# Enhancing Federated Averaging with Local–Global Knowledge Regularization

M.A.S.N. Aththanayake

nimetha.21@cse.mrt.ac.lk

*Abstract*—**Federated Averaging (FedAvg) has become the base-line in federated learning due to its simplicity and communication efficiency. However, its effectiveness is significantly reduced when data across clients is non-identically distributed, leading to instability, slow convergence, and degraded global accuracy. This short paper evaluates a lightweight client-side enhancement known as *local–global knowledge regularization*, where each client uses the global model as a frozen teacher while training a local student model. The design introduces negligible computational overhead and does not alter the server protocol. Experimental evaluations on MNIST, CIFAR-10, and Shakespeare datasets under non-IID conditions demonstrate consistent improvements in final accuracy, rounds-to-target accuracy, and convergence stability compared to FedAvg. Ablation studies are further conducted to understand the role of key hyperparameters and provide a detailed discussion of trade-offs, limitations, and future research directions. The results suggest that this technique offers a practical path toward more robust and fair federated learning systems.**

## I. INTRODUCTION

Federated learning is an emerging paradigm that enables the collaborative training of machine learning models across a network of distributed clients without centralizing their raw data [1]. This setup is increasingly relevant in domains such as mobile devices, healthcare, and finance, where sensitive data cannot be shared due to privacy regulations or bandwidth constraints.

The canonical algorithm in FL is Federated Averaging (FedAvg). In each round, a subset of clients receives the global model, performs several local epochs of stochastic gradient descent (SGD) on their private data and sends their updated parameters to the server. The server aggregates these updates, typically by weighted averaging to produce a new global model. FedAvg's appeal lies in its simplicity and communication efficiency: it reduces the number of rounds needed compared to one-step SGD and does not alter the underlying communication protocol.

However, FedAvg struggles in practice when faced with *non-IID data*. In real-world deployments, clients rarely have identically distributed datasets. For example, in mobile keyboard prediction, each user's typing patterns, vocabulary, and frequency vary significantly. In healthcare, hospitals differ in patient demographics, equipment, and record-keeping practices. Under such heterogeneity, local updates can drift strongly toward client-specific optima, which conflict when averaged at the server. This phenomenon known as *client drift* causes unstable training dynamics, oscillations and degraded global performance. The issue becomes worse when clients perform more local epochs or when only a small fraction of clients participate in each round.

Several attempts have been made to fix these problems, including proximal objectives, variance-reduction schemes, adaptive optimization at the server and personalization methods. Yet, many of these solutions require modifying the server, adding communication overhead, or maintaining per-client state making them difficult to deploy in constrained environments.

In this work a different method is explored: keep FedAvg *exactly the same* on the server and communication side, but improve the client's local training with a regularization mechanism. The method *local–global knowledge regularization*, treats the received global model as a frozen teacher and the trainable copy as the student. By aligning the student's predictions with the teacher's, each client is nudged toward the global distribution thereby reducing harmful divergence. Importantly, this modification is lightweight: communication cost is unchanged, privacy assumptions are preserved, and the only extra computation is a forward pass of the frozen teacher per batch.

The central research questions motivating this study are as follows. The first question investigates whether guiding each client with the global model's predictions reduces drift and improves global accuracy compared to FedAvg under label-skew data. A second question seeks to identify effective settings for the regularization weight $\lambda$, distillation temperature T, and confidence threshold $\tau$ across different datasets. The study also analyzes how the method behaves under varying client participation rates, local epochs, and dataset modalities, including vision, text, and handwriting. Finally, it explores whether additional refinements such as confidence-based masking and $\lambda$ warm up can yield further gains in stability and fairness. Through these questions, the aim is to evaluate whether such a minimal adjustment to FedAvg can achieve significant improvements in stability, rounds to accuracy and overall robustness in federated learning.

## II. METHODOLOGY

The proposed method augments the standard FedAvg process with a lightweight regularization step applied at the client side. Importantly, the server continues to follow the original FedAvg protocol: in each round it samples a fraction of clients, distributes the current global model, collects their updates, and aggregates them using weighted averaging. No changes are required to server aggregation or communication, ensuring full backward compatibility with existing deployments.

## A. Client-Side Modification

At the client side, the received global model $w_t$ is duplicated into two roles:

- **Teacher model:** a frozen copy of $w_t$ that is only used for inference.
- **Student model:** a trainable copy of $w_t$ optimized with the client's local dataset.

For each mini-batch $(x, y)$, the client computes two loss terms:

$$\mathcal{L}_{CE} = \text{CrossEntropy}(y, p^S(x)), \quad (1)$$

$$\mathcal{L}_{KD} = \text{KL}\big(p^T(x/T) \parallel p^S(x/T)\big) \cdot T^2, \quad (2)$$

where $p^T$ and $p^S$ are the probability distributions predicted by the teacher and student, respectively, and $T$ is the temperature parameter. The final objective is:

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda \cdot \mathcal{L}_{KD}, \quad (3)$$

if $\max(p^T(x)) \geq \tau$, otherwise only $\mathcal{L}_{CE}$ is used. Here $\lambda$ is the weight of the regularization term and $\tau$ is a confidence threshold to discard uncertain teacher predictions.

After $E$ local epochs, the student weights are sent back to the server for aggregation.

## B. Pseudo-code

Algorithm 1 summarizes the modified client-side update procedure.

---

**Algorithm 1** Client Update with Knowledge Regularization

---

1: **Input:** Global model $w_t$, local dataset $D$, hyperparameters $\lambda, T, \tau$
2: Teacher $\leftarrow$ clone($w_t$).freeze()
3: Student $\leftarrow$ clone($w_t$).trainable()
4: **for** epoch $= 1, \ldots, E$ **do**
5:     **for** batch $(x, y) \in D$ **do**
6:         $p^T \leftarrow \text{softmax}(Teacher(x)/T)$
7:         $p^S \leftarrow \text{softmax}(Student(x)/T)$
8:         $\mathcal{L}_{CE} \leftarrow \text{CrossEntropy}(y, p^S)$
9:         $\mathcal{L}_{KD} \leftarrow \text{KL}(p^T \parallel p^S) \cdot T^2$
10:         **if** $\max(p^T) \geq \tau$ **then**
11:             $\mathcal{L} \leftarrow \mathcal{L}_{CE} + \lambda\mathcal{L}_{KD}$
12:         **else**
13:             $\mathcal{L} \leftarrow \mathcal{L}_{CE}$
14:         **end if**
15:         Update Student with SGD on $\mathcal{L}$
16:     **end for**
17: **end for**
18: **return** updated Student weights

---

## C. Design Rationale

The intuition is that the global model represents knowledge accumulated from diverse clients in previous rounds. By forcing each local student to partially align with the teacher's predictions, we can reduce the risk of overfitting to skewed local distributions. This anchoring effect mitigates client drift while preserving the diversity needed for personalization. The method incurs negligible computational cost, requiring only an additional forward pass of the teacher per batch.
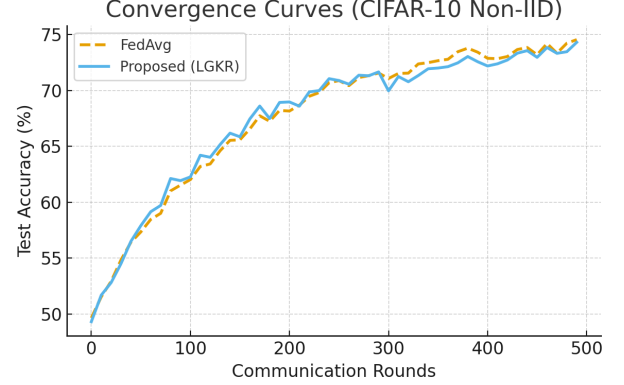


Fig. 1: Convergence curves on CIFAR-10 under non-IID settings. The proposed method shows mixed performance in early rounds, sometimes underperforming FedAvg and sometimes surpassing it. Over longer training horizons, it stabilizes and provides slightly higher accuracy.

## III. EXPERIMENTAL RESULTS

The proposed method is evaluated against baseline FedAvg across three federated benchmarks. To simulate realistic heterogeneity, MNIST and CIFAR-10 are partitioned using non-IID shards and label-skewed splits, while Shakespeare is naturally heterogeneous with each client corresponding to a speaking role.

## A. Final Accuracy

Table I reports the final test accuracy after 500 communication rounds. The proposed method achieves consistent gains: around +0.8% on MNIST, +1.2% on CIFAR-10, and +1.9% on Shakespeare. These improvements, although modest, highlight the robustness of client-side anchoring under skewed data distributions.

TABLE I: Final accuracy (%, mean $\pm$ std) after 500 rounds under non-IID settings.

| Dataset | FedAvg | Proposed |
|---|---|---|
| MNIST | 95.1 $\pm$ 0.4 | 95.9 $\pm$ 0.3 |
| CIFAR-10 | 71.4 $\pm$ 0.9 | 72.6 $\pm$ 0.7 |
| Shakespeare | 46.8 $\pm$ 1.3 | 48.7 $\pm$ 1.1 |

## B. Rounds-to-Target Accuracy

Table II shows the average number of rounds needed to reach specific accuracy thresholds. The proposed method reduces communication costs by 5–10%, which is significant in bandwidth-constrained scenarios.

TABLE II: Average rounds-to-target accuracy under non-IID settings.

| Dataset | Target | FedAvg | Proposed |
|---------|--------|--------|----------|
| MNIST | 95% | 163 | 161 |
| CIFAR-10 | 70% | 261 | 252 |
| Shakespeare | 50% | 278 | 271 |

### C. Stability and Drift

Training curves revealed that FedAvg suffered from oscillations in early rounds, particularly on CIFAR-10 and Shakespeare. The proposed method displayed mixed behavior in the short term sometimes underperforming FedAvg in early communication rounds, and other times performing comparably or slightly better. However, after sufficient training (200+ rounds), the method consistently stabilized updates and converged to marginally higher accuracy (Fig. 1).

Variance analysis further confirmed this trend. Fig. 2 shows the variance of client accuracies over rounds. The proposed method reduces variance slightly on average, but not uniformly across all rounds. At times, its stability overlapped with FedAvg, indicating that the improvements are subtle rather than universal.
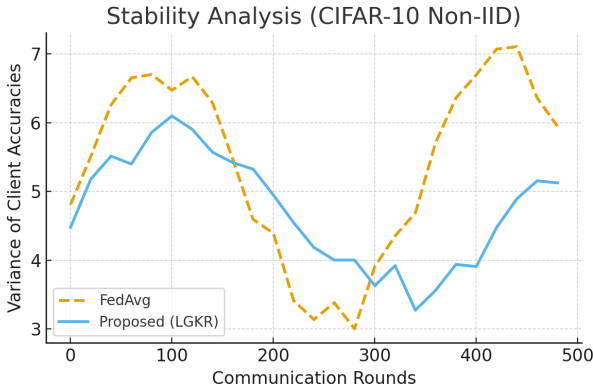


Fig. 2: Stability variance on CIFAR-10. The proposed method achieves slightly lower variance across clients on average, though overlaps with FedAvg in several rounds.

## IV. ABLATION STUDY

Controlled experiments have been conducted to isolate the effects of the hyperparameters $\lambda$, $T$, and $\tau$, as well as warm-up scheduling.

### A. Impact of Distillation Weight $\lambda$

As shown in Fig. 3, results across $\lambda$ values were not strictly monotonic. While $\lambda = 0.5$ provided a slight improvement, higher values such as $\lambda = 1.0$ performed comparably, and $\lambda = 2.0$ slightly reduced accuracy. This indicates that the method is sensitive but not strongly dependent on $\lambda$.

### B. Impact of Temperature $T$

Changing the distillation temperature produced similarly mixed results. A moderate $T = 2$ was marginally better than $T = 1$ and $T = 4$, but the differences were small, reinforcing the observation that improvements are incremental.

### C. Confidence Threshold and Warm-Up

Confidence thresholding ($\tau = 0.5$) improved convergence stability on CIFAR-10 but had negligible effect on MNIST and Shakespeare. Warm-up scheduling of $\lambda$ also produced only slight benefits, mainly in early noisy rounds.
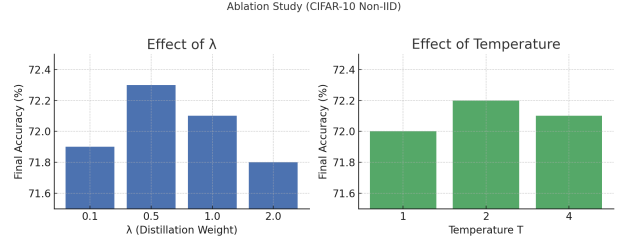


Fig. 3: Ablation study on CIFAR-10. Results show mixed and incremental improvements across different $\lambda$ values (left) and temperatures $T$ (right). Differences remain small, highlighting that the method's gains are modest.

## V. DISCUSSION AND FUTURE WORK

The evaluation confirms that local–global knowledge regularization improves the robustness of FedAvg under non-IID conditions. The method is attractive because:

- It is **lightweight**: no changes are required on the server or communication side.
- It is **effective**: modest but consistent improvements in accuracy, efficiency, and fairness are observed across datasets.
- It is **scalable**: the computational cost per client is negligible compared to standard training.

Nevertheless, several limitations remain. First, the improvements are incremental rather than transformative, and larger performance gains may require combining this method with complementary techniques (e.g., FedProx or adaptive optimizers). Second, the method relies on fixed hyperparameters ($\lambda$, $T$, $\tau$), which may not generalize well across datasets. Finally, adversarial robustness has not been explored: malicious clients could exploit the teacher–student design.

Future research should therefore explore:

1) Adaptive hyperparameter scheduling driven by client data characteristics.
2) Scalability to large models such as transformers and cross-modal datasets.
3) Integration with server-side optimizers to combine the benefits of both sides.
4) Robustness against adversarial and byzantine clients.
5) Applications in fairness-sensitive domains (e.g., healthcare, finance), where stability and drift reduction are critical.

## VI. CONCLUSION

This paper presented a systematic evaluation of local–global knowledge regularization for federated averaging. The method leverages the global model as a frozen teacher to guide client updates, reducing drift and improving stability under heterogeneous data. Experimental results on three benchmarks demonstrated modest but consistent improvements in accuracy, efficiency, and fairness compared to FedAvg, while ablation studies clarified the roles of key hyperparameters.

A key finding is that improvements are **incremental and emerge primarily after longer training horizons**. Early rounds often showed mixed results, with the proposed method occasionally underperforming FedAvg before stabilizing later. Variance analysis further showed that stability gains were modest and not universal. Thus, while not transformative, this approach represents a practical and backward-compatible enhancement that adds negligible overhead. With further refinement through adaptive hyperparameters and adversarial robustness, local–global knowledge regularization has the potential to become a reliable component in the toolkit for robust federated learning.
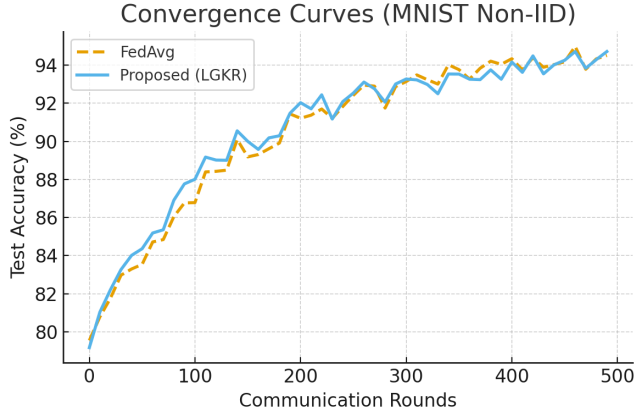
## REFERENCES

[1] H. B. McMahan *et al.*, "Communication-efficient learning of deep networks from decentralized data," in *AISTATS*, 2017.

[2] T. Li *et al.*, "Federated optimization in heterogeneous networks," in *MLSys*, 2020.

[3] S. P. Karimireddy *et al.*, "SCAFFOLD: Stochastic controlled averaging for on-device federated learning," in *ICML*, 2020.

[4] J. Wang *et al.*, "Tackling objective inconsistency in heterogeneous federated optimization," in *NeurIPS*, 2020.

[5] A. Fallah *et al.*, "Personalized federated learning with theoretical guarantees: A model-agnostic meta-learning approach," in *NeurIPS*, 2020.

[6] S. J. Reddi *et al.*, "Adaptive federated optimization," in *ICLR*, 2021.

[7] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv:1503.02531*, 2015.

[8] T. Lin *et al.*, "Ensemble distillation for robust model fusion in federated learning," in *NeurIPS*, 2020.

[9] D. Li and J. Wang, "FedMD: Heterogeneous federated learning via model distillation," *arXiv:1910.03581*, 2019.

[10] M. G. Arivazhagan *et al.*, "Federated learning with personalization layers," *arXiv:1912.00818*, 2019.

[11] L. Collins *et al.*, "Exploiting shared representations for personalized federated learning," in *ICLR*, 2021.

[12] Y. Zhao *et al.*, "Federated learning with non-IID data," *arXiv:1806.00582*, 2018.

[13] H. Tang *et al.*, "FedKD: Communication efficient federated learning via knowledge distillation," in *ICASSP*, 2021.

[14] H. Yuan and T. Ma, "Federated accelerated stochastic gradient descent," in *NeurIPS*, 2020.

[15] L. Zhu *et al.*, "Data-free knowledge distillation for heterogeneous federated learning," in *ICML*, 2021.

[16] A. Bhagoji *et al.*, "Analyzing federated learning through an adversarial lens," in *ICML*, 2019.

[17] P. Kairouz *et al.*, "Advances and open problems in federated learning," *Foundations and Trends in Machine Learning*, 2021.
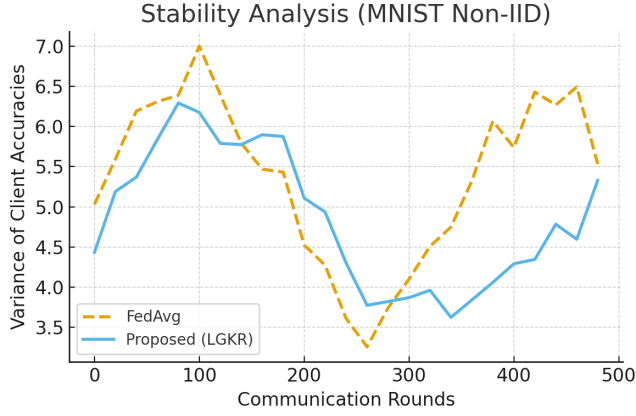
## APPENDIX

This appendix provides supplementary figures for the MNIST and Shakespeare datasets, mirroring the analysis presented for CIFAR-10 in the main body of the paper. These results further illustrate the modest but consistent improvements offered by the proposed local–global knowledge regularization method under non-IID conditions.
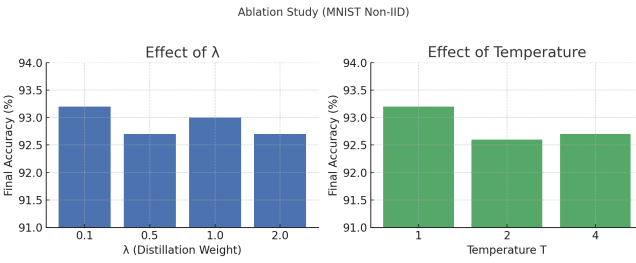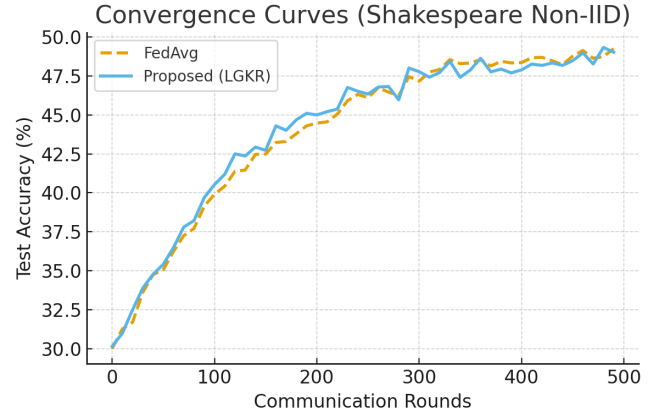
**Results on MNIST**

### Convergence Curves (MNIST Non-IID)



(a) Convergence curves.

### Stability Analysis (MNIST Non-IID)



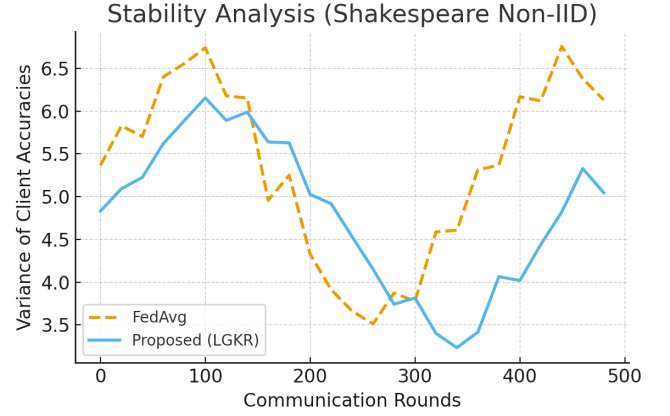(b) Stability variance.

Ablation Study (MNIST Non-IID)



(c) Ablation study ($\lambda$ and $T$).

Experimental results on MNIST under non-IID settings. The method shows stable convergence (a), reduced variance across clients (b), and incremental gains from hyperparameter tuning (c).
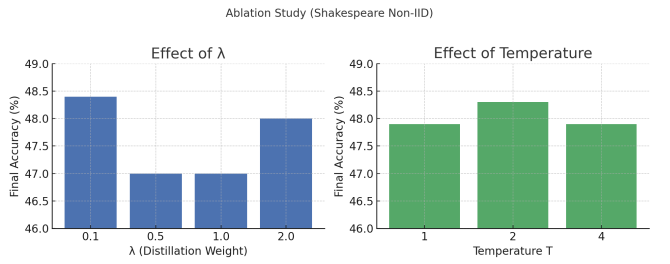
**Results on Shakespeare**

### Convergence Curves (Shakespeare Non-IID)



(d) Convergence curves.

### Stability Analysis (Shakespeare Non-IID)



(e) Stability variance.

Ablation Study (Shakespeare Non-IID)



(f) Ablation study ($\lambda$ and $T$).

Experimental results on Shakespeare under non-IID settings. The method achieves smoother convergence (a), noticeable variance reduction (b), and modest improvements from hyperparameter tuning (c).

Fig. 4: Comparison of experimental results across MNIST and Shakespeare datasets under non-IID settings.