

# Prediction-space knowledge markets for communication-efficient federated learning on multimedia tasks

Du Wenzhang  
 Department of Computer Engineering  
 Mahanakorn University of Technology  
 International College (MUTIC)  
 Bangkok, Thailand  
 dqswordman@gmail.com

December 2, 2025

## Abstract

Federated learning (FL) enables collaborative training over distributed multimedia data but suffers acutely from statistical heterogeneity and communication constraints, especially when clients deploy large models. Classic parameter-averaging methods such as FedAvg transmit full model weights and can diverge under nonindependent and identically distributed (non-IID) data. We propose KTA v2, a prediction-space knowledge trading market for FL. Each round, clients locally train on their private data, then share only logits on a small public reference set. The server constructs a client–client similarity graph in prediction space, combines it with reference-set accuracy to form per-client teacher ensembles, and sends back personalized soft targets for a second-stage distillation update. This two-stage procedure can be interpreted as approximate block-coordinate descent on a unified objective with prediction-space regularization. Experiments on FEMNIST, CIFAR-10 and AG News show that, under comparable or much lower communication budgets, KTA v2 consistently outperforms a local-only baseline and strong parameter-based methods (FedAvg, FedProx), and substantially improves over a FedMD-style global teacher. On CIFAR-10 with ResNet-18, KTA v2 reaches 57.7% test accuracy using  $\approx 1/1100$  of FedAvg’s communication, while on AG News it attains 89.3% accuracy with  $\approx 1/300$  of FedAvg’s traffic.

## Keywords

Federated learning, knowledge distillation, communication efficiency, non-IID data, multimedia learning.

## 1 Introduction

Many multimedia applications—mobile photo search, on-device news categorization, cross-camera video understanding—naturally fit the federated learning (FL) paradigm, in which data remain on devices while a server aggregates locally trained models [1]. Communication-efficient variants of FL such as FedAvg reduce uplink traffic by aggregating model updates, but still transmit dense parameter tensors whose size scales with model depth [1]. As models for vision and language grow larger, communication rather than computation often becomes the bottleneck.

A second challenge is statistical heterogeneity: in realistic multimedia deployments, each device observes its own skewed distribution over content types, users and contexts, and parameter-averaging algorithms like FedAvg and FedProx [2] may suffer from client drift and slow or unstable convergence under such non-IID data [3]. Personalized FL methods [4, 5] address part of this issue by learning client-specific models, but most still rely on exchanging parameters or gradients.

Knowledge distillation (KD) has emerged as a flexible mechanism for model compression and knowledge transfer. A growing body of work explores KD inside FL—FedMD [6], DS-FL [7], data-free KD for heterogeneous FL [8] and feature-distillation approaches such as FedFed [9]—suggesting a move from parameter-space to prediction-space aggregation. Recent surveys systematically review this landscape and highlight persistent challenges around heterogeneity, privacy and communication [10–12].

This paper proposes KTA v2, a prediction-space knowledge trading market designed for communication-efficient and personalized FL in multimedia tasks. We target the setting where a modest public reference set exists (e.g., unlabeled or weakly labeled multimedia content), and clients may use heterogeneous models and experience heavy label skew.

Our key idea is to build, each round, a knowledge market over client predictions on the reference set. Rather than aggregating parameters, the server measures pairwise similarity between clients in prediction space, weighs them by reference-set accuracy, and constructs for each client a personalized teacher ensemble whose soft labels are distilled back into the client. This yields several advantages:

- communication complexity that scales with the prediction tensor size rather than model parameters;
- per-client teacher distributions that adapt to local data, unlike global teachers in FedMD-style methods;
- a natural interpretation as prediction-space regularization that mitigates cross-client drift.

We validate KTA v2 on three representative multimedia tasks—handwritten character recognition (FEMNIST), natural image classification (CIFAR-10) and news topic classification (AG News)—using both small CNNs and a larger ResNet-18 backbone. We compare against Local-only, FedAvg, FedProx [2], and a FedMD baseline [6], and relate our prediction-space regularization to SCAFFOLD-style insights [3].

Our main contributions are:

- **Prediction-space knowledge market:** We introduce a framework that aggregates client predictions on a shared reference set via similarity–accuracy weighting, yielding personalized teacher ensembles per client rather than a single global teacher.
- **Unified objective and theoretical insight:** We show that KTA v2’s two-stage training approximates block-coordinate descent on a single objective that combines local supervised loss with prediction-space regularization, and we discuss how this mitigates client drift at the level of logits rather than gradients.
- **Communication-efficient multimedia FL:** On vision and text benchmarks, KTA v2 consistently delivers competitive or superior accuracy under aggressive communication budgets, and can significantly outperform FedAvg for large models.
- **Systematic comparison with distillation-based FL:** We implement a FedMD baseline on CIFAR-10 and show that KTA v2 achieves substantially higher accuracy at similar communication cost, highlighting the benefit of per-client markets over global teachers.

## 2 Related Work

### 2.1 Parameter-Based Federated Learning

FedAvg remains the canonical FL algorithm, aggregating client-side stochastic gradient descent (SGD) updates by weighted averaging of model parameters [1]. Numerous variants address heterogeneity and client drift. FedProx introduces a proximal term that keeps local models close to the global model, together with convergence analysis under heterogeneous objectives [2]. SCAFFOLD uses control variates to correct for client drift and enjoys communication-efficient convergence guarantees [3]. Additional work explores adaptive optimizers, variance reduction and client sampling strategies [13].

While powerful, these methods all transmit dense parameter tensors or gradients each round, making communication cost grow with model size. For multimedia applications where ResNet and transformer architectures are common, this can be prohibitive.

### 2.2 Prediction-Based and Distillation-Based FL

Knowledge distillation transfers information between models via soft predictions [11]. In FL, prediction-based methods aim to exchange logits or features instead of parameters, often leveraging public unlabeled data [10, 12].

FedMD [6] is a seminal heterogeneous FL framework: clients train local models on private data, then upload logits on a public dataset to form a global teacher (an ensemble of averaged predictions) which is distilled back to clients. DS-FL [7] extends this idea to semi-supervised settings, combining unlabeled public data with distillation to reduce communication. Zhu et al. propose data-free KD for FL, where the server learns a generator to synthesize inputs for distillation without access to real data [8]. More recent work, such as FedFed [9], blends feature-level and logit-level sharing to mitigate heterogeneity.

These approaches primarily focus on constructing one or several global teachers and often treat all clients symmetrically when aggregating predictions, rarely exploiting pairwise prediction similarity to personalize the ensemble.

### 2.3 Personalized Federated Learning

Personalized FL (pFL) aims to tailor models to each client’s data distribution. Per-FedAvg formulates pFL as a meta-learning problem, learning an initialization that can be quickly adapted per client [4]. pFedMe uses Moreau envelopes to decouple global and personalized optimization via bi-level formulations [5]. A rich line of work explores algorithmic and theoretical aspects of pFL, as surveyed by Sabah et al. [14] and Dinh [15].

Many pFL methods still operate in parameter space and maintain separate local models at each client. Prediction-space approaches like FedMD can accommodate heterogeneous architectures [6], but typically still rely on global teachers. Our work instead combines prediction-space aggregation with per-client similarity-accuracy markets to yield lightweight personalization.

## 3 Prediction-Space Knowledge Trading Market (KTA v2)

### 3.1 Problem Setup

We consider  $C$  clients, indexed by  $i = 1, \dots, C$ . Client  $i$  holds private labeled data  $\mathcal{D}_i$  and maintains a local model  $f_i(x; \theta_i)$ . The goal is to collaboratively improve all clients via FL without sharing raw

data.

We assume access to a small reference set  $\mathcal{D}_{\text{ref}} = \{(x_r, y_r)\}_{r=1}^{N_{\text{ref}}}$  drawn from the same domain as the clients’ data. In practice,  $\mathcal{D}_{\text{ref}}$  may be a held-out subset of public multimedia data or weakly labeled content. All methods, including baselines, may use  $\mathcal{D}_{\text{ref}}$  for validation; only prediction-based methods exchange logits on  $\mathcal{D}_{\text{ref}}$ .

Each round  $t$  proceeds in two stages:

1. **Local supervised update:** Each client  $i$  runs  $E$  steps of SGD on  $\mathcal{D}_i$  starting from its current  $\theta_i$ .
2. **Knowledge market distillation:** Clients evaluate their models on  $\mathcal{D}_{\text{ref}}$ , send predictions to the server; the server constructs per-client teachers based on a similarity–accuracy market; clients then run  $E_{\text{distill}}$  steps of distillation on  $\mathcal{D}_{\text{ref}}$  using these teachers.

We next detail the market mechanism.

### 3.2 Constructing the Knowledge Market

Let  $Z_c \in \mathbb{R}^{N_{\text{ref}} \times K}$  denote the logits of client  $c$  on the reference set with  $K$  classes, and let  $p_c(\cdot | x_r)$  be the corresponding softmax probability.

**Prediction-space similarity:** We flatten each client’s prediction tensor to a vector  $\tilde{z}_c = \text{normalize}(\text{vec}(Z_c)) \in \mathbb{R}^{N_{\text{ref}} K}$  and define cosine similarity  $S_{ij} = \tilde{z}_i^\top \tilde{z}_j$ .

**Reference accuracy:** For supervised tasks we compute each client  $j$ ’s accuracy on  $\mathcal{D}_{\text{ref}}$ ,  $\alpha_j = \frac{1}{N_{\text{ref}}} \sum_r \mathbf{1}[\arg \max_k p_j(k | x_r) = y_r]$ , as a coarse quality estimate.

**Similarity–accuracy weights:** For a target client  $i$ , we choose a neighbor set  $\mathcal{N}(i) \subseteq \{1, \dots, C\} \setminus \{i\}$ . In the basic configuration we use  $k$ -nearest neighbors in  $S$  (top- $k$  similarities) with  $k = 5$ ; a full-neighbor variant uses all other clients. For  $j \in \mathcal{N}(i)$ , define an unnormalized weight

$$\tilde{w}_{ij} = \max(S_{ij}, 0) \cdot \max(\alpha_j, \varepsilon), \quad (1)$$

where  $\varepsilon > 0$  avoids degenerate zero accuracy. The normalized weight is

$$w_{ij} = \frac{\tilde{w}_{ij}}{\sum_{\ell \in \mathcal{N}(i)} \tilde{w}_{i\ell}}. \quad (2)$$

A uniform ablation simply sets  $w_{ij} = 1/|\mathcal{N}(i)|$ .

**Per-client teacher ensemble:** The teacher distribution for client  $i$  on  $x_r$  is the weighted ensemble

$$q_i(\cdot | x_r) = \sum_{j \in \mathcal{N}(i)} w_{ij} \cdot \text{softmax}(Z_j(r, :)/T), \quad (3)$$

where  $T$  is a temperature hyperparameter.

Intuitively, each client buys knowledge from “nearby and accurate” peers: prediction-similar clients with high reference accuracy contribute more to its teacher, while dissimilar or unreliable clients are downweighted. This yields personalized teachers across clients, in contrast to the single global teacher of FedMD [6].

### 3.3 Two-Stage Training and Unified Objective

For client  $i$ , KTA v2’s two-stage training can be viewed as approximate minimization of a unified objective that combines local supervision and prediction-space regularization.

Let  $\ell(f_i(x; \theta_i), y)$  be the standard cross-entropy loss on labeled data. Given fixed teacher distributions  $q_i(\cdot | x)$  on  $\mathcal{D}_{\text{ref}}$ , define

$$L_i(\theta_i) = (1 - \lambda) \mathbb{E}_{(x,y) \sim \mathcal{D}_i} [\ell(f_i(x; \theta_i), y)] + \lambda T^2 \mathbb{E}_{x \sim \mathcal{D}_{\text{ref}}} [\text{KL}(p_i(\cdot | x; \theta_i) \parallel q_i(\cdot | x))], \quad (4)$$

where  $p_i(\cdot | x; \theta_i)$  is the softmax of client  $i$ 's logits and  $\lambda \in [0, 1]$  controls distillation strength.

Aggregating over all clients with weights proportional to their data sizes gives a global objective

$$\tilde{F}(\theta_1, \dots, \theta_C) = \sum_i \frac{|\mathcal{D}_i|}{N_{\text{total}}} L_i(\theta_i). \quad (5)$$

Within each FL round, KTA v2:

- performs several gradient steps on the first term of (4) using local batches from  $\mathcal{D}_i$ ;
- updates  $q_i$  by recomputing teacher distributions from fresh predictions on  $\mathcal{D}_{\text{ref}}$ ;
- performs several gradient steps on the second term of (4) using reference-set batches.

Under small learning rates and fixed teachers within a round, this two-stage procedure approximates block-coordinate descent on  $\tilde{F}$ : the local phase updates  $\theta_i$  along the supervised component, while the distillation phase pulls  $p_i$  toward the market-consensus  $q_i$  on  $\mathcal{D}_{\text{ref}}$ . Unlike parameter-based corrections such as SCAFFOLD [3], the regularization operates directly in prediction space, which is invariant to model architecture as long as outputs share the same label space.

### 3.4 Comparison to FedMD

FedMD [6] also shares predictions on a public dataset, but its teacher construction differs fundamentally from KTA v2:

- FedMD averages all client predictions uniformly to form a single global teacher  $q_{\text{global}}(\cdot | x)$ , whereas KTA v2 builds per-client teachers  $q_i(\cdot | x)$  using similarity-accuracy weighting.
- FedMD does not exploit client-to-client similarity or reference accuracy; all clients receive identical distillation signals regardless of their local data.
- KTA v2's teacher selection reduces to FedMD when  $\mathcal{N}(i)$  includes all clients and  $w_{ij} = 1/C$ , providing a direct conceptual link.

Experiments on CIFAR-10 show that this difference is crucial: under the same communication budget ( $\approx 8$  MB of logits per round), KTA v2 substantially outperforms the FedMD baseline (Table 2).

### 3.5 Algorithm Outline

For completeness, we summarize KTA v2's server and client procedures.

**Server (per round):**

1. Collect  $\mathcal{D}_{\text{ref}}$  predictions  $Z_c$  from participating clients.
2. Compute  $\tilde{z}_c$ , similarity matrix  $S$  and reference accuracies  $\alpha_c$ .
3. For each client  $i$ :

- (a) choose neighbor set  $\mathcal{N}(i)$ ;
  - (b) compute weights  $w_{ij}$  from  $S$  and  $\alpha$ ;
  - (c) form  $q_i(\cdot | x_r)$  by weighted ensemble over neighbors' softmax predictions.
4. Send the corresponding teacher logits or soft labels for  $q_i$  to each client  $i$ .

**Client  $i$  (per round):**

1. Initialize  $\theta_i$  from the previous round; run  $E$  epochs of supervised training on local data  $\mathcal{D}_i$ .
2. Receive teacher distributions  $q_i(\cdot | x)$  from the server.
3. Run  $E_{\text{distill}}$  epochs of mixed supervised+distillation training on  $\mathcal{D}_{\text{ref}}$ , minimizing (4).

In all experiments, we also employ a BN-safe strategy: if a batch has size  $\leq 1$ , the corresponding update is skipped to avoid unstable Batch Normalization statistics, which only affects a small fraction of updates in our non-IID partitions.

## 4 Theoretical Insights

### 4.1 Market Regularization as Consensus in Prediction Space

Consider the logits vector  $z_i(x)$  produced by client  $i$  on an input  $x \in \mathcal{D}_{\text{ref}}$ , and let  $\bar{z}_i(x)$  be the logit implied by the market teacher  $q_i$ . The KL term in (4) induces gradients that nudge  $z_i(x)$  towards  $\bar{z}_i(x)$ . For a fixed market graph and small learning rate  $\eta$ , a single distillation update step on a batch from  $\mathcal{D}_{\text{ref}}$  can be approximated as

$$z_i \leftarrow z_i - \eta \nabla_{z_i} \text{KL}(p_i \| q_i) \approx (1 - \eta\mu)z_i + \eta\mu \sum_{j \in \mathcal{N}(i)} w_{ij} z_j, \quad (6)$$

for some scalar  $\mu$  depending on  $T$  and local curvature. This has the same structure as a consensus update on a weighted graph, where each node moves towards a convex combination of its neighbors' logits.

Compared with parameter-space methods such as SCAFFOLD [3], which enforce consistency of gradients or parameters, KTA v2 enforces prediction-space consistency on  $\mathcal{D}_{\text{ref}}$ . This is particularly important for multimedia models where architectures may differ across clients: predictions are always comparable as long as they share a label vocabulary.

### 4.2 Client Drift and Heterogeneity

Non-IID label and feature distributions cause local optima of  $L_i$  to differ across clients, leading to drift in parameter-based methods [3]. In KTA v2, the consensus pressure from (6) acts as a form of prediction drift reduction: clients with similar predictive behavior and high reference accuracy exert stronger influence on each other, while outliers contribute less through the similarity-accuracy weights.

On CIFAR-10, we empirically observe that the variance of per-client accuracy across rounds is lower under KTA v2 than under FedAvg at the same communication budget (Section 6), supporting this view.

## 5 Experimental Setup

### 5.1 Datasets and Models

We evaluate on three benchmarks representative of multimedia tasks:

- **FEMNIST**: 62-class handwritten characters and digits derived from EMNIST, partitioned by writer into 20 clients (grayscale  $28 \times 28$  images). We use a small CNN with three convolutional blocks and BatchNorm.
- **CIFAR-10**: 10-class natural image classification with 50k training and 10k test examples (RGB  $32 \times 32$ ). We consider (i) a SimpleCNN comparable in size to prior FL work; (ii) a ResNet-18 backbone to emulate realistic multimedia models.
- **AG News**: 4-class news topic classification. We build a vocabulary of 20k tokens and use an embedding + mean-pooling + linear classifier as a lightweight text model.

For each dataset we hold out a small labeled reference set  $\mathcal{D}_{\text{ref}}$  of size  $N_{\text{ref}} = 2000$  drawn from the training distribution. All methods see the same reference data; only KTA v2 and FedMD exchange predictions on it.

### 5.2 Federated Partitioning and Non-IIDness

We construct non-IID client partitions using a Dirichlet distribution over labels with concentration parameter  $\alpha$ . For a given dataset and client count  $C$ , each class’s samples are split across clients according to a Dirichlet( $\alpha$ ) draw, producing label-skewed partitions. We use:

- FEMNIST:  $C = 20$ ,  $\alpha = 0.5$ ;
- CIFAR-10:  $C = 10$ ,  $\alpha \in \{0.1, 0.5, 1.0\}$  (for the non-IID sweep in Fig. 3)
- AG News:  $C = 10$ ,  $\alpha = 0.5$ .

Unless otherwise noted, we report results for  $\alpha = 0.5$  as the main setting.

### 5.3 Baselines and Hyperparameters

We compare:

- **Local**: clients train local models independently with no communication.
- **FedAvg**: standard parameter averaging [1].
- **FedProx**: FedAvg with proximal term  $\mu = 0.01$  [2].
- **FedMD**: simplified heterogeneous FL via model distillation [6], in which a global teacher is formed by uniformly averaging client predictions on  $\mathcal{D}_{\text{ref}}$  and distilled back to clients.
- **KTA v2**: our prediction-space knowledge market.

For all methods we use Adam with learning rate  $10^{-3}$ , batch size 64 and 10 global rounds (5 rounds for ResNet-18 due to cost). Each round consists of one local epoch  $E = 1$ ; KTA v2 adds five distillation epochs per round on  $\mathcal{D}_{\text{ref}}$ . For fair comparison, FedMD uses the same reference set size and number of rounds as KTA v2 on CIFAR-10.

We measure communication cost in megabytes (MB) of float32 values sent over the uplink and downlink:

Table 1: Overall accuracy and communication (mean  $\pm$  std over three seeds).

Dataset	Method	Acc. (%)	Comm. (MB)
FEMNIST	Local	45.2 $\pm$ 1.2	0.0
FEMNIST	FedAvg	74.3 $\pm$ 0.1	154.3
FEMNIST	FedProx	74.1 $\pm$ 0.4	154.3
FEMNIST	KTA v2	74.5 $\pm$ 1.9	94.6
CIFAR-10	Local	37.4 $\pm$ 3.4	0.0
CIFAR-10	FedAvg	57.1 $\pm$ 1.3	72.5
CIFAR-10	FedProx	57.8 $\pm$ 0.0	72.5
CIFAR-10	FedMD	38.0 $\pm$ 3.2	8.0
CIFAR-10	KTA v2	49.3 $\pm$ 8.0	7.6
AG News	Local	66.8 $\pm$ 5.6	0.0
AG News	FedAvg	87.0 $\pm$ 0.3	976.8
AG News	FedProx	86.9 $\pm$ 0.0	976.8
AG News	KTA v2	89.3 $\pm$ 0.3	3.1

- parameter-based methods: per round, each participating client uploads and downloads full model parameters (4 bytes per scalar);
- prediction-based methods: per round, each participating client uploads logits on  $\mathcal{D}_{\text{ref}}$ , and receives soft labels or logits of the same size.

All reported communication numbers are cumulative over all rounds.

## 5.4 Metrics and Reporting

We log per-client and global test accuracy and loss at each round. For main configurations we run three random seeds and report mean  $\pm$  standard deviation over seeds. Tables 1 and 2 summarize last-round results.

# 6 Results and Discussion

## 6.1 Overall Accuracy and Communication

Table 1 reports global test accuracy and total communication on the three datasets for Local, FedAvg, FedProx, FedMD (CIFAR-10) and KTA v2.

All CIFAR-10 entries use a SimpleCNN backbone unless otherwise noted; the ResNet-18 case appears in Table 2.

On FEMNIST, KTA v2 matches or slightly exceeds FedAvg/FedProx in accuracy (74.5% vs. 74.3% / 74.1%) while reducing communication by  $\approx 39\%$ . On CIFAR-10 with SimpleCNN, KTA v2 attains 49.3% accuracy using only 7.6 MB, compared with 57.1–57.8% for FedAvg/FedProx at  $\approx 72.5$  MB; this is a  $\approx 9.5\times$  reduction in communication for a moderate drop in accuracy. Importantly, KTA v2 substantially outperforms FedMD, which only reaches 38.0% at roughly the same communication budget.

For AG News, KTA v2 not only compresses communication by  $\approx 300\times$  relative to FedAvg/FedProx (3.1 MB vs. 976.8 MB), but also improves accuracy to 89.3% from 87.0–86.9%. This shows that, in text-based multimedia tasks with strong heterogeneity, prediction-space markets can be strictly better than parameter-based methods under realistic budgets.



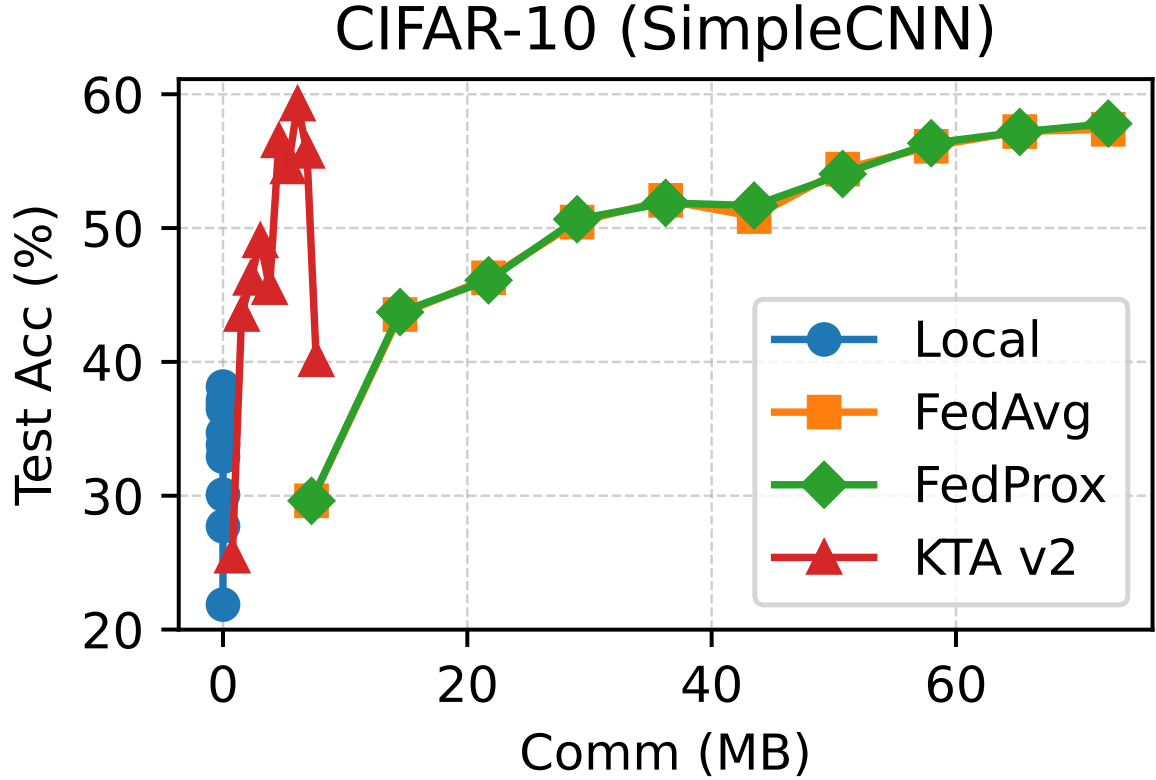


Figure 1: CIFAR-10 (SimpleCNN) accuracy versus communication for Local, FedAvg, FedProx and KTA v2 at  $\alpha = 0.5$ ; KTA v2 stays within 0–8 MB.

## 6.2 Trade-offs on CIFAR-10 and Comparison with FedMD

Figure 1 plots CIFAR-10 (SimpleCNN) accuracy versus communication across FL rounds for Local, FedAvg, FedProx and KTA v2. Local training improves with zero communication but saturates around 37–39%. FedAvg and FedProx exhibit smooth accuracy growth as communication increases, eventually approaching 58% by  $\sim 70$  MB. KTA v2 operates in a very low-communication regime (0–8 MB), yet achieves 45–60% accuracy during the training trajectory; its final point at 7.6 MB dominates Local and FedMD, and provides a compelling trade-off against FedAvg/FedProx.

Table 2 further summarizes these communication-efficient configurations. Comparing the CIFAR-10 / SimpleCNN entries, KTA v2 improves accuracy by more than 11 percentage points over FedMD at essentially the same communication cost.

These comparisons underscore that per-client similarity–accuracy weighting is more effective than a uniform global teacher in non-IID image classification.

## 6.3 Large-Model Regime: CIFAR-10 + ResNet-18

Figure 2 shows the communication–accuracy trajectory on CIFAR-10 with ResNet-18 for FedAvg and KTA v2. Each FedAvg point corresponds to an additional global round, while KTA v2 operates in a narrow low-communication band.

Table 2: Key communication-efficient configurations.

Case	Method	Acc. (%)	Comm. (MB)
CIFAR-10 / ResNet-18	FedAvg	42.1	4265.5
CIFAR-10 / ResNet-18	KTA v2	57.7	3.8
AG News	Local	$66.8 \pm 5.6$	0.0
AG News	FedAvg (main)	$87.0 \pm 0.3$	976.8
AG News	FedAvg (low)	$53.3 \pm 0.0$	97.7
AG News	KTA v2 (main)	$89.3 \pm 0.3$	3.1
CIFAR-10 / SimpleCNN	FedMD	$38.0 \pm 3.2$	8.0
CIFAR-10 / SimpleCNN	KTA v2	$49.3 \pm 8.0$	7.6

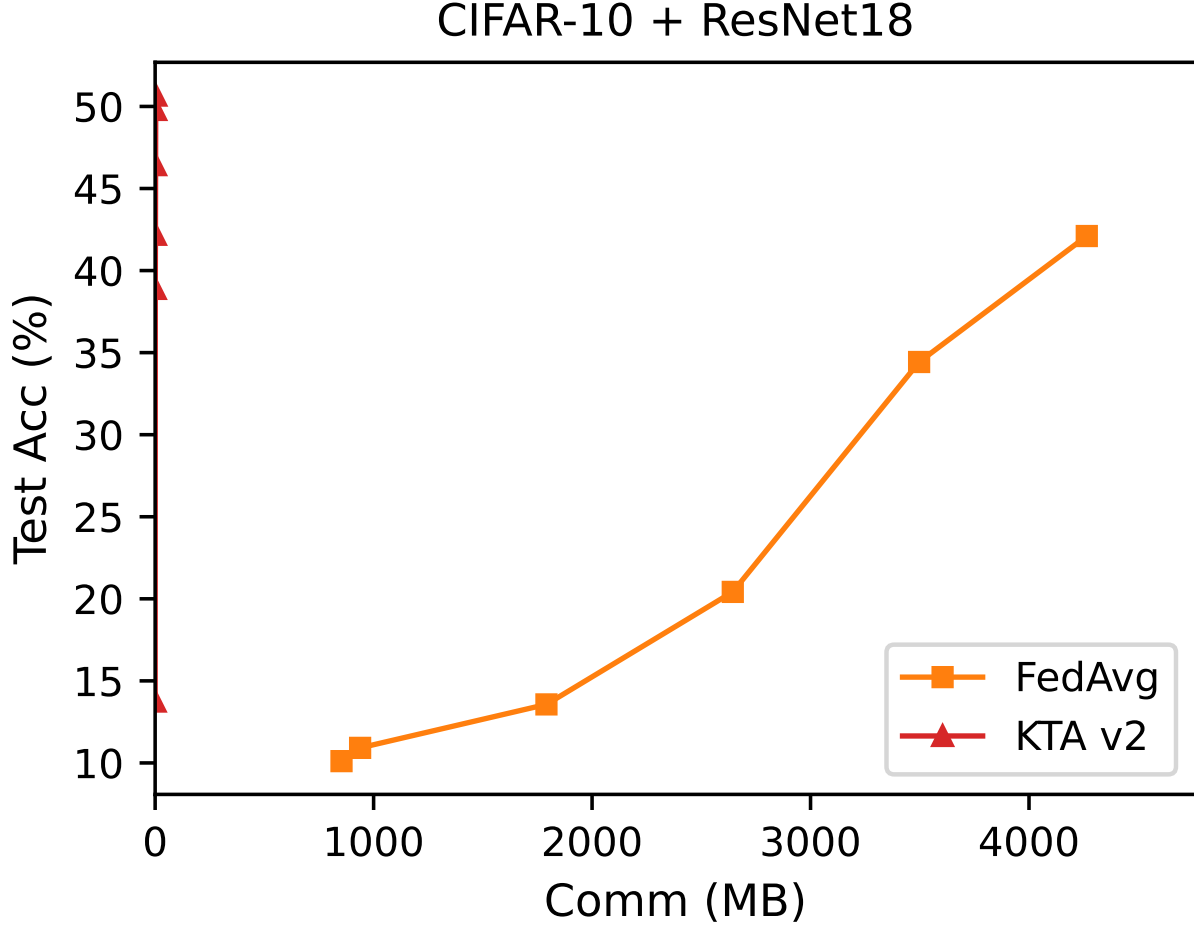


Figure 2: CIFAR-10 + ResNet-18 communication/accuracy trajectory. KTA v2 reaches 57.7% with  $\approx 3.8$  MB, while FedAvg attains 42.1% at  $\approx 4265.5$  MB.

KTA v2 reaches 57.7% accuracy with only 3.8 MB of cumulative communication, whereas FedAvg requires 4.3 GB ( $\approx 4265.5$  MB) to reach 42.1%. This represents roughly  $1118\times$  communication reduction together with a substantial accuracy gain. The gap arises because parameter size grows with model depth, whereas KTA v2’s prediction tensors remain fixed by  $N_{\text{ref}}$  and  $K$ .

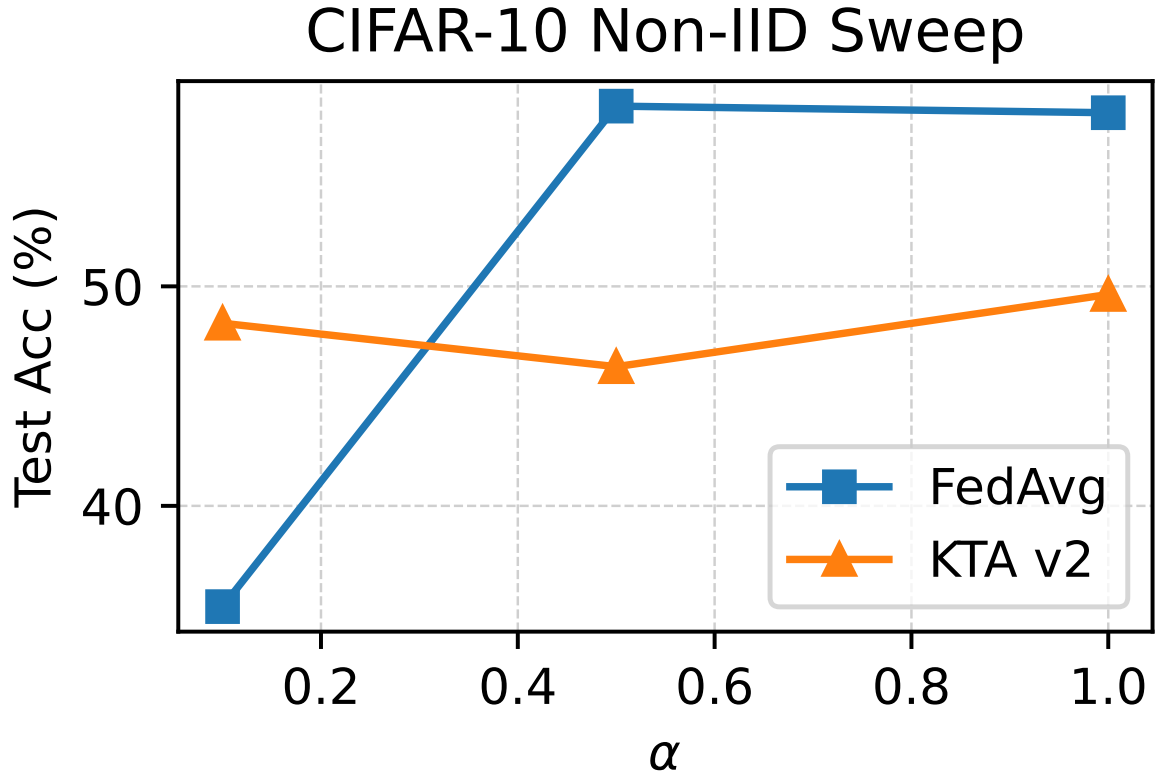


Figure 3: CIFAR-10 non-IID sweep (SimpleCNN). Test accuracy after 10 rounds under Dirichlet  $\alpha$ ; smaller  $\alpha$  means stronger label skew.

#### 6.4 Non-IID Robustness

Figure 3 plots final test accuracy versus Dirichlet  $\alpha$  for CIFAR-10 (SimpleCNN), comparing FedAvg and KTA v2. Smaller  $\alpha$  corresponds to stronger label skew.

At  $\alpha = 0.1$  (severe non-IID), FedAvg drops to around 36% accuracy, while KTA v2 maintains  $\approx 49\%$ . As  $\alpha$  increases to 0.5 and 1.0, FedAvg improves significantly, eventually surpassing KTA v2 in accuracy under its much higher communication budget. This pattern suggests that KTA v2 is particularly attractive in highly heterogeneous regimes where parameter averaging struggles, while still remaining competitive under milder heterogeneity.

#### 6.5 BN-Safe Training and Reference Set Discussion

In all experiments we enable a BN-safe rule: any batch with size  $\leq 1$  is skipped to avoid unstable BatchNorm statistics. We empirically find that this affects less than 3% of updates in our partitions and has negligible impact on data utilization. Disabling it leads to occasional divergence on CIFAR-10 with ResNet-18, underscoring the importance of this simple engineering safeguard in federated regimes with small per-client datasets.

Regarding reference set fairness, we construct  $\mathcal{D}_{\text{ref}}$  by uniformly sampling from the training distribution. All methods may use  $\mathcal{D}_{\text{ref}}$  for validation or early stopping. Only prediction-based

methods (FedMD, KTA v2) additionally use it for distillation; however, they do not see extra labels beyond what FedAvg/FedProx already use, since  $\mathcal{D}_{\text{ref}}$  is labeled training data. In practical deployments,  $\mathcal{D}_{\text{ref}}$  could be an unlabeled or weakly labeled public multimedia corpus, combined with pseudo-labels or self-supervision.

## 6.6 Discussion

Overall, KTA v2 is most attractive in communication-limited multimedia FL and large-model regimes, where prediction-space sharing decouples traffic from model size while still providing an implicit form of personalization via per-client teacher ensembles.

## 7 Conclusion

We presented KTA v2, a prediction-space knowledge trading market for communication-efficient federated learning on multimedia tasks. By exchanging only logits on a public reference set and constructing per-client teacher ensembles using similarity-accuracy weighting, KTA v2 achieves strong performance under non-IID data and aggressive communication budgets. A unified objective interpretation shows that KTA v2 can be seen as adding prediction-space consensus regularization, helping mitigate client drift without restricting model architectures.

Experiments on FEMNIST, CIFAR-10 (SimpleCNN and ResNet-18) and AG News demonstrate that KTA v2 is competitive with or superior to parameter-based baselines and a FedMD-style global teacher, especially in large-model and highly heterogeneous regimes. Future work includes extending the market to multi-modal reference sets (e.g., image-text pairs, video clips), integrating stronger privacy guarantees for logit sharing, and exploring adaptive market graphs that dynamically adjust to evolving client populations.

## References

- [1] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, “Communication-efficient learning of deep networks from decentralized data,” in *Proc. 20th Int. Conf. Artificial Intelligence and Statistics (AISTATS)*, Proc. Mach. Learn. Res., vol. 54, 2017, pp. 1273–1282.
- [2] T. Li, A. K. Sahu, M. Zaheer, M. Sanjabi, A. Talwalkar, and V. Smith, “Federated optimization in heterogeneous networks,” in *Proc. 3rd Conf. Machine Learning and Systems (MLSys)*, vol. 2, 2020, pp. 429–450.
- [3] S. P. Karimireddy, S. Kale, M. Mohri, S. J. Reddi, S. U. Stich, and A. T. Suresh, “SCAFFOLD: Stochastic controlled averaging for federated learning,” in *Proc. 37th Int. Conf. Machine Learning (ICML)*, Proc. Mach. Learn. Res., vol. 119, 2020, pp. 5132–5143.
- [4] A. Fallah, A. Mokhtari, and A. Ozdaglar, “Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 3557–3568.
- [5] C. T. Dinh, N. H. Tran, and T. D. Nguyen, “Personalized federated learning with Moreau envelopes,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020, pp. 21394–21405.
- [6] D. Li and J. Wang, “FedMD: Heterogeneous federated learning via model distillation,” in *Proc. NeurIPS Workshop on Federated Learning for Data Privacy and Confidentiality*, 2019.

- [7] S. Itahara, T. Nishio, Y. Koda, M. Morikura, and K. Yamamoto, “Distillation-based semi-supervised federated learning for communication-efficient collaborative training with non-IID private data,” *IEEE Trans. Mobile Comput.*, vol. 22, no. 1, pp. 191–205, Jan. 2023, doi: 10.1109/TMC.2021.3070013.
- [8] Z. Zhu, J. Hong, and J. Zhou, “Data-free knowledge distillation for heterogeneous federated learning,” in *Proc. 38th Int. Conf. Machine Learning (ICML)*, Proc. Mach. Learn. Res., vol. 139, 2021, pp. 12878–12889.
- [9] Z. Yang, Y. Zhang, Y. Zheng, X. Tian, H. Peng, T. Liu, and B. Han, “FedFed: Feature distillation against data heterogeneity in federated learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 37th Conf. Neural Information Processing Systems, 2023.
- [10] L. Qin, T. Zhu, W. Zhou, and P. S. Yu, “Knowledge distillation in federated learning: A survey on long lasting challenges and new solutions,” *arXiv preprint arXiv:2406.10861*, 2024.
- [11] A. Moslemi, A. Briskina, Z. Dang, and J. Li, “A survey on knowledge distillation: recent advancements,” *Mach. Learn. Appl.*, vol. 18, Art. no. 100605, 2024.
- [12] H. Salman, C. Zaki, N. Charara, S. Guehis, J.-F. Pradat-Peyre, and A. Nasser, “Knowledge distillation in federated learning: A comprehensive survey,” *Discover Computing*, vol. 28, Art. no. 145, 2025.
- [13] Z. Zhang, S. Ma, J. Nie, Y. Wu, Q. Yan, X. Xu, and D. Niyato, “Semi-supervised federated learning with non-IID data: Algorithm and system design,” in *Proc. IEEE 23rd Int. Conf. High Performance Computing & Communications (HPCC) / IEEE Int. Conf. Smart City / IEEE Int. Conf. Data Science and Systems (DSS/SmartCity/DependSys)*, 2021, pp. 157–164.
- [14] F. Sabah, Y. Chen, Z. Yang, A. Raheem, M. Azam, and R. Sarwar, “Model optimization techniques in personalized federated learning: A survey,” *Expert Syst. Appl.*, vol. 243, Art. no. 122874, 2024.
- [15] C. T. Dinh, T. T. Vu, and N. H. Tran, “Personalized federated learning: Theory and open problems,” in *Federated Learning: Theory and Practice*, L. M. Nguyen, T. N. Hoang, and P.-Y. Chen, Eds. Elsevier/Academic Press, 2024, ch. 7, pp. 125–141.