x}; \alpha_i \tau_i) = 0$  for  $\mathbf{x} \notin \mathcal{D}_i$  and  $i \in \mathbf{T}$ , and  $g_0(\mathbf{x}) = 0$  for  $\mathbf{x} \in \bigcup_{i \in \mathbf{T}} \mathcal{D}_i$ .

We now present the formal statements as well as the proofs of the main propositions.

**Proposition 1. (Information Loss in VED).** *Consider the VED procedure in the form of solving (3). Consider two device prototypes with a device capacity and solution dimension of  $Q^1, Q^2$  and  $W^1, W^2$ , respectively, and with associated eigenbases  $\mathcal{Q}^i, \mathcal{W}^i$ . Let the solution set of VED with prototype  $i$  as student be  $\hat{\Theta}_{VED}^i$  with  $\dim(\hat{\Theta}_{VED}^i) = W^{v_i}$  with eigenbasis  $\mathcal{W}^{v_i}$ . In addition, denote  $W^{s,t}$ ,  $s, t \in \{1, 2\}$  the dimension of the solution set for the student model trained on the data from the teacher device's ensembles. We assume that self-distillation is executed appropriately, e.g.,  $W^{1,1} = W^1$  and  $W^{2,2} = W^2$ .*

1. **Case 1:** Assume that  $Q^1 = Q^2$  and  $W^1 = W^2 = W^{1,2} = W^{2,1}$ . Then it holds that, in expectation,

$$\dim(\hat{\Theta}_{VED}^1 \cap [\mathcal{Q}^1 \setminus \mathcal{W}^1]) = O\left(\frac{(Q^1 - W^1)(W^1)!(Q^1 - W^{1,2})!}{Q^1!(W^1)!(Q^1 - W^1)! + Q^1!W^{1,2}!(Q^1 - W^{1,2})!}\right)$$

*This corresponds to the expected capacity of prototype 1 that is taken up for fitting logits that are not in the span of  $\mathcal{W}^1$ , that is, that do not fit the data corresponding to prototype 1.*

2. **Case 2:** Assume that  $Q^1 > Q^2$  and  $W^1 = W^{1,2} > W^2$ . Then the same quantity as for Case 1 holds. Moreover,

$$\dim(\hat{\Theta}_{VED}^1 \cap [\mathcal{Q}^1 \setminus (\mathcal{W}^1 \cup \mathcal{W}^{1,2})]) = O\left(\frac{(Q^1 - W^1)(W^1)!(W^{1,2} - W^2)!}{Q^1!W^{1,2}!(Q^1 - W^1)! + Q^1!W^2!(W^{1,2} - W^2)!}\right)$$

*This corresponds to capacity of client 1 that has been allocated but fits, in the model of prototype 1, neither the data of prototype 1, nor of the data of prototype 2.*

*Proof.* Formally,

$$\hat{\Theta}_{VED} := \operatorname{argmin}_{\theta \in \mathcal{Q}^1} \mathcal{L}_{ED} = \operatorname{argmin}_{\theta} \text{KL} \left[ \sum_{i=1,2} \sigma \left( f^i(\mathbf{x}, \theta_{avg}^i) \right), \sigma(\mathcal{S}(\mathbf{x})) \right]$$

Since by assumption  $\theta_{avg}^i$  solves the training problem on the data associated with device prototype  $i$ , the logit is accurate, and thus there is a map  $\mathcal{O}(i, j) : \mathcal{T}_i \rightarrow \mathcal{T}_j^i \subseteq \mathcal{W}^{i,j}$ . The self distillation, that is,  $\mathcal{S}_j$  defines a bijective map from  $\mathcal{W}^j$  to  $\mathcal{W}^j$  and thus does not affect the capacity allocation.

**Case 1:** In this case, generically (that is, measure zero on some non-atomic sampling on a distribution of operators)  $\mathcal{O}(i, j)$  is bijective. Now let us compute the expectation of the number of eigenvectors of, e.g.  $\mathcal{W}^{1,2}$  that are in the complement of the span of  $\mathcal{W}^1$ . Assuming, for simplicity, independence, this would correspond to counting the possible choices within the capacity of  $\mathcal{Q}^1 \setminus \mathcal{W}^1$  over the range of possible choices of filling the capacity of  $\mathcal{Q}^1$  with vectors in  $\mathcal{W}^1$  together with choices of filling it with vectors in  $\mathcal{W}^{1,2}$ :

$$\sum_{i=1}^{Q^1 - W^1} i \frac{Ch(Q^1 - W^1, i)}{Ch(Q^1, W^1) + Ch(Q^1, W^{1,2})}$$

For, e.g.,  $Q^1 = 4$  and  $W^1 = W^{1,2} = 2$  this is  $\frac{1}{3}$ .

To derive a scaling rate we can write:

$$\sum_i i \frac{\frac{(Q^1 - W^1)!}{i!(Q^1 - W^1 - i)!}}{\frac{Q^1!}{(W^1)!(Q^1 - W^1)!} + \frac{Q^1!}{W^{1,2}!(Q^1 - W^{1,2})!}} = O\left(\frac{(Q^1 - W^1)(W^1)!(Q^1 - W^{1,2})!}{Q^1!(W^1)!(Q^1 - W^1)! + Q^1!W^{1,2}!(Q^1 - W^{1,2})!}\right)$$

**Case 2:** In this case, it must be that, at best almost surely,  $\mathcal{O}(2, 1)$  is injective, but not surjective. This means that distilling from 2 to 1 does not fill the capacity of  $\mathcal{W}^{1,2}$ , and is thus a fundamentally

wasteful operation, that is  $|\mathcal{T}_i^j| = W^2 < W^{1,2}$ . Now let us compute the expectation of the number of eigenvectors of, e.g.  $\mathcal{W}^{1,2}$  that are in the complement of the span of  $\mathcal{W}^1$ . Since  $\mathcal{W}^{1,2}$  are being structurally allocated for fitting, the combinatorial expression is the same:

$$\sum_{i=1}^{Q^1-W^1} i \frac{Ch(Q^1-W^1, i)}{Ch(Q^1, W^1) + Ch(Q^1, W^{1,2})}$$

Thus for, e.g.,  $Q^1 = 4$  and  $W^1 = W^{1,2} = 2$  this is, again,  $\frac{1}{3}$ . The scaling in this case is

$$O\left(\frac{(Q^1-W^1)(W^1!)(Q^1-W^{1,2})!}{Q^1!W^1!(Q^1-W^1)! + Q^1!W^2!(Q^1-W^2)!}\right)$$

However, we observe that there are vectors in the range of  $\mathcal{W}^{1,2} \setminus \mathcal{O}(2, 1)(\mathcal{W}^2)$  that have been allocated by the VED but lie in neither  $\mathcal{W}^1$  nor in  $\mathcal{W}^{1,2}$ , that is, are garbage. We can compute those as the expected number of eigenvectors arising from allocating  $\mathcal{W}^{1,2} \setminus \mathcal{O}(2, 1)(\mathcal{W}^2)$  that intersect with  $\mathcal{Q}^1 \setminus \mathcal{W}^1$  (that is, the spare capacity not used for fitting data  $\mathcal{D}^1$ ). This is, using similar principles:

$$\sum_{i=1}^{W^{1,2}-W^1} i \frac{Ch(Q^1-W^1, i)}{Ch(Q^1, W^1) + Ch(Q^1, W^{1,2}-W^2)}$$

This is for, e.g.,  $Q^1 = 4$ ,  $W^{1,2} = 2$  and  $W^2 = 1$ , this would be  $\frac{3}{10}$

The scaling here is

$$O\left(\frac{(Q^1-W^1)(W^1!)(W^{1,2}-W^2)!}{Q^1!W^1!(Q^1-W^1)! + Q^1!W^2!(W^{1,2}-W^2)!}\right)$$

■

**Proposition 2. (Improve knowledge transfer with task arithmetic).** Consider the TAKFL procedure as in the form of computing (8). Consider two device prototypes with a device capacity and solution dimension of  $Q^1, Q^2$  and  $W^1, W^2$ , respectively, and with associated eigenbases  $\mathcal{Q}^i, \mathcal{W}^i$ . Let the solution set of TAKFL with prototype  $i$  as student be  $\hat{\Theta}_{TA}^i$  with  $\dim(\hat{\Theta}_{TA}^i) = W^i$  with eigenbasis  $\mathcal{W}^i$ . In addition, denote  $W^{s,t}$ ,  $s, t \in \{1, 2\}$  dimension of the solution set for the student model trained on the data from the teacher device's ensembles. . The following statements hold:

In the case that that  $Q^1 \geq Q^2$  and  $W^1 \geq W^2$ , it holds that the TAKFL preserves that the eigenbasis used to model the data  $\mathcal{D}^1$ 's accuracy for device prototype 1, that is for student 1

$$\dim(\mathcal{W}^v \cap [\mathcal{Q}^1 \setminus \mathcal{W}^1]) = 0$$

**Case 1:** Assume that  $Q^1 = Q^2$  and  $W^1 = W^2$ . Then it holds that,

$$\dim(\mathcal{W}^v \cap [\mathcal{Q}^1 \setminus (\mathcal{W}^1 \cup \mathcal{W}^{1,2})]) = 0$$

Moreover, it holds that,

$$\hat{\Theta}_{TA} \in \text{Span}(\mathcal{W}^1 \cap \mathcal{W}^{1,2})$$

Thus, with equal capacity, no information is lost in Task Arithmetic aided knowledge ensemble distillation and capacity is efficiently used to model the data from both prototype 1 and prototype 2.

**Case 2:** Assume that  $Q^1 > Q^2$  and  $W^1 > W^2$ . Then it again holds that,

$$\dim(\mathcal{W}^v \cap [\mathcal{Q}^1 \setminus (\mathcal{W}^1 \cup \mathcal{W}^{1,2})]) = 0$$

However, while  $\hat{\Theta}_{TA} \in \text{Span}(\mathcal{W}^1)$ , it holds that  $\dim(\mathcal{W}^v \cap \mathcal{W}^{1,2}) = W^2 < W^{1,2}$ .

*Proof.* We can see immediately from the weight disentanglement property of Task Arithmetic that,

$$f^1(\mathbf{x}; \theta_{avg}^1 + \alpha_1 \tau_1 + \alpha_2 \tau_2) = g^{1,1}(\mathbf{x}; \alpha_1 \tau_1) + g^{1,2}(\mathbf{x}; \alpha_2 \tau_2) + g^{1,0}(\mathbf{x})$$

with  $g^{1,1}(\mathbf{x}; \alpha_1 \tau_1)$  for  $\mathbf{x} \notin \mathcal{D}^1$ ,  $g^{1,2}(\mathbf{x}; \alpha_2 \tau_2)$  for  $\mathbf{x} \notin \mathcal{D}^2$  and  $g^{1,0}(\mathbf{x}) = 0$  for  $\mathbf{x} \in \mathcal{D}^1 \cup \mathcal{D}^2$ . From this, we can immediately conclude the first statement of the Proposition as well as the expression

$$\dim(\mathcal{W}^v \cap [\mathcal{Q}^1 \setminus (\mathcal{W}^1 \cup \mathcal{W}^{1,2})]) = 0$$

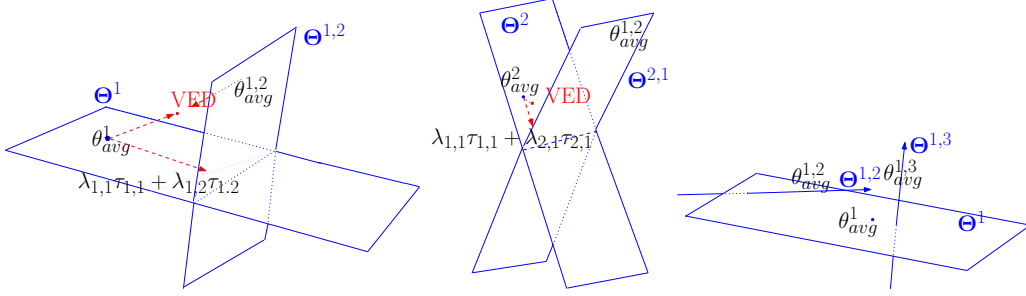


Figure 5: **Illustration of Geometric Intuition.** Each panel presents a different case example. The left and center panels present the geometric intuition of KD for vanilla ensemble distillation (VED) and TAKFL in the case where two different large device prototypes performing knowledge transfer. The planes represent the solution subspaces. The right panel presents a circumstance by which two small device prototypes (2, and 3) serve as teacher for transferring knowledge to a larger device prototype 1.

and also, in the case of  $W^1 = W^{1,2} = W^2$  implies

$$\hat{\Theta}_{TA} \in \text{Span}(\mathcal{W}^1 \cap \mathcal{W}^{1,2})$$

For the last statement we observe again as in the second Case in the Proposition describing VED,  $\dim(\mathcal{O}(2, 1)(\mathcal{W}^2)) < W^{1,2}$  from which we can conclude that, generically

$$\dim(\{v : v \in \mathcal{O}(2, 1)(\mathcal{W}^2) \subseteq \mathcal{W}^v, \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}^2} l(f^1(\mathbf{x}; v), y) > 0\}) = W^{1,2} - W^2$$

proving the final statement. ■

We observe that a key mechanism of the proof is the dimension of the target space of the teaching operator  $\mathcal{O}(i, j)$ . As an informative model, we can consider coefficients  $\lambda_j$  of task vectors as restricting the rank, relative to other teachers. For instance, in the previous Proposition, if  $W^{1,2} = 2$  and  $W^2 = 1$ , then  $\lambda_2 = 1/2$ , so as to enforce one vector of  $\mathcal{W}^{1,2}$  is a target for the map  $\tilde{\mathcal{O}}(2, 1)$ , would be appropriately sensible.

## C.2 Geometric Intuition

In this section we aim to provide geometric intuition for the mechanism of VED and Task Arithmetic KD on three different cases. Figure 5 presents the geometry illustration for three different cases. We discuss each case in the following.

**Case I: KD between two large prototypes with different data distributions.** Consider Figure 5 left panel. This panel corresponds to a setting where two large device prototypes with similar total capacity, i.e.  $Q^1 = Q^2 = 3$  perform knowledge transfer. We consider the solution dimensions of both prototypes to be the same, i.e.  $W^1 = W^2 = 2$ . These would correspond to planes in the ambient space. Therefore, one plane corresponds to the solution subspace of the prototype 1 trained on its own data, i.e.  $\Theta^1$  subspace in the panel, and the other corresponds to the (theoretical) solution subspace of this prototype trained on prototype 2's data, i.e.  $\Theta^{(1,2)}$  in the panel. In this case, since the data distributions of the prototypes are fairly disparate, this has resulted into near orthogonal subspaces corresponding to these solutions. As we can see from the panel, VED will lead to point which is far away from either of the planes corresponding to optimal solution subspaces, and far from the optimal set of parameters, which is their intersection, suggesting a loss of knowledge. By contrast, the TAKFL approach, by customizing the merging coefficients and putting each to half, i.e.  $\lambda_{1,1}\lambda_{1,2} = 0.5$ , can traverse in the tangent space of zero loss surface and get into the intersection subspace which is exactly the optimal solution ( $\theta_{merged}^* = \theta_{avg}^1 + \lambda_{1,1}\tau_{1,1} + \lambda_{1,2}\tau_{1,2}$ ).

**Case II: KD between two large prototypes with similar data distributions.** Consider Figure 5 center panel. Similar to Case I, this panel corresponds to a setting where two large device prototypes with similar total capacity, i.e.  $Q^1 = Q^2 = 3$  performing knowledge transfer. We consider the solution dimensions of both prototypes to be the same, i.e.  $W^1 = W^2 = 2$ . These would correspond



to planes in the ambient space. Therefore, one plane corresponds to the solution subspace of the prototype 1 trained on its own data, i.e.  $\Theta^1$  subspace in the panel, and the other corresponds to the (theoretical) solution subspace of this prototype trained on prototype 2's data, i.e.  $\Theta^{(1,2)}$  in the panel. In this case, since the data distributions of the prototypes are fairly close, this has resulted into non-orthogonal solution subspaces. As we can see from the panel, while VED could still lead to some information loss, by and large we expect straightforward KD. Our task arithmetic (TA) approach, again by customizing the merging coefficients  $\lambda_{1,1}\lambda_{1,2} = 0.5$ , can traverse in the tangent space of zero loss surface on each plane and get into the intersection subspace, corresponding to the most efficient allocation of device prototype capacity for fitting simultaneously the logits corresponding to accurate modeling of device prototype 1 as well as device prototype 2's data distribution.

**Case III: KD between two small prototypes and one large prototype.** Now consider the right panel in Figure 5. This panel corresponds to a setting where two small prototypes serve as teachers and one large prototype is the student. The  $\theta_{avg}^1$  plane corresponds to the solution subspace of the large prototype 1 on its own data,  $\mathcal{D}^1$ . The line  $\theta_{avg}^{1,2}$  corresponds to the subspace of solutions in prototype 1's parameter space projected into the capacity of the information transferred from device prototype 2. Finally, the line labelled  $\theta_{avg}^{1,3}$  corresponds to the subspace of solutions in prototype 1's parameter space projected into the capacity of the information transferred from device prototype 3. Here, we can see from the relative angle of the lines with respect to the plane that the distribution  $\mathcal{D}^1$  is closer to the distribution  $\mathcal{D}^2$  than to  $\mathcal{D}^3$ . Comparing this case to the previous cases,  $\theta_{avg}^{1,3}$  is like case I and  $\theta_{avg}^{1,2}$  is like case II. We can apply the same conclusions here as well regarding the performance of vanilla ensemble distillation and our adaptive task arithmetic approach. We can see from the geometric visualization that knowledge distillation towards  $\theta_{avg}^{1,3}$  has more margin of error for prototype 1. Therefore, with the TAKFL approach large prototype 1 can strategically select which prototype to learn more from, and since  $\theta_{avg}^{1,2}$  has closer data distribution to prototype 1, TAKFL will prioritize this by putting a larger merging coefficient, i.e.  $\lambda_{1,2} > \lambda_{1,3}$ . By contrast, VED lacks this customization and results in sub-optimal knowledge distillation.

The geometric intuition discussed here is consistent with our detailed experimental analysis in D.1.1.

### C.3 Analytical Properties of Learning Dynamics

Here we provide additional insights from the literature as to the nature and properties of learning as it takes place on overparametrized models. Specifically, we comment on literature in the area of Stochastic Differential Equation (SDE) models for SGD training dynamics, and its correspondence to the results above. Overparametrization has been conjectured to be a significant factor in contributing to the (unexpected, by classical bias-variance tradeoffs) generalization ability of deep neural networks, from a number of perspectives [13].

Consider the diffusion model of SGD training for overparametrized NNs provided in [29]. Their analysis relies on the following two assumptions. For our purposes  $L$  is shorthand for a client group's loss,  $L(\theta) = \sum_{k \in \mathbb{C}^j} \mathbb{E}_{(x,y) \sim \mathbb{D}_k^j} [\ell(f^j(x; \theta), y)]$  for some  $j$ , which will be identified from the context.

**Assumption 2.**  $L : \mathbb{R}^Q \rightarrow \mathbb{R}$  is  $C^3$  and the solution set  $\Gamma$  is a  $W$ -dimensional  $C^2$ -submanifold of  $\mathbb{R}^D$  for  $0 \leq W \leq D$  and  $\text{rank}(\nabla^2 L(\theta)) = Q - W$

**Assumption 3.** Assume that  $U$  is an open neighborhood of  $\Gamma$  satisfying that gradient flow starting in  $U$  converges to some point in  $\Gamma$

From these, [29] derives Theorem 4.6. This Theorem decomposes the random process of the parameter weights driven by SGD after it has reached the solution manifold, e.g., the diffusive random walk of  $\theta_{avg}^1$  in Figure 5 along its respective solution manifold.

**Theorem 1.** Given  $L$ ,  $\Gamma$  and  $\theta_\mu(0) = \theta(0) \in U$  as by Assumptions 2 and 3 the SDE modeling the optimization of  $F$  by SGD, that is, defining  $\Phi(X)$  to be the gradient flow applied to the state random variable  $X$ , then for  $T$  as long as  $\mathbb{P}[Y(t) \in U, \forall 0 \leq t \leq T] = 1$ ,  $\theta_\eta(\lfloor T/\eta^2 \rfloor)$  converges in distribution to the stochastic process  $Y(T)$  as  $\eta \rightarrow 0$ , with  $Y(t)$  given as,

$$dY(t) = \Sigma_{\parallel}^{1/2}(Y) dW(t) - \frac{1}{2} \nabla^2 L(Y)^\dagger \partial^2 (\nabla L)(Y) [\Sigma_{\parallel}(Y)] dt \\ - \frac{1}{2} \partial \Phi(Y) (2\partial^2 (\nabla L)(Y) [\nabla^2 L(Y)^\dagger \Sigma_{\perp, \parallel}(Y)] + \partial^2 (\nabla L)(Y) [\mathcal{L}_{\nabla^2 L}^{-1}(\Sigma_{\perp}(Y))]) dt \quad (11)$$

where  $\Sigma \equiv \sigma\sigma^T$  and  $\Sigma_{\parallel}, \Sigma_{\perp}, \Sigma_{\parallel,\perp}$  are given as,

$$\begin{aligned}\Sigma_{\parallel} &= \partial\Phi\Sigma\partial\Phi, \Sigma_{\perp} = (I_D - \partial\Phi)\Sigma(I_D - \partial\Phi), \\ \Sigma_{\parallel,\perp} &= \partial\Phi\Sigma(I_D - \partial\Phi)\end{aligned}\tag{12}$$

This theorem indicates that the asymptotic flow of SGD on the client training can be decomposed into a covariance-driven random walk in the tangent space, drift to preserve the flow into the tangent plane, the tangent-normal portion of the noise covariance and noise in the normal direction.

This analytical expression provides the probabilistic foundations for the more higher level theoretical results above. In particular, local gradient dynamics, as employed by individual device prototypes  $j$  using FedAvg on its local clients, yields a flow for the stochastic process defined by the weights. At this point of learning, the weights are traversing the solution set, with noise predominantly in the tangent directions. Thus knowledge distillation which preserves this noise structure is going to be more effective as far as preserving accuracy across data.

## D Detailed Experimental Results

In this section we present a more detailed version of experimental results presented in the main paper Section 7. Additional experimental results are also presented here.

### D.1 Main Experimental Results on Computer Vision (CV) Task

**The experiments in this section complements the main experimental results in the main paper Section 7.2.**

**Experimental Setup.** For the evaluation on CV task, we employ CIFAR-10 and CIFAR-100 [24] datasets. For CIFAR-10, we use CIFAR-100 as the unlabeled public dataset, while ImageNet-100, a subset of ImageNet [10] with 100 classes (see Appendix F.1.1), is used for CIFAR-100. We distribute the training dataset among the device prototypes in a ratio of 1:3:6 for S, M, and L, respectively. Each device prototype’s data portion is further distributed among its clients using a Dirichlet distribution. We apply two levels of data heterogeneity for a comprehensive evaluation: low heterogeneity, i.e. Dir(0.3), and high heterogeneity, i.e. Dir(0.1). Additionally, we configure the number of clients and their sampling rates as follows: 100 clients for S, 20 for M, and 4 for L, with sampling rates set at 0.1, 0.2, and 0.5 respectively. To comprehensively evaluate, we use two distinct architectural settings: the “*homo-family*” setting, where all device prototypes’ architectures are from the same family—employing ResNet8, ResNet14, and ResNet18 [17] for S, M, and L, respectively; and the “*hetero-family*” setting, which diverse architectures are used—ViT-S [12] for S, ResNet14 for M, and VGG-16 [42] for L. All models are initialized from scratch, and the communication round is set at 60 rounds. Further details regarding hyper-parameters can be found in Table 12.

**Overview of Performance Results.** Table 4 presents the performance of TAKFL across diverse architecture settings on the CIFAR-10 and CIFAR-100 datasets. TAKFL consistently improves all device prototypes of different sizes in various cases by a significant margin compared to the base-lines, achieving SOTA performance. Notably, in the homo-family architecture setting with Dir(0.3) on CIFAR-10, TAKFL improves average performance across all prototypes by 8%, and by 4% on CIFAR-100. In the hetero-family settings with Dir(0.1) on CIFAR-10 and Dir(0.3) on CIFAR-100, TAKFL enhances performance by  $\sim 3\%$  and 1%, respectively. Furthermore, we observe that our self-regularization technique has successfully mitigated issues associated with the noisy and unsupervised ensemble distillation process, thereby enhancing performance. Generally, the performance gains from self-regularization are more pronounced in low data heterogeneity cases, where prototypes’ models perform better and possess higher quality self-knowledge. Thus, self-regularization proves more effective as it preserves this higher quality self-knowledge.

#### D.1.1 Consistency Analysis: Experimental and Theoretical Correlations

In this part, we elaborate on our key experimental observations and their alignment with our theoretical findings.

**Insight 1:** From Table 4, it is evident that prior KD-based methods show inconsistent performance across various device prototypes, particularly for the large (L) prototype. For instance, in the CIFAR-10 homo-family setting with Dir(0.3), while small (S) and medium (M) prototypes see performance

Table 4: **Performance Results for CV task on CIFAR-10 and CIFAR-100.** Training data is distributed among S, M, and L device prototypes in a 1:3:6 ratio, subdivided among clients using Dirichlet distribution. Public datasets are CIFAR-100 for CIFAR-10 and ImageNet-100 for CIFAR-100. Client configurations include 100, 20, and 4 clients for S, M, and L, with sampling rates of 0.1, 0.2, and 0.5. In homo-family settings, architectures are ResNet8, ResNet14, and ResNet18; in hetero-family settings, they are ViT-S, ResNet14, and VGG-16. All models are trained from scratch for 60 rounds. See Appendix F.1 for more details.

Dataset	Baseline	Homo-family Architecture Setting							
		Low Data Heterogeneity				High Data Heterogeneity			
		S	M	L	Average	S	M	L	Average
CIFAR-10	FedAvg	36.21 $\pm$ 2.24	46.41 $\pm$ 2.33	59.46 $\pm$ 6.17	47.36	22.01 $\pm$ 0.78	25.26 $\pm$ 3.89	51.51 $\pm$ 3.52	32.93
	FedDF	49.31 $\pm$ 0.15	50.63 $\pm$ 0.73	49.82 $\pm$ 0.98	49.92	34.71 $\pm$ 1.48	35.27 $\pm$ 4.74	51.08 $\pm$ 4.04	40.35
	FedET	49.21 $\pm$ 0.72	55.01 $\pm$ 1.81	53.60 $\pm$ 6.47	52.61	29.58 $\pm$ 3.00	30.96 $\pm$ 4.70	45.53 $\pm$ 6.46	35.36
	TAKFL	55.90 $\pm$ 1.70	57.93 $\pm$ 3.49	60.58 $\pm$ 2.35	58.14	37.40 $\pm$ 1.68	38.96 $\pm$ 0.17	51.49 $\pm$ 6.15	42.61
	TAKFL+Reg	<b>56.37<math>\pm</math>0.46</b>	<b>58.60<math>\pm</math>0.43</b>	<b>65.69<math>\pm</math>1.28</b>	<b>60.22</b>	<b>40.51<math>\pm</math>1.05</b>	<b>40.12<math>\pm</math>1.24</b>	<b>53.24<math>\pm</math>2.51</b>	<b>44.62</b>
CIFAR-100	FedAvg	13.22 $\pm$ 0.14	21.39 $\pm$ 1.11	29.47 $\pm$ 0.86	21.36	11.86 $\pm$ 0.08	14.63 $\pm$ 0.65	26.25 $\pm$ 1.64	17.58
	FedDF	19.54 $\pm$ 0.20	24.32 $\pm$ 0.45	29.29 $\pm$ 1.45	24.38	16.09 $\pm$ 0.32	19.80 $\pm$ 0.17	26.59 $\pm$ 0.25	20.83
	FedET	19.67 $\pm$ 0.35	25.27 $\pm$ 0.66	31.10 $\pm$ 1.53	25.35	11.18 $\pm$ 1.68	18.22 $\pm$ 0.35	26.40 $\pm$ 0.65	18.60
	TAKFL	24.48 $\pm$ 0.42	27.60 $\pm$ 0.25	29.84 $\pm$ 0.94	27.31	<b>22.90<math>\pm</math>0.18</b>	23.63 $\pm$ 0.72	26.98 $\pm$ 0.13	24.50
	TAKFL+Reg	<b>27.18<math>\pm</math>0.27</b>	<b>29.14<math>\pm</math>0.20</b>	<b>31.15<math>\pm</math>0.97</b>	<b>29.16</b>	22.88 $\pm$ 0.37	<b>23.92<math>\pm</math>0.57</b>	<b>28.01<math>\pm</math>0.34</b>	<b>24.94</b>
Dataset	Baseline	Hetero-family Architecture Setting							
		Low Data Heterogeneity				High Data Heterogeneity			
		S	M	L	Average	S	M	L	Average
CIFAR-10	FedAvg	27.53 $\pm$ 0.83	47.30 $\pm$ 3.17	55.10 $\pm$ 8.60	43.31	20.93 $\pm$ 1.54	25.62 $\pm$ 6.04	36.80 $\pm$ 5.47	27.78
	FedDF	34.15 $\pm$ 0.87	54.06 $\pm$ 1.06	69.07 $\pm$ 4.99	52.43	24.20 $\pm$ 0.74	34.07 $\pm$ 3.08	39.81 $\pm$ 5.45	32.69
	FedET	33.24 $\pm$ 1.27	58.86 $\pm$ 0.94	65.56 $\pm$ 3.49	52.55	24.37 $\pm$ 1.26	37.77 $\pm$ 4.71	43.64 $\pm$ 3.36	35.26
	TAKFL	33.29 $\pm$ 0.15	57.64 $\pm$ 0.19	68.44 $\pm$ 0.66	53.12	24.92 $\pm$ 1.32	38.07 $\pm$ 3.19	48.01 $\pm$ 3.99	37.00
	TAKFL+Reg	<b>33.34<math>\pm</math>3.36</b>	<b>59.01<math>\pm</math>3.12</b>	<b>70.22<math>\pm</math>4.40</b>	<b>54.19</b>	<b>25.10<math>\pm</math>1.87</b>	<b>38.81<math>\pm</math>5.36</b>	<b>50.26<math>\pm</math>6.42</b>	<b>38.06</b>
CIFAR-100	FedAvg	8.51 $\pm$ 0.37	22.11 $\pm$ 0.58	37.91 $\pm$ 2.60	22.84	7.01 $\pm$ 0.47	14.94 $\pm$ 0.96	28.51 $\pm$ 1.46	16.82
	FedDF	10.46 $\pm$ 0.17	23.46 $\pm$ 0.65	36.81 $\pm$ 0.82	23.58	7.76 $\pm$ 0.40	18.92 $\pm$ 0.39	29.81 $\pm$ 1.09	18.83
	FedET	11.16 $\pm$ 0.18	25.40 $\pm$ 0.30	37.38 $\pm$ 0.60	24.65	8.20 $\pm$ 0.54	20.66 $\pm$ 0.50	28.95 $\pm$ 1.79	19.27
	TAKFL	10.29 $\pm$ 0.11	27.14 $\pm$ 0.89	<b>39.15<math>\pm</math>0.88</b>	25.53	7.88 $\pm$ 0.68	21.41 $\pm$ 0.37	31.31 $\pm$ 0.66	20.20
	TAKFL+Reg	<b>11.25<math>\pm</math>0.37</b>	<b>27.86<math>\pm</math>0.86</b>	38.68 $\pm$ 0.45	<b>25.93</b>	<b>8.45<math>\pm</math>0.20</b>	<b>22.16<math>\pm</math>0.87</b>	<b>31.95<math>\pm</math>1.13</b>	<b>20.85</b>

gains, the L prototype experiences up to a  $\sim 10\%$  performance decline compared to vanilla FedAvg, which lacks server-side knowledge distillation. This trend is consistent across other settings, such as CIFAR-10 Dir(0.1) homo-family and CIFAR-100 Dir(0.3) homo-family. These outcomes underline the dilution problem inherent in existing methods, where the valuable insights from larger, more capable device prototypes are overshadowed by less informative outputs from smaller devices, thereby degrading the performance of L prototypes. These empirical findings are supported by our theoretical insights as discussed in Remark 1. Specifically, Proposition 1 illustrates that vanilla ensemble distillation (VED) leads to knowledge dilution and inaccuracies due to misaligned device capacity allocations. Moreover, this issue becomes more significant when the smaller device prototype serve as teacher.

**Insight 2:** From Table 4, the suboptimality of existing KD-based methods is evident from the significant performance improvements of our method, especially for S and M prototypes across various settings. This underscores the ineffectiveness of the one-size-fits-all approach these methods employ, where a single averaged logits distillation target is used for all device sizes, proving to be sub-optimal. Our experimental observations regarding the shortcomings of vanilla ensemble distillation methods align with our theoretical findings, as substantiated in Remark 1 and 2. It becomes evident that an efficient knowledge distillation process must allocate capacity in a manner that appropriately corresponds to the information value of the teacher ensemble prototypes.

**Insight 3:** Our experiments, detailed in Table 4, demonstrate TAKFL’s adept handling of knowledge from various device prototypes under different data heterogeneity conditions. We observed consistent performance gains for small (S) and medium (M) prototypes across both low and high data heterogeneity, compared to vanilla FedAvg. However, in high heterogeneity settings, large (L) prototypes show less improvement, prompting the question: *What can smaller device prototypes offer to larger ones?*

In low heterogeneity scenarios, large prototypes significantly benefit from the collective knowledge, showing enhanced performance. Conversely, in conditions of extreme heterogeneity, where smaller models contribute less effectively, the performance improvements for larger devices are notably reduced. This pattern highlights TAKFL’s ability to intelligently manage and utilize the available knowledge, selectively distilling information based on the intrinsic capacity and contributions of

Table 5: **Performance Results for CV task on TinyImageNet, STL-10, and CINIC-10 using pre-trained models.**

Private	Public	Baseline	S	M	L	Average
TinyImageNet	STL-10	FedAvg	8.97	13.03	15.12	12.37
		FedMH	15.08	17.10	17.83	16.67
		FedET	10.60	16.39	17.62	14.87
		TAKFL	16.10	17.60	19.03	17.58
		TAKFL+Reg	<b>16.55</b>	<b>17.98</b>	<b>19.74</b>	<b>18.09</b>
STL-10	CIFAR-100	FedAvg	26.01	34.47	42.88	34.45
		FedMH	28.64	34.55	39.25	34.15
		FedET	29.87	33.00	38.26	33.71
		TAKFL	29.57	37.57	42.53	36.56
		TAKFL+Reg	<b>30.78</b>	<b>37.89</b>	<b>43.38</b>	<b>37.35</b>
CINIC-10	CIFAR-100	FedAvg	44.87	55.49	51.33	50.56
		FedMH	45.52	55.75	53.48	51.58
		FedET	46.31	57.43	53.01	52.25
		TAKFL	<b>48.21</b>	<b>57.81</b>	52.74	<b>52.92</b>
		TAKFL+Reg	47.66	57.54	<b>53.25</b>	52.82

each prototype, and integrating only the most valuable knowledge from smaller devices when beneficial.

By contrast, our analysis of existing KD-based methods shows their failure to effectively discern and utilize the most informative knowledge across prototypes. These methods often overload the capacity of larger prototypes with suboptimal or irrelevant information, particularly in high heterogeneity environments, leading to not just stagnation but an accumulation of inefficiencies. These experimental observations align with our theoretical insights, as outlined in Remark 2, which emphasizes the crucial combinatorial constraint of capacity and diverse information. This further confirms the superiority of TAKFL’s adaptive approach to knowledge distillation in diverse federated learning environments.

## D.2 Additional Experimental Results on CV Task

**Experimental Setup.** For additional evaluation of the CV task, we conducted experiments on TinyImageNet [25], STL-10 [7], and CINIC-10 [9] datasets using pre-trained models. For TinyImageNet [25], we utilized STL-10 [7] as the unlabeled public dataset. STL-10 [7] and CINIC-10 [9] both employ CIFAR-100 [24] as their respective public datasets. We distributed the training datasets among the device prototypes in a 2:3:5 ratio for small (S), medium (M), and large (L) prototypes, respectively. The data portion for each prototype was further subdivided among its clients using a Dirichlet distribution: Dir(1.0) for TinyImageNet and Dir(0.3) for both STL-10 and CINIC-10. Client configurations were set with 4, 3, and 2 clients for S, M, and L, respectively, all with a sampling rate of 1.0. The architectures employed were MobileNetV3-Large [19] for S, MobileViTV2 [34] for M, and ResNet-34 for L, sourced from the TIMM library.<sup>1</sup> The local training was conducted over 10 epochs using an Adam [23] optimizer with a learning rate of 1e-3 and weight decay of 1e-5. For server-side distillation, the epoch count was 10 for TinyImageNet and 1 for STL-10 and CINIC-10, with a batch size of 128, employing an Adam optimizer with a learning rate of 1e-5 and weight decay of 1e-5. For TinyImageNet, public dataset images from STL-10 were resized to 64×64, while for STL-10, images were resized to 32×32. No data augmentation was used. The communication rounds is fixed to 40. These experiments were conducted by only 1 trial. Table 14 details the configurations.

**Performance Results.** Table 5 presents the results. The superiority of TAKFL’s performance across these challenging datasets using pre-trained models is evident here as well.

## D.3 Additional Experimental Results on Natural Language Processing (NLP) Task

The experiments in this section complements the main experimental results in the main paper Section 7.2.

<sup>1</sup><https://github.com/huggingface/pytorch-image-models>

Table 6: **Performance Results for NLP Task on 4 Datasets.** Training data is distributed among S, M, and L device prototypes in a 1:3:6 ratio, subdivided among clients using Dir(0.5). Client configurations are 8, 4, and 2 clients for S, M, and L, with sample rates of 0.3, 0.5, and 1.0, respectively. Architectures include Bert-Tiny, Bert-Mini, and Bert-Small for S, M, and L, initialized from pre-trained parameters and fine-tuned for 20 communication rounds. See Appendix F.2 for more details.

Private	Public	Baseline	S	M	L	Average
MNLI	SNLI	FedAvg	36.15 $\pm$ 0.46	54.47 $\pm$ 2.48	57.51 $\pm$ 2.79	49.37
		FedDF	54.21 $\pm$ 0.15	60.44 $\pm$ 1.91	66.71 $\pm$ 1.09	60.45
		FedET	48.03 $\pm$ 6.32	50.33 $\pm$ 7.87	53.80 $\pm$ 6.18	50.72
		TAKFL	57.43 $\pm$ 0.21	63.58 $\pm$ 0.31	68.74 $\pm$ 0.12	63.25
		TAKFL+Reg	<b>57.61<math>\pm</math>0.89</b>	<b>63.91<math>\pm</math>1.05</b>	<b>68.96<math>\pm</math>1.10</b>	<b>63.49</b>
SST2	Sent140	FedAvg	54.98 $\pm$ 1.81	74.71 $\pm$ 8.22	86.69 $\pm$ 0.06	72.13
		FedDF	74.41 $\pm$ 2.62	80.71 $\pm$ 1.63	84.35 $\pm$ 1.66	79.82
		FedET	66.63 $\pm$ 9.14	65.89 $\pm$ 16.35	70.05 $\pm$ 15.83	67.52
		TAKFL	74.73 $\pm$ 0.55	82.17 $\pm$ 0.31	86.93 $\pm$ 0.42	81.28
		TAKFL+Reg	<b>74.88<math>\pm</math>0.43</b>	<b>82.40<math>\pm</math>0.83</b>	<b>87.33<math>\pm</math>0.63</b>	<b>81.54</b>
MARC	Yelp	FedAvg	33.76 $\pm$ 1.13	49.08 $\pm$ 1.28	59.26 $\pm$ 1.43	47.36
		FedDF	53.01 $\pm$ 1.24	55.37 $\pm$ 0.87	56.81 $\pm$ 0.99	55.06
		FedET	52.63 $\pm$ 2.29	54.28 $\pm$ 2.31	56.11 $\pm$ 2.61	54.34
		TAKFL	55.70 $\pm$ 2.08	58.64 $\pm$ 1.75	59.39 $\pm$ 1.16	57.91
		TAKFL+Reg	<b>55.96<math>\pm</math>1.66</b>	<b>59.18<math>\pm</math>1.13</b>	<b>59.61<math>\pm</math>1.89</b>	<b>58.25</b>
AG-News	DBPedia	FedAvg	83.64 $\pm$ 3.51	83.47 $\pm$ 2.35	91.48 $\pm$ 2.22	86.20
		FedDF	85.97 $\pm$ 2.45	89.10 $\pm$ 1.85	91.37 $\pm$ 1.10	88.81
		FedET	75.27 $\pm$ 3.85	81.13 $\pm$ 3.21	83.19 $\pm$ 4.58	79.86
		TAKFL	87.37 $\pm$ 1.31	90.11 $\pm$ 1.56	92.48 $\pm$ 1.12	89.99
		TAKFL+Reg	<b>87.66<math>\pm</math>1.83</b>	<b>90.30<math>\pm</math>2.05</b>	<b>92.61<math>\pm</math>1.72</b>	<b>90.19</b>

**Experimental Setup.** For the evaluation of NLP tasks, we utilize four datasets: MNLI [50], SST-2 [43], MARC [22], and AG-news [58]. The corresponding unlabeled public datasets are SNLI [2] for MNLI, Sentiment140 [14] for SST-2, Yelp [59] for Amazon, and DBPedia [57] for AG-News. The training data is distributed among the device prototypes in a ratio of 1:3:6 for small (S), medium (M), and large (L) categories, respectively, with each portion further subdivided among its clients using a Dirichlet distribution (Dir(0.5)). The client configurations and their sampling rates are set as follows: 8, 4, and 2 clients for S, M, and L categories, respectively, with sampling rates of 0.3, 0.5, and 1.0. The architectures employed for each prototype size are BERT [45] -Tiny, -Small, and -Mini, respectively, each initialized from pre-trained parameters and tested over 20 communication rounds. Additional details regarding hyper-parameters and datasets are presented in Appendix F.2 and Table 16.

**Performance on NLP Task.** Table 6 presents the results on four different datasets: MNLI, SST-2, MARC, and AG-News. Similar to the CV task, TAKFL has consistently improved performance across all device prototypes of varying sizes, achieving state-of-the-art results. On MNLI, it has enhanced average performance across all prototypes by 3%, on SST-2 by  $\sim$ 2%, on MARC by 3%, and on AG-News by  $\sim$ 1.50%. As observed in the CV task, the suboptimality of existing KD-based methods is also evident here. Notably, FedET exhibits very poor performance compared vanilla FedAvg, failing to achieve satisfactory results on all datasets except for the MARC dataset. Particularly, the performance of the L prototype has consistently decreased across all datasets compared to vanilla FedAvg. This behavior can be attributed to FedET’s reliance on neural network confidence scores for uncertainty estimates in its uncertainty-weighted distillation. However, neural networks, especially pretrained language models (PLMs), are often poorly calibrated and prone to overconfidence, which compromises their ability to provide reliable uncertainty estimates [48, 15, 5, 53].

#### D.4 Scalability Evaluation

This section complements the experimental results in the main paper Section 7.3.

**Experimental Setup.** To evaluate the effectiveness and scalability of our method across a broad spectrum of device prototypes, ranging from very small to very large sizes, we conduct experiments involving 3 to 7 different prototypes. Our objective is to assess how effectively our method adapts from a uniform array of small-size prototypes (3 device prototypes) to a diverse mix that



includes prototypes ranging from extremely small (XXS) to extremely large (XXL) (7 device prototypes). These experiments involve training image classification models from scratch on the CINIC-10 dataset, using CIFAR-100 as the unlabeled public dataset. We randomly distribute the dataset among prototypes with dataset ratios set to 1:2:3:4:5:6:7 from XXS to XXL. Each dataset portion is further distributed among clients using a Dirichlet distribution ( $\text{Dir}(0.5)$ ). The number of clients ranges from 35 to 5 from XXS to XXL, respectively. Client sample rates are set at 0.1, 0.1, 0.15, 0.15, 0.2, 0.3, and 0.6 from XXS to XXL. We use a series of ResNet architectures—ResNet10-XXS, ResNet10-XS, ResNet10-S, ResNet10-M, ResNet10, ResNet18, and ResNet50—scaled appropriately for each prototype. The local training epochs are set at 2, 2, 2, 5, 10, 10, and 20 from XXS to XXL to account for resource constraints, with fewer epochs assigned to smaller devices. We employ the Adam optimizer with a learning rate of  $1e-3$  and a weight decay of  $5e-5$  for local training. For XL and XXL, a step learning rate scheduler reduces the learning rate by a factor of 0.1 at half epoch. Server-side distillation employs a fixed batch size of 128, using the Adam optimizer with learning rate of  $1e-3$  and weight decay of  $5e-5$ . The softmax temperature is set at 3 for ensemble distillation and 20 for self-regularization. The number of communication rounds is fixed at 30. These experiments are conducted over 3 trials with different random seeds, and the average performance with standard deviation is reported. The entire device prototypes configurations are given in Table 15.

The detailed results are presented in Tables 7, 8, and 9.

Table 7: **Scalability Evaluation.** Detailed performance results for 7 device prototypes case.

Baseline	XXS	XS	S	M	L	XL	XXL	Average
FedAvg	23.17 $\pm$ 1.26	30.66 $\pm$ 0.14	32.81 $\pm$ 0.21	31.77 $\pm$ 0.21	37.69 $\pm$ 0.08	41.78 $\pm$ 0.05	50.52 $\pm$ 0.01	35.49
FedDF	27.98 $\pm$ 0.66	<b>37.47<math>\pm</math>0.33</b>	40.61 $\pm$ 0.01	40.26 $\pm$ 0.18	43.83 $\pm$ 0.22	45.58 $\pm$ 0.18	52.18 $\pm$ 0.12	41.13
FedET	26.75 $\pm$ 0.98	36.99 $\pm$ 0.31	40.51 $\pm$ 0.19	41.60 $\pm$ 0.16	46.12 $\pm$ 0.31	48.39 $\pm$ 0.11	52.71 $\pm$ 0.09	41.87
TAKFL	27.30 $\pm$ 0.08	36.93 $\pm$ 0.16	43.31 $\pm$ 0.42	40.88 $\pm$ 0.01	48.52 $\pm$ 0.15	50.95 $\pm$ 0.04	54.27 $\pm$ 0.43	43.17
TAKFL+Reg	<b>29.28<math>\pm</math>0.16</b>	37.10 $\pm$ 0.45	<b>43.96<math>\pm</math>1.65</b>	<b>41.83<math>\pm</math>0.73</b>	<b>48.77<math>\pm</math>0.37</b>	<b>51.43<math>\pm</math>0.46</b>	<b>54.63<math>\pm</math>0.84</b>	<b>43.86</b>

Table 8: **Scalability Evaluation.** Detailed performance results for 5 device prototypes case.

Baseline	XXS	XS	S	M	XL	Average
FedAvg	24.19 $\pm$ 1.03	21.04 $\pm$ 0.76	33.62 $\pm$ 0.88	38.91 $\pm$ 0.74	46.93 $\pm$ 0.05	32.94
FedDF	<b>28.31<math>\pm</math>0.61</b>	34.66 $\pm$ 0.00	39.91 $\pm$ 0.07	38.24 $\pm$ 0.36	46.81 $\pm$ 0.11	37.59
FedET	26.88 $\pm$ 0.95	34.11 $\pm$ 0.27	<b>41.15<math>\pm</math>0.29</b>	40.81 $\pm$ 0.87	48.14 $\pm$ 0.06	38.22
TAKFL	27.91 $\pm$ 0.12	37.09 $\pm$ 0.11	40.46 $\pm$ 0.34	41.06 $\pm$ 0.02	49.02 $\pm$ 0.35	39.11
TAKFL+Reg	28.24 $\pm$ 0.46	<b>37.30<math>\pm</math>1.10</b>	40.76 $\pm$ 0.94	<b>43.09<math>\pm</math>0.27</b>	<b>50.86<math>\pm</math>0.22</b>	<b>40.05</b>

Table 9: **Scalability Evaluation.** Detailed performance results for 3 device prototypes case.

Baseline	XXS	S	M	Average
FedAvg	24.19 $\pm$ 1.03	33.62 $\pm$ 0.88	<b>38.91<math>\pm</math>0.74</b>	32.24
FedDF	27.85 $\pm$ 0.10	<b>37.83<math>\pm</math>0.12</b>	37.74 $\pm$ 0.41	34.47
FedET	26.04 $\pm$ 0.67	36.87 $\pm$ 0.68	37.66 $\pm$ 0.09	33.52
TAKFL	26.62 $\pm$ 0.16	37.32 $\pm$ 0.40	38.13 $\pm$ 0.58	34.02
TAKFL+Reg	<b>27.90<math>\pm</math>0.98</b>	37.63 $\pm$ 0.87	38.20 $\pm$ 0.91	<b>34.58</b>

## E Ablation Studies

### E.1 Understanding Merging Coefficient

In this section, we conduct an ablation study to further understand how TAKFL customizes knowledge integration and understand how the merging coefficients  $\lambda_i$  are achieving this. This experiment aims to further understand the trade-offs between customized knowledge integration approach from the one-size-fits-all strategy employed in vanilla ensemble distillation and prior works.

**Experimental Setup.** Our experimentation focuses on two device prototypes: XXS and XXL, selected from the scalability evaluation detailed in Section 7.3, Appendix D.4, and Table 15. We employ the image classification task on the CINIC-10 [9] dataset, starting from scratch. Each prototype receives a randomly selected, non-overlapping subset of the training dataset—3.57% for XXS and 25% for XXL—distributed among their clients in a non-i.i.d. manner using  $\text{Dir}(0.5)$ . Both prototypes have three clients each. The architectures used are ResNet10-XXS for the XXS prototype and ResNet-50 for the XXL prototype. To focus solely on the evaluation of the server-side distillation process and its evolution with varying  $\lambda$ , we pre-train each prototype using standard FedAvg

for 10 communication rounds, with a sample rate of 1.0. Local training involves 20 epochs for XXS and 20 epochs for XXL using an Adam optimizer with a learning rate of 1e-3 and weight decay of 5e-5. The XXL prototype employs a step learning rate scheduler that reduces the rate by a factor of 0.1 at local epoch 10. For server-side distillation, we utilize a batch size of 128 and an Adam optimizer with a learning rate of 1e-5 and weight decay of 5e-5. CIFAR-100 [24] serves as the unlabeled public dataset. We save the final updated client and server models from both prototypes for further experimentation, focusing on the impact of merging coefficients without self-regularization in TAKFL. The merging coefficient  $\lambda$  is varied linearly from 0 to 1 in increments of 0.05. For simplicity, the XXS prototype is referred to as the small (S) prototype and XXL as the large (L) prototype.

**Discussion.** Figure 6 illustrates the significant impact of customized knowledge integration on the performance of both small and large device prototypes compared to the one-size-fits-all approach typical of vanilla ensemble distillation in the prior works, at different distillation epochs. Here, TAKFL adeptly manages customization for both small and large prototypes by controlling the merging coefficient  $\lambda$ . The merged model for both the small and large student prototypes is obtained using the formula  $\theta_{merged} = \theta_{avg} + ((1 - \lambda)\tau_S + \lambda\tau_L)$ . Notably, the performance is benchmarked at  $\lambda \approx 0.5$  in all cases, reflecting similar results to vanilla ensemble distillation (FedDF), where no customization in knowledge transfer occurs. This baseline performance is critical for understanding the effects of further customization.

In small distillation epochs ( $I_{distill} < 10$ ), minimal benefit is observed from customized knowledge integration, as both small and large prototypes achieve optimal performance at the non-customized  $\lambda \approx 0.5$ . However, as the distillation process progresses beyond 10 epochs, the influence of  $\lambda$  becomes increasingly pronounced. For  $\lambda > 0.5$ , the knowledge from the large prototype’s ensembles predominates, enhancing their impact, while for  $\lambda < 0.5$ , integration is more influenced by the small prototype’s ensembles. This pattern suggests that increased distillation epochs enable more effective distillation of each prototype’s unique knowledge for extreme cases of extremely small and large prototypes, thereby making the customization benefits evident. In scenarios with small distillation epochs, the absence of significant unique knowledge results in optimal performance at  $\lambda \approx 0.5$ . Conversely, as the number of distillation epochs rises ( $I_{distill} \geq 20$ ), the one-size-fits-all strategy proves suboptimal, underscoring the importance of tailored knowledge integration strategies. Optimal performance increasingly occurs at  $\lambda > 0.5$ , indicating effective leveraging of each prototype’s strengths to maximize overall performance. These findings confirm the necessity for customized knowledge integration in environments with significant prototype size variations and support our theoretical insights as detailed in Remark 1 and 2.

## E.2 Impact of Public Dataset

In this section, we explore the influence of the public dataset on the performance of TAKFL and existing KD-based methods when the public dataset used for server-side distillation is less similar to the private dataset, which is the actual learning objective. For this analysis we employ the same experimental setup previously outlined in Section 7.2 and Appendix D.1, using the CIFAR-10 homo-family architecture. To measure dataset similarity, we compute cosine similarity between the averaged features of datasets, extracted using an off-the-shelf pre-trained CLIP model [40] (CLIP ViT-B/32) available from the official GitHub repository.<sup>2</sup>

**Discussion.** Table 10 presents our results, highlighting a significant observation: the performance of existing methods drastically deteriorates as the similarity between the public dataset and private datasets decreases. In contrast, TAKFL exhibits robustness, suffering much less performance degradation under the same conditions. This demonstrates TAKFL’s practical utility in real-world scenarios where the server typically lacks knowledge of the private datasets to select a closely aligned public dataset for distillation. Notably, FedET underperforms significantly when using a less similar public dataset, performing worse than vanilla FedAvg in both low and high data heterogeneity scenarios. A similar pattern was observed with FedET in the NLP tasks discussed in Section 7.2 and Appendix D.3. This issue is likely due to FedET’s dependence on the overconfident and poorly calibrated confidence scores from neural networks [15, 53] for uncertainty estimates in its uncertainty-weighted distillation approach.

<sup>2</sup><https://github.com/OpenAI/CLIP>

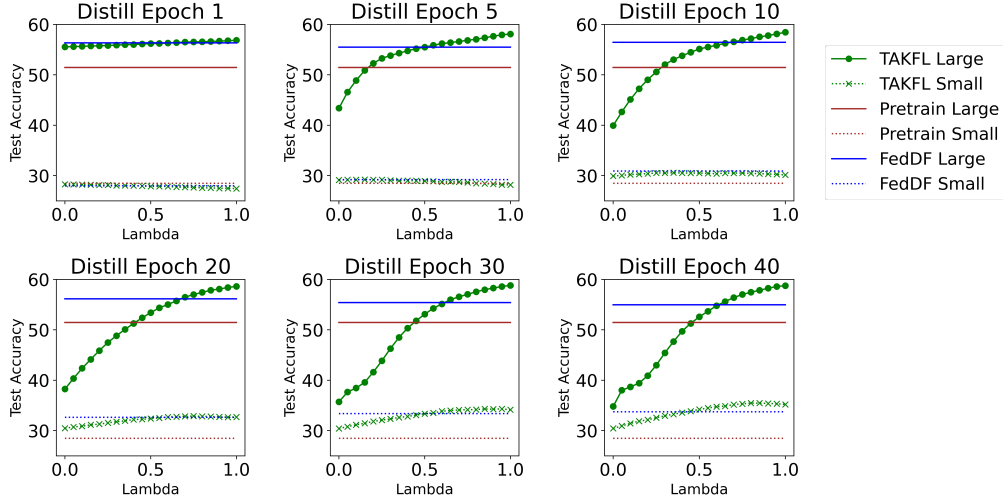


Figure 6: **Understanding the Impact of Merging Coefficients.** This figure showcases server-side knowledge distillation between two device prototypes, XXS and XXL, referred to as small and large, respectively, utilizing CIFAR-100 as the unlabeled public dataset. Both prototypes were pre-trained from scratch using standard FedAvg for 10 communication rounds. The CINIC-10 dataset was distributed between the small and large prototypes in ratios of 3.57% and 25%, respectively, and further subdivided non-i.i.d. among the clients using Dir(0.5). Each prototype has three clients with a sample rate of 1.0. The small prototype utilizes a ResNet10-XXS architecture, while the large prototype employs a ResNet-50.

Table 10: **Impact of Public Dataset on performance results.** Same experimental setting described in Section 7.2 and Appendix D.1 on CIFAR-10 homo-family setting is used for this experiment. The numbers in parentheses represent the similarity scores between private and public datasets, obtained using a pre-trained CLIP ViT-B/32 model.

Public Dataset	Baseline	Low Data Heterogeneity (Dir(0.3))				High Data Heterogeneity (Dir(0.1))			
		S	M	L	Average	S	M	L	Average
—	FedAvg	36.21 $\pm$ 2.24	46.41 $\pm$ 2.33	59.46 $\pm$ 6.17	47.36	22.01 $\pm$ 0.78	25.26 $\pm$ 3.89	51.51 $\pm$ 3.52	32.93
	FedDF	49.31 $\pm$ 0.15	50.63 $\pm$ 0.73	49.82 $\pm$ 0.98	49.92	34.71 $\pm$ 1.48	35.27 $\pm$ 4.74	51.08 $\pm$ 4.04	40.35
	FedET	49.21 $\pm$ 0.72	55.01 $\pm$ 1.81	53.60 $\pm$ 6.47	52.61	29.58 $\pm$ 3.00	30.96 $\pm$ 4.70	45.53 $\pm$ 6.46	35.36
	TAKFL	55.90 $\pm$ 1.70	57.93 $\pm$ 3.49	60.58 $\pm$ 2.35	58.14	37.40 $\pm$ 1.68	38.96 $\pm$ 0.17	51.49 $\pm$ 6.15	42.62
	TAKFL+Reg	<b>56.37<math>\pm</math>0.46</b>	<b>58.60<math>\pm</math>0.43</b>	<b>65.69<math>\pm</math>1.28</b>	<b>60.22</b>	<b>40.51<math>\pm</math>1.05</b>	<b>40.12<math>\pm</math>1.24</b>	<b>53.24<math>\pm</math>2.51</b>	<b>44.62</b>
TinyImagenet (0.92)	FedDF	49.37 $\pm$ 1.58	49.41 $\pm$ 4.21	55.06 $\pm$ 6.71	51.28	31.41 $\pm$ 6.61	30.73 $\pm$ 7.77	39.82 $\pm$ 5.16	33.99
	FedET	33.95 $\pm$ 0.92	37.26 $\pm$ 1.64	39.77 $\pm$ 3.44	36.99	24.12 $\pm$ 1.84	24.58 $\pm$ 2.13	28.91 $\pm$ 1.09	25.87
	TAKFL	55.20 $\pm$ 0.07	56.36 $\pm$ 0.40	60.71 $\pm$ 0.22	57.42	40.08 $\pm$ 0.19	40.26 $\pm$ 0.04	43.56 $\pm$ 1.10	41.30
	TAKFL+Reg	<b>56.28<math>\pm</math>0.09</b>	<b>57.14<math>\pm</math>0.03</b>	<b>60.90<math>\pm</math>0.22</b>	<b>58.11</b>	<b>40.88<math>\pm</math>0.11</b>	<b>41.10<math>\pm</math>1.15</b>	<b>46.25<math>\pm</math>5.95</b>	<b>42.74</b>
Celeb-A (0.77)	FedDF	48.99 $\pm$ 0.37	50.06 $\pm$ 0.43	55.12 $\pm$ 4.95	51.39	29.80 $\pm$ 0.39	32.28 $\pm$ 4.41	44.04 $\pm$ 4.60	35.36
	FedET	28.56 $\pm$ 3.00	28.80 $\pm$ 1.00	37.20 $\pm$ 2.78	31.52	15.28 $\pm$ 1.75	19.00 $\pm$ 3.43	23.29 $\pm$ 5.04	19.19
	TAKFL	45.65 $\pm$ 2.72	54.53 $\pm$ 1.72	58.13 $\pm$ 0.13	52.77	<b>31.02<math>\pm</math>0.68</b>	<b>36.76<math>\pm</math>1.58</b>	48.33 $\pm$ 0.53	38.70
	TAKFL+Reg	46.93 $\pm$ 0.67	<b>56.67<math>\pm</math>1.26</b>	<b>60.13<math>\pm</math>1.38</b>	<b>54.58</b>	30.88 $\pm$ 3.51	35.95 $\pm$ 5.40	<b>52.68<math>\pm</math>1.90</b>	<b>39.84</b>

## F Hyper-parameters and Implementation

In this section we bring the details of the hyper-parameters we used and our implementation. We implement our entire code in PyTorch [38] and release it at <https://github.com/MMorafah/TAKFL>. We use two NVIDIA RTX 3090 gpus to conduct the entire experimentation in this paper.

### F.1 Computer Vision (CV) Task

For comprehensive evaluation of our method, we consider federated learning image classification training from scratch.

#### F.1.1 Datasets

**Datasets.** We experiment with several image classification datasets: CIFAR-10, CIFAR-100, CINIC-10, TinyImageNet, STL-10. The details of each dataset is the following:

- **CIFAR-10** [24]: CIFAR-10 consists of 60,000 images of size  $32 \times 32$  RGB across 10 classes, with each class containing 6,000 images.
- **CIFAR-100** [24]: CIFAR-100 comprises 60,000 images of size  $32 \times 32$  RGB distributed across 100 classes, with 500 images per class.
- **CINIC-10** [9]: CINIC-10 has 270,000 images of size  $32 \times 32$  RGB across 10 classes, each class containing 27,000 images.
- **TinyImageNet** [25]: TinyImageNet contains 100,000 images of size  $64 \times 64$  RGB across 200 classes.
- **STL-10** [7]: STL-10 has 100,000 unlabeled images and 13,000 labeled images of size  $96 \times 96$  RGB across 10 classes.

The ImageNet-100 as the unlabeled public dataset is constructed by randomly selecting 100 classes from the ImageNet [10] dataset.

#### F.1.2 Architectures

**Experiments in Section 7.2 and Appendix D.1.** The architecture that we use are two distinct architectural settings: the “*homo-family*” setting, where all device prototypes’ architectures are from the same family, and the “*hetero-family*” setting, where architectures do not necessarily belong to the same family. For the homo-family scenario, we employ ResNet-8 for S, ResNet-14 for M, and ResNet-18 for L. For the hetero-family scenario, we use ViT-S for S, ResNet-14 for M, and VGG-16 for L. All models are initialized from random initialization.

For the ResNet architecture configuration, we utilize the standard ‘BasicBlock’ as the building block. This consists of a convolutional block, followed by four residual block stages, an adaptive average pooling layer, and a classifier layer. The models within this family differ in terms of the number of repetitions of the residual block and the number of filters in each stage. The configurations for different capacities are detailed below:

- *ResNet-18* is configured with [64, 128, 256, 512] filters and repeats the ‘BasicBlock’ [2, 2, 2, 2] times.
- *ResNet-14* is configured with [64, 128, 256, 512] filters and repeats the ‘BasicBlock’ [1, 2, 2, 1] times.
- *ResNet-8* is configured with [64, 128, 256] filters, corresponding to the first three stages, and repeats the ‘BasicBlock’ [1, 1, 1] times.

For VGG-16 [42], we use the standard architecture which includes convolutional layers followed by max-pooling layers. The configuration of filters for each layer is as follows:

*VGG-16*: [64, 64, ‘M’, 128, 128, ‘M’, 256, 256, 256, ‘M’, 512, 512, 512, ‘M’, 512, 512, 512, ‘M’]

The final classification head consists of two linear layers with a hidden size of 512, followed by a ReLU activation, and a final linear classifier layer.

For ViT-S, we adopt the standard Vision Transformer [12] architecture implementation from Github.<sup>3</sup> We configure ViT-S with 6 attention blocks, each with 16 heads and a hidden dimension of 64. The final MLP dimension is set to 256. In our experiments with ViT-S, we set the patch size to 4, and the input image size is  $32 \times 32$ .

**Experiment in Appendix D.2.** The pre-trained MobileNetV3-Large [19], MobileViTV2 [34], and ResNet34 [49] were instantiated using the TIMM library.<sup>4</sup>

**Scalability Experiments in Section 7.3.** The architectures in these experiments are inspired by [20]. The details of architectures are as following:

- *ResNet10-XXS* is configured with [8, 8, 16, 16] filters and repeats the ‘BasicBlock’ [1, 1, 1, 1] times.
- *ResNet10-XS* is configured with [8, 16, 32, 64] filters and repeats the ‘BasicBlock’ [1, 1, 1, 1] times.
- *ResNet10-S* is configured with [16, 32, 64, 128] filters and repeats the ‘BasicBlock’ [1, 1, 1, 1] times.
- *ResNet10-M* is configured with [8, 16, 32, 64] filters and repeats the ‘BasicBlock’ [1, 1, 1, 1] times.
- *ResNet10* is configured with [64, 128, 256, 512] filters and repeats the ‘BasicBlock’ [1, 1, 1, 1] times.
- *ResNet18* is configured with [64, 128, 256, 512] filters and repeats the ‘BasicBlock’ [1, 1, 1, 1] times.
- *ResNet50* is configured with [64, 128, 256, 512] filters and repeats the ‘BasicBlock’ [3,4,6,3] times.

### F.1.3 FL configuration and Hyper-parameters

**Base Hyper-parameters.** The following hyperparameter values apply to all CV experiments unless stated otherwise. We set the diversity regularizer coefficient of FedET to 0.1 for our entire experimentation per the original paper [6]. We use the Adam optimizer with a learning rate of  $1e-5$ , weight decay value of  $5e-5$ , and a batch size of 128 for distillation. The softmax distillation temperature is set to 3, the distillation epoch to 1, and the self-regularizer softmax temperature to 20 for both CV and NLP experiments. Table 12 details the hyper-parameters.

**Experiments in Section 7.2 and Appendix D.1.** For tables 1 and 4, there are 100 clients with the S device prototype, 20 clients with the M device prototype, and 4 clients with the L device prototype. For each round, 10, 4, and 2 clients from each S, M, and L prototype are randomly sampled respectively for participation. 10% of the data goes to the S prototype, 30% to M, and 60% to L. The data is distributed to each client among each prototype in a Dirichlet distribution. Table 13 details the FL configuration.

**Experiments in Appendix D.2 and Appendix E.2.** For tables 5 and 10, there are 4, 3, and 2 clients for S, M, and L device prototypes, respectively. Each round, every client participates in FL. 20% of the data is distributed to prototype S, 30% to prototype M, and 50% to prototype L. Table 14 details the FL configuration.

**Scalability Experiments in Section 7.3 and D.4.** For tables 7, 8, and 9, there are 35, 30, 25, 20, 15, 10, and 5 clients for the prototypes XXS, XS, S, M, L, XL, and XXL, respectively. The sample rate is set to 0.1, 0.1, 0.15, 0.15, 0.2, 0.3, and 0.6 from XXS to XXL. The data is distributed for each prototype in the ratio 1:2:3:4:5:6:7 from XXS to XXL. Table 15 details the hyper-parameters and configuration.

**Validation Set.** For TAKFL, the validation set used for the heuristic method (see F.3) is 5% of the training dataset. The validation set and the private dataset does not overlap.

## F.2 Natural Language Processing (NLP) Task

For the NLP task, we fine-tune federated learning text classification task using pretrained models.

<sup>3</sup><https://github.com/lucidrains/vit-pytorch>

<sup>4</sup><https://github.com/huggingface/pytorch-image-models>

### F.2.1 Datasets.

All NLP datasets were provided by Hugging Face.<sup>5</sup>

- **MNLI** [50]: MNLI contains 433K sentence pairs, each sentence pair labeled as one of ‘entailment,’ ‘neutral,’ and ‘contradiction.’
- **SNLI** [2]: SNLI is similar to MNLI, with 570K sentence pairs each labeled one of 3 labels.
- **SST2** [43]: SST2 consists of 67K phrases, each labeled as sentiment ‘positive’ or ‘negative.’
- **Sentiment140** [14]: Sentiment140 is a dataset of 1.6M Twitter messages each labeled with one of 2 sentiment values.
- **MARC** [22]: MARC (Multilingual Amazon Reviews Corpus) is a dataset with online reviews in multiple languages from the Amazon delivery service website. Each review has a label which is one of 1-5 stars. We only use the English reviews from this dataset, which results in 260,000 English reviews total.
- **Yelp** [59]: The Yelp reviews dataset contains 700K reviews each labeled 1-5 stars from the Yelp service which publishes public reviews of businesses.
- **AG News** [58]: AG News contains 127,600 news article titles. Each article is one of four classifications of news articles.
- **DBpedia** [57]: The DBpedia dataset consists of 630K DBpedia article summaries each labeled one of 14 categorizations.

### F.2.2 Architectures

**Experiments in Section 7.2 and Appendix D.3.** We use three variations of the BERT architecture: BERT-Tiny, BERT-Mini, and BERT-Small from [45]. The weights were pre-trained on the BookCorpus dataset and extracted text from Wikipedia. Further details regarding each model are described extensively on Github.<sup>6</sup> The tokenizer used for these transformer models are the same ones provided by the authors of [45].

- *BERT-Tiny* contains 2 transformer layers and an embedding size of 128.
- *BERT-Mini* contains 4 transformer layers and an embedding size of 256.
- *BERT-Small* contains 4 transformer layers and an embedding size of 512.

### F.2.3 FL configuration and hyper-parameters

**Base Hyper-parameters.** For distillation, we use the Adam optimizer with a learning rate of  $3e-5$ , no weight decay, and batch size of 32. The distillation epoch is set to 1, the ensemble distillation softmax temperature to 3, and the self-regularizer softmax temperature to 20 for all NLP experiments. Table 16 details the hyper-parameters.

**Experiments in Section 7.2 and Appendix D.3.** For tables 2 and 6, we limit the private dataset to 100,000 samples, randomly sampled from the original dataset i.i.d. The public dataset is limited to 30,000 examples sampled i.i.d as well. There are 8, 4, and 2 clients for the S, M, and L prototypes. The private data is split across each prototype in the following proportions: 0.1, 0.3, 0.6. Table 17 details the FL configuration.

**Validation Set.** The validation dataset used for TAKFL is 5,000 samples taken from the original training dataset that does not overlap with the 100,000 private dataset.

## F.3 Hyper-parameters of TAKFL

**Merging Coefficients.** We conducted extensive experiments with different merging coefficients on the main 3-device prototype setting of small (S), medium (M), and large (L) discussed in Section 7.2 and Appendix D.1. We empirically observed that the small (S) prototype typically achieves the best performance using a uniformly increasing merging coefficient, where the larger the prototype, the

<sup>5</sup><https://github.com/huggingface/datasets>

<sup>6</sup><https://github.com/google-research/bert>



larger the merging coefficient, i.e.,  $\lambda_S \leq \lambda_M \leq \lambda_L$ . As we move towards larger prototypes, they benefit more from increasingly skewed merging coefficients towards the larger ones. In the extreme case of the large (L) prototype, highly skewed merging coefficients generally led to better performance, i.e.,  $\lambda_S \ll \lambda_M \ll \lambda_L$ . This pattern is intuitive as small prototypes can benefit from everyone while gaining more from the larger, more informative prototypes. However, larger prototypes benefit less from smaller ones, as they typically offer less information, especially in high data heterogeneity cases. Notably, in high data heterogeneity cases, more skewed merging coefficients seemed to be more advantageous as the smaller prototypes (S and M) possess lower quality knowledge.

Based on these observations, we designed a simple and cost-effective heuristic method that randomly instantiates merging coefficients following this intuition. Our heuristic method, presented in F.3, leverages these observations by generating candidate merging coefficients that incorporate both uniformly increasing and different degrees of skewed merging coefficients. This dual approach enables us to explore a wide range of merging strategies and identify the most effective configurations for different prototypes. The optimal merging coefficient candidate is determined using the performance on the held-out validation set.

```
import numpy as np
def heuristic(num_devices=3, n_candidates=10):
    candidates = [[1/num_devices for _ in range(num_devices)]]
    for exponent in [1, 5, 10]:
        for i in range(n_candidates):
            candidate = np.random.beta(a=1, b=100, size=num_devices)
            candidate = candidate ** exponent
            candidate = np.sort(candidate)
            candidate = candidate / np.sum(candidate)
            candidates.append(candidate)
    return candidates
```

Listing 1: Implementation of the heuristic method for merging coefficients in Python. The exponent term controls the degree of skewness or peaking in the merging coefficients.

Furthermore, we experiment with manually determining the merging coefficients and fixating them throughout the federation. We achieved similar results with this approach compared to adaptively finding the coefficients using the heuristic method and a small held-out validation set. We present the merging coefficient candidates that performed reasonably well during our experiments in Table 11.

Table 11: Details of the experimentally determined merging coefficients for the 3-device prototype setting discussed in Section 7.2 and Appendix D.1. Coefficients are ordered as  $[\lambda_S, \lambda_M, \lambda_L]$ .

Merging Coefficient Candidate	Small Prototype	Medium Prototype	Large Prototype
1	[0.2, 0.3, 0.5]	[0.1, 0.2, 0.7]	[0.1, 0.2, 0.7]
2	[0.3, 0.3, 0.4]	[0.05, 0.15, 0.8]	[0.01, 0.09, 0.99]
3	[0.2, 0.3, 0.5]	[0.1, 0.2, 0.7]	[0.05, 0.2, 0.75]
4	[0.05, 0.1, 0.85]	[0.01, 0.19, 0.8]	[0.01, 0.09, 0.90]
5	[0.1, 0.15, 0.75]	[0.05, 0.15, 0.8]	[0.01, 0.09, 0.90]
6	[0.05, 0.1, 0.85]	[0.05, 0.05, 0.9]	[0.001, 0.009, 0.99]
7	[0.05, 0.15, 0.80]	[0.05, 0.2, 0.75]	[0.001, 0.009, 0.99]
8	[0.05, 0.15, 0.80]	[0.05, 0.1, 0.85]	[0.001, 0.009, 0.99]
9	[0.3, 0.35, 0.35]	[0.2, 0.3, 0.5]	[0.1, 0.2, 0.7]

**Self-Regularization Coefficient.** Extensive experiments were conducted on the self-regulation coefficients for different device prototypes and settings. Although no consistent pattern emerged, we experimentally determined that optimal performance for the small prototype was achieved with self-regulation coefficients  $\gamma_S \in 0.1, 0.01, 0.001$ . For the medium prototype, the coefficients were  $\gamma_M \in 0.5, 0.1, 0.01, 0.001, 0.0001$ , and for the large prototype,  $\gamma_L \in 1.0, 0.8, 0.5, 0.1, 0.01, 0.001, 0.0001$  yielded the best results.

Table 12: Details of hyper-parameters for CV task in Section 7.2, Appendix D.1, and Appendix E.2.

Local/Server	Hyperparameter	Small Prototype	Medium Prototype	Large Prototype
Local Training	Training epochs	20	80	100
	Batch Size	64	64	64
	Optimizer	Adam	Adam	Adam
	Learning Rate	1e-3	1e-3	1e-3
	Weight Decay	5e-5	5e-5	5e-5
	LR scheduler	None	None	StepLR(step_size = 10, gamma = 0.1)
Server KD Training	Optimizer	Adam	Adam	Adam
	Learning Rate	1e-5	1e-5	1e-5
	Weight Decay	5e-5	5e-5	5e-5
	Batch Size	128	128	128
	Training Epochs	1	1	1
	Ensemble Distillation Softmax Temperature	3	3	3
	Self-Regularizer Softmax Temperature	20	20	20

Table 13: Details of Architecture parameters and FL configuration for CV task in Section 7.2 and Appendix D.1.

Architecture Setting	Device Prototype	Architecture	CIFAR-10 Parameters	CIFAR-100 Parameters	Dataset Portion	Clients	Sample Rate
Homo-Family	Prototype S	ResNet8	1.23M	1.25M	0.1	100	0.1
	Prototype M	ResNet14	6.38M	6.43M	0.3	20	0.2
	Prototype L	ResNet18	11.17M	11.22M	0.6	4	0.5
Hetero-Family	Prototype S	ViT-S	1.78M	1.79M	0.1	100	0.1
	Prototype M	ResNet14	6.38M	6.43M	0.3	20	0.2
	Prototype L	VGG16	15.25M	15.30M	0.6	4	0.5

Table 14: Details of Architecture parameters and FL configuration for CV task experiment using pre-trained models in Appendix D.2.

Device Prototype	Architecture	STL-10/CINIC-10 Parameters	TinyImageNet Parameters	Dataset Portion	Clients	Sample Rate
Prototype S	mobilenetv3-large-100	4.21M	4.45M	0.2	4	1.0
Prototype M	mobilevitv2-175	13.36M	13.53M	0.3	3	1.0
Prototype L	ResNet34	21.28M	21.38M	0.5	2	1.0

Table 15: Details of Architecture parameters for Scalability Section 7.3, and Appendix D.4.

Device Prototype	CINIC-10					
	Architecture	Parameters	Dataset Portion	Clients	Sample Rate	Local Epochs
Prototype XXS	ResNet10-XXS	11K	0.0357	35	0.1	2
Prototype XS	ResNet10-XS	78K	0.0714	30	0.1	2
Prototype S	ResNet10-S	309K	0.1071	25	0.15	2
Prototype M	ResNet10-M	1.2M	0.1428	20	0.15	5
Prototype L	ResNet10	4.9M	0.1785	15	0.2	10
Prototype XL	ResNet18	11M	0.2142	10	0.3	10
Prototype XXL	ResNet50	24M	0.25	5	0.6	20

Table 16: Details of hyper-parameters for NLP task experiments in Section 7.2 and Appendix D.3.

Local/Server	Hyperparameter	Prototype S	Prototype M	Prototype L
Local Training	training epochs	1	1	1
	batch size	32	32	32
	optimizer	Adam	Adam	Adam
	Learning rate	3e-5	3e-5	3e-5
	Weight decay	0	0	0
	lr scheduler	None	None	None
Server KD Training	optimizer	Adam	Adam	Adam
	learning rate	3e-5	3e-5	3e-5
	weight decay	3e-5	3e-5	3e-5
	batch size	32	32	32
	training epochs	1	1	1
	Ensemble distillation softmax temperature	3	3	3
	Self-regularizer softmax temperature	20	20	20

Table 17: Details of Architecture parameters and FL configuration for NLP task experiment using pre-trained models in Section 7.2 and Appendix D.3.

Device Prototype	Architecture	Parameters	Clients	Dataset Portion	Sample Rate
Prototype S	BERT-Tiny	4.39M	8	0.1	0.4
Prototype M	BERT-Mini	11.17M	4	0.3	0.5
Prototype L	BERT-Small	28.77M	2	0.6	1.0

## G Acknowledgement

We would like to express our sincere gratitude to Ang Li, Matias Mendieta, and Guangyu Sun for their invaluable feedback and insightful discussions, which significantly contributed to the development and refinement of this work. Their thoughtful suggestions and careful review of earlier drafts were instrumental in enhancing the quality of this paper.

This research was partially supported by a grant from Cisco Systems, Inc. We also gratefully acknowledge the use of the computational infrastructure provided by the OP VVV funded project CZ.02.1.01/0.0/0.0/16\_019/0000765, “Research Center for Informatics,” which enabled us to conduct the experiments presented in this work.