

# Decentralized Learning with Multi-Headed Distillation

Andrey Zhmoginov   Mark Sandler   Nolan Miller   Gus Kristiansen   Max Vladymyrov  
Google Research

{azhmogin,sandler,namiller,gusatb,mxv}@google.com

## Abstract

*Decentralized learning with private data is a central problem in machine learning. We propose a novel distillation-based decentralized learning technique that allows multiple agents with private non-iid data to learn from each other, without having to share their data, weights or weight updates. Our approach is communication efficient, utilizes an unlabeled public dataset and uses multiple auxiliary heads for each client, greatly improving training efficiency in the case of heterogeneous data. This approach allows individual models to preserve and enhance performance on their private tasks while also dramatically improving their performance on the global aggregated data distribution. We study the effects of data and model architecture heterogeneity and the impact of the underlying communication graph topology on learning efficiency and show that our agents can significantly improve their performance compared to learning in isolation.*

## 1. Introduction

Supervised training of large models historically relied on access to massive amounts of labeled data. Unfortunately, since data collection and labeling are very time-consuming, curating new high-quality datasets remains expensive and practitioners are frequently forced to get by with a limited set of available labeled datasets. Recently it has been proposed to circumvent this issue by utilizing the existence of large amounts of siloed private information. Algorithms capable of training models on the entire available data without having a direct access to private information have been developed with Federated Learning approaches [24] taking the leading role.

While very effective in large-scale distributed environments, more canonical techniques based on federated averaging, have several noticeable drawbacks. First, gradient aggregation requires individual models to have fully compatible weight spaces and thus iden-

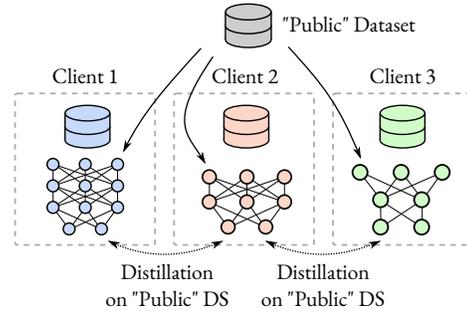


Figure 1. Conceptual diagram of a distillation in a distributed system. Clients use a public dataset to distill knowledge from other clients, each having their primary private dataset. Individual clients may have different architectures and different objective functions.

tical architectures. While this condition may not be difficult to satisfy for sufficiently small models trained across devices with compatible hardware limitations, this restriction may be disadvantageous in a more general setting, where some participant hardware can be significantly more powerful than the others. Secondly, federated averaging methods are generally trained in a centralized fashion. Among other things, this prohibits the use of complex distributed communication patterns and implies that different groups of clients cannot generally be trained in isolation from each other for prolonged periods of time.

Another branch of learning methods suitable for distributed model training on private data are those based on distillation [3,6,15]. Instead of synchronizing the inner states of the models, such methods use outputs or intermediate representations of the models to exchange the information. The source of data for computing exchanged model predictions is generally assumed to be provided in the form of publicly available datasets [12] that do not have to be annotated since the source of annotation can come from other models in the ensemble (see Figure 1). One interesting interpretation of model distillation is to view it as a way of using queries

from the public dataset to indirectly gather information about the weights of the network (see Appendix A). Unlike canonical federated-based techniques, where the entire model state update is communicated, distillation only reveals *activations* on specific samples, thus potentially reducing the amount of communicated bits of information. By the data processing inequality, such reduction, also translates into additional insulation of the private data used to train the model from adversaries. However, it is worth noting that there exists multiple secure aggregation protocols including SecAgg [5] that provide data privacy guarantees for different Federated Learning techniques.

The family of approaches based on distillation is less restrictive than canonical federated-based approaches with respect to the communication pattern, supporting fully distributed knowledge exchange. It also permits different models to have entirely different architectures as long as their outputs or representations are compatible with each other. It even allows different models to use various data modalities and be optimizing different objectives, for example mixing supervised and self-supervised tasks within the same domain. Finally, notice that the distillation approaches can and frequently are used in conjunction with weight aggregation [21, 30, 31, 37], where some of the participating clients may in fact be entire ensemble of models with identical architectures continuously synchronized using federated aggregation (see Figure 8 in Supplementary).

**Our contributions.** In this paper, we propose and empirically study a novel distillation-based technique that we call Multi-Headed Distillation (MHD) for distributed learning on a large-scale ImageNet [9] dataset. Our approach is based on two ideas: (a) inspired by self-distillation [2, 10, 38] we utilize multiple model heads distilling to each other (see Figure 2) and (b) during training we simultaneously distill client model predictions and intermediate network embeddings to those of a target model. These techniques allow individual clients to effectively absorb more knowledge from other participants, achieving a much higher accuracy on a set of all available client tasks compared with the naive distillation method.

In our experiments, we explore several key properties of the proposed model including those that are specific to decentralized distillation-based techniques. First, we analyse the effects of data heterogeneity, studying two scenarios in which individual client tasks are either identical or very dissimilar. We then investigate the effects of working with nontrivial communication graphs and using heterogeneous model architectures. Studying complex communication patterns,

we discover that even if two clients in the ensemble cannot communicate directly, they can still learn from each other via a chain of interconnected clients. This “transitive” property relies in large part on utilization of multiple auxiliary heads in our method. We also conduct experiments with multi-client systems consisting of both ResNet-18 and ResNet-34 models [14] and demonstrate that: (a) smaller models benefit from having large models in the ensemble, (b) large models learning from a collection of small models can reach higher accuracies than those achievable with small models only.

## 2. Related Work

**Personalized Federated Learning.** While many early canonical Federated Learning approaches trained a single global model for all clients [24], it has been quickly realized that non-IID nature of private data in real systems may pose a problem and requires personalized approaches [20]. Since then many Personalized Federated Learning approaches have been developed, many covered in the surveys [18, 33].

**Federated Distillation.** Emergence of Federated Distillation was motivated by the need to perform learning across ensembles of heterogeneous models<sup>1</sup>, reducing communication costs and improving performance on non-IID data. Existing distillation-based approaches can be categorized based on the system setup and the types of the messages passed between participants. A number of approaches including [8, 12, 21, 23, 30, 31, 37, 40] combine aggregation of weight updates with model distillation. They are typically centralized and frequently involve client-side distillation, which may restrict the size of the aggregated model. A different body of work is concentrated on centralized systems, where only model predictions are communicated between the clients and the server [11, 13, 16, 19, 26, 29, 32, 39]. Another related family of approaches is based on communicating embedding prototypes [34], or using embeddings for distillation directly [1, 26]. In this paper, we concentrate on a more general decentralized setup, where there is not single central authority and all clients exchange knowledge via distillation [4].

## 3. Model

### 3.1. Setup

We consider a system of  $K$  clients  $C = \{C_1, \dots, C_K\}$ . Each client  $C_i$  is assumed to possess

<sup>1</sup>note that multiple existing approaches like [28, 35] allow using FedAvg for training heterogeneous model ensembles

their own *private dataset*  $\mathcal{D}_i$  while training a *private model*  $\mathcal{M}_i$  that solves a corresponding task  $\mathcal{T}_i$ . In the following, we assume that all tasks  $\mathcal{T}_i$  are supervised.

While using their local dataset  $\mathcal{D}_i$  to train the private model, each client can also communicate with other clients to learn from them. At each global training step  $t$ , we define a local directed graph  $\mathcal{G}_t$  that determines the pattern of this communication. While the set of nodes of  $\mathcal{G}_t$  is fixed to be the set of all clients, the set of edges  $\mathcal{E}_t$  with the corresponding incidence function can be dynamic and change every training step.

The local datasets  $\mathcal{D}_i$  are not directly exchanged between the clients, instead the information exchange occurs via a shared *public* source of unlabeled data  $\mathcal{D}_*$ . We assume that at training step  $t$ , each client  $C_i$  can perform inference on a set of public samples and request the results of a similar computation on the same samples from other clients that are incident to it by directed edges of  $\mathcal{G}_t$ . In other words, each client  $C_i$  is optimizing a local objective  $\mathcal{L}_i$  defined as:

$$\mathcal{L}_i(t) = \mathcal{L}_{i,\text{CE}} + \sum_{\alpha} \mathbb{E}_{x \sim \mathcal{D}_*} \mathcal{L}_{\text{dist}}^{\alpha}(\psi_i^{\alpha}(x), \Phi_{t,i}^{\alpha}(x)), \quad (1)$$

where  $\mathcal{L}_{i,\text{CE}} \equiv \mathbb{E}_{(x,y) \sim \mathcal{D}_i} \mathcal{L}_{\text{CE}}(x,y)$  and  $\mathcal{L}_{\text{CE}}$  is a cross-entropy loss optimized locally by each client on their private data  $\mathcal{D}_i$ ,  $\mathcal{L}_{\text{dist}}^{\alpha}$  is a collection of different *distillation losses* enumerated by  $\alpha$  that use some *local* computation result  $\psi_i^{\alpha}$  and a *remote* results  $\Phi_{t,i}^{\alpha}(x) \equiv \{\phi_j^{\alpha}(x) | j \in e_t(i)\}$  computed on the same sample and  $e_t(i)$  is a set of clients connected to  $i$  via a set of outgoing edges (from  $\mathcal{G}_t$ ).

Notice that in contrast with Federated Learning, here we do not require different models  $\mathcal{M}_i$  to have compatible architectures, but instead optimize local and remote sample representations  $\psi_i(x)$  and  $\phi_j(x)$  to be compatible. In the next section, we discuss several potential choices of the distillation losses.

In this paper, we are interested in evaluating the impact that the communication and cross-learning between the clients has on (a) how well these models can be suited for their original private tasks and (b) how much of the knowledge gets shared and distributed to the other tasks over time. Notice that if each client has a sufficiently simple model and enough training data (making the model underfit), the communication between individual models is not expected to improve their private task performance, but can only enhance their learned representations making them more suitable for adapting to other client’s tasks. However, if the private training data is scarce (making the model overfit), the model communication could improve generalization and ultimately improve client performance on their private tasks.

## 3.2. Distillation Losses

**Embedding distillation.** We utilize the embedding regularization loss [1, 26] in our experiments. If  $\xi_i(x)$  is an intermediate embedding produced for a sample  $x$  coming from the shared public dataset by the model  $\mathcal{M}_i$ , then we can choose  $\psi_i^{\text{emb}}(x) \equiv \xi_i(x)$ ,  $\phi_j^{\text{emb}}(x) \equiv \xi_j(x)$  and define  $\mathcal{L}_{\text{dist}}^{\text{emb}}(\psi_i^{\text{emb}}(x), \Phi_{t,i}^{\text{emb}}(x))$  as

$$\nu_{\text{emb}} \sum_{j \in e_t(i)} \rho(\|\psi_i^{\text{emb}}(x) - \phi_j^{\text{emb}}(x)\|), \quad (2)$$

or simply  $\nu_{\text{emb}} \sum_{j \in e_t(i)} \rho(\|\xi_i(x) - \xi_j(x)\|)$ , where  $\nu_{\text{emb}}$  is the weighting constant and  $\rho(x) \in C^{\infty}$  is some monotonically growing function. The choice of this distillation loss forces compatibility between sample embeddings across the ensemble. In practice, we noticed that the embedding norms of different models frequently diverge during training, and to adapt to that we use normalized embeddings preserving regularization consistency across the entire duration of training:  $\psi_i^{\text{norm}}(x) \equiv \xi_i(x) / \|\xi_i(x)\|$ .

**Prediction distillation.** Ability to predict on classes that are rarely present in private data can be improved by utilizing prediction vector as an additional distillation target. However, since  $\mathcal{M}_i$  is tasked with fitting ground truth on a particular dataset  $\mathcal{D}_i$ , distilling this prediction to labels relevant for another client may be damaging for the model performance on  $\mathcal{T}_i$ . Instead, we choose to add another single prediction head to  $\mathcal{M}_i$  that is distilled from all existing tasks thus (a) not polluting the main prediction head of the model  $\mathcal{M}_i$ , but (b) at the same time forcing the intermediate representation  $\xi_i(x)$  to contain information relevant for solving all existing tasks  $\{\mathcal{T}_j | j \in 1, \dots, K\}$ .

Let  $\mathbf{h}_i(\xi_i(x))$  be the main head of the model  $\mathcal{M}_i$  used for computing  $\mathcal{L}_{\text{CE}}$  and  $\mathbf{h}_i^{\text{aux}}(\xi_i(x))$  be the auxiliary head. Then, the naïve prediction distillation loss takes the following form:

$$\mathcal{L}_{\text{dist}}^{\text{aux}}[\mathbf{h}^{\text{aux}}, \mathbf{h}] \equiv -\nu_{\text{aux}} \sum_{j \in e_t(i)} \mathbf{h}_j \log \mathbf{h}_i^{\text{aux}}(x), \quad (3)$$

where  $\nu_{\text{aux}}$  is the auxiliary loss weight. Here all the distillation targets from  $e_t(i)$  are essentially treated the same irrespective of their confidence in their prediction. One way of integrating the knowledge of the distillation target quality is to use some *confidence* metric for their prediction on  $x$ . For example, we could consider the following modification of the loss (3):

$$-\nu_{\text{aux}} \sum_{j \in e_t(i) \cup \{i\}} Q[\Lambda(\mathbf{h}_j); H[\mathbf{h}]] \times \mathbf{h}_j \log \mathbf{h}_i^{\text{aux}}(x), \quad (4)$$

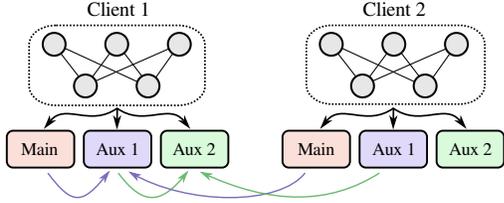


Figure 2. A pattern used for distilling multiple auxiliary heads. Here multiple auxiliary heads of “Client 1” are distilled from other auxiliary heads of the same model and from auxiliary heads of other clients (here “Client 2”). Auxiliary head *Aux 1* is distilled from the main heads, auxiliary head *Aux 2* is distilled from auxiliary heads *Aux 1* and so on.

where  $\Lambda(\mathbf{h}(x))$  is the confidence of the classifier prediction,  $Q$  is some function of the client confidence and  $H[\mathbf{h}] \equiv \{\Lambda(\mathbf{h}_j) | j \in e_t(i) \cup \{i\}\}$  is the information about confidence of all possible distillation targets including the  $i^{\text{th}}$  client itself. We considered perhaps the simplest choice for  $\Lambda$  defining it as  $\arg \max_k h_k(x)$ . This measure of the model confidence that we end up using in our method is, of course, not reliable (see Appendix A) and using a separate per-client density model  $\rho_i(x)$  for detecting in-distribution and out-of-distribution samples could potentially improve model performance (for an alternative approach see [22]). For  $Q$ , we only considered perhaps the most obvious choice of  $Q[\Lambda(\mathbf{h}_j)] = 1$  if  $j^{\text{th}}$  client has the largest confidence from  $H$  and 0 otherwise, effectively selecting the most confident client and using it as the distillation target (see Appendix A for a detailed discussion).

### Self-distillation with multiple auxiliary heads.

Self-distillation is a well-known technique that improves model performance by repeatedly using the previous iteration of the model as the distillation target for itself [2, 10, 25, 38]. The most direct application of this technique to training an ensemble of models is to perform multiple cycles of self-distillation across all available networks. Here, however, we propose a different approach, where we modify a conventional training procedure by equipping each classifier with a collection of multiple auxiliary heads  $\{\mathbf{h}^{\text{aux},1}, \dots, \mathbf{h}^{\text{aux},m}\}$ . These auxiliary heads distill from each other by optimizing the following loss:

$$\mathcal{L}_{\text{dist}}^{\text{aux}}[\mathbf{h}^{\text{aux},1}, \mathbf{h}] + \sum_{k=2}^m \mathcal{L}_{\text{dist}}^{\text{aux}}[\mathbf{h}^{\text{aux},k}, \mathbf{h}^{\text{aux},k-1}], \quad (5)$$

where  $\mathcal{L}_{\text{dist}}^{\text{aux}}[\mathbf{h}^{(a)}, \mathbf{h}^{(b)}]$  is defined according to Eq. (4). In other words,  $\mathbf{h}^{\text{aux},1}$  distills from  $\mathbf{h}$  and  $\mathbf{h}^{\text{aux},k}$  distills from  $\mathbf{h}^{\text{aux},k-1}$  for all  $1 < k \leq m$ . This approach

illustrated in Figure 2 is one of the core contributions of our paper.

### Communication efficiency.

In terms of communication efficiency, this approach could suffer from ineffective communication when the distillation targets are frequently a poor source of knowledge for a particular sample class. This problem would ideally require client awareness of the label distribution on each client that it communicates with. However, since in practice, prediction distillation (embedding distillation is more costly) only requires a transmission of several highest-confidence predictions for each sample, each step with batch size of 512 would require a communication of only a few thousand floating point numbers (assuming that shared public set images could be uniquely identified with a small hash). At the same time, a single back-and-forth round of FedAvg communication of a ResNet-34 model would require more than 100 million floating-point parameters, which would be equivalent to around 50k prediction distillation steps.

### 3.3. Dataset

In this work, we study distributed learning in systems with varying degrees of data heterogeneity: from those where the distribution of data is the same across all clients, to more extreme cases where each client specializes on its own unique task. We simulate these scenarios using an underlying labeled dataset  $\mathcal{D}$ . Let  $S$  be the set of all samples from  $\mathcal{D}$ . Some fraction of samples  $\gamma_{\text{pub}}$  (typically around 10%) is treated as a set of unlabeled *public* samples. The remaining samples are treated as the source of private data and are distributed without repetition across all of  $K$  clients as discussed below.

### Label assignment.

Each client  $C_i$  is assigned a subset  $\ell_i$  of all labels, which are treated as *primary labels* for  $C_i$ . Remaining labels from  $\mathcal{D}$  not belonging to  $\ell_i$  are treated as *secondary labels* for  $C_i$ . For each label  $l$ , we take all available samples and randomly distribute them across all clients. The probability of assigning a sample with label  $l$  to a client  $C_i$  is chosen to be  $1 + s$  times higher for clients that have  $l$  as their primary label. We call the parameter  $s$  *dataset skewness*. As a result, in the iid case with  $s = 0$  all samples are equally likely to be assigned to any one of the clients. However, in the non-iid case in the limit of  $s \rightarrow \infty$ , all samples for label  $l$  are only distributed across clients for which  $l$  is primary.

We considered two choices for selecting the primary label sets for the clients. One choice (we refer to as *even*) is to subdivide the set of all labels in such a

way that each label has exactly  $m$  corresponding primary clients. Another choice (we refer to as *random*) is to randomly assign each client  $C_i$  a random fixed-size subset of all labels. This choice creates a variation in the number of primary clients for different labels, making it a less idealized and more realistic setup even in the limit of  $s \rightarrow \infty$ . For example, for ImageNet with 1000 classes, if it is subdivided between 8 clients each receiving 250 random labels: (a) around 100 labels will be distributed evenly across all clients (no primary clients), (b) around 270 labels will have a single primary client, (c) around 310 labels will have two primary clients, (d) around 210 labels will have three primary clients and (e) around 110 remaining labels will have 4 or more primary clients.

## 4. Experiments

### 4.1. Experimental Framework

In most of our experiments, we used ImageNet dataset with samples distributed across multiple clients as discussed in Section 3.3. The public dataset used for distillation was chosen by selecting  $\gamma_{\text{pub}} = 10\%$  of all available training samples and the remaining 90% were distributed across clients as private labeled data. We used both *random* and *even* label distribution strategies and considered two cases of  $s = 0$  and  $s = 100$  corresponding to homogeneous and heterogeneous task distributions correspondingly. In most of our experiments, unless indicated otherwise, we used ResNet-34 models as individual clients, trained 8 clients and each was assigned 250 primary labels at *random*. The models were typically trained for 60 000 or 120 000 steps with SGD with momentum, batch size of 512, cosine learning rate decay and the initial learning rate of 0.1 and momentum 0.9.

Our experimental platform was based on distillation losses outlined in Section 3.2. However, being restricted by computational efficiency needed to run numerous experiments, we made several implementation choices that deviated from the general formulation of Section 3.2. Most importantly, individual clients do not directly exchange their predictions on the public dataset, but instead each client  $C_i$  keeps a rolling pool  $\mathcal{P}_i$  of  $N_{\mathcal{P}}$  model checkpoints. In most of our experiments,  $N_{\mathcal{P}}$  was chosen to be equal to the total number of clients in the system. Every step, each client  $C_i$  picks a  $\Delta$  random checkpoints from  $\mathcal{P}_i$  and uses them for performing a distillation step on a new batch. Each pool  $\mathcal{P}_i$  is updated every  $S_{\mathcal{P}}$  steps, when a new checkpoint for one of the other clients is added into the pool (replacing another random checkpoint). In most of our experiments, we used a single distillation

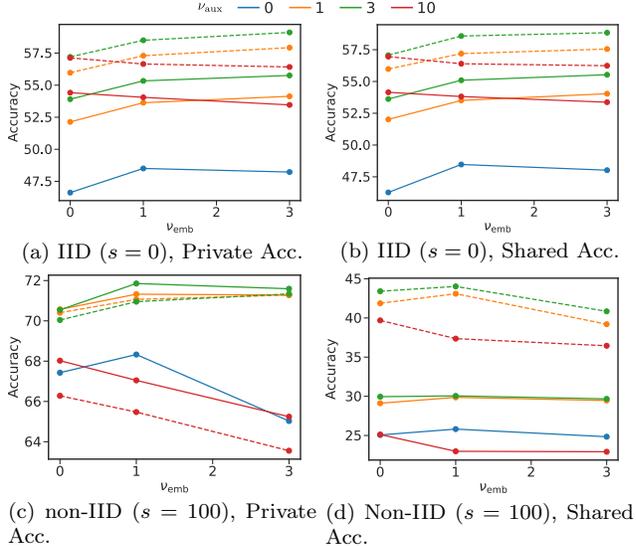


Figure 3. Comparison of *private* (on the client’s dataset) and *shared* accuracies (on a uniform class distribution) for models trained on datasets with iid and non-iid distributions (see Sec. 3.3) (a) with  $s = 0$  and (b)  $s = 100$ . Both the main head (solid) and the auxiliary head accuracies (dashed) are shown. Four values of  $\nu_{\text{aux}}$  are shown: 0.0 (blue), 1.0 (orange), 3.0 (green), 10.0 (red). The accuracies are seen to peak for  $\nu_{\text{aux}} = 3$  and  $\nu_{\text{emb}} = 3$  for  $s = 0$  and  $\nu_{\text{emb}} = 1$  for  $s = 100$ .

client on every step, i.e.,  $\Delta = 1$  and  $e_t(i)$  defined in Sec. 3.1 contains a single element every step  $t$ . However, a separate exploration of the parameter  $\Delta$  was also performed. Also, since in most of our experiments we used  $S_{\mathcal{P}} = 200$ , infrequent pool updates would typically introduce a time lag causing the model to distill knowledge from somewhat outdated checkpoints.

### 4.2. Embedding and Multi-Headed Distillation

In this section we start exploring distillation technique in the simplest scenario with identical model architectures and a complete graph connectivity, where each model can distill knowledge from any other existing client.

#### 4.2.1 Evaluating Basic Distillation Approaches

Consider a set of models with identical ResNet-based architectures learning on their private subsets of ImageNet and distilling the knowledge from each other assuming a complete connectivity of the communication graph. Here we compare the efficiency of knowledge transfer for different distillation approaches: (a) distilling sample embeddings preceding the final log-

its layer (*embedding distillation*) and (b) distilling actual model predictions (*prediction distillation*) (see Sec. 3.2). Specifically, we consider two extreme cases of an iid ( $s = 0$ ) and non-iid ( $s = 100$ ) distributed ImageNet datasets and study the final performance of individual agents while varying the strengths of the embedding and the prediction distillation losses,  $\nu_{\text{emb}}$  and  $\nu_{\text{aux}}$  correspondingly.

In our experiments, we study the performance of primary and auxiliary model heads on two data distributions: (a) *private dataset* defining the primary problem that the client is tasked with and (b) *shared dataset* reflecting the uniform label distribution averaged across all clients. Any technique improving the private dataset accuracy  $\beta_{\text{priv}}$  can be viewed as successful at learning from other clients and translating the acquired knowledge into better performance on their own task. On the other hand, a technique improving the shared dataset accuracy  $\beta_{\text{sh}}$  is successful at learning a more robust representation that can be easily adapted to solving other possible tasks (seen by other clients). Both of these potential capabilities can be viewed as positive outcomes of cross-client communication and learning, but their utility may be application specific.

Figure 3 summarizes our empirical results (see Appendix B for raw numbers) showing the measurements of the average private accuracy  $\beta_{\text{priv}}$ , that is the accuracy of each client on their respective dataset  $\mathcal{D}_i$ , and the averaged shared accuracy  $\beta_{\text{sh}}$  measured on a dataset with a uniform label distribution identical to that of the original ImageNet. While  $\beta_{\text{priv}}$  measures how well a particular client performs on their own task,  $\beta_{\text{sh}}$  is a reflection of the world knowledge (some may be irrelevant for the private task) that the client learns from other participants.

Figure 3 contains several interesting findings: (a) while both regularization techniques are useful for improving model performance, there is a threshold beyond which they start deteriorating both accuracies; (b) taken alone prediction distillation seems to have a stronger positive effect than the embedding distillation, while embedding distillation is more effective in the  $s = 0$  case; (c) however, the best results are obtained by combining both distillation techniques. Furthermore, we see that the distillation techniques generally improve both  $\beta_{\text{priv}}$  and  $\beta_{\text{sh}}$  simultaneously. Notice that the positive effect of  $\nu_{\text{aux}}$  suggests that training a separate auxiliary head has an effect on the model embedding that leads to an improved performance on the main head trained with the client’s private dataset alone. Another interesting observation is that for uniform datasets with a small  $s$ , the auxiliary head ends up having better performance on both the private and

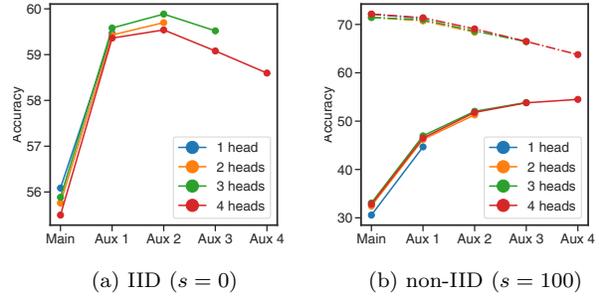


Figure 4. Private (dot-dashed) and shared (solid) dataset accuracies of main and auxiliary heads in ensembles trained with different numbers of auxiliary heads: 1 aux head (blue), 2 heads (orange), 3 heads (green) and 4 heads (red). For the IID case the private and shared performance match.

shared tasks (identical for  $s = 0$ ). At the same time, in a non-iid dataset with  $s = 100$ , auxiliary head performs much better on the shared dataset, but lags behind on the private task since it is not trained on it directly.

#### 4.2.2 Improving Distillation Efficiency

While Figure 3 shows a clear evidence that distillation techniques can be useful for distributed learning even in the case of heterogeneous client data, there is a room for further improvement.

**Ignoring poor distillation targets.** In some cases, agents can be distilling knowledge about particular categories from agents that themselves do not possess accurate information. It is even possible that the agent’s auxiliary head is already “more knowledgeable” about the class than the main head of another agent that it is trying to distill from. As a result, the performance of the auxiliary head may degrade. One approach that we study here is to skip distillation on a sample if the auxiliary head confidence is already higher than that of the head it is trying to distill from. In our experiments, we observed that that this simple idea had virtually no effect for  $s = 0$ , but allowed us to improve the performance of the auxiliary head for heterogeneous data distributions with  $s = 100$ . Specifically, for 8 clients and  $s = 100$ , this technique improved auxiliary head  $\beta_{\text{sh}}$  from 44.7% to 46.5%, while having virtually no effect on the private dataset accuracy  $\beta_{\text{priv}}$  of the main model head, which stayed at 72.2%. While effective for single auxiliary head, this technique did not improve results in multiple auxiliary heads scenario (see Appendix B) that we will discuss next.

$s = 0$	Accuracy	$s = 100$	Accuracy
Separate	46.3%	Separate	25.1%
MHD (Ours)	59.9%	MHD (Ours)	54.5%
MHD+ (Ours)	68.6%	MHD+ (Ours)	63.4%
FA, $u = 200$	70.5%	FA, $u = 200$	68.0%
FA, $u = 1000$	69.1%	FA, $u = 1000$	65.7%
Supervised	68.9%	–	–

Table 1. Comparison of the shared accuracies  $\beta_{\text{sh}}$  for our technique and two “upper-bound” baselines trained for 60k steps on 90% of ImageNet: (a) *supervised* and (b) trained with Federated Averaging (FA) performed every  $u$  steps. *MHD+* experiments were conducted with 180k steps and used the entire ImageNet as a public dataset (regime of plentiful public data). *Separate* corresponds to shared dataset performance for clients trained independently on their own private data. FA accuracy being higher than the supervised could be explained by a much larger number of samples being effectively processed during training ( $\times 8$ ).

**Multiple auxiliary heads.** Here we empirically study the multi-head approach inspired by self-distillation and described in detail in Section 3.2. Guided by earlier results from Section 4.2.1, we choose  $\nu_{\text{emb}} = 1$  and  $\nu_{\text{aux}} = 3$ . We then train an ensemble of 8 models, each with 250 primary labels and two choices of dataset skew:  $s = 0$  and  $s = 100$ . For each choice of parameters, we independently trained models with 1 to 4 auxiliary heads and then measured the performance of the main and every auxiliary head on the client’s private dataset and a shared test set with a uniform label distribution. The results of our experiments are presented in Figure 4 (see Appendix B for raw numbers). For a uniform data distribution, i.e.,  $s = 0$ , we see that distilling multiple auxiliary heads has a positive impact on all model heads for up to 3 auxiliary heads, after which performance starts to degrade. Among the heads themselves, the peak performance is seen to be attained by the 2<sup>nd</sup> auxiliary head. However, we hypothesize that with the increase of the number of training steps, the final head will end up having the highest accuracy.

In the case of a non-iid distribution with  $s = 100$ , we observed that increasing the number of auxiliary heads has a very profound positive affect on the shared dataset performance  $\beta_{\text{sh}}$  of the final auxiliary head. However, it is the main head that achieves the highest private dataset accuracy  $\beta_{\text{priv}}$ . All consecutive auxiliary heads appear to loose their private dataset performance  $\beta_{\text{priv}}$  by specializing on capturing the overall data distribution.

**Dependence on the number of distillation targets  $\Delta$ .** We studied the effect of using multiple distillation targets  $\Delta$  at every training step by considering a typical 8-client setup with  $s = 100$ , 4 auxiliary heads,  $\nu_{\text{emb}} = 1$  and  $\nu_{\text{aux}} = 3$ . While increasing  $\Delta$  from 1 to 3 had virtually no effect on the main head private accuracy  $\beta_{\text{priv}}$ , the shared dataset accuracy  $\beta_{\text{sh}}$  for the last auxiliary head improved from 54.5% to 56.1% and then to 56.4% as we increased  $\Delta$  from 1 to 3. At  $\Delta = 4$ ,  $\beta_{\text{sh}}$  appeared to saturate and fell to 56.2% (within the statistical error of about 0.2%). Overall, earlier auxiliary heads appeared to be affected by  $\Delta$  more strongly.

**Choice of the confidence measure.** The choice of the confidence  $\Lambda(\mathbf{h}(x))$  is central to the distillation technique. We compared our current choice based on selecting the most confident head, with a random selection of the distillation target. In our experiments with 8 clients each with 250 random primary labels,  $\nu_{\text{emb}} = 1$ ,  $\nu_{\text{aux}} = 3$ ,  $s = 0$  and 3 auxiliary heads, we observed that randomizing confidence caused the main head  $\beta_{\text{priv}}$  degradation from 56% to 55.2% and the last auxiliary head  $\beta_{\text{sh}}$  went down from 59.5% to 58.4%. The degradation of model performance is more significant in the case of heterogeneous client data. In experiments with  $s = 100$  and 4 auxiliary heads, we observed the main head  $\beta_{\text{priv}}$  degraded from 72.1% to 71.3% and the last auxiliary head  $\beta_{\text{sh}}$  decreased from 54.5% to 49%.

**Dependence on the technique efficiency on the public dataset size.** The efficiency of model distillation depends on the amount of data used for performing this distillation, in our case, on the size of the public dataset. In our experiments outlined in Appendix B.2, increasing the size of the public dataset while fixing the amount of private training data has a positive impact on the final model performance.

In practice, since unlabeled data is more abundant, one can expect that the public dataset size will be comparable or even larger than the total amount of labeled data available to clients. Being constrained by the ImageNet size and attempting to keep the amount of private training data unaffected, we simulate the abundance of public data by reusing the entirety of the ImageNet dataset as an unlabeled public dataset. This, of course, is not realistic and somewhat biased given that we reuse the same samples as labeled and unlabeled, but it allows us to explore the limits of the distributed training efficiency with distillation.

MHD Base	MHD	FedMD Base	FedMD
60.6%	57.0% / 0.6%	56.5%	50.2% / 2.7%

Table 2. Comparison of mean test accuracies (first number) and their deviations (second number after /) across 10 clients for our method and FedMD as reported in Ref. [19]. Baselines (Base) are obtained by training clients with all available private data.

### 4.3. Baseline Comparisons

Before comparing our technique with a similar distillation-based method, we compared its performance with two strong “upper-bound” baselines (see Table 1): supervised training on all ImageNet and FedAvg algorithm implemented within our framework. A large performance gap between shared dataset accuracies obtained using our method and the strong baselines can be viewed as a price paid for learning via distillation in a decentralized multi-agent system. At the same time, we see that increasing the public dataset size and training for a longer period of time, allowing the information to propagate across all clients (*Our* results), brings us close to the supervised model performance. Notice that like many other distillation-based techniques [19,39], our method reaches higher accuracy in the homogeneous data scenario.

We compared our method with FedMD [19] a similar, but *centralized* distillation-based methods. This comparison was carried out by replicating the dataset and 10 model architectures from the publicly available implementation. The dataset is based on CIFAR-100 [17] and makes use of 20 coarse labels, while the public dataset is chosen to be CIFAR-10. Due to the differences in the training process, our baseline results with individual models trained on all private data pooled together was higher than that reported in [19]. At the same time, we observed a much smaller gap in performance between this upper baseline and the results obtained using our method than the gap reported in [19] (see Table 2). Interestingly, we also observe a much smaller performance spread across all 10 models trained with our technique (deviation of 0.6% compared to 2.7% for FedMD).

### 4.4. Communication Topology Effects

In order to explore how our approach might scale to larger systems in which pairwise connections between all agents are not feasible, we aim to evaluate how the communication topology affects performance. In particular we are interested in the question of whether “transitive distillation” is possible with our approach – that is whether two agents that are not directly con-

nected to one-another can still learn from each-other through an intermediary.

To evaluate this and determine how auxiliary heads play a role in the performance we ran a training sweep with 4 agents arranged in 3 different topologies (Figure 5) with 3 auxiliary heads each. In all cases we trained for 120k steps, with 250 primary labels per agent with  $s = 100$ . We observe (Figure 6) that performance on the shared dataset improves significantly between island and cycle topology, with the baseline performance matching closely the cycle performance. Without transitive distillation we would expect island and cycle performance to match closely so this provides strong evidence for transitive distillation. Also note that this behavior is only present on auxiliary heads and is more pronounced for later heads.

We further analyze the performance of each agent on other agents’ private data. Predictably we observe that island topologies perform well on in-island other agents, and poorly on agents from outside their island. Cycle topology agents perform best on their di-

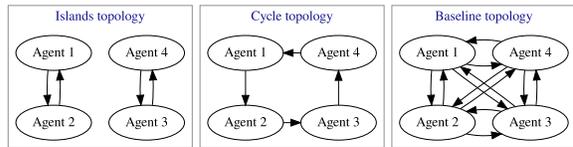


Figure 5. Topologies compared to validate transitive distillation.

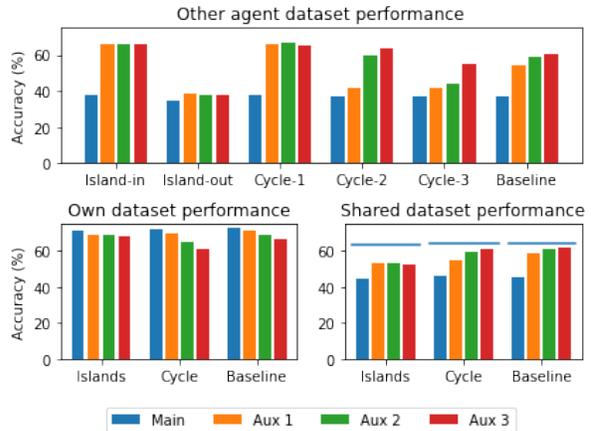


Figure 6. The performance by topology and distance between distillation teacher and student. On the shared dataset the blue horizontal lines indicate the upper bound per the embedding quality – computed by fine tuning a head on the frozen model embeddings. Note that island embedding accuracy for “Islands” is still worse than “Cycle”. Best viewed in color.

rect teacher (*Cycle-1*), but auxiliary heads 2 and 3 perform well on the “1-hop” transitive teacher (*Cycle-2*), and auxiliary head 3 has markedly improved performance on the “2-hop” transitive teacher (*Cycle-3*). We take this as strong evidence that auxiliary heads enable transitive distillation, and that additional heads make learning across additional degrees of separation more efficient.

#### 4.5. Learning in Heterogeneous Systems

In Section 4.2, we conducted experiments with homogeneous ensembles of models. However, in many realistic scenarios of distributed deep learning, client devices may have different hardware-defined limitations and it may be desirable to train smaller models on some clients, while allowing other devices to utilize much larger networks. While model distillation allows one to achieve this, it is reasonable to ask why would this even be desirable? What do we expect to gain from having much larger models in the ensemble? Here we show two positive effects emerging from having larger models in an ensemble of smaller clients: (a) informally speaking, small models benefit from having stronger teachers and (b) large models can gain complex knowledge by distilling from smaller and simpler models.

Our ImageNet experiments were conducted with 4 clients each assigned 500 primary labels with one client being a ResNet34 model and the remaining clients being ResNet18. Primary label assignment was random across clients and we trained the model for 240k steps.

First, we observed that the presence of a larger model improved the accuracy of smaller clients suggesting that they benefited from seeing a stronger teacher holding some of the relevant data. Specifically, we observed that the presence of a ResNet34 model instead of ResNet18 in the ensemble led to an increase in the average shared accuracy  $\beta_{sh}$  of ResNet18 models from 66.2% to 66.7%.

Secondly, if small models achieve high performance on their limited personalized domains, a large model distilling from such an ensemble can potentially learn a much more complex picture of the entire dataset than would otherwise be accessible to any individual small learner. This observation has already inspired centralized distillation-based methods like [13]. In our experiments, we witnessed this by observing that ResNet34 trained in conjunction with 3 ResNet18 clients reached the shared accuracy  $\beta_{sh}$  of 68.6%, which exceeds the 67.7% accuracy of an ensemble of 4 ResNet18 models trained with FedAvg or 66.0% if trained with our approach (both with 200 steps between updates). Notice that if the ResNet34 model is isolated from ResNet18 models, it only reaches  $\beta_{sh}$  of 39.4%.

## 5. Discussion and Conclusions

In this paper, we proposed a novel distributed machine learning technique based on model distillation. The core idea of our approach lies in using a hierarchy of multiple auxiliary heads distilling knowledge from each other and across the ensemble. We show that this technique is much more effective than naive distillation and allows us to get close to the supervised accuracy on a large ImageNet dataset given a large public dataset and longer training time necessary for information to spread across the system. We also study two key capabilities of a distributed distillation-based learning technique. Specifically, we demonstrate that in systems where direct communication between the clients is limited, multiple auxiliary heads allow information exchange across clients that are not directly connected. We also demonstrate two positive effects of adding larger models into the system of small models: (a) small models benefit from seeing larger teachers and that (b) large models learning from a collection of small models can reach higher accuracies than those achievable with small models only.

## References

- [1] Gustavo Aguilar, Yuan Ling, Yu Zhang, Benjamin Yao, Xing Fan, and Chenlei Guo. Knowledge distillation from internal representations. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 7350–7357. AAAI Press, 2020. 2, 3
- [2] Sungsoo Ahn, Shell Xu Hu, Andreas C. Damianou, Neil D. Lawrence, and Zhenwen Dai. Variational information distillation for knowledge transfer. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 9163–9171. Computer Vision Foundation / IEEE, 2019. 2, 4
- [3] Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? In Zoubin Ghahramani, Max Welling, Corinna Cortes, Neil D. Lawrence, and Kilian Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, pages 2654–2662, 2014. 1
- [4] Ilai Bistriz, Ariana J. Mann, and Nicholas Bambos. Distributed distillation for on-device learning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Pro-*

- cessing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual, 2020. [2](#)
- [5] Kallista A. Bonawitz, Vladimir Ivanov, Ben Kreuter, Antonio Marcedone, H. Brendan McMahan, Sarvar Patel, Daniel Ramage, Aaron Segal, and Karn Seth. Practical secure aggregation for privacy-preserving machine learning. In Bhavani Thuraisingham, David Evans, Tal Malkin, and Dongyan Xu, editors, *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, CCS 2017, Dallas, TX, USA, October 30 - November 03, 2017*, pages 1175–1191. ACM, 2017. [2](#)
- [6] Cristian Bucila, Rich Caruana, and Alexandru Niculescu-Mizil. Model compression. In Tina Eliassi-Rad, Lyle H. Ungar, Mark Craven, and Dimitrios Gunopulos, editors, *Proceedings of the Twelfth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, Philadelphia, PA, USA, August 20-23, 2006*, pages 535–541. ACM, 2006. [1](#)
- [7] George Cazenavette, Tongzhou Wang, Antonio Torralba, Alexei A. Efros, and Jun-Yan Zhu. Dataset distillation by matching training trajectories. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 10708–10717. IEEE, 2022. [15](#)
- [8] Hong-You Chen and Wei-Lun Chao. Fedbe: Making bayesian model ensemble applicable to federated learning. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021. [2](#)
- [9] 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. Ieee, 2009. [2](#)
- [10] Tommaso Furlanello, Zachary Chase Lipton, Michael Tschannen, Laurent Itti, and Anima Anandkumar. Born-again neural networks. In Jennifer G. Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholm, Sweden, July 10-15, 2018*, volume 80 of *Proceedings of Machine Learning Research*, pages 1602–1611. PMLR, 2018. [2](#), [4](#)
- [11] Xuan Gong, Abhishek Sharma, Srikrishna Karanam, Ziyang Wu, Terrence Chen, David S. Doermann, and Arun Innanje. Preserving privacy in federated learning with ensemble cross-domain knowledge distillation. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 11891–11899. AAAI Press, 2022. [2](#)
- [12] Neel Guha, Ameet Talwalkar, and Virginia Smith. One-shot federated learning. *CoRR*, abs/1902.11175, 2019. [1](#), [2](#)
- [13] Chaoyang He, Murali Annavaram, and Salman Avestimehr. Group knowledge transfer: Federated learning of large cnns at the edge. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. [2](#), [9](#)
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 770–778. IEEE Computer Society, 2016. [2](#)
- [15] Geoffrey E. Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. *CoRR*, abs/1503.02531, 2015. [1](#)
- [16] Sohei Itahara, Takayuki Nishio, Yusuke Koda, Masahiro Morikura, and Koji Yamamoto. Distillation-based semi-supervised federated learning for communication-efficient collaborative training with non-iid private data. *CoRR*, abs/2008.06180, 2020. [2](#)
- [17] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. [8](#)
- [18] Viraj Kulkarni, Milind Kulkarni, and Aniruddha Pant. Survey of personalization techniques for federated learning. *CoRR*, abs/2003.08673, 2020. [2](#)
- [19] Daliang Li and Junpu Wang. Fedmd: Heterogeneous federated learning via model distillation. *CoRR*, abs/1910.03581, 2019. [2](#), [8](#)
- [20] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In Inderjit S. Dhillon, Dimitris S. Papailiopoulos, and Vivienne Sze, editors, *Proceedings of Machine Learning and Systems 2020, MLSys 2020, Austin, TX, USA, March 2-4, 2020*. mlsys.org, 2020. [2](#)
- [21] Tao Lin, Lingjing Kong, Sebastian U. Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. [2](#)
- [22] Jiaxin Ma, Ryo Yonetani, and Zahid Iqbal. Adaptive distillation for decentralized learning from heterogeneous clients. In *25th International Conference on Pattern Recognition, ICPR 2020, Virtual Event / Milan, Italy, January 10-15, 2021*, pages 7486–7492. IEEE, 2020. [4](#)
- [23] Disha Makhija, Xing Han, Nhat Ho, and Joydeep Ghosh. Architecture agnostic federated learning for neural networks. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and

- Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 14860–14870. PMLR, 2022. [2](#)
- [24] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-efficient learning of deep networks from decentralized data. In Aarti Singh and Xiaojin (Jerry) Zhu, editors, *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, AISTATS 2017, 20-22 April 2017, Fort Lauderdale, FL, USA*, volume 54 of *Proceedings of Machine Learning Research*, pages 1273–1282. PMLR, 2017. [1](#), [2](#)
- [25] Hossein Mobahi, Mehrdad Farajtabar, and Peter L. Bartlett. Self-distillation amplifies regularization in hilbert space. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. [4](#)
- [26] Minh N. H. Nguyen, Huy Q. Le, Shashi Raj Pandey, and Choong Seon Hong. CDKT-FL: cross-device knowledge transfer using proxy dataset in federated learning. *CoRR*, abs/2204.01542, 2022. [2](#), [3](#)
- [27] Timothy Nguyen, Roman Novak, Lechao Xiao, and Jaehoon Lee. Dataset distillation with infinitely wide convolutional networks. *CoRR*, abs/2107.13034, 2021. [15](#)
- [28] Krishna Pillutla, Kshitiz Malik, Abdelrahman Mohamed, Michael Rabbat, Maziar Sanjabi, and Lin Xiao. Federated learning with partial model personalization. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pages 17716–17758. PMLR, 2022. [2](#)
- [29] Felix Sattler, Arturo Marbán, Roman Rischke, and Wojciech Samek. Communication-efficient federated distillation. *CoRR*, abs/2012.00632, 2020. [2](#)
- [30] Tao Shen, Jie Zhang, Xinkang Jia, Fengda Zhang, Gang Huang, Pan Zhou, Fei Wu, and Chao Wu. Federated mutual learning. *CoRR*, abs/2006.16765, 2020. [2](#)
- [31] Stefán Páll Sturluson, Samuel Trew, Luis Muñoz-González, Matei Grama, Jonathan Passerat-Palmbach, Daniel Rueckert, and Amir Alansary. Fedrad: Federated robust adaptive distillation. *CoRR*, abs/2112.01405, 2021. [2](#)
- [32] Lichao Sun and Lingjuan Lyu. Federated model distillation with noise-free differential privacy. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJ-CAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pages 1563–1570. ijcai.org, 2021. [2](#)
- [33] Alysa Ziyang Tan, Han Yu, Lizhen Cui, and Qiang Yang. Towards personalized federated learning. *CoRR*, abs/2103.00710, 2021. [2](#)
- [34] Yue Tan, Guodong Long, Lu Liu, Tianyi Zhou, Qinghua Lu, Jing Jiang, and Chengqi Zhang. Fedproto: Federated prototype learning across heterogeneous clients. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 8432–8440. AAAI Press, 2022. [2](#)
- [35] Tianchun Wan, Wei Cheng, Dongsheng Luo, Wenchao Yu, Jingchao Ni, Liang Tong, Haifeng Chen, and Xiang Zhang. Personalized federated learning via heterogeneous modular networks. *CoRR*, abs/2210.14830, 2022. [2](#)
- [36] Tongzhou Wang, Jun-Yan Zhu, Antonio Torralba, and Alexei A. Efros. Dataset distillation. *CoRR*, abs/1811.10959, 2018. [15](#)
- [37] Chuhan Wu, Fangzhao Wu, Ruixuan Liu, Lingjuan Lyu, Yongfeng Huang, and Xing Xie. Fedkd: Communication efficient federated learning via knowledge distillation. *CoRR*, abs/2108.13323, 2021. [2](#)
- [38] Chenglin Yang, Lingxi Xie, Siyuan Qiao, and Alan L. Yuille. Training deep neural networks in generations: A more tolerant teacher educates better students. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 5628–5635. AAAI Press, 2019. [2](#), [4](#)
- [39] Jie Zhang, Song Guo, Xiaosong Ma, Haozhao Wang, Wenchao Xu, and Feijie Wu. Parameterized knowledge transfer for personalized federated learning. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 10092–10104, 2021. [2](#), [8](#)
- [40] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 12878–12889. PMLR, 2021. [2](#)

## A. Analysis and Discussion of Our Method

### A.1. Analysis of Multi-Headed Distillation

As described in the main text, the multi-headed distillation involves simultaneous training of multiple model heads that communicate with each other. Rigorous theoretical analysis of this process in the most general case is very complicated. However, by making several assumptions, here we find an approximate solution for the weights of the heads of rank  $k$  being given the weights of the heads of rank  $k - 1$ . In our future work, we hope to study this model for different prediction aggregation techniques and compare conclusions obtained from this simple model with those obtained empirically in realistic systems.

Let  $X$  be the input space and  $L = \mathbb{R}^d$  be the logit space, where  $d$  is the number of classes. The logits  $f(x)$  for a model  $f : X \rightarrow L$  are then converted to label assignment probabilities  $p(y|x)$  via softmax, i.e.,  $p(y|x) = \text{softmax}(f(x))$ .

Consider a single client  $h_i : X \rightarrow L$  distilling information from some model head  $h : X \rightarrow L$ . The corresponding distillation loss  $\mathcal{L}[h_i; h]$  admits many possible choices, but we will assume that

$$\mathcal{L} \equiv \mathbb{E}_{x \sim \mathcal{D}} D[h_i(y|x; \psi_i) \parallel p_h(y|x)]$$

with  $\mathcal{D}$  being the shared (proxy) dataset and  $D$  being some divergence (more general  $f$ -divergence or KL-divergence as some examples). The distillation is then carried out by performing optimization of  $\mathcal{L}$ , for example via gradient descent:

$$\Delta \psi_i = -\gamma \frac{\partial \mathcal{L}}{\partial \psi_i}.$$

Notice that the components of  $\psi_i$  corresponding to the model backbone may receive updates from multiple heads reusing the same model embedding.

In our system, we assume that there is a set of heads  $\{h_i^{(1)}, h_i^{(2)}, \dots, h_i^{(n)}\}$  for each client  $i$ . For simplicity, let us first consider distillation procedure independent of prediction confidence. In this case, the loss for head  $k$  of the client  $i$  may look like:

$$\mathcal{L}_i^{(k)} \equiv \sum_{j=1}^N \rho_{ij} \Gamma[h_i^{(k)} \parallel h_j^{(k-1)}],$$

where

$$\Gamma[h_i^{(k)} \parallel h_j^{(k-1)}] \equiv \mathbb{E}_{x \sim \mathcal{D}} D[p_i^{(k)}(y|x) \parallel p_j^{(k-1)}(y|x)],$$

$p_i^{(k)}(y|x)$  is a shorter notation for  $p_{h_i^{(k)}}(y|x)$  and  $\rho_{ij}$  is some distribution defining the probability of picking a

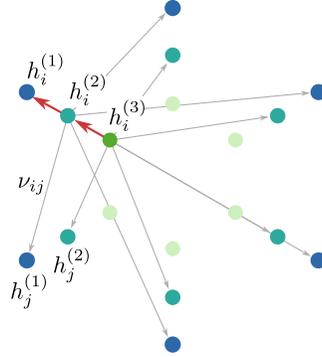


Figure 7. Illustration of multi-head distillation as discussed in Appendix A.1. Large red arrow shows strong distillation of  $h_i^{(k)}$  to  $h_i^{(k-1)}$  and smaller gray arrows indicate attraction towards  $h_j^{(k-1)}$  with effective “strength”  $\nu_{ij}$ .

particular client for distillation. Again, here we assume that  $\rho_{ij}$  does not depend on the sample confidence and is simply fixed.

While we talked about  $h^{(k)}$  distilling to  $h^{(k-1)}$ , we have not yet discussed the “main head”  $h^{(1)}$ . This head is normally trained locally on the client’s private data. For simplicity, in the following we thus assume that its behavior is known, i.e.,  $h_i^{(1)}$  is a specified function of the training step. Furthermore, in the following, we start analyzing the problem by assuming that all  $h_i^{(1)}$  already converged and are all generally different due to some differences in the client’s private data. The behavior of all other heads is then defined by the losses outlined above.

Let us first consider the simplest case of  $\rho_{ij} = \delta_{ij}$ . In other words, each head only distills from the same client’s “prior” head. The choice of  $h_i^{(n)} = \dots = h_i^{(1)}$  would obviously minimize all losses  $\mathcal{L}_i^{(k)}$  since all corresponding  $\Gamma[\cdot]$  values vanish. But as soon as we introduce a small correction  $\rho_{ij} = \delta_{ij} + \nu_{ij}$  with  $\sum_j \nu_{ij} = 0$ , this trivial solution is no longer optimal. Instead, each client’s head is now optimizing:

$$\mathcal{L}_i^{(k)} = \Gamma[h_i^{(k)} \parallel h_i^{(k-1)}] + \sum_{j=1}^N \nu_{ij} \Gamma[h_i^{(k)} \parallel h_j^{(k-1)}].$$

Notice that if  $\Gamma$  was a metric in the  $h$  space, we could interpret this optimization objective geometrically as a minimization of the head’s distance to its lower-order state (towards  $h_i^{(k-1)}$ ) coupled with a weak ( $\sim \nu$ ) attraction towards a number of other heads ( $h_j^{(k-1)}$ ). See Figure 7 for illustration.

Here we have to make yet another simplifying assumption and consider a prescribed model backbone (and corresponding embedding) that we are not optimizing or updating with backpropagated gradients.

Doing so, we disentangle individual heads and can treat their optimization as independent tasks. For sufficiently small  $\nu_{ij}$  it will hold that  $p_i^{(k)} = p_i^{(k-1)} + O(\nu)$  and we can therefore write:

$$\mathcal{L}_i^{(k)} = \mathbb{E}_{x \sim \mathcal{D}} \left\{ D \left[ p_{h_i^{(k-1)} + \kappa_i^{(k)}} \parallel p_{h_i^{(k-1)}} \right] + \sum_{j=1}^N \nu_{ij} D \left[ p_{h_i^{(k-1)} + \kappa_i^{(k)}} \parallel p_{h_j^{(k-1)}} \right] \right\},$$

where  $h_i^{(k)} \equiv h_i^{(k-1)} + \kappa_i^{(k)}$  and  $\kappa_i^{(k)} \sim O(\nu)$ . Introducing  $\delta_i^{(k)} \equiv p_i^{(k)} - p_i^{(k-1)}$ , we obtain:

$$\mathcal{L}_i^{(k)} = \mathbb{E}_{x \sim \mathcal{D}} \left\{ D \left[ p_i^{(k-1)} + \delta_i^{(k)} \parallel p_i^{(k-1)} \right] + \sum_j \nu_{ij} D \left[ p_i^{(k-1)} + \delta_i^{(k)} \parallel p_j^{(k-1)} \right] \right\}.$$

Noticing that the first term needs to be decomposed near the global minimum and the second term permits linear expansion, we obtain:

$$\mathcal{L}_i^{(k)} \approx \mathbb{E}_{x \sim \mathcal{D}} \left\{ D'' \left[ p_i^{(k-1)} \parallel p_i^{(k-1)} \right] \frac{\delta_i^{(k)} \delta_i^{(k)}}{2} + \sum_j \nu_{ij} D' \left[ p_i^{(k-1)} \parallel p_j^{(k-1)} \right] \delta_i^{(k)} \right\},$$

where  $D''$  and  $D'$  are the derivatives of  $D$  with respect to the first argument. Recalling that  $\delta_i^{(k)} \in \mathbb{R}^d$  we can rewrite the loss function as:

$$\mathcal{L}_i^{(k)} \approx \mathbb{E}_{x \sim \mathcal{D}} \left[ \delta^\top A \delta + b^\top \delta \right],$$

where  $\delta \equiv \delta_i^{(k)}$  for brevity,

$$A \equiv D'' \left[ p_i^{(k-1)} \parallel p_i^{(k-1)} \right] / 2$$

is effectively a matrix and

$$b \equiv \sum_j \nu_{ij} D' \left[ p_i^{(k-1)} \parallel p_j^{(k-1)} \right] \in \mathbb{R}^d$$

can be thought of as a column vector.

At this point we can connect the probability distribution perturbation  $\delta$  to the logit perturbation  $\kappa \equiv \kappa_i^{(k)}$  using the fact that  $p_m \equiv e^{h_m} / Z$ , where  $Z \equiv \sum_k e^{h_k}$  (we omit this simple calculation here):

$$\begin{aligned} p_i^{(k)} &= p_{h_i^{(k-1)} + \kappa_i^{(k)}} = p_i^{(k-1)} + \delta = \\ &= p_i^{(k-1)} + \kappa * p_i^{(k-1)} - (\kappa \cdot p_i^{(k-1)}) p_i^{(k-1)}, \end{aligned}$$

where  $\mathbf{a} * \mathbf{b}$  is an element-wise product of two vectors and therefore:

$$\delta = \kappa * p_i^{(k-1)} - (\kappa \cdot p_i^{(k-1)}) p_i^{(k-1)} \equiv C \kappa, \quad (6)$$

where  $C$  is a matrix constructed from the components of  $p_i^{(k-1)}(x) \in \mathbb{R}^d$ . Notice that  $\sum_m \delta_m = 0$ , which agrees with  $\delta$  being the perturbation of the normalized probability distribution.

Finally, remember that  $\kappa$  itself is a perturbation of model logits. Given the sample embedding  $\xi_i(x) \in \mathbb{R}^t$ , the sample logits are constructed as  $W_i \xi_i(x)$  with  $W_i$  being a  $d \times t$  matrix. The perturbation  $\kappa$  transforming  $W_i^{(k-1)} \xi_i$  into  $W_i^{(k)} \xi_i$  can thus be characterized by the logit weight perturbation  $\mu \equiv \mu_i^{(k)} := W_i^{(k)} - W_i^{(k-1)}$  and we get  $\kappa = \mu \xi(x)$ . Combining everything together, we see that the loss function transforms to:

$$\mathcal{L}_i^{(k)} \approx \mathbb{E}_{x \sim \mathcal{D}} \left[ \xi(x)^\top \mu^\top C^\top A C \mu \xi(x) + b^\top C \mu \xi(x) \right], \quad (7)$$

where  $A$ ,  $C$  and  $b \sim \nu$  all depend on the sample  $x$  via  $p_i^{(k-1)}(x)$  and  $\xi$  is a function of  $x$ , while  $\mu$  is effectively an unknown sample-independent matrix that we need to tune with the goal of minimizing  $\mathcal{L}_i^{(k)}$ . The optimum can be identified by taking a derivative with respect to  $\mu_{\alpha\beta}$  and setting it to 0:

$$\mathbb{E}_{x \sim \mathcal{D}} \left[ 2(\xi(x)^\top \mu^\top C^\top A C)_\alpha \xi_\beta(x) + (b^\top C)_\alpha \xi_\beta(x) \right] = 0.$$

This is a linear equation on  $\mu \sim \nu$  that can be solved in a closed form to give us a logit weight perturbation  $\mu$  as a complex nonlinear function of  $\nu_{ij}$  and  $\{p_\ell^{(k-1)}\}$ .

Note that since  $\mu$  is only a small perturbation, we can introduce  $W_i^{(k)}$  as a function of a *continuous* parameter  $k$  and approximate  $dW_i^{(k)}/dk$  with a finite difference  $W_i^{(k)} - W_i^{(k-1)} = \mu$  leaving us with a differential equation (the approximation is valid in the first order in  $\nu$ ):

$$\frac{dW_i(k)}{dk} = G[\nu, \{W_\ell(k)\}]$$

with  $G$  being a linear function with respect to  $\nu$ , but very complex nonlinear function with respect to  $\{W_\ell\}$ . If  $\nu_{ij}$  is localized around  $i = j$  (which would be the case for communication patterns with partial connectivity, like in the case of long chains), this differential equation resembles a complex nonlinear diffusion equation defining the spread of information across the clients as we look at deeper and deeper heads (with the head rank  $k$  essentially playing the role of time).

It is also worth noting here that if  $\nu$  was not fixed, but was itself a function of model confidence (while still remaining small), our conclusions would not change except that  $\nu$  itself would now itself be a complex nonlinear function of  $\{W_\ell(k)\}$  and  $x$ . In our future work, we

hope to study the effect that this confidence-dependent aggregation has on head dynamics and the final stationary state.

Finally, let us look at the stationary state of system dynamics. Equation (7) suggests that  $\mu = 0$  is a local optimum when  $b^\top C = 0$ , or

$$\sum_{i,j} \nu_{ij} D' [p_i \| p_j] C_{ik} = 0,$$

or after noticing that  $D' [p_i \| p_j] = p_j / p_i$  and recalling that  $C$  is defined by Eq. (6) we obtain for every  $k$ :

$$\mathbb{E}_{x \sim \mathcal{D}} \left[ \sum_{i,j} \nu_{ij} p_j (\delta_{ik} - p_k) \right] = 0. \quad (8)$$

Since  $\sum_j \nu_{ij} = 0$ , the trivial solution of this system of equations is the case of identical models, i.e.,  $p_1 = \dots = p_n$ , but since generally the models might have different embeddings and cannot be made identical, the solution of Eq. (8) restricts the system stationary state.

## A.2. Value of $p(y|x)$ as Classifier Confidence

In our model distillation approach, we need to combine predictions of multiple different model heads. If all predictions  $p_k(y|x)$  (by heads  $\{h_k\}$ ) come with reliable error estimates, this information can be taken into account. For example, if we know that for the true distribution  $p(y|x)$  and every prediction  $p_k(y|x)$  it holds that  $D[p_k(y|x) \| p(y|x)] \leq e_k(x)$  with  $D$  being some divergence, the true  $p(y|x)$  belongs to the intersection of “balls”<sup>2</sup>  $\mathcal{B}_k \equiv \{p' | D[p' \| p] \leq e_k\}$ . We can then choose any point in this set and compute a prediction error as a maximum distance from a chosen distribution to any point in the intersection. Unfortunately, however such reliable measures of classifier error are not generally available and even approximating them can be quite difficult.

Instead we choose a very simple approach based on estimating classifier confidence and picking the most confident model, effectively ignoring other predictions. The confidence measure itself is chosen as a value of the largest component of the classifier prediction  $o(x) \equiv \text{softmax}(f(x; \theta))$  with  $f(x; \theta) = W\xi(x; \theta)$  and  $\xi$  being the embedding vector.

**Why choose  $\max_k o_k$ .** This value has a number of trivial properties that can actually make it a useful measure of classifier uncertainty. First is that after seeing a supervised training sample  $(x, y)$ , the value of  $o_y(x)$  is increased. Second is that if the class prototypes in the embedding space are nearly orthogonal

<sup>2</sup>note that  $D$  is not generally a metric

for different classes, then updates for samples of different classes would not “interfere” with each other and high-confidence predictions would not generally be disrupted by observing unrelated samples with different labels. For a simple logits layer  $W\xi(x)$  trained with cross-entropy loss, both of these properties trivially follow from the following expression for  $\Delta o_k(x')$  after training on a sample  $(x, y)$ :

$$\begin{aligned} \Delta o_k(x') &= \lambda o_k(x') [\xi(x) \cdot \xi(x')] \times \\ &\quad \times \sum_i (\delta_{k,i} - o_i(x')) (\delta_{y,i} - o_i(x)). \end{aligned}$$

**Drawbacks.** But while  $\max_k o_k(x)$  has these useful properties, it is not guaranteed to be a reliable measure of classifier confidence for out-of-distribution samples and the training objective never explicitly optimizes for that<sup>3</sup>. A density model  $\rho(x)$  would allow detecting such out-of-distribution samples, but could also reveal information about the client samples in their private dataset. Combining classification models with locally-trained density models, or adopting other existing similar techniques could be a logical extension of our present work.

## A.3. Distillation as Revelation of Some Information about Model Weights

The canonical version of FedAvg combines the knowledge of individual clients by periodically aggregating their weight snapshots. Distillation-based techniques are instead based on communicating model predictions on datasets accessible to all participants. While these two approaches appear to be different, communication in distillation-based methods can of course also be viewed as a way of revealing incomplete information about model weights.

The amount of revealed information can be defined as follows. Assuming the knowledge of the prior  $p(\theta)$  on the model weights and model predictions  $(y_1, \dots, y_n)$  on a public dataset  $\mathcal{D}_* = (x_1, \dots, x_n)$ , one can compare the difference of entropies for the original  $p(\theta)$  and  $p(\theta|y_1, \dots, y_n)$  with

$$p(\theta|y_1, \dots, y_n) = \frac{p(y_1, \dots, y_n | \theta) p(\theta)}{\int d\theta p(y_1, \dots, y_n | \theta) p(\theta)}.$$

While generally intractable, it might be possible to obtain the lower bound on the amount of the revealed information by training a model that predicts the weights  $\theta$  from  $(y_1, \dots, y_n)$ .

<sup>3</sup>Contrast this to the energy-based models, for example, where the energy update and the MCMC sampling are explicitly contributing to the model “awareness” of what in-distribution samples are and are not.

Deeper understanding of this question can have an impact on the optimal choice of the public dataset  $\mathcal{D}_*$  that would allow us to retrieve the knowledge of interest from a trained model using only a small number of samples. Ongoing research on dataset distillation [7, 27, 36] is very closely related to this question.

## B. Additional Experiments and Experimental Data

### B.1. Effect of Distilling to Self and Same-Level Heads

In Section 4.2.2 we reported that including a head into the list of its own distillation targets (“self”) improved the model accuracy, but the gain was still smaller than that of a model with multiple auxiliary heads. Here we explore what happens if we use the head as its own potential distillation target, while also using a number of auxiliary heads. Furthermore, what if we modify our method to include distillations to other heads of the same rank (see Figure 9)?

We conducted a set of experiments with a heterogeneous dataset with  $s = 100$ ,  $\nu_{\text{emb}} = 1$ ,  $\nu_{\text{aux}} = 3$ , four auxiliary heads and 250 randomized labels per each of 8 clients. The results of experiments using different combinations of distillation targets and both  $\Delta = 1$  and  $\Delta = 2$  (choosing two other clients at a time as potential distillation targets) are presented in Table 3. We observed that using same-level heads and “self” targets *separately* provides noticeable benefit only for earlier heads. But when used together, these two techniques result in  $\sim 1\%$  accuracy improvement and this improvement is realized for the 2<sup>nd</sup> auxiliary head. Also, not unexpectedly, using two clients to distill to ( $\Delta = 2$ ) instead of one, leads to a noticeable 1.5% accuracy improvement. Combined together, all these techniques, in conjunction with using the entire ImageNet as the public dataset improve the accuracy to 59.4% if trained for 60k steps, or 65.7% if trained for 180k steps.

### B.2. Dependence on the Public Dataset Size

In a separate set of experiments, we trained 8 clients with 4 auxiliary heads,  $s = 100$ ,  $\nu_{\text{emb}} = 1$ ,  $\nu_{\text{aux}} = 3$  and 250 randomly assigned “private” labels and “private” samples drawn from 70% of the ImageNet training data. The remaining 30% of ImageNet samples were fully or partly used as a public dataset, i.e.,  $\gamma_{\text{pub}} \leq 30\%$ . As one would expect, increasing the size of the “public” dataset while fixing the amount of “private” training data has a positive impact on the final model performance (see Table 4).

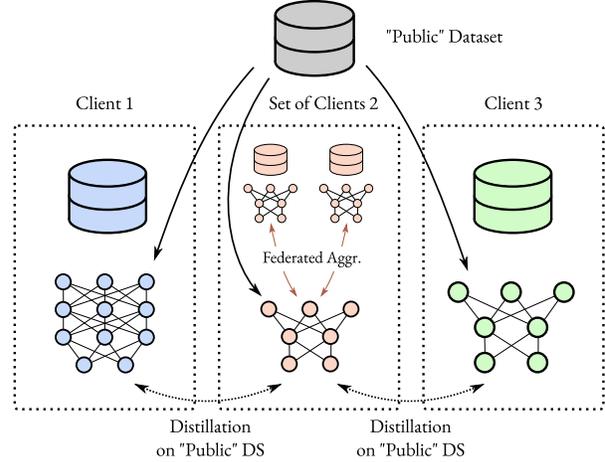


Figure 8. Conceptual diagram of a distillation in a distributed system. Clients use a “public” dataset to distill knowledge from other clients, each having their primary private dataset. Individual clients may have different architectures and different objective functions. Furthermore, some of the “clients” may themselves be collections of models trained using federated learning.

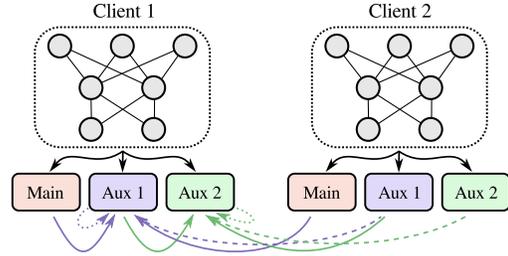


Figure 9. A pattern used for distilling multiple auxiliary heads with two additional types of distillation targets: (a) distilling to heads of the same “rank” (dashed), (b) distilling to “self” (dotted). Here distilling to the same “rank” means, for example, that *Aux1* head is distilled to the most confident of *Main* heads, or *Aux1* heads of adjacent clients. Distilling to “self” means that the samples on which the distilled head is already most confident will effectively be ignored.

## C. Additional Tables and Figures

### C.1. Raw Experimental Data

Tables 5 and 6 contain raw values used for producing Figure 3, while Tables 7 and 8 complement Figure 4.

Experiment	$\beta_{\text{priv}}^{\text{Main}}$	$\beta_{\text{sh}}^{\text{Aux1}}$	$\beta_{\text{sh}}^{\text{Aux2}}$	$\beta_{\text{sh}}^{\text{Aux3}}$	$\beta_{\text{sh}}^{\text{Aux4}}$
Base	70.9%	46.7%	51.8%	53.9%	<b>54.6%</b>
$\Delta = 2$	71.1%	50.9%	55.1%	<b>56.1%</b>	<b>56.0%</b>
SL	70.8%	48.6%	53.6%	<b>54.7%</b>	<b>54.7%</b>
SF	71.3%	48.1%	53.4%	<b>54.9%</b>	<b>54.8%</b>
SL+SF	70.3%	53.0%	<b>55.5%</b>	53.9%	52.4%
All	70.8%	53.5%	<b>55.8%</b>	54.5%	52.9%
All+	72.7%	56.5%	<b>59.4%</b>	57.9%	56.1%
All+, 180k steps	76.2%	62.3%	<b>65.7%</b>	65.0%	64.0%

Table 3. Experimental results exploring the usage of different distillation heads trained for 60k steps. Here “*Base*” is the original experiment with  $\Delta = 1$  and conventional heads as described in Sec. 4.2.2; “*SL*” adds same-level heads to distillation targets; “*SF*” adds the distilled head (“self”) as a potential target; “*All*” combines same-level and “self” heads and  $\Delta = 2$  (each step distilling to two other clients), “*All+*” is the same as *All*, but also uses the entire ImageNet as the public dataset.

Public DS fraction	10%	20%	30%	All
main head $\beta_{\text{priv}}$	70.1%	71.1%	70.9%	71.9%
last aux head $\beta_{\text{sh}}$	52.4%	53.9%	54.1%	55.3%

Table 4. The dependence of the main head “private” accuracy  $\beta_{\text{priv}}$  and the “shared” accuracy of the 4<sup>th</sup> auxiliary head on the size of the public dataset (fraction of ImageNet training set). Experiments were conducted for a system of 8 clients with 4 auxiliary heads,  $s = 100$ ,  $\nu_{\text{emb}} = 1$ ,  $\nu_{\text{aux}} = 3$  and 250 randomly assigned “private” labels. Private training samples were drawn from 70% of the ImageNet training set in all experiments. “*All*” column shows the accuracy attained by using the entire ImageNet training set as a public dataset (while still using only 70% of it as private data).

$\nu_{\text{emb}}$	$\nu_{\text{aux}}$	$\beta_{\text{priv}}^{(m)}$	$\beta_{\text{sh}}^{(m)}$	$\beta_{\text{priv}}^{(\text{aux})}$	$\beta_{\text{sh}}^{(\text{aux})}$
0.0	0.0	46.3%	46.3%	0.1%	0.1%
	1.0	52.2%	52.0%	56.0%	56.0%
	3.0	54.1%	53.5%	57.3%	57.0%
	10.0	54.3%	54.1%	57.1%	57.1%
1.0	0.0	48.5%	48.5%	0.1%	0.1%
	1.0	53.6%	53.5%	57.3%	57.2%
	3.0	55.4%	55.2%	58.6%	58.5%
	10.0	54.1%	53.6%	56.9%	56.3%
3.0	0.0	48.3%	48.0%	0.1%	0.1%
	1.0	54.3%	54.0%	58.2%	57.6%
	3.0	55.7%	55.5%	59.3%	58.8%
	10.0	53.3%	53.4%	56.5%	56.3%

Table 5. Results for 8-client experiments with 250 random classes per client,  $s = 0$  and a varying values of  $\nu_{\text{emb}}$  and  $\nu_{\text{aux}}$ .

$\nu_{\text{emb}}$	$\nu_{\text{aux}}$	$\beta_{\text{priv}}^{(m)}$	$\beta_{\text{sh}}^{(m)}$	$\beta_{\text{priv}}^{(\text{aux})}$	$\beta_{\text{sh}}^{(\text{aux})}$
0.0	0.0	68.0%	25.2%	0.1%	0.1%
	1.0	70.6%	29.1%	70.5%	42.0%
	3.0	70.9%	30.0%	70.1%	43.3%
	10.0	68.0%	25.3%	66.0%	39.7%
1.0	0.0	69.0%	26.0%	0.1%	0.1%
	1.0	71.8%	29.8%	71.5%	43.0%
	3.0	72.0%	29.9%	71.0%	44.1%
	10.0	66.8%	23.0%	65.1%	37.5%
3.0	0.0	65.2%	24.9%	0.1%	0.1%
	1.0	71.7%	29.7%	72.1%	39.1%
	3.0	71.8%	29.7%	71.9%	40.8%
	10.0	65.4%	23.1%	63.4%	36.4%

Table 6. Results for 8-client experiments with 250 random classes per client,  $s = 100$  and a varying values of  $\nu_{\text{emb}}$  and  $\nu_{\text{aux}}$ .

Heads	1	2	3	4
$\beta_{\text{priv}}^{(m)}$	56.2%	56.1%	55.8%	55.9%
$\beta_{\text{sh}}^{(m)}$	56.1%	55.8%	55.8%	55.5%
$\beta_{\text{priv}}^{(1)}$	59.6%	59.6%	59.4%	59.4%
$\beta_{\text{sh}}^{(1)}$	59.4%	59.5%	59.6%	59.4%
$\beta_{\text{priv}}^{(2)}$		60.0%	59.7%	59.7%
$\beta_{\text{sh}}^{(2)}$		59.7%	59.9%	59.5%
$\beta_{\text{priv}}^{(3)}$			59.5%	59.1%
$\beta_{\text{sh}}^{(3)}$			59.5%	59.1%
$\beta_{\text{priv}}^{(4)}$				58.7%
$\beta_{\text{sh}}^{(4)}$				58.6%

Table 7. Results for 8-client experiments with 250 random classes per client,  $s = 0$ ,  $\nu_{\text{emb}} = 1$ ,  $\nu_{\text{aux}} = 3$  and a varying number of auxiliary heads (separate columns).

Heads	1	2	3	4
$\beta_{\text{priv}}^{(\text{m})}$	72.5%	71.6%	71.1%	72.5%
$\beta_{\text{sh}}^{(\text{m})}$	30.5%	32.5%	33.1%	32.7%
$\beta_{\text{priv}}^{(1)}$	71.4%	70.6%	70.7%	71.4%
$\beta_{\text{sh}}^{(1)}$	44.7%	46.6%	46.9%	46.4%
$\beta_{\text{priv}}^{(2)}$		68.5%	68.1%	68.7%
$\beta_{\text{sh}}^{(2)}$		51.6%	52.0%	51.6%
$\beta_{\text{priv}}^{(3)}$			66.1%	66.1%
$\beta_{\text{sh}}^{(3)}$			53.8%	53.6%
$\beta_{\text{priv}}^{(4)}$				63.4%
$\beta_{\text{sh}}^{(4)}$				54.5%

Table 8. Results for 8-client experiments with 250 random classes per client,  $s = 100$ ,  $\nu_{\text{emb}} = 1$ ,  $\nu_{\text{aux}} = 3$  and a varying number of auxiliary heads (separate columns).