

Beyond the Teacher-Student Paradigm: The Role of Network Structure in Knowledge Distillation

Nasif Fahmid Prangon and Jie Wu

Department of Computer and Information Sciences, Temple University, USA

Email: {nasifprangon, jiewu}@temple.edu

Abstract—Traditional knowledge distillation relies on one-shot, static teacher–student interaction, limiting its effectiveness in large-scale, heterogeneous, or resource-constrained environments. This work presents a multi-round, peer-to-peer distillation framework in which student models repeatedly exchange knowledge over structured network topologies. We analyze how the classical graph-theoretic metrics—diameter, node degree, and bisection bandwidth well as social network measures, including eccentricity and closeness centrality, impact the speed and uniformity of knowledge propagation. In order to minimize the dependence on continuous teacher availability, we propose a dynamic reinforcement strategy that selectively selects the most influential students for teacher supervision, allowing them to serve as anchors in accelerating the processes of diffusion across the topology. Experiments on line, ring, and torus topologies confirm that lower diameter topologies result in faster convergence, higher student saturation, and more consistent learning behavior. When combined with adaptive reinforcement, these findings indicate that topology-aware, dynamically guided distillation provides more scalable and stable performance than traditional single-teacher approaches.

Index Terms—Distributed machine learning, dynamic reinforcement, knowledge distillation, network topology, peer-to-peer learning.

I. INTRODUCTION

Knowledge distillation enables a teacher model to transfer its learned representations to one or more students [1]. Traditional KD is typically a static, one-shot teacher–student interaction [2], which limits its ability to support repeated or distributed knowledge transfer. As systems scale to large, heterogeneous, or resource-constrained environments, researchers have increasingly explored multi-round and dynamic forms of interaction that better capture these settings [1]. Here, ‘single contact’ refers to a one-shot teacher–student update, while ‘multiple’ and ‘dynamic’ contacts denote repeated or adaptive interactions across rounds.

In such settings, the teacher may be intermittently reachable due to bandwidth limits, power constraints, or the high computational cost of running the full model, and hence there is a greater need for adaptive and collaborative learning approaches [3]. Consequently, research on repeated interactions, dynamic reinforcement, and varied network topologies has become important both practically and theoretically [4]. These factors significantly impact knowledge propagation efficiency and effectiveness. This study has explored in detail how more efficient knowledge distillation can be accomplished by using numerous teacher-student interactions, adaptive interaction

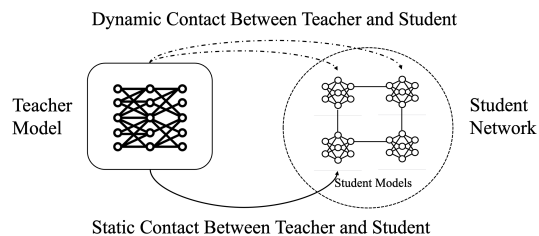


Fig. 1: Illustration of the teacher-student distillation with different contact strategies.

patterns, and well-organized network topologies. We offer an array of experiments that demonstrate the performance improvement and scalability advantages of these general frameworks and thus constitute a justification for their general acceptance in real-life settings.

A single teacher–student contact remains appealing in constrained environments due to its simplicity, low overhead, with no requirement for persistent access to the teacher. This is the case when the teacher is only intermittently reachable due to bandwidth, power, or computing limitations. However, relying solely on one contact also restricts the student’s capability of refining knowledge over time, which motivated the need for peer-to-peer propagation where students can continue improving even when the teacher is unavailable [5], [6]. Such an iterative diffusion can keep learning progress going despite network disruptions in edge environments.

Dynamic contacts introduce some flexibility by allowing nodes to reassess and adapt their roles, making it more suitable for a scalable system. The result may be that a node, having increased its performance substantially through initial contact, becomes essentially a “teacher” to the others subsequently and thus propels the overall network in the spread of top-quality information. This reflects real-world learning environments where knowledge is seldom passed on unidirectionally or with one sole authority figure.

Because multi-round interactions depend on how fast information can move between students, the underlying communication topology becomes a key factor. The network topology plays an important role in defining the efficiency of knowledge diffusion. Ring or fully connected topologies reduce the diameter of the network and will allow faster and more uniform flows of knowledge transfer. Redundancy and fault tolerance may also be achieved by minimizing path length, ensuring that even resource-constrained nodes keep

fine-tuning their models without being isolated [7].

For systematic analysis of these effects, we examine five metrics: diameter, node degree, and bisection bandwidth of graph theory, and eccentricity and closeness centrality of social networks. These first three encode structural characteristics that are general to all topologies, whereas the last two are significant in the asymmetric configurations (e.g., line, star), in which a core node speeds up diffusion. Our simulations mostly estimate the diameter, linearly proportional to node degree and bisection bandwidth, but we report all five metrics to provide the overall topological characterization. The challenge is how to propagate knowledge efficiently when teacher supervision is limited and communication is constrained by the underlying topology. To manage scarce teacher supervision, the system must decide which student should be reinforced and how topology influences that choice.

This work consequently explores dynamic reinforcement in peer-to-peer knowledge distillation and the prospect of controlling topology structure design and reinforcement decisions with these measures. We propose a topology-aware framework with multi-round transactions and a dynamic strategy in which the teacher selectively trains influential nodes. Experiments on line, ring, and torus demonstrate that training structurally central nodes overwhelmingly enhances convergence and learning homogeneity.

II. TOPOLOGY-AWARE DISTILLATION FRAMEWORK

A. Topology and Graph Metrics

The student nodes are connected through a predefined network topology $G = (\mathcal{S}, E)$, where each edge represents an allowable peer-to-peer interaction. The structure of this graph influences the rate and coverage of knowledge diffusion. We focus on the following key metrics:

Diameter (D): Defined as the longest shortest-path distance between any two nodes in the graph. This serves as an upper bound on the number of rounds needed for complete knowledge diffusion across the network.

Eccentricity ($\epsilon(S_i)$): For a given node S_i , this is the maximum shortest-path distance from S_i to any other node. It determines the minimum number of rounds required for knowledge originating from S_i to reach the entire network.

Closeness Centrality ($C(S_i)$): A measure of how centrally located node S_i is, defined as the inverse of the average shortest-path distance from S_i to all other nodes. Nodes with high $C(S_i)$ can disseminate or receive knowledge more efficiently, making them strong candidates for reinforcement.

We also note two related measures: the *node degree*, which reflects direct connectivity, and the *bisection bandwidth*, which quantifies the minimum communication cut across the network. Although our simulations primarily use diameter (proportional to both), we include all five measures for completeness.

B. Student Update Rule

The key symbols used throughout our model are summarized in Table I. Each student node S_i maintains a local model

TABLE I: Key Symbols and Descriptions

Symbol	Description
S_i	Student node i in the network
T	Teacher node
$\mathbf{w}_i^{(t)}$	Parameters of student S_i at round t
$f(\mathbf{w})$	Output logits from model with parameters \mathbf{w}
y	Ground-truth label
α	Weight balancing supervised and distillation losses
$\gamma_{ij}^{(t)}$	Peer influence weight from student S_j to S_i
$\gamma_{iT}^{(t)}$	Teacher influence weight on student S_i
$g_i^{(t)}$	Distillation target for student S_i
$A_i^{(t)}$	Accuracy of student S_i at round t
$R_i^{(t)}$	Reinforcement score for student S_i after round t
$\epsilon(S_i)$	Eccentricity of node S_i
D	Graph diameter
$C(S_i)$	Closeness centrality of node S_i

with parameters $\mathbf{w}_i^{(t)}$ at training round t . The model is trained using a weighted combination of two objectives:

- A supervised learning loss based on ground-truth labels.
- A distillation loss that aligns the student’s output with a target derived from its peers and, if reinforced, the teacher.

The overall loss function minimized by each student at round t is:

$$\mathcal{L}_i^{(t)} = \alpha \text{CE} \left(f(\mathbf{w}_i^{(t)}), y \right) + (1 - \alpha) \text{KL} \left(f(\mathbf{w}_i^{(t)}), g_i^{(t)} \right) \quad (1)$$

Here, CE denotes the cross-entropy loss between the student’s prediction and the ground-truth label y , while KL represents the Kullback–Leibler divergence between the student’s output and the distillation target $g_i^{(t)}$. The parameter $\alpha \in [0, 1]$ controls the trade-off between the supervised loss and the distillation loss, allowing the model to balance learning from labeled data and from its peers or the teacher.

The distillation target $g_i^{(t)}$ is a weighted average of the outputs (logits) of the student’s neighbors and, optionally, the teacher:

$$g_i^{(t)} = \sum_{j \in \mathcal{N}_i} \gamma_{ij}^{(t)} f(\mathbf{w}_j^{(t)}) + \gamma_{iT}^{(t)} f(\mathbf{w}_T) \quad (2)$$

In this expression, \mathcal{N}_i denotes the set of peer student nodes connected to S_i in the communication graph. The term $\gamma_{ij}^{(t)}$ represents the weight assigned to the influence of peer S_j on S_i , while $\gamma_{iT}^{(t)}$ is the weight associated with teacher supervision, which is greater than zero only if S_i is selected for reinforcement in round t . These weights are normalized such that the total influence sums to one: $\sum_{j \in \mathcal{N}_i} \gamma_{ij}^{(t)} + \gamma_{iT}^{(t)} = 1$.

This setup allows each student to learn not only from its own data but also from surrounding peers and, occasionally, the teacher. When $\gamma_{iT}^{(t)} = 0$, the student is fully peer-supervised. When $\gamma_{iT}^{(t)} > 0$, the student benefits from direct reinforcement, which typically accelerates convergence and improves its influence in the network during rounds.

C. Reinforcement Objective

The teacher can supervise only a limited number of students simultaneously. Thus, it must select students to reinforce each

round. To guide this decision, we define a reinforcement score $R_i^{(t)}$ for each student S_i , which considers four factors: (i) the potential for improvement, favoring students with lower current accuracy; (ii) the influence of the student, measured by the number of its neighbors; and (iii) the communication cost, penalizing students that are farther from the teacher. The reinforcement score is defined as:

$$R_i^{(t)} = W_1(1 - A_i^{(t)}) + W_2(N_i/N_{\max}) - W_3\text{Dist}_i \quad (3)$$

Here, $A_i^{(t)}$ is the accuracy of student S_i at round t , N_i is the number of neighbors of S_i , N_{\max} is the maximum number of neighbors any student has in the network, Dist_i represents the average communication distance between S_i and its neighbors. The terms W_1 , W_2 , and W_3 are tunable hyperparameters that control the relative importance of accuracy improvement, peer influence, communication cost, and structural resilience in the reinforcement decision. Specifically, W_1 prioritizes low-accuracy students, W_2 favors highly connected ones, and W_3 penalizes nodes with larger communication distance.

The student selected for reinforcement in the next round is the one with the highest reinforcement score:

$$S_{\text{reinforced}} = \arg \max_{S_i} R_i^{(t)} \quad (4)$$

Only this student will receive direct teacher supervision, reflected by setting $\gamma_{iT}^{(t+1)} > 0$ in the subsequent round's update rule.

D. Constraints and Optimization Objective

Our objective is to maximize the total accuracy of all student models at the final round of training expressed as:

$$\max \sum_{i=1}^N A_i^{(t_{\text{last}})} \quad (5)$$

Here, $A_i^{(t_{\text{last}})}$ represents the final accuracy of student S_i . After a fixed number of peer-to-peer rounds, reinforcement is assigned to the student with the highest computed reinforcement score by the teacher.

To achieve this goal under practical constraints, the following conditions must be satisfied at each round t :

1. Accuracy must not decrease over time, ensuring that each student's performance improves or remains the same: $A_i^{(t+1)} \geq A_i^{(t)}$, $\forall i, t$.

2. The number of students reinforced by the teacher in any round cannot exceed a predefined limit: $|\mathcal{R}^{(t)}| \leq R_{\max}$, $\forall t$.

For our experiments, we fix $R_{\max} = 1$, constraining that one student only is directly supervised by the teacher per round. The constraint mimics a resource-constrained scenario, whereby the capacity of the teachers (e.g., computation resources, bandwidth, or power) is greatly constrained. Although the framework supports greater values for R_{\max} , examining the single-reinforcement scenario enables us to control and assess the effect of reinforcement techniques and topological aspects more effectively.

3. Reinforcement may be subject to a communication budget, where the total cost of reinforcing selected students must not exceed a global limit: $\sum_{S_i \in \mathcal{R}^{(t)}} \kappa_i \leq C_{\max}$, $\forall t$.

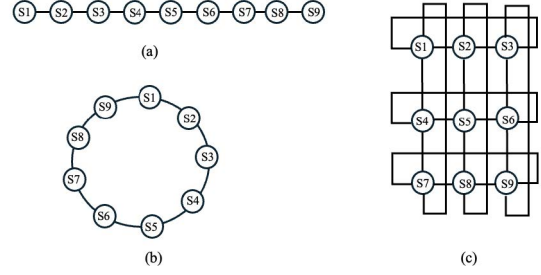


Fig. 2: Students arranged in different topologies: (a) Line ($D = 8$), (b) Ring ($D = 4$), and (c) Torus ($D = 2$).

Here, $\mathcal{R}^{(t)} \subseteq \mathcal{S}$ represents the set of students reinforced by the teacher at round t , κ_i denotes the cost of reinforcing student S_i , and C_{\max} total communication budget per round.

Reinforcement in our experiments is given following a fixed number of peer-to-peer rounds, regardless of any dynamic cost modeling. Therefore, we assume C_{\max} to be so high that it does not limit reinforcement. This makes evaluation easy by maintaining focus on the impact of the reinforcement strategy and network structure alone.

The overall goal, as defined in Equation (5), is to determine the optimal reinforcement choice each round that guides which students receive teacher supervision, thereby maximizing all the students' performance at the final round while respecting structural and resource constraints.

III. TOPOLOGY-AWARE DISTILLATION DESIGN

A. Topologies Used

We evaluate our framework across three structured network topologies: Line, Ring, and Torus, depicted in Figure 2. These define how students are connected and influence how quickly teacher-refined knowledge spreads through the network. Each topology impacts diffusion efficiency based on its diameter and eccentricity. The significant characteristics of these topologies, as shown in Figure 2, are discussed below:

Line Topology: Students form a 1D chain, and the largest diameter ($D = 8$) is obtained along with a high average communication cost. Centrality reinforcement towards a middle node (e.g., S_5) reduces the eccentricity (4) but is slow and inefficient in the diffusion of knowledge to edge nodes.

Ring Topology: By linking the endpoints of the line, the ring minimizes the diameter to 4 and provides bidirectional flow of knowledge. Although eccentricities are comparable to the line, better connectivity speeds up convergence and improves general performance.

Torus Topology: A wraparound linked 2D grid (e.g., 3×3) provides each node with four neighbors and lowers the diameter to 2. Since all nodes are now available in two hops, peer updates and reinforcement disseminate rapidly and in parallel, and balanced and efficient distillation is achieved.

B. Training and Learning Setup

We compare our model on CIFAR-10 and SVHN, two popular image classification datasets. CIFAR-10 comprises

60,000 colored images from 10 classes (50,000 training, 10,000 testing), while SVHN contains more than 70,000 digit images from real-world street scenes. A pre-trained ResNet-50 is employed as the teacher in both cases, while the student is initialized with a light-weight ResNet-18 to resolve the tradeoff between representational strength and communication efficiency in serving resource-constrained environments.

Training occurs in three stages: (1) distillation from the teacher to the student node in the center, (2) propagation among peers in the topological structure, and (3) reinforcement of selected student nodes to improve global accuracy.

Peer-to-Peer Propagation Mechanism: After the first interaction with the teacher, knowledge is distributed exclusively through peer-to-peer interactions on the network graph. Every student updates itself using a weighted combination of its supervised loss and the distillation loss from surrounding nodes. The empowered node has the maximum influence, and students directly connected to it put more weight on its output. Subsequently, nodes further away emphasize peer outputs along the shortest distance to the center node.

Reinforcement Strategy Details: To further enhance learning, we implement a dynamic reinforcement strategy where only one student is selected after a few rounds to receive direct supervision from the teacher. This selection is guided by a reinforcement score based on the student’s performance, influence on the network, and communication cost. Reinforcement accelerates knowledge diffusion by periodically updating strategic nodes within the topology.

IV. TOPOLOGICAL DISTILLATION PROPERTIES

From standard graph theory, the time for information to reach all nodes from a source is upper-bounded by the source’s eccentricity, which is itself bounded by the graph’s diameter. We restate this property only to justify why diameter is an essential predictor of diffusion speed in our setting.

Property 1. *The number of rounds required for teacher-refined knowledge to reach a student is equal to the shortest-path distance from the reinforced node to that student.*

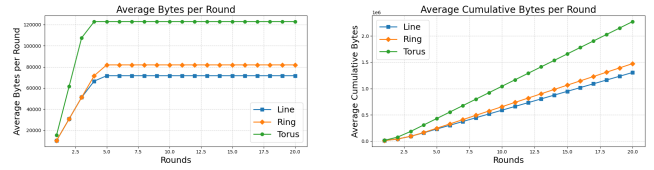
Property 2. *The maximum number of rounds for full network-wide diffusion equals the eccentricity of the reinforced node—i.e., the largest shortest-path distance from it to any other node.*

Theorem 1. *Let S_c be the student reinforced by the teacher. Then the number of rounds T required for complete diffusion satisfies:*

$$T \geq \epsilon(S_c) \leq D,$$

where $\epsilon(S_c)$ is the eccentricity of S_c , and D is the graph diameter.

Proof. Since each round propagates knowledge one hop, a node k hops from S_c must wait k rounds. Thus, at least $\epsilon(S_c)$ rounds are required for the farthest node. By definition, $\epsilon(S_c) \leq D$, completing the bound. \square



(a) Average Bytes per Round

(b) Average Cumulative Bytes

Fig. 3: Bandwidth usage comparison for Line, Ring, and 2D Torus topologies.

V. SIMULATION RESULTS

We evaluate our reinforcement strategies under line, ring, and torus topologies. Traditional baselines like FedAvg or Deep Mutual Learning assume full teacher availability or synchronous global aggregation. Because our setting restricts reinforcement to a single teacher under explicit topology constraints, we compare against fixed and rotating reinforcement strategies to ensure a fair and consistent evaluation. Each reinforcement round is preceded by 5 peer-to-peer diffusion steps across the student topology. The study analyzes topology variations, parallelization effects, reinforcement strategies, and dynamic reinforcement decision-making.

A. Impact of Topology

The topological arrangement of student models has a major influence on distillation peer-to-peer dynamics. We consider three topologies—*Line*, *Ring*, and *Torus*—with different diameters and connectivity.

Table II shows the number of students who achieved 35% accuracy and the average epochs taken. On CIFAR-10, both *Line* and *Ring* saturated 7 out of 9 students in 1.71 epochs, whereas *Torus* saturated all 9 in only 1.56 epochs. On SVHN, *Line* and *Ring* saturated 6 in 2.00 epochs, and *Torus* saturated 8 in 1.78 epochs. Figure 3 shows the bytes transmitted for three topologies. This shows that while line minimizes communication, the torus provides a better accuracy–efficiency tradeoff by enabling faster diffusion without significantly increasing bandwidth. From the figure, we see that line topology uses the least bandwidth; however, ring and torus topologies use almost the same bandwidth. This highlights a clear trade-off: the line minimizes communication, whereas the torus incurs slightly higher cost but offers significantly faster and more uniform convergence. However, Figure 4 shows accuracy progression across epochs, where we clearly see the gain from torus topology over line and ring. This highlights that the torus topology consistently achieves the best accuracy–speed balance due to its lower diameter. In the *Line* topology (Figures 4a, 4d), central students converge quickly, but peripheral nodes lag due to longer paths. The

TABLE II: Student Saturation Analysis Across Topologies

Dataset	Topology	# Saturated Students	Avg. Epochs
CIFAR	Line	7 / 9	1.71
	Ring	7 / 9	1.71
	Torus	9 / 9	1.56
SVHN	Line	6 / 9	2.00
	Ring	6 / 9	2.00
	Torus	8 / 9	1.78

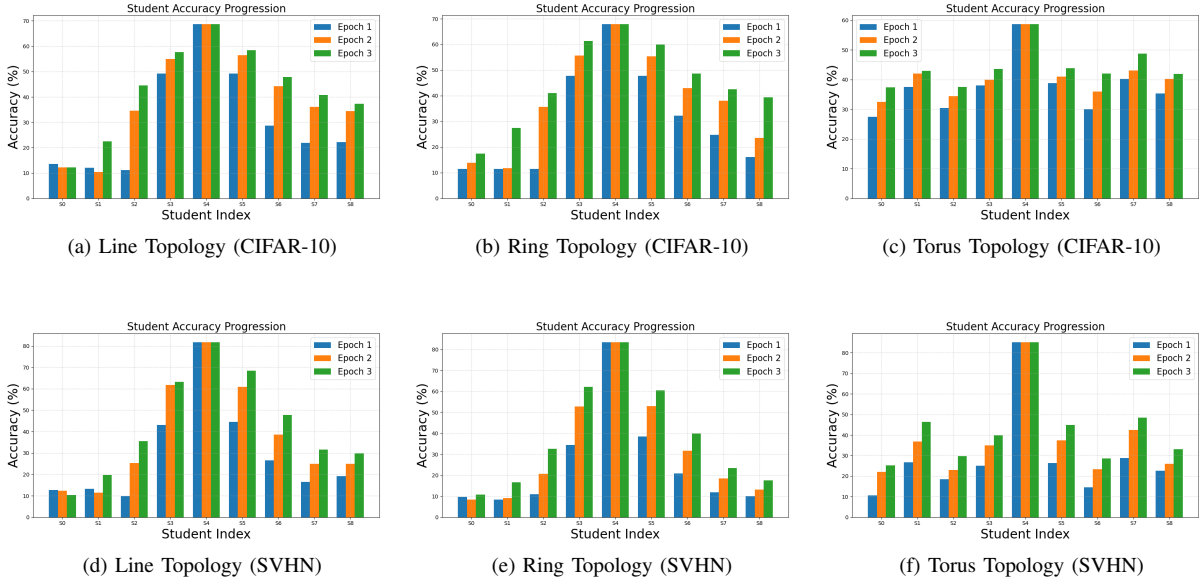


Fig. 4: Student accuracy progression across epochs for different topologies.

Ring (Figures 4b, 4e) allows more balanced updates through bi-directional flow, though early-round asymmetries persist. The *Torus* (Figures 4c, 4f) achieves the most rapid and uniform convergence, with nearly all students exceeding 40% accuracy by epoch 3, thanks to its low diameter and high interconnectivity.

Collectively, these findings show how node degree and diameter directly impact the distributed learning efficiency. Richer topologies like the torus invariably produce more rapid convergence and a smaller performance gap among learners.

B. Impact of Parallelization and Diameter

We evaluate how parallelism and structural adjustments affect learning efficiency on the ring topology. In the first strategy, a single central student S_4 is reinforced across all rounds. The second strategy introduces parallelism by reinforcing two distant students (e.g., S_3 and S_7) simultaneously, enabling bidirectional diffusion. The third strategy virtually divides the ring into two halves by strengthening students such as S_0 and S_6 , effectively reducing the diameter of the topology. Figures 5a and 5b show that greater parallelism accelerates convergence and improves final accuracy: single-teacher supervision diffuses knowledge slowly, two-teacher supervision creates faster dual-front propagation, and the splitting strategy performs best due to shorter paths and wider early coverage. These results highlight that parallel reinforcement and diameter reduction significantly boost transmission efficiency. However, additional teachers are not always beneficial, as noted in [8], and their number should be tuned according to the behavior of the data set.

C. Impact of Reinforcement Strategy

To assess the effectiveness of reinforcement allocation under single-node teacher supervision, we compare three strategies—Fixed, Rotating, and Dynamic reinforcement—on a line

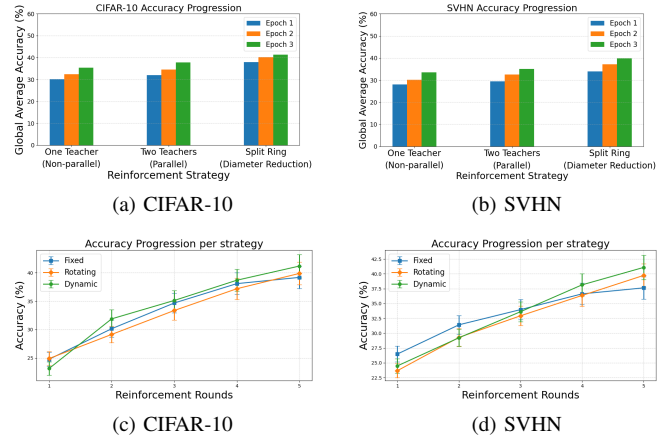


Fig. 5: Accuracy trends: (a–b) show the effect of decreasing diameter on ring topology; (c–d) compare fixed, rotating, and dynamic reinforcement on a line topology.

topology as summarized in Table III. In all cases, only one student is reinforced per round. Reinforcement in the *Fixed* strategy targets the same student (e.g., S_5) repeatedly. In the *Rotating* strategy, reinforcement cycles sequentially through the students, giving each student a chance for direct supervision. Finally, the *Dynamic* strategy (proposed) adaptively selects the student with the highest reinforcement score based on current accuracy, peer connectivity, and cost. We note that the absolute accuracies stay below 50% mainly because we use lightweight students and only five reinforcement rounds; running more rounds or larger datasets (e.g., CIFAR-100 or ImageNet) would naturally increase these values and is part of our planned future work.

Figures 5c and 5d illustrate the progression of global average accuracy across five reinforcement rounds for CIFAR-10 and SVHN. The Fixed strategy provides stable gains but plateaus early due to localized knowledge propagation. The

TABLE III: Global Accuracy (%) per Round Across Strategies

Data	Strategy	R1	R2	R3	R4	R5
CIFAR	Fixed	24.76	30.15	34.68	38.10	39.20
	Rotating	24.89	29.14	33.36	37.20	39.90
	Dynamic	23.13	31.89	35.13	38.70	41.20
SVHN	Fixed	26.48	31.41	33.97	36.63	37.64
	Rotating	23.67	29.29	32.93	36.36	39.74
	Dynamic	24.48	29.22	33.62	38.15	41.09

rotating strategy improves fairness but reduces consistency since each student receives limited reinforcement, whereas dynamic reinforcement consistently outperforms both strategies.

These results highlight that when only a single student can be reinforced per round, the strategy of selection significantly influences convergence. Dynamic reinforcement with context-aware decision-making maximizes learning gains under constrained supervision and proves to be the most efficient. Together, these results demonstrate that structural properties such as diameter and reinforcement choice have a measurable and consistent effect on convergence speed and accuracy.

VI. RELATED WORKS

Hinton et al. [2] proposed knowledge distillation by transferring softened class probabilities from the teacher to the smaller student model. Buciluă et al. [9] previously exemplified that with low-capacity models, high accuracy is possible. FitNets [10] took that further by copying intermediate representations, with compression possible in resource-constrained environments like edge cloud. To address the shortcomings of fixed one-shot distillation, such as Deep Mutual Learning [11] and federated learning [12], allow the student models to share and hone knowledge iteratively, and this improves robustness in decentralized or dynamic environments. Dynamic teacher-student roles, e.g., during meta-teaching [13] and curriculum learning [14], enable top-performing learners to assume teaching roles temporarily, supporting adaptive and peer-driven learning. Decentralized methods like federated averaging [15] and sparse DGD [16] tackle communication overhead through partial connectivity. Graph-based models, such as GNNs [17], further optimize peer-to-peer transmission through structure-aware message passing. Lightweight distillation methods have addressed edge and IoT settings, emphasizing energy-conscious teacher-student design [18] and offloading of tasks under bandwidth and computationally constrained environments [19].

Considering the diversity of previous works, there is still unexplored work in the integrated framework that brings together iterative interaction, dynamic role switching, and topology-awareness. This work fills that gap by consolidating these aspects in order to support scalable and robust knowledge transfer in structured networks.

VII. CONCLUSION

In this work, we showed how topology, diameter, and parallelism affect peer-to-peer distillation of knowledge in decentralized student networks. Experiments in line, ring, and torus topologies brought to light the structural effect on diffusion rates and equality. With single-student reinforcement,

our dynamic approach met with efficient convergence, though extensions to general topology and multi-student reinforcement are left as open avenues. Scalability to bigger datasets (e.g., CIFAR-100, Tiny-ImageNet, and ImageNet), fairness under node heterogeneity, and robustness against dynamic failures are directions for future exploration. This work can also extend to FedKD, edge-centric model sharing, and other decentralized distillation applications.

REFERENCES

- [1] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.
- [2] G. Hinton, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [3] G. Wu and S. Gong, "Peer collaborative learning for online knowledge distillation," in *Proceedings of the AAAI Conference on artificial intelligence*, vol. 35, no. 12, 2021, pp. 10302–10310.
- [4] J. Jankowski, "Habituation effect in social networks as a potential factor silently crushing influence maximisation efforts," *Scientific reports*, vol. 11, no. 1, p. 19055, 2021.
- [5] Q. Lu, E. Xun, and G. Tang, "Mta4dpr: Multi-teaching-assistants based iterative knowledge distillation for dense passage retrieval," in *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024, pp. 5871–5883.
- [6] D. Chen, J.-P. Mei, C. Wang, Y. Feng, and C. Chen, "Online knowledge distillation with diverse peers," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 04, 2020, pp. 3430–3437.
- [7] A. Bellet, A.-M. Kermarrec, and E. Lavoie, "D-cliques: Compensating for data heterogeneity with topology in decentralized federated learning," in *2022 41st International Symposium on Reliable Distributed Systems (SRDS)*. IEEE, 2022, pp. 1–11.
- [8] S. I. Mirzadeh, M. Farajtabar, A. Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh, "Improved knowledge distillation via teacher assistant," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 34, no. 04, 2020, pp. 5191–5198.
- [9] C. Buciluă, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2006, pp. 535–541.
- [10] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *arXiv preprint arXiv:1412.6550*, 2014.
- [11] Y. Zhang, T. Xiang, T. M. Hospedales, and H. Lu, "Deep mutual learning," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4320–4328.
- [12] H. B. McMahan, F. Yu, P. Richtarik, A. Suresh, D. Bacon et al., "Federated learning: Strategies for improving communication efficiency," in *Proceedings of the 29th Conference on Neural Information Processing Systems (NIPS)*, Barcelona, Spain, 2016, pp. 5–10.
- [13] H. Wang, C. T. Yip, and B. Li, "Dynamics of meta-learning representation in the teacher-student scenario," *arXiv preprint arXiv:2408.12545*, 2024.
- [14] P. Soviany, R. T. Ionescu, P. Rota, and N. Sebe, "Curriculum learning: A survey," *International Journal of Computer Vision*, vol. 130, no. 6, pp. 1526–1565, 2022.
- [15] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.
- [16] A. F. Aji and K. Heafield, "Sparse communication for distributed gradient descent," *arXiv preprint arXiv:1704.05021*, 2017.
- [17] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl, "Message passing neural networks," *Machine learning meets quantum physics*, pp. 199–214, 2020.
- [18] Z. Wu, S. Sun, Y. Wang, M. Liu, X. Jiang, and R. Li, "Survey of knowledge distillation in federated edge learning," *arXiv preprint arXiv:2301.05849*, 2023.
- [19] M. Abouelyazid, "Machine learning algorithms for dynamic resource allocation in cloud computing: Optimization techniques and real-world applications," *Journal of AI-Assisted Scientific Discovery*, vol. 1, no. 2, pp. 1–58, 2021.