

REAL: Efficient Distributed Training of Dynamic GNNs with Reuse-Aware Load Balancing Scheduling

ANONYMOUS AUTHOR

Dynamic Graph Neural Networks (DGNNs) have emerged as powerful models for learning from evolving graphs by jointly capturing structural and temporal dependencies. However, training DGNNs at scale is challenging due to three conflicting requirements: maximizing resource utilization, minimizing data transfers, and maximizing data reuse. Existing systems often prioritize one objective at the expense of others, resulting in either communication bottlenecks from fine-grained partitioning or load imbalance and redundant computations caused by temporal variations. To break this trilemma, we present REAL, an efficient distributed system designed to simultaneously optimize these conflicting objectives. REAL minimizes data transfers via snapshot-group scheduling, and leverages a novel incremental aggregation mechanism to decouple computation from structural updates, effectively pruning redundant I/O and FLOPS via snapshot- and group-level reuse. To orchestrate these optimizations, we formulate the workload distribution as a makespan minimization problem and propose a cross-group scheduling algorithm that co-optimizes data reuse and resource utilization within a triple-stream execution pipeline. Furthermore, we design two complementary scheduling policies: REAL-L, which leverages ILP for globally optimal planning, and REAL-G, a greedy fallback capable of scaling to massive graphs. Extensive experiments on six dynamic graph datasets and four DGNN models show that REAL-L and REAL-G achieve up to 3.26 \times and 3.16 \times speedup over state-of-the-art baselines, respectively.

1 INTRODUCTION

Dynamic graphs are graphs whose structures and attributes change continuously over time. Dynamic Graph Neural Networks (DGNNs) have emerged as state-of-the-art methods for processing dynamic graphs, exhibiting a strong ability to capture both structural and temporal dependencies [1, 4, 17, 29, 31, 36, 43, 50, 51, 53, 56, 57, 59, 60, 62, 65, 68, 69]. Depending on the event model, dynamic graphs are categorized into discrete-time dynamic graphs (DTDGs) and continuous-time dynamic graphs (CTDGs). Discrete-time dynamic graphs evolve via snapshot-based updates, whereas continuous-time dynamic graphs are driven by continuous event streams. DGNNs are categorized similarly according to the dynamic graphs they process. In this work, we focus on DGNNs for DTDGs, which model temporal dynamics via discrete snapshots and are widely adopted in both real-world applications [7, 60, 63] and mainstream graph learning frameworks [9, 10, 54].

While significant efforts have been made to improve DGNN training efficiency [13, 14, 20, 28, 39, 48, 52], the escalating scale of graph datasets renders distributed training across multiple GPUs increasingly important [3, 5, 11, 46]. Distributed DGNN training requires simultaneously (1) *maximizing resource utilization*, (2) *minimizing data transfer overhead* (both cross-GPU and CPU-GPU), and (3) *maximizing data reuse*, which constitute three often conflicting objectives. Existing distributed systems navigate these trade-offs with limited success. ESDG [3] distributes adjacent snapshots across GPUs; this necessitates expensive hidden state transfers and, due to RNN sequential dependencies, forces GPUs to idle while awaiting the previous hidden states, degrading utilization. BLAD [11] circumvents such transfer overhead by processing each snapshot group on the same GPU while assigning different groups on different GPUs. However, variations across groups lead to load imbalance, which also results in inefficient resource utilization. Similarly, DGC [5] employs chunk-based partitioning to jointly optimize load balance and communication; however, prioritizing partition quality frequently exacerbates load imbalance due to the intrinsic conflict between these objectives. DynaHB [46] leverages CPU vertex caching to balance workloads and eliminate inter-GPU communication, yet this incurs prohibitive CPU-GPU data transfer overhead, shifting the bottleneck to the PCIe bus. Moreover, these distributed frameworks overlook the opportunities for

Author's Contact Information: Anonymous author.

data reuse. While prior studies [20, 21, 45, 52] demonstrate that leveraging topological similarity between adjacent snapshots significantly reduces redundant computation and I/O, exploiting this reuse introduces computational irregularity due to varying redundancy rates across snapshots. This irregularity exacerbates load imbalance in distributed settings, yet no existing system co-optimizes such data reuse with workload distribution.

To solve this challenging problem, we design REAL, a distributed DGNN training system that jointly optimizes the three objectives above and balances the trade-offs among them. Our key insight is that *different snapshot groups are treated as independent samples in DGNN training*, which allows us to flexibly combine groups and schedule them in an arbitrary order. **First**, REAL adopts *snapshot groups* as the basic scheduling unit, which naturally avoids cross-GPU data transfers. **Second**, REAL enables topology-aware data reuse at both snapshot and group levels. REAL designs new operators that explicitly exploit structural similarities between snapshots and overlapping snapshots across groups, significantly reducing redundant computation and data loading overhead. In addition, REAL overlaps data transfer with computation across heterogeneous resources in a pipelined manner, further improving resource utilization. **Third**, REAL performs cross-group combination to balance workload across GPUs. The cross-group scheduling process jointly considers inter-group data reuse and load balance to improve overall resource utilization. Based on these designs, we formulate an optimal scheduling problem to minimize per-epoch time. Depending on the scenario, REAL solves the problem using either Integer Linear Programming (ILP) or a greedy heuristic, resulting in two variants: REAL-L and REAL-G.

To validate the system's efficacy, we evaluate REAL on physical clusters equipped with NVIDIA H200 GPUs, scaling from single-node (8 GPUs) to distributed (32 GPUs) setups. REAL-L and REAL-G achieve $1.11\times$ – $3.26\times$ and $1.10\times$ – $3.16\times$ speedups over state-of-the-art methods, respectively, driven by substantial reductions in data transfer volume and redundant FLOPs. We further present ablation studies showing that all design components of REAL are necessary. In simulations, as GPU scale increases, REAL-G shows good scalability, scaling to 512 GPUs with 95% parallel efficiency.

We make the following main contributions.

- We identify key inefficiencies in existing distributed DGNN training methods, including insufficient resource utilization, redundant data transfers, and limited data reuse.
- We propose REAL, a system that integrates snapshot- and group-level reuse, and a cross-group scheduling algorithm to jointly optimize resource utilization and data reuse without incurring extra transfer overhead.
- We introduce two variants of the scheduling algorithm: REAL-L and REAL-G, based on ILP and a greedy heuristic, respectively. The system first attempts REAL-L and falls back to REAL-G when ILP solving exceeds the time limit.
- We evaluate REAL-L and REAL-G on six dynamic graph datasets and four DGNN models. REAL-L and REAL-G achieve up to $3.26\times$ and $3.16\times$ higher throughput compared to state-of-the-art method, respectively.

2 BACKGROUND AND MOTIVATION

2.1 Dynamic Graph Neural Networks

Dynamic Graph Neural Networks (DGNNs) capture the evolution of graph data by stacking multiple spatio-temporal blocks. Regardless of specific architectural variations, a typical DGNN layer processes a snapshot \mathcal{G}_t through two logical phases: *structural encoding* and *temporal evolution*. First, a structural encoder (e.g., GCN, GAT) aggregates neighborhood information to capture spatial dependencies. Formally, for a node v , the feature aggregation is defined as:

$$\mathcal{H}_t(v) = \text{SpatioAggr}_v(\mathcal{W}_{\text{gnn}}, \{\mathcal{X}_t(u) \mid u \in \mathcal{N}(v)\}, \mathcal{X}_t(v)) \quad (1)$$

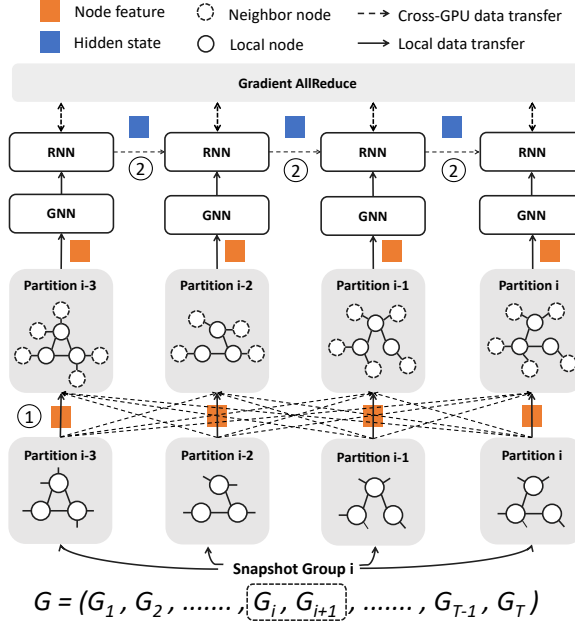


Fig. 1. **Workflow of distributed training for DGNNs.** The dynamic graph is processed in snapshot groups, with partitions distributed across GPUs (gray boxes). Solid and dashed circles represent local and cross-partition neighbors, respectively. The workflow highlights two data transfer overheads: ① **Neighbor Feature Transfer** (dashed arrows) exchanges boundary-node features (orange squares) prior to GNN aggregation; ② **Hidden State Transfer** synchronizes boundary-node states (blue squares) to maintain temporal consistency across snapshots.

where $\mathcal{N}(v)$ denotes the neighborhood of v . Second, the aggregated features are fed into a temporal encoder (e.g., RNN, LSTM) to update the node's hidden state based on historical contexts:

$$h_t = \text{TemporalAggr}(\mathcal{W}_{\text{rnn}}, \mathcal{H}_t, h_{t-1}), \quad h_0 = 0 \quad (2)$$

This establishes a sequential dependency where h_t strictly relies on the previous state h_{t-1} .

Representative models follow this paradigm primarily by integrating spatial aggregation with temporal recurrence. Canonical examples like WD-GCN [31], TGCN [4], and GAT-LSTM [59] combine spatial encoders (e.g., GCN or GAT) with recurrent units (e.g., LSTM or GRU) to jointly capture structural and temporal dependencies in node embeddings. Even parameter-centric variants like EvolveGCN [36], which evolve model weights rather than node embeddings, strictly adhere to this fundamental snapshot-based execution pattern.

2.2 Workflow of Distributed DGNN Training

Training DGNNs typically involves partitioning the dynamic graph sequence $\mathcal{G} = \{G_1, \dots, G_T\}$ into *snapshot groups*, which serve as the basic training units. In a distributed setting, mapping the spatio-temporal dependencies defined in §2.1 to multiple GPUs introduces distinct execution phases and communication patterns, as illustrated in Figure 1.

Category	Dimension	AliGraph [70]	DistDGL [66]	ESDG [3]	DGC [5]	BLAD [11]	DynaHB [46]	REAL
Data Transfer	No Neighbor Feature Transfer	×	×	✓	×	✓	✓	✓
	No Hidden State Transfer	✓	✓	×	×	✓	✓	✓
	No Extra CPU-GPU Transfer	✓	✓	✓	✓	✓	×	✓
Resource Util.	Resource-Level Parallelism	×	×	×	×	✓	×	✓
	Load Balance	✓	✓	×	✓	×	✓	✓
Data Reuse	Topology-Aware Reuse	×	×	×	×	×	×	✓

Table 1. **Comparison of existing solutions across six dimensions in three categories:** *data transfer* (neighbor feature, hidden state, CPU-GPU transfer), *resource utilization* (resource-level parallelism, load balance), and *data reuse* (topology-aware reuse).

Data Loading. The workflow begins by transferring the assigned graph partition and input features from host to device (HtoD). This step prepares the necessary data for GPU computation but imposes significant pressure on PCIe bandwidth, particularly for large-scale graphs.

Spatial Aggregation & Neighbor Feature Transfer. To execute the spatial aggregation (Eq. 1), a GPU must access the features of all neighbors $u \in \mathcal{N}(v)$. If the graph is partitioned, the neighborhood of boundary nodes often spans multiple devices. This necessitates *Neighbor Feature Transfer* (① in Figure 1), where feature tensors must be fetched from remote GPUs before computation can proceed. This irregular, topology-dependent communication constitutes a major performance bottleneck for dense graphs.

Temporal Evolution & Hidden State Transfer. Following spatial aggregation, the temporal update (Eq. 2) requires the previous hidden state h_{t-1} . When consecutive snapshots of a node (or its relevant boundary dependencies) are assigned to different devices, the system must perform *Hidden State Transfer* (②) to synchronize states. Unlike neighbor aggregation, this introduces a strict sequential synchronization barrier, often forcing GPUs to idle while awaiting remote states.

Parameter Synchronization. Finally, a global *Gradient AllReduce* synchronizes model parameters across all participating GPUs after the backward pass, ensuring consistency for the next iteration.

2.3 Trilemma in Distributed Training of DGNNs

Efficient distributed DGNN training requires jointly optimizing three objectives: high resource utilization, minimal data transfer overhead, and maximal data reuse. Existing approaches exhibit a range of strengths and weaknesses, as shown in Table 1; yet none of them address all key objectives. We next analyze these limitations according to the partitioning strategies employed.

Vertex-Based Partitioning. Vertex-based partitioning (e.g., DistDGL) adopt vertex-level partitioning to achieve fine-grained load balancing across GPUs by evenly distributing nodes and their associated computations. However, this strategy fragments the graph structure within each snapshot, resulting in frequent neighbor feature exchanges across GPUs. Consequently, it incurs substantial inter-GPU data transfer during message passing. For example, the profiling of the forward pass of WD-GCN [31] on three real-world datasets shows that DistDGL incurs notable data transfer overhead, accounting for an average of 34% of the total forward pass time (Figure 2).

Snapshot-Based Partitioning. To reduce neighbor feature data transfer, ESDG performs snapshot-level scheduling by assigning each snapshot to a single GPU, ensuring that neighbor aggregation is local. However, it introduces hidden state exchange across GPUs when modeling temporal dependencies, and variations in snapshot sizes cause inter-GPU load imbalance. Moreover, the sequential dependency of RNNs exacerbates resource underutilization, as GPUs often idle while waiting for hidden states to be transferred. Using WD-GCN [31] on three real-world datasets, we observe that inter-GPU data transfer (including idle time) accounts for an average of 81% of forward

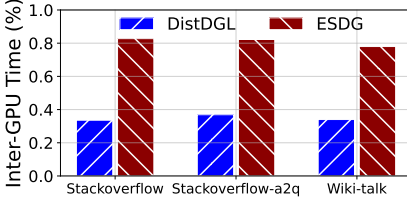


Fig. 2. Inter-GPU data transfer time during forward pass for DistDGL and ESDG.

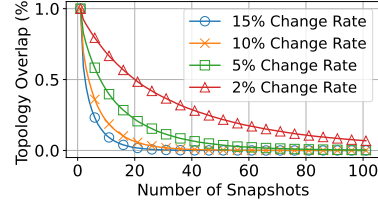


Fig. 3. Effect of change rate and snapshot count on graph topology overlap.

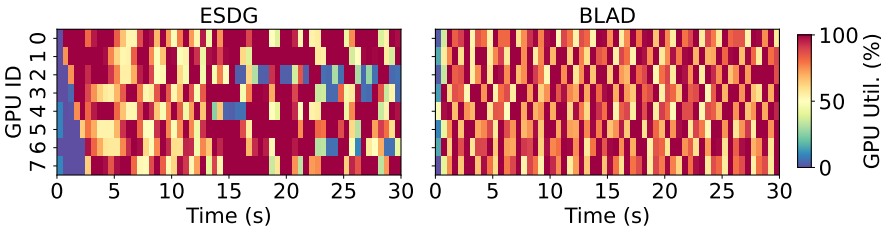


Fig. 4. GPU utilization of ESDG and BLAD on the Stackoverflow dataset.

pass time (Figure 2). In parallel, GPU utilization on Stackoverflow [47] dataset (Figure 4) also reveals severe underutilization of compute resources.

Snapshot Group-Based Partitioning. BLAD extends snapshot-level scheduling to snapshot groups to improve training efficiency. This approach reduces both neighbor feature transfer and hidden state data transfer, as snapshots within the same group are processed locally and sequentially. However, grouping multiple different snapshots together amplifies the imbalance problem: variations in graph size and structure accumulate within a group, making it harder to evenly distribute workload across GPUs. We also visualize GPU utilization under BLAD in Figure 4, which shows uneven utilization caused by imbalance.

Ineffective Utilization of Topological Similarity. Current approaches struggle to leverage topological similarity to accelerate distributed DGNN training, ranging from ignoring structural evolution to suffering from diminishing reuse potential and limited model applicability. For instance, DGC [5] reuses graph embeddings across epochs, but ignores topological similarity between graphs. In the single-GPU setting, PiPAD [52] adopts a global structured reuse scheme by identifying the global intersection across snapshots. However, its scalability is limited: as the number of snapshots grows, the globally shared intersection quickly shrinks. With fixed edge change rates (5%, 10% and 15%), the intersection ratio drops below 20% after 30 snapshots (Figure 3), diminishing its reuse potential. Similarly, reINC [21] employs incremental aggregation to prune redundant computations. Moreover, neither PiPAD nor reINC extends to attention-based GNNs such as GATs, as efficiently maintaining the global non-linear softmax normalization using only incremental information remains challenging. Furthermore, in the distributed context, reINC also fails to co-optimize reuse with load balancing. It ignores the computational heterogeneity induced by varying reuse rates across snapshots, which exacerbates load imbalance and straggler effects among GPUs.

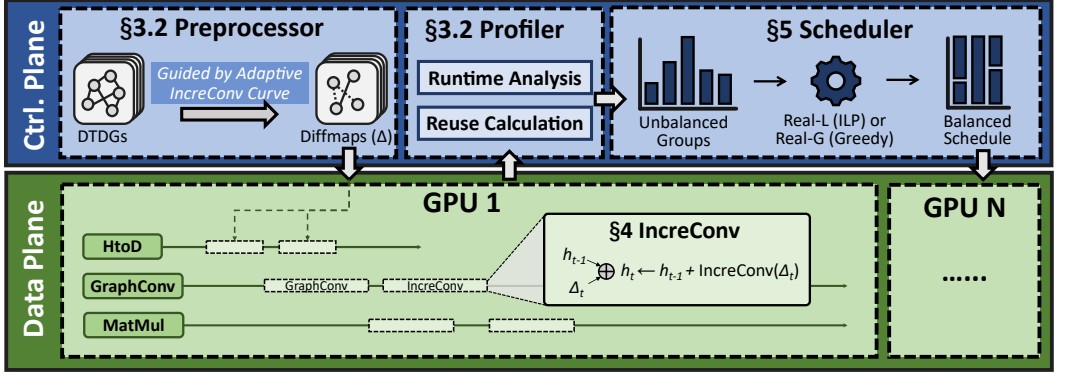


Fig. 5. **Overview of the REAL.** The system decouples the Control Plane (CPU) for preprocessing and scheduling from the Data Plane (GPU) for parallel training. Guided by runtime profiles, the Scheduler selects between the ILP solver (REAL-L) and the greedy heuristic (REAL-G) to dispatch reuse-optimized, load-balanced workloads.

3 SYSTEM OVERVIEW

In this section, we first introduce our design principles (§3.1), then present the architecture of REAL in detail (§3.2), and detail REAL’s pipelined execution design (§3.3).

3.1 Design Principles

Minimize Data Transfer. REAL employs a *snapshot group-based partitioning strategy* that limits cross-GPU data transfer to all-reduce gradient only. In addition, REAL proactively identifies *reusable graph structures* across snapshots to reduce redundant CPU-GPU data transfer.

Maximize Data Reuse. REAL proactively considers *both snapshot-level and group-level topology similarity*, optimizing data reuse and reducing redundant computations.

Maximize Resource Utilization. REAL improves resource utilization through two techniques: (1) a *triple-stream pipeline* that overlaps heterogeneous operations, and (2) *cross-group scheduling* that balances workloads across GPUs.

3.2 REAL Architecture

Figure 5 illustrates the architecture of REAL. To minimize resource contention, REAL adopts a two-plane design that decouples coordination from execution. The control plane (CPU) orchestrates graph preprocessing, profiling, and scheduling, while the data plane (GPU) focuses exclusively on distributed DGNN training. A detailed workflow is provided in Appendix A. The core components are described below.

Preprocessor. The preprocessor analyzes the DTDG topology to guide efficient scheduling and execution. As an offline profiling step, the preprocessor fits degree–runtime cost curves on synthetic graphs only once during system initialization; the curves are reused across different datasets as the basis for operator selection. Then, for each snapshot \mathcal{G}_i ($i > 1$), it computes a sparse difference map (dmap) by diffing against its predecessor \mathcal{G}_{i-1} , which captures localized structural changes. Finally, depending on the estimated cost from dmap and the fitted curves, the preprocessor selects for each snapshot an appropriate form, either the full graph or its dmap. For example, a four-size group can be constructed as $\mathcal{SG}_i = [\mathcal{G}_i, \text{dmap}_{i+1}, \mathcal{G}_{i+2}, \text{dmap}_{i+3}]$. Details are provided in §4.1.

Profiler. The profiler quantifies training costs and reuse opportunities to drive the scheduler's decision-making.. It first collects fine-grained training time details, including graph-level metrics such as HtoD, Conv, and MatMul times, as well as group-level RNN and backward times. To ensure reliable statistics, the profiler analyzes multiple training epochs and eliminates outliers caused by runtime noise. Based on these profiled values, it further estimates the potential reuse time $\mathcal{R}_{i,j}$ when two groups \mathcal{SG}_i and \mathcal{SG}_j are scheduled to the same GPU in the same iteration. This reuse time captures the overlapping data transfer and computation in HtoD and GNN phases. Both the profiled timeline and $\mathcal{R}_{i,j}$ table are later used by the scheduler to simulate the total epoch time.

Scheduler. The scheduler estimates the training time of an epoch by leveraging profiling statistics, incorporating both the standalone execution cost of each group and the reuse $\mathcal{R}_{i,j}$ when groups with overlapping snapshots are co-located on the same GPU. Based on this performance model, REAL supports two scheduling backends: an ILP-based solver (REAL-L) and a greedy heuristic (REAL-G), described in detail in §5. To balance scheduling quality and runtime efficiency, REAL adopts a dynamic selection strategy. At the beginning of training, REAL first attempts to solve the scheduling problem using REAL-L, and a timer is started upon invocation. If the solver fails to return a feasible schedule within a predefined threshold, the scheduler falls back to REAL-G.

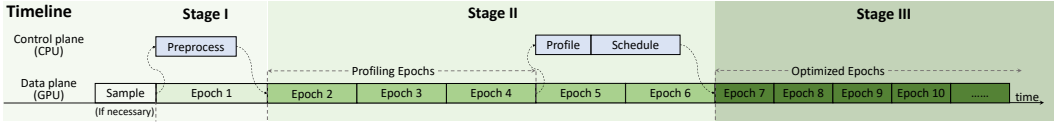


Fig. 6. **Pipelined execution timeline.** *REAL overlaps control plane overheads (preprocessing and scheduling) with data plane training to hide latency, enabling a non-blocking transition to the optimized schedule in Stage III.*

3.3 Overlapping Control Plane Overhead

Figure 6 shows the training timeline, illustrating how REAL overlaps control plane operations with data plane training. We decompose the training into three phases: Cold-start, Profiling, and Optimized Execution (visualized as distinct stages in Figure 6). During the first stage, the data plane initiates training without any optimization, while the control plane concurrently performs preprocessing on the input graph. Once preprocessing completes, REAL enters the second stage, where type I reuse (*i.e.*, snapshot-level reuse in §4.1) is enabled based on the constructed difference maps, allowing the data plane to benefit from reduced computation without waiting for full scheduling decisions. Concurrently, the control plane collects runtime statistics across several profiling epochs. These statistics drive the profiling and scheduling modules to compute an optimized execution plan for the remaining training. In the third stage, the data plane adopts the optimized schedule, which incorporates type II reuse (*i.e.*, group-level reuse in §4.2) and cross-group scheduling (§5), to further accelerate training. This pipelined design prevents CPU coordination from blocking GPU execution, sustaining high utilization.

4 TOPOLOGY-AWARE DATA REUSE

In this section, we present our design to exploit the topology reuse opportunities in DGNN training, including both snapshot-level and group-level data reuse.

4.1 Type I: Snapshot-Level Data Reuse

In DTDGs, consecutive snapshots are usually highly similar: most nodes and edges remain unchanged, with only a small fraction updated. However, existing distributed DGNN systems treat

updates locally:

$$w_{ij} = \begin{cases} +\frac{1}{\sqrt{d_i^{(t)} d_j^{(t)}}} & \text{dmap}(i, j) = +1 \\ -\frac{1}{\sqrt{d_i^{(t-1)} d_j^{(t-1)}}} & \text{dmap}(i, j) = -1 \end{cases} \quad (4)$$

Substituting this signed weight w_{ij} into Eq. (3) yields the final incremental embedding update. Crucially, this design significantly reduces memory overhead by eliminating the need to maintain a global edge weight table.

Softmax-Aware InCreConv for GAT. In contrast to GCN, where normalized edge weights only depend on node degrees and can be incrementally updated from the dmap with degree tracking, GAT poses an additional challenge: due to the softmax operation, each attention coefficient α_{ij} depends on the *entire* neighborhood of node j . Thus, even a single edge update alters the denominator shared by all neighbors, making it impossible to derive updated attention weights solely from the dmap. To incrementally update GAT aggregation, we maintain for each node j the softmax denominator $\mathcal{D}_j^{(t-1)}$ and the aggregation result $h_j^{(t-1)}$ from snapshot $t - 1$:

$$\mathcal{D}_j^{(t-1)} = \sum_{i \in \mathcal{N}(j, t-1)} \exp(e_{ij}^{(t-1)}) \quad (5)$$

$$h_j^{(t-1)} = \sum_{i \in \mathcal{N}(j, t-1)} \frac{\exp(e_{ij}^{(t-1)})}{\mathcal{D}_j^{(t-1)}} \mathcal{W} h_i \quad (6)$$

where $e_{ij}^{(t-1)}$ is the unnormalized attention score on edge (i, j) , and \mathcal{W} is the linear transformation weight. When snapshot t introduces an update set $\Delta\mathcal{E}_j$ of incident edges (with sign +1 for addition and -1 for deletion in dmap), we first compute the incremental change of the denominator:

$$\mathcal{D}_j^{(t)} = \mathcal{D}_j^{(t-1)} + \sum_{i \in \Delta\mathcal{E}_j} \text{dmap}(i, j) \cdot \exp(e_{ij}^{(t)}) \quad (7)$$

The previous aggregation is then rescaled and updated with the contributions of changed edges:

$$h_j^{(t)} = h_j^{(t-1)} \cdot \frac{\mathcal{D}_j^{(t-1)}}{\mathcal{D}_j^{(t)}} + \sum_{i \in \Delta\mathcal{E}_j} \text{dmap}(i, j) \cdot \frac{\exp(e_{ij}^{(t)})}{\mathcal{D}_j^{(t)}} \mathcal{W} h_i \quad (8)$$

This design allows unchanged neighbors to be updated in $O(1)$ by denominator rescaling, while only the affected edges incur $O(|\Delta\mathcal{E}_j|)$ cost. Overall, the per-snapshot complexity is reduced from $O(|\mathcal{E}|)$ to $O(|\Delta\mathcal{E}|)$, where $\Delta\mathcal{E}$ is the number of changed edges.

Adaptive InCreConv Execution. Selecting the most efficient aggregation target is non-trivial due to the non-monotonic performance characteristics of the GNN operator. Profiling indicates that GraphConv latency exhibits a U-shaped correlation with node density. Consequently, processing a dmap, despite having fewer edges, can paradoxically degrade performance compared to the original graph G in certain sparsity regimes. This stems from the irregularity of sparse structures, where the overhead of uncoalesced memory accesses outweighs the theoretical reduction in computational FLOPs. To reconcile this, we implement a self-adaptive selection mechanism based on a fitted degree-latency cost curve (see Appendix B). At runtime, the system estimates the execution costs for both \mathcal{G} and dmap based on their current average degrees and dispatches the computation to the path with the lowest predicted cost.

4.1.2 Transformation Phase Reuse. The applicability of transformation reuse depends on the dependency order between the MatMul and GraphConv. *Case 1: Pre-aggregation Transformation.* In architectures where feature transformation precedes aggregation, the input node features are *time-invariant* within a snapshot group. Consequently, the projection is computed only once for the first snapshot and shared across all subsequent snapshots, eliminating redundant computations entirely. *Case 2: Post-aggregation Transformation.* When transformation follows aggregation, the input to MatMul depends on the temporally evolving topology. Here, we employ *selective execution*: only nodes with topologically updated neighborhoods (identified via dmap) trigger new MatMul operations, while unaltered nodes directly reuse cached results from the preceding snapshot.

4.1.3 Triple-Stream Pipeline. To further improve end-to-end efficiency, we extend our execution to three parallel CUDA streams: (i) an HtoD stream for host-to-device data transfer, (ii) a Conv stream for Conv operations, and (iii) a MatMul stream for feature transformation. These three stages stress different hardware resources: PCIe bandwidth, HBM bandwidth, and GPU compute, respectively, and thus can be effectively overlapped. The streams are coordinated via `cudaEvent` to preserve inter-operator dependencies while maximizing concurrency. Compared to conventional sequential execution, our triple-stream pipeline enables fine-grained inter-snapshot interleaving. As illustrated in Figure 7, HtoD, Conv, and MatMul operations from different snapshots can overlap, effectively hiding latency and improving pipeline throughput. When transformation precedes GraphConv (e.g., GAT), the dependency becomes $\text{HtoD} \rightarrow \text{MatMul} \rightarrow \text{Conv}$. Here, MatMul is computed once and reused across snapshots, so the pipeline effectively overlaps HtoD and Conv.

4.2 Type II: Group-Level Data Reuse

Group-level data reuse leverages the property that model parameters remain frozen within a single distributed training iteration. Consequently, when consecutive snapshot groups share overlapping graph, their associated data can persist in resident GPU memory across group boundaries. This mechanism effectively eliminates redundant HtoD transfers and expensive GraphConv computations for the overlapping portions. To fully exploit this opportunity, the REAL scheduler adopts a similarity-driven scheduling strategy. Instead of adhering to a naive order that results in disjoint pairs (e.g., processing $[\mathcal{G}_1, \mathcal{G}_2]$ followed by $[\mathcal{G}_3, \mathcal{G}_4]$), the scheduler proactively reorders the execution sequence to prioritize snapshot overlap across consecutive groups (e.g., chaining $[\mathcal{G}_1, \mathcal{G}_2]$ with $[\mathcal{G}_2, \mathcal{G}_3]$). This similarity-aware rescheduling constructs a reuse-centric execution chain that significantly reduces I/O and compute overhead while maintaining workload balance. The concrete *cross-group scheduling* algorithm is detailed in § 5.

4.3 Reconstruction-Based Execution for Higher GNN Layers

Both Type I (snapshot-level) and Type II (group-level) data reuse mechanisms are applied exclusively to the first GNN layer in our design. This restriction stems from the temporal dependency introduced by the RNN module. After the first GNN layer, node representations are updated by the RNN, causing the features of the same node to diverge across snapshots even if the underlying topology remains unchanged. As a result, feature-level equivalence across snapshots no longer holds, rendering reuse across snapshots invalid for subsequent GNN layers. For higher GNN layers, we therefore adopt a different strategy. We reconstruct the full graph directly on the GPU by incrementally applying the dmap to the base snapshot of each group, whose first snapshot is always materialized as a full graph. This dmap-based reconstruction is performed entirely on the GPU and enables standard GraphConv execution, while avoiding host-to-device transfers of full graphs. We note that DGNNs typically adopt shallow architectures with two layers being the most common [24, 26, 27, 42, 64], making first-layer reuse sufficient to yield significant gains, as shown in our evaluation.

5 CROSS-GROUP SCHEDULING

In this section, we present our cross-group scheduling algorithm, which jointly considers inter-group data reuse and load balance to improve performance.

5.1 Optimization Objectives

We consider a set of snapshot groups with standalone execution times $\mathcal{SG} = \{t_1, \dots, t_n\}$, a set of available GPUs (devices) $\mathcal{D} = \{1, \dots, D\}$, and a sequence of training iterations $\mathcal{I} = \{1, \dots, m\}$. The primary objective of REAL is to minimize the total per-epoch training time $T = \sum_{i \in \mathcal{I}} T_i$, where T_i denotes the duration (makespan) of iteration i . We define binary decision variables $x_{k,i,j} \in \{0, 1\}$ to indicate whether snapshot group $k \in \mathcal{SG}$ is assigned to GPU $j \in \mathcal{D}$ in iteration $i \in \mathcal{I}$. Let t_k be the standalone execution time of group k , and \mathcal{R}_{k_1,k_2} be the reuse benefit (time saving) if groups k_1 and k_2 are co-located. Considering device memory limits and solver efficiency, we set a capacity cap of two groups per GPU per iteration. Consequently, the actual processing time (load) of GPU j in iteration i , denoted as $\mathcal{L}_{i,j}$, accounts for both computation cost and reuse reduction:

$$\mathcal{L}_{i,j} = \underbrace{\sum_{t_k \in \mathcal{SG}} t_k \cdot x_{k,i,j}}_{\text{Base Computation}} - \underbrace{\sum_{k_1 < k_2} \mathcal{R}_{k_1,k_2} \cdot \mathbb{I}(x_{k_1,i,j} = 1 \wedge x_{k_2,i,j} = 1)}_{\text{Reuse Benefit}}, \quad (9)$$

where $\mathbb{I}(\cdot)$ is the indicator function. The duration of iteration i is determined by the bottleneck GPU plus synchronization overhead α :

$$T_i = \max_{j \in \mathcal{D}} \{\mathcal{L}_{i,j}\} + \alpha. \quad (10)$$

Minimizing T requires jointly optimizing the assignment x to balance loads and maximize reuse under the capacity constraint $L = 2$. This problem generalizes the classical makespan minimization on parallel machines [15] and is NP-hard.

5.2 REAL-L: ILP-based Scheduling

To solve the formulation efficiently, we propose REAL-L, which models the problem as an Integer Linear Programming (ILP) task. Standard formulations of Eq. (10) often involve the max operator and logical AND conditions, which are non-linear. We linearize these components as follows:

Assignment and Capacity Constraints. We strictly enforce that every snapshot group is scheduled exactly once, and no GPU exceeds its capacity limit L :

$$\sum_{i \in \mathcal{I}} \sum_{j \in \mathcal{D}} x_{k,i,j} = 1, \quad \forall k \in \mathcal{SG}, \quad (11)$$

$$\sum_{k \in \mathcal{SG}} x_{k,i,j} \leq L, \quad \forall i \in \mathcal{I}, \forall j \in \mathcal{D}. \quad (12)$$

Reuse Linearization (McCormick Envelopes). To handle the conditional reuse term, we introduce auxiliary binary variables $y_{k_1,k_2,i,j}$ for all pairs $k_1 < k_2$. We enforce $y_{k_1,k_2,i,j} = 1 \iff x_{k_1,i,j} = 1 \wedge x_{k_2,i,j} = 1$ using standard McCormick envelope constraints:

$$\begin{aligned} y_{k_1,k_2,i,j} &\leq x_{k_1,i,j}, \\ y_{k_1,k_2,i,j} &\leq x_{k_2,i,j}, \\ y_{k_1,k_2,i,j} &\geq x_{k_1,i,j} + x_{k_2,i,j} - 1. \end{aligned} \quad \forall i \in \mathcal{I}, \forall j \in \mathcal{D}, \forall k_1 < k_2. \quad (13)$$

Min-Max Formulation via Epigraph Transformation. A direct linearization of the max operator in Eq. (10) typically requires Big-M constraints, which loosen the linear relaxation and hinder solver

performance. Instead, we treat T_i as a continuous decision variable and apply an epigraph-style formulation. We constrain T_i to be lower-bounded by the load of every GPU in iteration i :

$$T_i \geq \underbrace{\sum_{k \in \mathcal{SG}} t_k \cdot x_{k,i,j}}_{\text{Total Cost}} - \underbrace{\sum_{k_1 < k_2} \mathcal{R}_{k_1,k_2} \cdot y_{k_1,k_2,i,j}}_{\text{Reuse Benefit}} + \alpha, \quad \forall j \in \mathcal{D}. \quad (14)$$

Since the objective is to minimize $\sum T_i$, the solver naturally tightens T_i to the maximum load among GPUs ($\max_{j \in \mathcal{D}} L_{i,j} + \alpha$) without requiring auxiliary binary indicators. This formulation significantly tightens the relaxation bound, pruning the branch-and-bound search space efficiently. We implement REAL-L using the Gurobi solver [23]. Given the NP-hardness, we set a termination criterion based on a preset optimality gap (e.g., 10%) to balance schedule quality and solving time.

5.3 REAL-G: Greedy-based Scheduling

For large-scale scenarios, we employ REAL-G (Algorithm 1) as a scalable fallback upon ILP timeout. This heuristic iteratively generates candidate schedules and selects compatible groups to co-optimize reuse and load balance, thereby minimizing execution bubbles (or “wasted time”).

Input. The algorithm takes as input the GPU count D , the reuse matrix \mathcal{R} capturing type II reuse savings, and the group list $\mathcal{SG} = [t_1, t_2, \dots, t_n]$ with t_i being the execution time of group i .

Group preprocessing. Line 1 initializes preprocessing by calling Preprocess (lines 14–16), which inserts a *zero group* ($t_0 = 0$) as a dummy option. This ensures that any group can either pair with another group or with the zero group (i.e., run alone). All groups are then sorted by execution time t_i in ascending order to facilitate efficient matching.

Candidate target generation. Line 3 generates a set of potential makespan targets, denoted as *targets*, using GetCandidateTargets (lines 16–20). This procedure pairs the group with the longest remaining execution time (denoted as $\mathcal{SG}[n]$ after sorting) with other groups to project potential bottleneck durations for the current iteration.

Group selection and waste time calculation. Lines 5–9 evaluate each candidate target $\mathcal{T} \in \text{targets}$. For each target \mathcal{T} , GetPairs (lines 21–31) performs a two-point search to efficiently identify the group combination that best matches \mathcal{T} . The algorithm then calculates the wasted time W_i (lines 8–9) to quantify the quality of the schedule:

$$W_i = \underbrace{D \cdot (T_i - \alpha)}_{\text{Allocated Capacity}} - \underbrace{\sum_{j \in \mathcal{D}} \mathcal{L}_{i,j}}_{\text{Total Effective Load}}. \quad (15)$$

The algorithm selects the schedule that yields the minimum waste W_i , ensuring efficient group scheduling for the current iteration.

Complexity analysis. The overall time complexity of Algorithm 1 is dominated by the main loop while running at $O(n^3)$, making it computationally feasible for large-scale scenarios. In practice, when n reaches the $O(10^3)$ scale, the solving time still remains within a few seconds.

5.4 Scheduling Strategy Selection

To balance schedule quality and solving overhead, REAL employs a dynamic timeout mechanism bounded by the estimated epoch time T_{epoch} derived from profiling (Stage II in Figure 6). Specifically, we set the solver time limit to $\tau = 2 \times T_{epoch}$. If the ILP solver (REAL-L) does not converge within this threshold, the system seamlessly falls back to the greedy heuristic (REAL-G). This strategy ensures that the scheduling overhead is strictly bounded while maximizing the number of epochs that

Algorithm 1: REAL-G**Input** : GPU count D , Reuse matrix \mathcal{R} , Snapshot group times $\mathcal{SG} = [t_1, t_2, \dots, t_n]$ **Output**: Schedule strategy \mathbb{X} 1 Initialize $\mathcal{SG} \leftarrow \text{Preprocess}(\mathcal{SG})$, $n \leftarrow |\mathcal{SG}|$, $\mathbb{X} \leftarrow \emptyset$;2 **while** $|\mathcal{SG}| > D$ **do**3 $\text{targets} \leftarrow \text{GetCandidateTargets}(\mathcal{SG}, \mathcal{R}, n)$;4 $W_{\min} \leftarrow \infty$;5 **foreach** $\mathcal{T} \in \text{targets}$ **do**6 $\text{pairs} \leftarrow \text{GetPairs}(\mathcal{T}, \mathcal{R}, \mathcal{SG}, D, n)$;7 $W \leftarrow D \cdot \max_j (\text{pairs}[j].T, \mathcal{T})$;8 $W \leftarrow W - \sum_{j=1}^{D-1} \text{pairs}[j].T$;9 **if** $W < W_{\min}$ **then** $W_{\min} \leftarrow W$;10 Update \mathcal{SG}, \mathbb{X} ;11 Record rest groups in \mathbb{X} ;12 **return** \mathbb{X} ;13 **Function** $\text{Preprocess}(\mathcal{SG})$:14 Add $t_0 = 0$ to \mathcal{SG} and sort ascending;15 **return** \mathcal{SG} ;16 **Function** $\text{GetCandidateTargets}(\mathcal{SG}, \mathcal{R}, n)$:17 Initialize $\text{targets} \leftarrow \emptyset$;18 **for** $i = 0$ **to** $n - 1$ **do**19 Add $\{\mathcal{SG}[i] + \mathcal{SG}[n] - \mathcal{R}[i][n]\}$ to targets ;20 **return** targets ;21 **Function** $\text{GetPairs}(\mathcal{T}, \mathcal{R}, \mathcal{SG}, D, n)$:22 Initialize $\text{pairs} \leftarrow \emptyset$;23 **while** $|\text{pairs}| < D - 1$ **do**24 Initialize $\text{head}, \text{tail}, \text{best_pair}$;25 **while** $\text{head} < \text{tail}$ **do**26 $T \leftarrow \mathcal{SG}[\text{head}] + \mathcal{SG}[\text{tail}] - \mathcal{R}[\text{head}][\text{tail}]$;27 **if** $|T - \mathcal{T}| < |\text{best_pair}.T - \mathcal{T}|$ **then**28 $\text{best_pair} \leftarrow [T, \text{head}, \text{tail}]$;29 $T < \mathcal{T} ? \text{head}++ : \text{tail}--$ 30 Add best_pair to pairs ;31 **return** pairs ;

benefit from optimized execution, whether via the globally optimal plan from ILP or the efficient approximation from the greedy fallback.

6 IMPLEMENTATION

We implement REAL on PyTorch [37] (training backend) and DGL [54] (graph backend). REAL introduces system-level extensions in two places: the GNN operator and the data-loading runtime.

Component	Cluster A (4×8)	Cluster B (2×8)
CPU	192-core Intel Xeon Platinum 8558 @ 4.0 GHz	192-core Intel Xeon Platinum 8575C @ 4.0 GHz
GPU	8× NVIDIA H200	8× NVIDIA H20
Intra-node GPU BW	900 GB/s NVLink	900 GB/s NVLink
Inter-node Network	3.2 Tbps RoCE (8× ConnectX-7)	1.6 Tbps RoCE (4× ConnectX-7)

Table 2. **Hardware configuration of our testbeds.** *The main evaluation is conducted on Cluster A, while results for Cluster B are provided in the Appendix C.*

Operator extensions. We override DGL’s GraphConv and GATConv to expose intermediate aggregation states across snapshots and to support incremental execution. The operators take diffmaps as auxiliary inputs and selectively update the affected nodes while keeping full compatibility with DGL’s operator interfaces.

Scheduling runtime. We replace DGL’s default dataloader with a custom REALDataLoader. It treats snapshot groups as training units, materializes snapshots or diffmaps on demand and dispatches groups according to the REAL-L or REAL-G schedule. REALDataLoader integrates with PyTorch’s training loop through standard iterator semantics and performs asynchronous prefetching to overlap data preparation with GPU execution.

7 EVALUATION

In this section we first present our experimental setup, including the testbed, models, datasets, and baselines (§7.1). We evaluate REAL for training throughput, data transfer time, GraphConv FLOPs, and load imbalance (§7.2). We conduct ablation studies (§7.3), and confirm its accuracy (§7.4). We also assess scalability via simulations (§7.5).

7.1 Methodology

Testbed. We evaluate Real on two testbeds summarized in Table 2: a flagship H200-based cluster and a lower-tier H20-based cluster. Cluster A contains four servers, each equipped with 192-core Intel Xeon Platinum 8558 CPUs and eight NVIDIA H200 GPUs (141 GB HBM3e per GPU) interconnected via full-bandwidth NVLink (900 GB/s). Each server also integrates eight ConnectX-7 NICs, providing up to 3.2 Tbps RoCE bandwidth in multi-server experiments. Cluster B comprises two servers with 192-core Intel Xeon Platinum 8575C CPUs and eight NVIDIA H20 GPUs (141 GB HBM3e per GPU) per server, using the same NVLink (900 GB/s) topology and four ConnectX-7 NICs (1.6 Tbps RoCE) per server. We report results on Cluster A in the main text and include Cluster B results in the Appendix C. All experiments are conducted within the DGL NGC Container (version 24.07-py3) [34], which provides DGL v2.4 [54] for scalable graph processing, PyTorch v2.4.0 [37] as the deep learning backend, and CUDA 12.5 for GPU acceleration.

Datasets. We evaluate on six DTDG datasets (Table 3), including both real-world temporal networks and synthetic dynamic graphs derived from static datasets. The real-world datasets are Stackoverflow and Stackoverflow-a2q [47], which record user interactions on the Stack Exchange website, and Wiki-talk [35], which captures edits among Wikipedia users. We also construct synthetic dynamic graphs from three static graph datasets: Arxiv [25], Products [25], and Reddit [41]. Following BLAD [11], we create snapshots by randomly deleting some of the edges of the static graph. The evolution pattern of the number of edges in these snapshots mirrors the trend observed in the Stackoverflow dataset. For all datasets, the snapshot group size is set to four.

Benchmark DGNN models. We use four representative DGNNs: EvolveGCN [36], WD-GCN [31], TGCN [4], and GAT-LSTM [59]. The first three models are GCN-based DGNNs, while GAT-LSTM is a GAT-based model. These models are widely used due to their effectiveness in dynamic graph

Dataset	$ \overline{\mathcal{V}} $	$ \overline{\mathcal{E}} $	$ \overline{\mathcal{V}_d} $	$ \overline{\mathcal{E}_d} $	d_v	β	γ
Arxiv [25]	169.3k	1.8M	88.7k	242.6k	128	10.6	100
Products [25]	2.4M	36.7M	1.5M	5.9M	100	15.0	100
Reddit [41]	232.9k	31.9M	211.6k	5.5M	602	137.3	100
Stackoverflow [47]	2.6M	23.1M	790.3k	3.1M	128	8.9	100
Stackoverflow-a2q [47]	2.5M	10.8M	554.3k	1.5M	128	4.4	100
Wiki-talk [35]	1.1M	6.3M	182.8k	439.4k	128	5.5	100
Arxiv [†]	169.3k	1.8M	89.3k	245.9k	128	10.6	300
Products [†]	2.2M	18.6M	1.6M	3.6M	100	7.6	300

Table 3. **Attributes of the six datasets.** $|\overline{\mathcal{V}}|$ and $|\overline{\mathcal{E}}|$ are the average nodes and edges per snapshot; $|\overline{\mathcal{V}_d}|$ and $|\overline{\mathcal{E}_d}|$ are the average nodes and edges per dmap; d_v is the node feature dimension; β and γ are the average degree and snapshot count. The first six datasets are used for single-server (1x8 GPUs) experiments, while the $\gamma=300$ variants (Arxiv[†] and Products[†]) are constructed for multi-server evaluation.

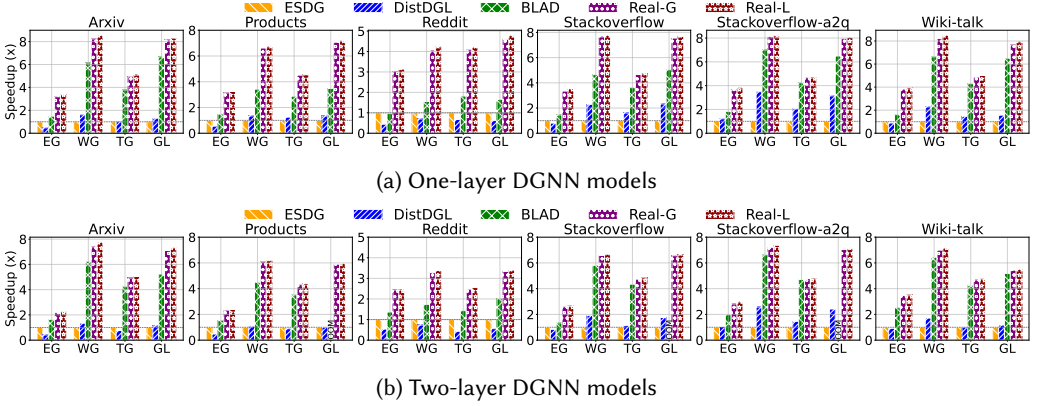


Fig. 8. **Training speedup of different systems across DGNN models and datasets on a single server (8 GPUs) of Cluster A.** Speedup is normalized to ESDG [3]. EG, WG, TG, and GL denote EvolveGCN, WD-GCN, TGCN, and GAT-LSTM, respectively.

learning. For each DGNN, we evaluate both one-layer and two-layer variants, where a layer consists of a spatial aggregation module (GNN) followed by a temporal update module (RNN). The hidden dimension is set to 64 across all model configurations.

Baselines. We compare REAL-G and REAL-L with existing state-of-the-art distributed DGNN training methods, including ESDG [3], DistDGL [66] and BLAD [11]. In ESDG, snapshots in a snapshot group are evenly distributed across GPUs based on their temporal intervals. DistDGL partitions each snapshot into subgraphs and distributes them across GPUs. BLAD utilizes a two-stage pipeline to collaboratively train two consecutive snapshot groups. In contrast, REAL-G and REAL-L execute two scheduled groups sequentially. DGC [5] is excluded due to its *lack of open-source code*, and DynaHB [46] is an *asynchronous* training framework, thus not included in the comparison.

7.2 Experimental results

7.2.1 Overall Performance. We first compare the training speedups of all methods on a single server in Cluster A. We prioritize this setting to evaluate the intrinsic system performance, preventing

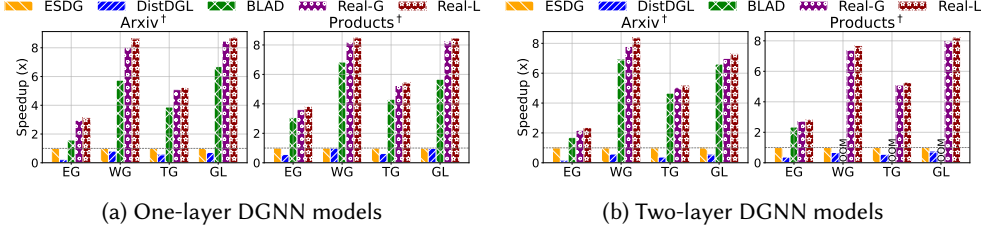


Fig. 9. **Training speedup of different systems across DGNN models and datasets on four servers (32 GPUs) of cluster A.** Speedup is also normalized to ESDG [3].

the heavy inter-GPU communication overhead of baselines from dominating the execution time due to limited inter-node bandwidth. As shown in Figure 8, for one-layer DGNN models, REAL-L delivers substantial gains over all baselines, achieving a $3.14\times$ – $8.56\times$ speedup over ESDG (avg. $5.73\times$), $2.26\times$ – $7.76\times$ over DistDGL (avg. $4.68\times$), and $1.11\times$ – $3.26\times$ over BLAD (avg. $1.83\times$). REAL-G also delivers consistent improvements, with a $3.05\times$ – $8.24\times$ speedup over ESDG (avg. $5.57\times$), $2.23\times$ – $7.42\times$ over DistDGL (avg. $4.52\times$), and $1.10\times$ – $3.16\times$ over BLAD (avg. $1.77\times$). For two-layer models, REAL also maintains its performance advantage. Specifically, REAL-L achieves speedups of $2.28\times$ – $7.77\times$ over ESDG (avg. $4.88\times$), $2.79\times$ – $6.57\times$ over DistDGL (avg. $4.63\times$), and $1.02\times$ – $1.99\times$ over BLAD (avg. $1.39\times$). Similarly, REAL-G achieves a $2.18\times$ – $7.39\times$ speedup over ESDG (avg. $4.78\times$), $2.65\times$ – $6.42\times$ over DistDGL (avg. $4.53\times$), and $1.02\times$ – $1.94\times$ over BLAD (avg. $1.36\times$). The gains of REAL are particularly pronounced in WD-GCN, TGCN, and GAT-LSTM, where all reuse mechanisms can be exploited. In contrast, for EvolveGCN, the weights in the MatMul stage change with each timestamp, preventing reuse in this phase. ESDG suffers from hidden state transfers between GPUs. The cost is minor for EvolveGCN with small hidden states, but becomes a major bottleneck for models with larger hidden states. DistDGL is further hindered by frequent neighbor aggregation, which limits its ability to scale throughput despite balanced workloads. BLAD, as the current SOTA, outperforms both ESDG and DistDGL mainly by reducing inter-GPU communication; however, its lack of data reuse and load balancing considerations still caps its maximum performance gains. Moreover, BLAD requires concurrently loading two groups of data into GPU memory to support its pipelining. This high peak memory footprint leads to Out-Of-Memory (OOM) errors in memory-intensive scenarios, preventing it from running on large datasets. Additionally, we provide performance results for three-layer DGNN models in Appendix D.1.

We further report results on 4 servers (32 GPUs) in Figure 9. While the speedups of REAL over ESDG and BLAD remain comparable to the single-server setting, the performance gap with DistDGL widens significantly—with REAL-L achieving average speedups of $9.93\times$ (One-layer) and $12.43\times$ (Two-layer). DistDGL is more adversely affected in the multi-machine setting, as its vertex-level partitioning necessitates cross-machine neighbor aggregation, incurring heavy communication overhead due to limited inter-node bandwidth. In contrast, ESDG avoids cross-machine hidden state transfers when the group size matches the per-server GPU count (i.e., four), effectively confining hidden states within each server and enabling efficient inter-group data parallelism.

7.2.2 Data transfer time breakdown. Figure 10 illustrates the breakdown of per-epoch data transfer time on a single server in Cluster A. In ESDG, inter-GPU hidden state transfer is the dominant cost for most models except EvolveGCN, and is particularly expensive when processing large single-graph datasets such as Stackoverflow. DistDGL suffers from neighbor aggregation overhead, especially on dense graphs or with high feature dimensions, where many neighbors are on different GPUs. For example, in the Reddit dataset, with an average degree of 137 and feature dimension of

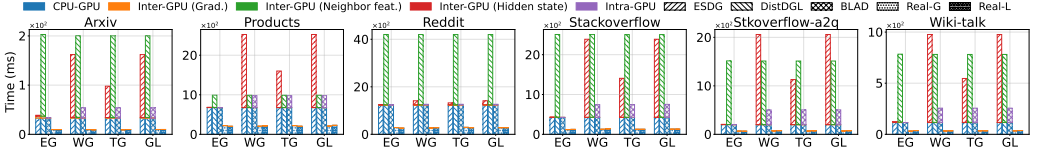


Fig. 10. **Breakdown of data transfer time on a single server in Cluster A.** The time is decomposed into CPU-GPU transfer, inter-GPU communication (gradients, neighbors, hidden states), and intra-GPU transfer.

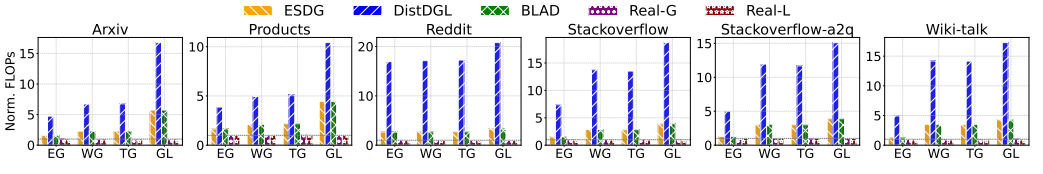


Fig. 11. **Normalized GraphConv FLOPs on a single server in Cluster A.** FLOPs is normalized to REAL-L.

602, this results in a particularly high transfer cost. BLAD avoids both types of cross-GPU data transfer but still loads the full graph, making HtoD transfer a bottleneck. Its pipeline further involves hidden state transfer across processes on the same GPU, which becomes significant on datasets such as Stackoverflow and Wiki. REAL also avoids both types of cross-GPU data transfer and, within each GPU, leverages DiffMap-based reuse to greatly reduce HtoD volume.

7.2.3 GraphConv FLOPs. We profile the total GraphConv floating-point operations (FLOPs) within one training epoch on a single server in Cluster A. As shown in Figure 11, REAL significantly reduces computational redundancy. Compared to DistDGL, REAL-G achieves a FLOPs reduction of $3.88\times$ – $21.70\times$ (avg. $11.92\times$), while REAL-L achieves $3.87\times$ – $20.80\times$ (avg. $11.64\times$). Against ESDG and BLAD, which perform full-graph computation and thus incur identical FLOPs, REAL-G reduces operations by $1.27\times$ – $6.45\times$ (avg. $2.95\times$), and REAL-L by $1.22\times$ – $5.74\times$ (avg. $2.87\times$). These savings stem from skipping full GraphConv on structurally unchanged nodes. DistDGL’s vertex-level partitioning still applies complete GraphConv to imported cross-GPU vertices whose outputs are never used in subsequent computation, incurring substantial redundant FLOPs. ESDG and BLAD overlook structural redundancy and perform full-graph convolution for every snapshot. In contrast, REAL leverages structural reuse to strictly limit computation to the incremental updates, effectively pruning these redundant operations.

7.2.4 GPU Utilization and Load Balance. We profile the utilization of all GPUs during the first 30 seconds of training. Due to space constraints, we only present the GPU utilization heatmaps of REAL-L in Figure 12, while the results for other baselines are provided in the Appendix D.2. The heatmaps demonstrate that REAL-L maintains a consistently high overall GPU utilization, mainly driven by the *Triple-Stream Pipeline* in Type I reuse. By effectively overlapping HtoD transfers with concurrent Conv and MatMul computation streams, the system successfully masks data movement latency and eliminates execution bubbles. Furthermore, REAL-L also achieves a remarkably balanced workload distribution across GPUs. This is primarily due to its *cross-group* scheduling, which optimizes group assignment to reduce inter-GPU variance and alleviate bottlenecks.

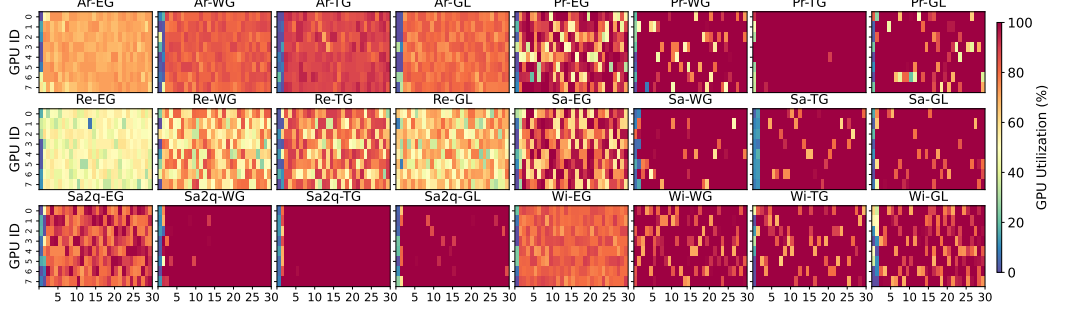


Fig. 12. GPU utilization heatmaps of REAL-L across different DGNN models and datasets on a single server in Cluster A. Ar, Pr, Re, Sa, Sa2q, and Wi denote the datasets Arxiv, Products, Reddit, Stackoverflow, Stackoverflow-a2q, and Wiki-talk, respectively.

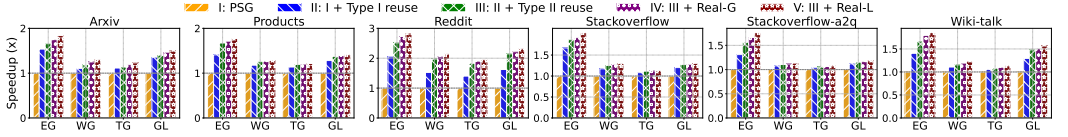


Fig. 13. Speedup ablation across different models and datasets. Speedup is normalized to PSG.

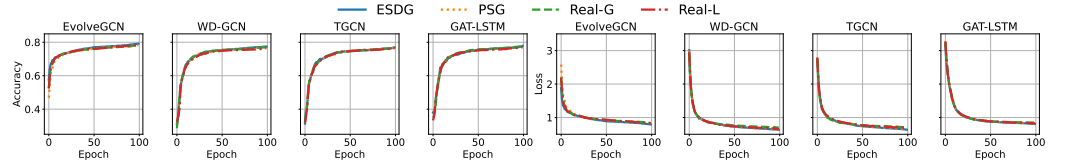


Fig. 14. Accuracy and loss curves on the Arxiv dataset. The trajectories of REAL align closely with the baseline (ESDG) across various models, confirming statistical correctness.

7.3 Ablation Study

To evaluate the individual contribution of each component in REAL, we perform a stepwise ablation study normalized to the *Partition-by-Snapshot-Group* (PSG) baseline, where each GPU sequentially trains one complete snapshot group per iteration without reuse or load balancing. Figure 13 reports cumulative speedups across models and datasets.

The PSG baseline achieves $1.00\times$ by definition, serving as the starting point. First, adding **type I reuse**—via IncrConv, selective MatMul, and the triple-stream pipeline—targets intra-group redundancy. This improves performance to $1.04\times$ – $2.05\times$ by effectively eliminating redundant computations and overlapping data transfers between snapshots. Next, incorporating **type II reuse**, which by default pairs two consecutive snapshot groups to maximize inter-group locality, further raises speedups to $1.05\times$ – $2.52\times$. This gain stems from removing redundant HtoD transfers and GraphConv operations across group boundaries. To address the potential imbalance introduced by reuse, the greedy cross-group scheduler (**REAL-G**) is introduced; it balances workloads while preserving reuse benefits, achieving $1.06\times$ – $2.73\times$. Finally, replacing the greedy scheduler with the ILP solver (**REAL-L**) delivers the highest cumulative speedups of $1.08\times$ – $2.81\times$, benefiting from a global view that optimizes the entire epoch’s schedule rather than making local decisions.

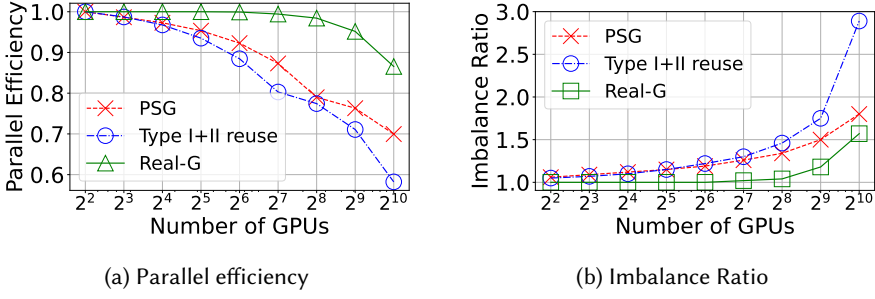


Fig. 15. Simulated parallel efficiency and imbalance ratio on clusters with 4 to 1024 GPUs.

7.4 Convergence Analysis

Compared to ESDG, REAL effectively increases the batch size per iteration by a factor of K . Given that the convergence dynamics of large-batch training are well-established [16, 18, 32], REAL achieve comparable accuracy to the baseline by proportionally scaling the learning rate. To formally validate this behavior, we analyze the convergence properties under standard assumptions (e.g., L -smoothness and bounded variance). Our analysis confirms that by aggregating gradients across N GPUs and K snapshot groups, the gradient variance is reduced by a factor of $\frac{1}{NK}$, ensuring convergence to a first-order stationary point at the standard rate of $\mathcal{O}(1/\sqrt{T})$ (detailed proofs in Appendix F). We also empirically verify this consistency in Figure 14 using the Arxiv dataset, where the training trajectories (accuracy and loss) of REAL closely align with the baseline. Additional evaluations on other datasets are provided in Appendix E.

7.5 Scalability

Due to limited hardware resources, we extend the evaluation to larger clusters through simulation. The fidelity of our simulation rests on two observations: **First**, the training time of a specific snapshot group on a single GPU exhibits *workload independence*, meaning it depends solely on the GPU's compute capability and the group's intrinsic characteristics (e.g., topology and feature size). Consequently, the execution time remains invariant regardless of the cluster scale, allowing us to leverage profiled real-world execution times as reliable ground truth. **Second**, the methods under comparison (PSG, Reuse-only, and REAL-G) incur *minimal communication overhead*, as the unique cross-GPU data transfer involved is the gradient AllReduce. Since DGNN models are typically compact, this overhead is negligible, accounting for less than 0.1% of the total training time in our measurements. Nevertheless, to ensure rigor, we explicitly model this overhead (α in Eq. (10)) following the communication modeling methodology from MG-WFBP [44].

To evaluate performance at larger scales, we conduct simulations using 10,000 snapshots generated by DyGraph [33]. In this scenario, solving with ILP-based methods is time-consuming, so we evaluate REAL-G, which can be solved quickly even at large scales, against baselines that only incorporate data reuse without load balancing and the naive PSG method. The scheduling strategies from different methods are applied to Eq. (10) for time simulation; detailed simulation settings are provided in the Appendix G. Figure 15 reports both *parallel efficiency*, defined as the ratio of per-GPU throughput at scale to the single-GPU throughput, and the *imbalance ratio*, defined as the training time of the most heavily loaded GPU divided by that of the least loaded one (excluding synchronization wait time). Figure 15 (left) shows that while baselines suffer from declining efficiency as the GPU count increases—particularly beyond 64 GPUs, REAL-G sustains 95% efficiency

at 512 GPUs and over 85% at 1024 GPUs. Figure 15 (right) links this scalability to lower imbalance ratios, where REAL-G consistently outperforms baselines.

8 RELATED WORK

Since distributed DGNN training has been discussed in the Introduction, we here discuss other relevant research landscapes.

Distributed GNN Training System. Distributed training systems for static graphs have laid the foundation for scalability by optimizing communication and resource utilization. Systems like P3 [12] and NeutronStar [55] minimize transfer overheads through pipelined push-pull strategies and hybrid dependency management, respectively, while Sancus [38] employs decentralized staleness-aware communication to bypass synchronization barriers. Addressing sampling bottlenecks, GNNLab [61] utilizes pre-sampling caching to optimize GPU memory, and DSP [2] dynamically adapts sampling policies for load balancing. At the kernel level, MGG [58] enables fine-grained intra-kernel pipelining to overlap sparse-dense operations. XGNN [49] abstracts unified global memory spaces for billion-scale datasets, F²CGT [30] leverages two-level feature compression to alleviate processing bottlenecks, and LeapGNN [6] proposes a feature-centric paradigm that migrates lightweight models to data locations to minimize retrieval overhead. However, these systems are designed for static graphs and do not account for temporal dependency or snapshot-level redundancy, making them orthogonal to distributed DGNN training.

Single-GPU DGNN Training Acceleration. DGNN training acceleration within a single GPU focuses on maximizing memory bandwidth efficiency and eliminating redundant computations. PiPAD [52] introduces a holistic pipeline that reconstructs the training paradigm to overlap graph loading, transfer, and computation, effectively hiding memory access latency. SEIGN [40] decouples spatial aggregation and temporal evolution to enable efficient mini-batch training without neighborhood sampling. ETC [13] further optimizes data access by enforcing temporal locality through a coherent execution plan, grouping contiguous snapshots to minimize cache misses. WinGNN [71] reduces model complexity by eliminating temporal encoders, instead employing a random gradient aggregation window to capture temporal dependencies. While effective for small-scale graphs, these single-GPU solutions hit the memory wall with growing datasets and fail to address the distributed challenges: cross-device communication planning and global load balancing.

CTDG Training Optimization. While REAL targets DTDGs, significant progress continues in optimizing CTDGs. TGL [67] establishes a baseline for billion-scale CTDG learning with specialized temporal data structures and parallel sampling. SIMPLE [14] optimizes CPU-GPU transfers through a dynamic data placement strategy that adapts to temporal access patterns. D3-GNN [22] employs a distributed dataflow architecture to efficiently manage cascading updates. CTAN [19] leverages ODEs to capture long-range dependencies, while SIG [8] integrates causal inference for interpretable link prediction. However, these systems assume continuous-time execution and are not directly applicable to snapshot-based DGNN training.

9 CONCLUSION

In this paper, we introduce REAL, an efficient distributed training system for DGNNs. REAL directly targets the three fundamental challenges of distributed DGNN training: maximize resource utilization, minimize data transfer, and maximize data reuse. REAL integrates snapshot-level and group-level reuse with cross-group scheduling, and provides two variants: REAL-L, which employs ILP to achieve globally optimized schedules, and REAL-G, a greedy fallback used when ILP solving times out. Extensive evaluation on six DTDG datasets and four DGNN models demonstrates that REAL-L and REAL-G deliver $1.11\times$ – $3.26\times$ and $1.10\times$ – $3.16\times$ speedup over SOTA baseline, respectively.

References

- [1] Guangji Bai, Chen Ling, and Liang Zhao. 2022. Temporal domain generalization with drift-aware dynamic neural networks. *arXiv preprint arXiv:2205.10664* (2022).
- [2] Zhenkun Cai, Qihui Zhou, Xiao Yan, Da Zheng, Xiang Song, Chenguang Zheng, James Cheng, and George Karypis. 2023. DSP: Efficient GNN Training with Multiple GPUs. In *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming* (Montreal, QC, Canada) (PPoPP '23). Association for Computing Machinery, New York, NY, USA, 392–404. <https://doi.org/10.1145/3572848.3577528>
- [3] Venkatesan T. Chakaravarthy, Shivmaran S. Pandian, Saurabh Raj, Yogish Sabharwal, Toyotaro Suzumura, and Shashanka Ubaru. 2021. Efficient scaling of dynamic graph neural networks (SC '21). Association for Computing Machinery, New York, NY, USA, Article 77, 15 pages. <https://doi.org/10.1145/3458817.3480858>
- [4] Bo Chen, Wei Guo, Ruiming Tang, Xin Xin, Yue Ding, Xiuqiang He, and Dong Wang. 2020. TGCN: Tag graph convolutional network for tag-aware recommendation. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. 155–164.
- [5] Fahao Chen, Peng Li, and Celimuge Wu. 2023. DGC: Training Dynamic Graphs with Spatio-Temporal Non-Uniformity using Graph Partitioning by Chunks. *Proc. ACM Manag. Data* 1, 4, Article 237 (dec 2023), 25 pages. <https://doi.org/10.1145/3626724>
- [6] Weijian Chen, Shuibing He, Haoyang Qu, and Xuechen Zhang. 2025. LeapGNN: Accelerating Distributed GNN Training Leveraging Feature-Centric Model Migration. In *23rd USENIX Conference on File and Storage Technologies (FAST 25)*. USENIX Association, Santa Clara, CA, 255–270. <https://www.usenix.org/conference/fast25/presentation/chen-weijian-leap>
- [7] Songgaojun Deng, Huzefa Rangwala, and Yue Ning. 2019. Learning Dynamic Context Graphs for Predicting Social Events. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Anchorage, AK, USA) (KDD '19). Association for Computing Machinery, New York, NY, USA, 1007–1016. <https://doi.org/10.1145/3292500.3330919>
- [8] Lanting Fang, Yulian Yang, Kai Wang, Shanshan Feng, Kaiyu Feng, Jie Gui, Shuliang Wang, and Yew-Soon Ong. 2024. SIG: Efficient Self-Interpretable Graph Neural Network for Continuous-time Dynamic Graphs. *arXiv:2405.19062* [cs.LG] <https://arxiv.org/abs/2405.19062>
- [9] Matthias Fey and Jan E. Lenssen. 2019. Fast Graph Representation Learning with PyTorch Geometric. In *ICLR Workshop on Representation Learning on Graphs and Manifolds*.
- [10] Matthias Fey, Jinu Sunil, Akihiro Nitta, Rishi Puri, Manan Shah, Blaž Stojanovič, Ramona Bendias, Alexandria Barghi, Vid Kocijan, Zecheng Zhang, Xinwei He, Jan Eric Lenssen, and Jure Leskovec. 2025. PyG 2.0: Scalable Learning on Real World Graphs. *arXiv:2507.16991* [cs.LG] <https://arxiv.org/abs/2507.16991>
- [11] Kaihua Fu, Quan Chen, Yuzhuo Yang, Jiuchen Shi, Chao Li, and Minyi Guo. 2023. BLAD: Adaptive Load Balanced Scheduling and Operator Overlap Pipeline For Accelerating The Dynamic GNN Training. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis* (Denver, CO, USA) (SC '23). Association for Computing Machinery, New York, NY, USA, Article 37, 13 pages. <https://doi.org/10.1145/3581784.3607040>
- [12] Swapnil Gandhi and Anand Padmanabha Iyer. 2021. P3: Distributed Deep Graph Learning at Scale. In *15th USENIX Symposium on Operating Systems Design and Implementation (OSDI 21)*. USENIX Association, 551–568. <https://www.usenix.org/conference/osdi21/presentation/gandhi>
- [13] Shihong Gao, Yiming Li, Yanyan Shen, Yingxia Shao, and Lei Chen. 2024. ETC: Efficient Training of Temporal Graph Neural Networks over Large-scale Dynamic Graphs. *Proceedings of the VLDB Endowment* 17, 5 (2024), 1060–1072.
- [14] Shihong Gao, Yiming Li, Xin Zhang, Yanyan Shen, Yingxia Shao, and Lei Chen. 2024. SIMPLE: Efficient Temporal Graph Neural Network Training at Scale with Dynamic Data Placement. *Proceedings of the ACM on Management of Data* 2, 3 (2024), 1–25.
- [15] Michael R. Garey and David S. Johnson. 1990. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., USA.
- [16] Saeed Ghadimi and Guanghui Lan. 2013. Stochastic First- and Zeroth-order Methods for Nonconvex Stochastic Programming. *arXiv:1309.5549* [math.OC] <https://arxiv.org/abs/1309.5549>
- [17] Palash Goyal, Sujit Rokka Chhetri, and Arquimedes Canedo. 2020. dyngraph2vec: Capturing network dynamics using dynamic graph representation learning. *Knowledge-Based Systems* 187 (2020), 104816.
- [18] Priya Goyal, Piotr Dollár, Ross Girshick, Pieter Noordhuis, Lukasz Wesolowski, Aapo Kyrola, Andrew Tulloch, Yangqing Jia, and Kaiming He. 2018. Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour. *arXiv:1706.02677* [cs.CV] <https://arxiv.org/abs/1706.02677>
- [19] Alessio Gravina, Giulio Lovisotto, Claudio Gallicchio, Davide Bacciu, and Claas Grohnfeldt. 2024. Long Range Propagation on Continuous-Time Dynamic Graphs. *arXiv:2406.02740* [cs.LG] <https://arxiv.org/abs/2406.02740>

- [20] Mingyu Guan, Anand Padmanabha Iyer, and Taesoo Kim. 2022. DynaGraph: dynamic graph neural networks at scale. In *Proceedings of the 5th ACM SIGMOD Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA)*. 1–10.
- [21] Mingyu Guan, Saumia Singhal, Taesoo Kim, and Anand Padmanabha Iyer. 2025. ReInc: Scaling Training of Dynamic Graph Neural Networks. arXiv:2501.15348 [cs.LG] <https://arxiv.org/abs/2501.15348>
- [22] Rustam Guliyev, Aparajita Haldar, and Hakan Ferhatosmanoglu. 2024. D3-GNN: Dynamic Distributed Dataflow for Streaming Graph Neural Networks. *Proceedings of the VLDB Endowment* 17, 11 (July 2024), 2764–2777. <https://doi.org/10.14778/3681954.3681961>
- [23] Gurobi Optimization, LLC. 2022. Gurobi - The Fastest Solver. <https://www.gurobi.com> Accessed: 2022-01-01.
- [24] William L. Hamilton, Rex Ying, and Jure Leskovec. 2018. Inductive Representation Learning on Large Graphs. arXiv:1706.02216 [cs.SI] <https://arxiv.org/abs/1706.02216>
- [25] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, and Jure Leskovec. 2020. Open graph benchmark: Datasets for machine learning on graphs. *Advances in neural information processing systems* 33 (2020), 22118–22133.
- [26] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. arXiv:1609.02907 [cs.LG] <https://arxiv.org/abs/1609.02907>
- [27] Chao Li, Runshuo Liu, Jinhu Fu, Zhongying Zhao, Hua Duan, and Qingtian Zeng. 2024. DGNN-MN: Dynamic Graph Neural Network via memory regenerate and neighbor propagation. *Applied Intelligence* 54, 19 (July 2024), 9253–9268. <https://doi.org/10.1007/s10489-024-05500-3>
- [28] Haoyang Li and Lei Chen. 2021. Cache-based gnn system for dynamic graphs. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*. 937–946.
- [29] Hongxi Li, Zuxuan Zhang, Dengzhe Liang, and Yuncheng Jiang. 2024. K-Truss Based Temporal Graph Convolutional Network for Dynamic Graphs. In *Asian Conference on Machine Learning*. PMLR, 739–754.
- [30] Yuxin Ma, Ping Gong, Tianming Wu, Jiawei Yi, Chengru Yang, Cheng Li, Qirong Peng, Guiming Xie, Yongcheng Bao, Haifeng Liu, and Yinlong Xu. 2024. Eliminating Data Processing Bottlenecks in GNN Training over Large Graphs via Two-level Feature Compression. *Proc. VLDB Endow.* 17, 11 (July 2024), 2854–2866. <https://doi.org/10.14778/3681954.3681968>
- [31] Franco Manessi, Alessandro Rozza, and Mario Manzo. 2020. Dynamic graph convolutional networks. *Pattern Recognition* 97 (2020), 107000.
- [32] Sam McCandlish, Jared Kaplan, Dario Amodei, and OpenAI Dota Team. 2018. An Empirical Model of Large-Batch Training. arXiv:1812.06162 [cs.LG] <https://arxiv.org/abs/1812.06162>
- [33] Andrew McCrabb, Hellina Nigatu, Absalat Getachew, and Valeria Bertacco. 2022. DyGraph: a dynamic graph generator and benchmark suite. In *Proceedings of the 5th ACM SIGMOD Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA)* (Philadelphia, Pennsylvania) (GRADES-NDA '22). Association for Computing Machinery, New York, NY, USA, Article 7, 8 pages. <https://doi.org/10.1145/3534540.3534692>
- [34] NVIDIA. 2024. Deep Graph Library (DGL) Container. <https://catalog.ngc.nvidia.com/orgs/nvidia/containers/dgl>. Accessed: 2025-08-20.
- [35] Ashwin Paranjape, Austin R. Benson, and Jure Leskovec. 2017. Motifs in Temporal Networks. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining* (Cambridge, United Kingdom) (WSDM '17). Association for Computing Machinery, New York, NY, USA, 601–610. <https://doi.org/10.1145/3018661.3018731>
- [36] Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao Schardl, and Charles Leiserson. 2020. Evolvegcn: Evolving graph convolutional networks for dynamic graphs. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 34. 5363–5370.
- [37] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Yang, Zach DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. *PyTorch: an imperative style, high-performance deep learning library*. Curran Associates Inc., Red Hook, NY, USA.
- [38] Jingshu Peng, Zhao Chen, Yingxia Shao, Yanyan Shen, Lei Chen, and Jiannong Cao. 2022. Sancus: staleness-aware communication-avoiding full-graph decentralized training in large-scale graph neural networks. *Proc. VLDB Endow.* 15, 9 (May 2022), 1937–1950. <https://doi.org/10.14778/3538598.3538614>
- [39] Xiao Qin, Nasrullah Sheikh, Chuan Lei, Berthold Reinwald, and Giacomo Domeniconi. 2023. Seign: A simple and efficient graph neural network for large dynamic graphs. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. IEEE, 2850–2863.
- [40] Xiao Qin, Nasrullah Sheikh, Chuan Lei, Berthold Reinwald, and Giacomo Domeniconi. 2023. SEIGN: A Simple and Efficient Graph Neural Network for Large Dynamic Graphs. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*. 2850–2863. <https://doi.org/10.1109/ICDE55515.2023.00218>
- [41] Reddit. n.d.. Reddit Dataset. <https://www.reddit.com/>.

- [42] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. 2020. Temporal Graph Networks for Deep Learning on Dynamic Graphs. *arXiv:2006.10637* [cs.LG] <https://arxiv.org/abs/2006.10637>
- [43] Aravind Sankar, Yanhong Wu, Liang Gou, Wei Zhang, and Hao Yang. 2020. Dysat: Deep neural representation learning on dynamic graphs via self-attention networks. In *Proceedings of the 13th international conference on web search and data mining*. 519–527.
- [44] Shaohuai Shi, Xiaowen Chu, and Bo Li. 2021. MG-WFBP: Merging Gradients Wisely for Efficient Communication in Distributed Deep Learning. *IEEE Transactions on Parallel and Distributed Systems* 32, 8 (2021), 1903–1917. <https://doi.org/10.1109/TPDS.2021.3052862>
- [45] Xinkai Song, Tian Zhi, Zhe Fan, Zhenxing Zhang, Xi Zeng, Wei Li, Xing Hu, Zidong Du, Qi Guo, and Yunji Chen. 2021. Cambricon-G: A polyvalent energy-efficient accelerator for dynamic graph neural networks. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 41, 1 (2021), 116–128.
- [46] Zhen Song, Yu Gu, Qing Sun, Tianyi Li, Yanfeng Zhang, Yushuai Li, Christian S Jensen, and Ge Yu. 2024. DynaHB: A Communication-Avoiding Asynchronous Distributed Framework with Hybrid Batches for Dynamic GNN Training. *Proceedings of the VLDB Endowment* 17, 11 (2024), 3388–3401.
- [47] Stack-Overflow. 2023. Stack-Overflow Dataset. <https://snap.stanford.edu/data/sx-stackoverflow.html>. Accessed: [Insert the date you accessed the dataset].
- [48] Junwei Su, Difan Zou, and Chuan Wu. 2024. PRES: Toward Scalable Memory-Based Dynamic Graph Neural Networks. *arXiv preprint arXiv:2402.04284* (2024).
- [49] Dahai Tang, Jiali Wang, Rong Chen, Lei Wang, Wenyuan Yu, Jingren Zhou, and Kenli Li. 2024. XGNN: Boosting Multi-GPU GNN Training via Global GNN Memory Store. *Proc. VLDB Endow.* 17, 5 (Jan. 2024), 1105–1118. <https://doi.org/10.14778/3641204.3641219>
- [50] Yuxing Tian, Yiyan Qi, and Fan Guo. 2023. FreeDyG: Frequency Enhanced Continuous-Time Dynamic Graph Model for Link Prediction. In *The Twelfth International Conference on Learning Representations*.
- [51] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. 2019. Dyrep: Learning representations over dynamic graphs. In *International conference on learning representations*.
- [52] Chunyang Wang, Desen Sun, and Yuebin Bai. 2023. PiPAD: pipelined and parallel dynamic GNN training on GPUs. In *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming*. 405–418.
- [53] Junshan Wang, Wenhao Zhu, Guojie Song, and Liang Wang. 2022. Streaming graph neural networks with generative replay. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 1878–1888.
- [54] Minjie Yu Wang. 2019. Deep graph library: Towards efficient and scalable deep learning on graphs. In *ICLR workshop on representation learning on graphs and manifolds*.
- [55] Qiange Wang, Yanfeng Zhang, Hao Wang, Chaoyi Chen, Xiaodong Zhang, and Ge Yu. 2022. NeutronStar: Distributed GNN Training with Hybrid Dependency Management. In *Proceedings of the 2022 International Conference on Management of Data* (Philadelphia, PA, USA) (*SIGMOD '22*). Association for Computing Machinery, New York, NY, USA, 1301–1315. <https://doi.org/10.1145/3514221.3526134>
- [56] Xuhong Wang, Ding Lyu, Mengjian Li, Yang Xia, Qi Yang, Xinwen Wang, Xinguang Wang, Ping Cui, Yupu Yang, Bowen Sun, and Zhenyu Guo. 2021. APAN: Asynchronous Propagation Attention Network for Real-time Temporal Graph Embedding. In *Proceedings of the 2021 International Conference on Management of Data* (Virtual Event, China) (*SIGMOD '21*). Association for Computing Machinery, New York, NY, USA, 2628–2638. <https://doi.org/10.1145/3448016.3457564>
- [57] Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. 2021. Inductive Representation Learning in Temporal Networks via Causal Anonymous Walks. In *International Conference on Learning Representations (ICLR)*.
- [58] Yuke Wang, Boyuan Feng, Zheng Wang, Tong Geng, Kevin Barker, Ang Li, and Yufei Ding. 2023. MGG: Accelerating Graph Neural Networks with Fine-Grained Intra-Kernel Communication-Computation Pipelining on Multi-GPU Platforms. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*. USENIX Association, Boston, MA, 779–795. <https://www.usenix.org/conference/osdi23/presentation/wang-yuke>
- [59] Tianlong Wu, Feng Chen, and Yun Wan. 2018. Graph attention LSTM network: A new model for traffic flow forecasting. In *2018 5th international conference on information science and control engineering (ICISCE)*. IEEE, 241–245.
- [60] Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. 2020. Inductive representation learning on temporal graphs. *arXiv preprint arXiv:2002.07962* (2020).
- [61] Jianbang Yang, Dahai Tang, Xiaoni Song, Lei Wang, Qiang Yin, Rong Chen, Wenyuan Yu, and Jingren Zhou. 2022. GNNLab: a factored system for sample-based GNN training over GPUs. In *Proceedings of the Seventeenth European Conference on Computer Systems* (Rennes, France) (*EuroSys '22*). Association for Computing Machinery, New York, NY, USA, 417–434. <https://doi.org/10.1145/3492321.3519557>
- [62] Jiaxuan You, Tianyu Du, and Jure Leskovec. 2022. ROLAND: graph learning framework for dynamic graphs. In *Proceedings of the 28th ACM SIGKDD conference on knowledge discovery and data mining*. 2358–2366.

- [63] Bing Yu, Haoteng Yin, and Zhanxing Zhu. 2018. Spatio-Temporal Graph Convolutional Networks: A Deep Learning Framework for Traffic Forecasting. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*. International Joint Conferences on Artificial Intelligence Organization, 3634–3640. <https://doi.org/10.24963/ijcai.2018/505>
- [64] Haonan Yuan, Qingyun Sun, Xingcheng Fu, Cheng Ji, and Jianxin Li. 2024. Dynamic Graph Information Bottleneck. arXiv:2402.06716 [cs.LG] <https://arxiv.org/abs/2402.06716>
- [65] Zeyang Zhang, Xin Wang, Ziwei Zhang, Zhou Qin, Weigao Wen, Hui Xue, Haoyang Li, and Wenwu Zhu. 2024. Spectral invariant learning for dynamic graphs under distribution shifts. *Advances in Neural Information Processing Systems* 36 (2024).
- [66] Da Zheng, Chao Ma, Minjie Wang, Jinjing Zhou, Qidong Su, Xiang Song, Quan Gan, Zheng Zhang, and George Karypis. 2021. DistDGL: Distributed Graph Neural Network Training for Billion-Scale Graphs. arXiv:2010.05337 [cs.LG] <https://arxiv.org/abs/2010.05337>
- [67] Hongkuan Zhou, Da Zheng, Israt Nisa, Vasileios Ioannidis, Xiang Song, and George Karypis. 2022. TGL: a general framework for temporal GNN training on billion-scale graphs. *Proceedings of the VLDB Endowment* 15, 8 (2022), 1572–1580.
- [68] Lekui Zhou, Yang Yang, Xiang Ren, Fei Wu, and Yueting Zhuang. 2018. Dynamic network embedding by modeling triadic closure process. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.
- [69] Linhong Zhu, Dong Guo, Junming Yin, Greg Ver Steeg, and Aram Galstyan. 2016. Scalable temporal latent space inference for link prediction in dynamic social networks. *IEEE Transactions on Knowledge and Data Engineering* 28, 10 (2016), 2765–2777.
- [70] Rong Zhu, Kun Zhao, Hongxia Yang, Wei Lin, Chang Zhou, Baole Ai, Yong Li, and Jingren Zhou. 2019. AliGraph: A Comprehensive Graph Neural Network Platform. arXiv:1902.08730 [cs.DC] <https://arxiv.org/abs/1902.08730>
- [71] Yifan Zhu, Fangpeng Cong, Dan Zhang, Wenwen Gong, Qika Lin, Wenzheng Feng, Yuxiao Dong, and Jie Tang. 2023. WinGNN: Dynamic Graph Neural Networks with Random Gradient Aggregation Window. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (Long Beach, CA, USA) (KDD '23)*. Association for Computing Machinery, New York, NY, USA, 3650–3662. <https://doi.org/10.1145/3580305.3599551>

Appendix

A Detailed System Workflow

Figure 16 illustrates the end-to-end execution pipeline of REAL, which decouples logical coordination (Control Plane) from parallel computation (Data Plane). The workflow proceeds in three logical phases: preprocessing, profiling, and scheduling.

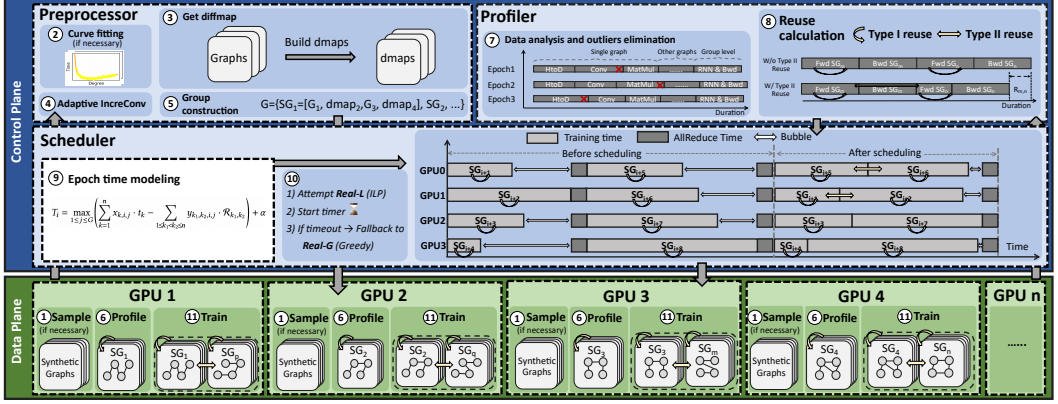


Fig. 16. **The detailed execution workflow of REAL.** The system coordinates the Control Plane (preprocessing, profiling, scheduling) and Data Plane (sampling, training) to optimize resource utilization and data reuse.

Initialization and Preprocessing. The workflow initiates with cost modeling and data preparation. The Data Plane first samples synthetic graphs (Step ①) to capture hardware-specific runtime characteristics. Using these samples, the Preprocessor fits degree-latency curves (Step ②) to guide the Adaptive InceConv operator. Simultaneously, it computes sparse difference maps (dmaps) from raw snapshots (Step ③) and determines the optimal execution path—either full-graph convolution or incremental aggregation—based on the fitted curves (Step ④). Finally, snapshots and dmaps are assembled into snapshot groups (Step ⑤), serving as the atomic units for scheduling.

Profiling and Reuse Modeling To drive the scheduler, REAL collects precise runtime statistics. The Data Plane executes a limited number of profiling epochs (Step ⑥), recording fine-grained kernel durations and transfer latencies. The Profiler cleans this data by eliminating runtime noise and outliers (Step ⑦). Crucially, it quantifies the potential latency reduction \mathcal{R}_{k_1, k_2} derived from both snapshot-level (Type I) and group-level (Type II) reuse (Step ⑧), mapping the theoretical reuse opportunities to concrete time savings.

Scheduling and Execution In the final phase, the Scheduler formulates the workload distribution as a makespan minimization problem using the profiled cost model (Step ⑨). It employs a hybrid solving strategy (Step ⑩): it first attempts to solve the Integer Linear Programming (ILP) formulation (REAL-L); if the solver exceeds the strict time budget τ , it seamlessly falls back to the greedy heuristic (REAL-G). The optimized schedule X , which defines the group-to-GPU assignment and execution order, is then dispatched to the Data Plane for high-throughput distributed training (Step ⑪).

B Performance Modeling for InceConv

Figure 17 profiles the runtime of the GraphConv operator across varying graph scales (100k to 5M nodes) and average degrees. The x-axis represents the average node degree on a logarithmic scale, while the y-axis reports execution time (ms). A key observation is that GraphConv exhibits a characteristic **U-shaped latency curve**: runtime initially decreases as the graph transitions from

extremely sparse to moderately connected, but subsequently rises in high-degree regimes due to increased aggregation overheads. We capture this non-monotonic behavior using the following regression model:

$$T_{\text{graphconv}}(N, d) = a_0 \cdot N + a_1 \cdot d + \frac{a_2}{d} + a_3,$$

where N and d denote the node count and average degree, respectively.

This fitted model serves as the oracle for our adaptive execution strategy. At runtime, we evaluate the predicted latency for two potential execution paths: (i) applying GraphConv on the full baseline graph \mathcal{G} , and (ii) applying GraphConv on the differential dmap. The system then dynamically dispatches the computation to the target that yields the minimized predicted cost, ensuring optimal efficiency across varying density regimes.

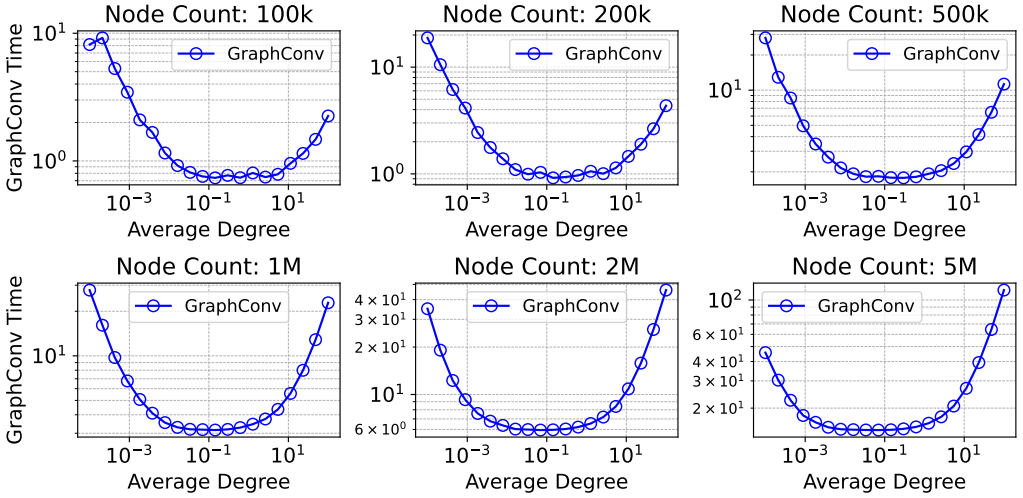


Fig. 17. GraphConv runtime exhibits a U-shaped trend with respect to node degree.

C Supplementary Evaluation on Cluster B

C.1 Experimental results

C.1.1 Overall Performance. Figure 18a illustrates the training speedups on Cluster B. In the single-server setting, REAL-L consistently outperforms all baselines, achieving $2.71\times$ – $9.28\times$ speedup over ESDG (avg. $5.98\times$), $1.13\times$ – $4.53\times$ over BLAD (avg. $2.03\times$), and $1.94\times$ – $8.92\times$ over DistDGL (avg. $4.23\times$). REAL-G also delivers consistent improvements, with a $2.65\times$ – $8.99\times$ speedup over ESDG (avg. $5.80\times$), $1.11\times$ – $3.99\times$ over BLAD (avg. $1.95\times$), and $1.86\times$ – $8.59\times$ over DistDGL (avg. $4.07\times$). The gains of REAL are particularly pronounced in WD-GCN, TGCN, and GAT-LSTM, where all reuse mechanisms can be exploited. In contrast, for EvolveGCN, the weights in the MatMul stage change with each timestamp, preventing reuse in this phase. ESDG suffers from hidden state transfers between GPUs. The cost is minor for EvolveGCN with small hidden states, but becomes a major bottleneck for models with larger hidden states. DistDGL is further hindered by frequent neighbor aggregation, which limits its ability to scale throughput despite balanced workloads. BLAD, as the current SOTA, outperforms both ESDG and DistDGL mainly by reducing inter-GPU communication; however, its lack of data reuse and load balancing considerations still caps its maximum performance gains.

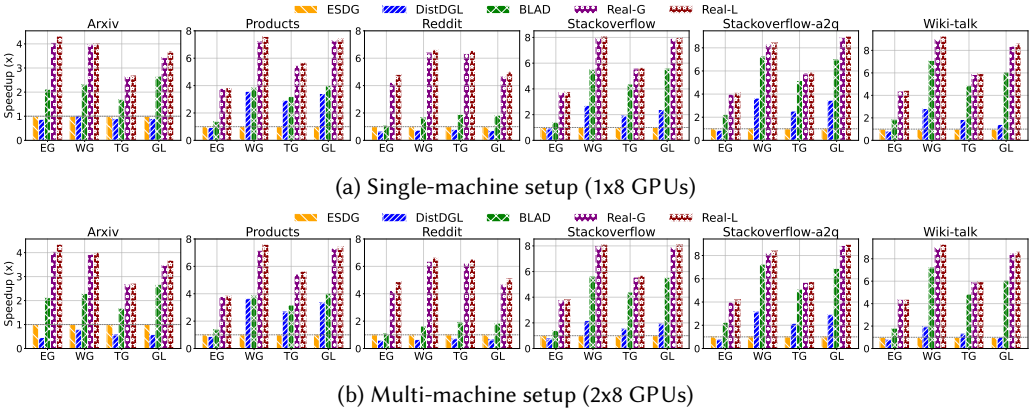


Fig. 18. **Training speedup of different systems across DGNN models and datasets.** Speedup is normalized to ESDG [3]. EG, WG, TG, and GL denote EvolveGCN, WD-GCN, TGCN, and GAT-LSTM, respectively.

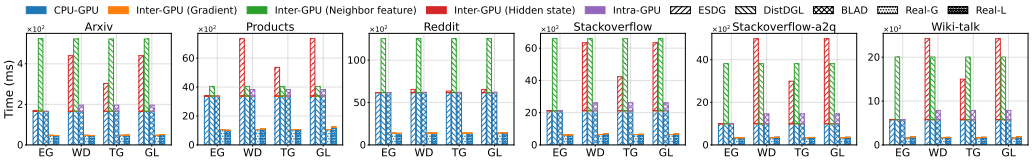


Fig. 19. **Breakdown of data transfer time per epoch on the testbed A.** Each bar is decomposed into CPU-GPU transfer, inter-GPU transfer (gradients, neighbors, hidden states), and intra-GPU transfer.

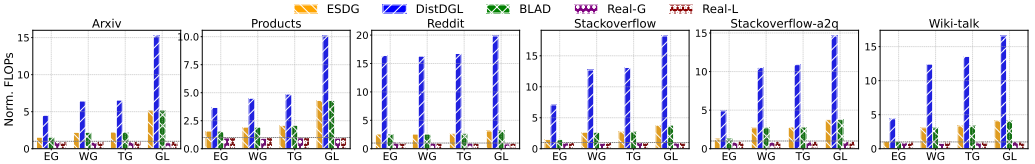


Fig. 20. **Normalized GraphConv FLOPs per epoch on the testbed A.** FLOPs is normalized to REAL-L.

Figure 18b extends the evaluation to the multi-server setting (Cluster B, 2×8 H20 GPUs). The speedup trends of REAL remain consistent with the single-server results for most baselines. The notable exception is DistDGL, which suffers severe degradation in the multi-machine environment; its vertex-level partitioning necessitates cross-node neighbor aggregation, inducing substantial network communication overhead. Conversely, ESDG avoids cross-machine hidden state transfers (with a group size of four), as dependencies are effectively localized within each GPU group, thereby preserving inter-group data parallelism.

C.1.2 Transfer Time Breakdown. Figure 19 decomposes the per-epoch data transfer latency. For ESDG, inter-GPU hidden state synchronization dominates the overhead (except for EvolveGCN), becoming particularly prohibitive on large datasets like Stackoverflow. DistDGL is constrained by neighbor aggregation overhead, which scales poorly on dense or high-dimensional graphs due to cross-partition edges. For instance, on Reddit (avg. degree 137, feature dim. 602), this induces substantial communication costs. While BLAD eliminates cross-GPU traffic, it remains bottlenecked

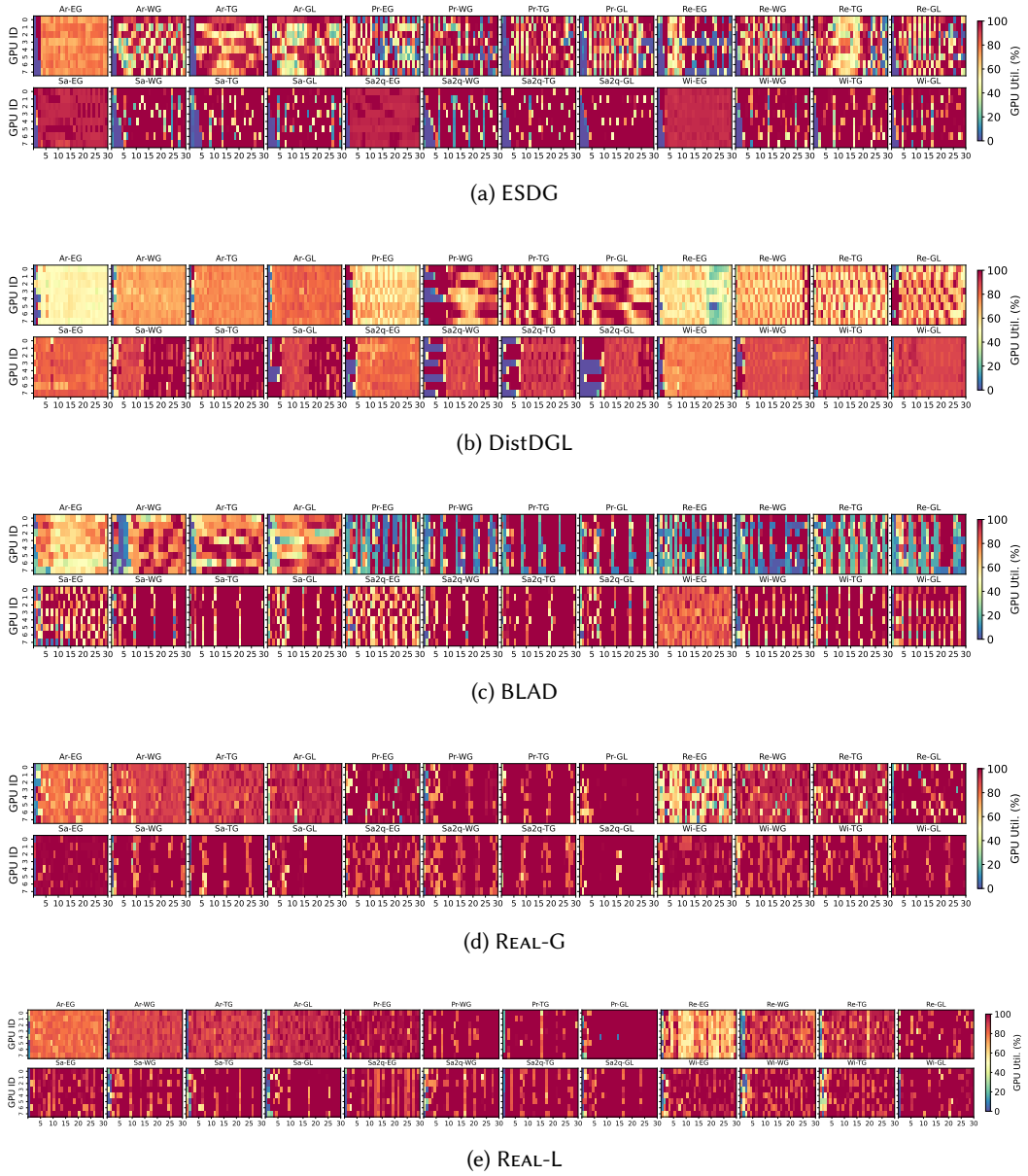


Fig. 21. GPU utilization heatmaps across different dynamic GNN models and datasets. We compare (a) ESDG, (b) DistDGL, (c) BLAD, (d) REAL-G, and (e) REAL-L. Labels Ar, Pr, Re, Sa, Sa2q, and Wi denote the datasets Arxiv, Products, Reddit, Stackoverflow, Stackoverflow-a2q, and Wiki-talk, respectively.

by HtoD bandwidth due to full-graph loading. Furthermore, its inter-process state transfer incurs non-negligible overhead on datasets such as Stackoverflow and Wiki. REAL similarly eliminates inter-GPU communication but goes further by leveraging DiffMap-based reuse to minimize HtoD data volume.

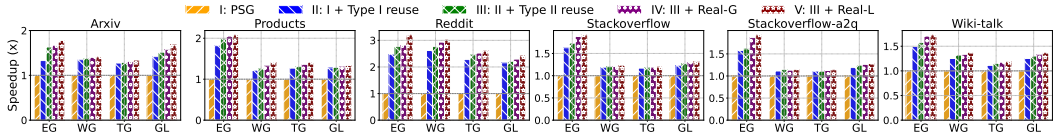


Fig. 22. **Speedup ablation across different models and datasets.** *Speedup is normalized to PSG.*

C.1.3 GraphConv FLOPs Reduction. Figure 20 profiles the total GraphConv FLOPs per training epoch. Compared to DistDGL, REAL-G and REAL-L eliminate substantial redundancy, achieving reduction factors of $3.81\times$ – $23.23\times$ and $3.68\times$ – $19.85\times$, respectively. Relative to ESDG and BLAD (which perform exhaustive computation and thus incur identical FLOPs), REAL delivers reductions of $1.09\times$ – $6.43\times$. These efficiency gains stem from our topology-aware pruning. DistDGL incurs high overhead by performing superfluous computations on *halo nodes* (imported cross-GPU vertices) whose features are irrelevant to local updates. Similarly, ESDG and BLAD remain agnostic to structural sparsity, forcing full-graph convolution on every snapshot. In contrast, REAL exploits both Type I and Type II reuse to effectively bypass aggregation and transformation for topologically invariant nodes.

C.1.4 GPU Utilization and Load Balance. We profile the utilization of all GPUs during the first 30 seconds of training. We present the GPU utilization heatmaps REAL and all baselines in Figure 21. The heatmaps show that REAL-L maintains a high overall GPU utilization, mainly because the *Triple-Stream Pipeline* in Type I reuse overlaps HtoD transfers with Conv and MatMul operations, removing idle GPU time. Furthermore, REAL-L also achieves balanced workload distribution across GPUs, primarily attributed to its *cross-group* scheduling.

C.2 Ablation Study

To evaluate the individual contribution of each component in REAL, we perform a stepwise ablation study normalized to the *Partition-by-Snapshot-Group* (PSG) baseline, where each GPU sequentially trains one complete snapshot group per iteration without reuse or load balancing. Figure 22 reports cumulative speedups across models and datasets.

The PSG baseline achieves $1.00\times$ by definition, serving as the starting point. Adding **type I reuse**, via IncreConv, selective MatMul, and the triple-stream pipeline, improves performance to $1.11\times$ – $2.62\times$ by eliminating redundant computations and data transfers between snapshots. Incorporating **type II reuse**, which by default pairs two consecutive snapshot groups to maximize reuse, further raises speedups to $1.11\times$ – $2.77\times$ by removing redundant HtoD transfers and GraphConv across groups. Introducing the greedy cross-group scheduler (**REAL-G**) balances workloads while preserving reuse benefits, achieving $1.12\times$ – $2.93\times$. Finally, replacing the greedy scheduler with the ILP solver (**REAL-L**) delivers the highest cumulative speedups of $1.14\times$ – $3.21\times$, benefiting from a global view for optimal scheduling.

D Additional Experiment Result on Cluster A

D.1 Performance on Three-Layer DGNN Models

We further evaluate the training performance on deeper, three-layer DGNN models to assess system scalability. As illustrated in Figure 23, REAL maintains substantial performance advantages consistent with the shallow model results. Specifically, REAL-L achieves speedups of $2.30\times$ – $10.02\times$ over ESDG (avg. $5.49\times$), $2.61\times$ – $7.88\times$ over DistDGL (avg. $5.02\times$), and $1.03\times$ – $1.95\times$ over BLAD (avg. $1.43\times$). REAL-G exhibits similar efficiency, delivering average speedups of $5.33\times$ over ESDG,

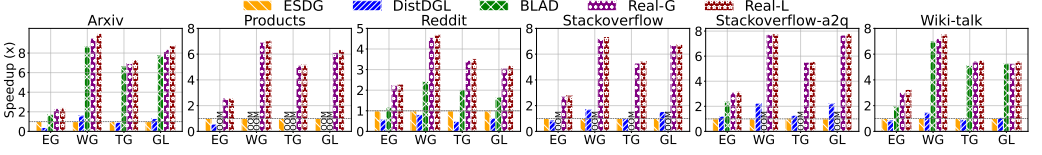
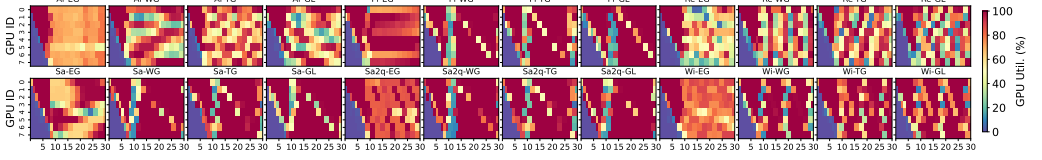
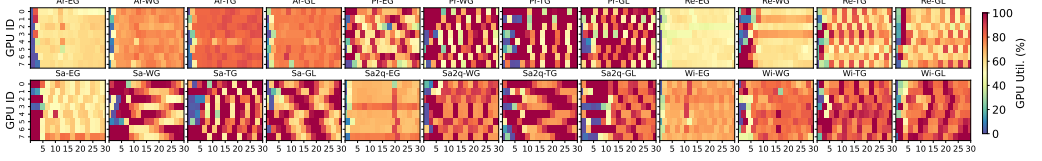


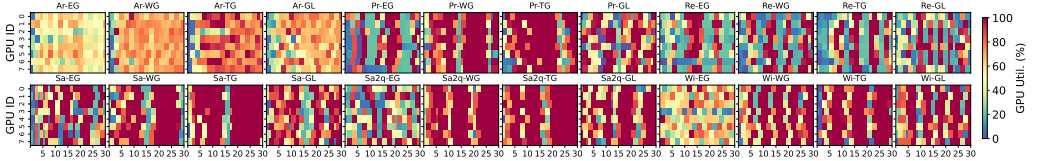
Fig. 23. Training speedup of different systems on three-layer DGNN models. The results are measured on a single server in Cluster A, with speedup normalized to ESDG.



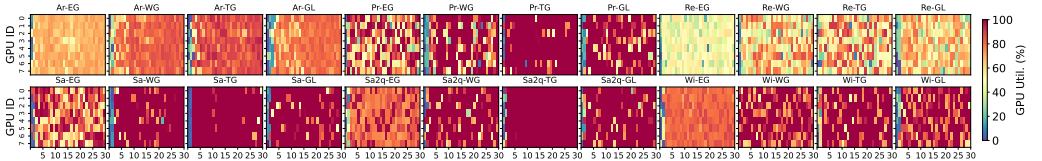
(a) ESDG



(b) DistDGL



(c) BLAD



(d) REAL-G

Fig. 24. GPU utilization heatmaps across different dynamic GNN models and datasets on Cluster A.

4.86 \times over DistDGL, and 1.38 \times over BLAD. These results confirm that the benefits of reuse-aware scheduling and load balancing persist as model depth increases.

D.2 GPU Utilization and Load Balance.

We also visualize the per-GPU utilization heatmaps for the remaining four baselines in Figures 24a–24d, measured during the first 30 seconds of training. For ESDG (Figure 24a), each non-initial rank must wait for the RNN hidden state from other ranks, which frequently causes GPU idleness. This imbalance is mild for EvoLveGCN due to its small hidden state size, but becomes

more pronounced for other DGNN models. **DistDGL** (Figure 24b) mitigates imbalance by partitioning graphs to equalize workload sizes across ranks, enabling fine-grained load control. However, its lack of overlap between computation and data transfer results in uniformly low utilization. **BLAD** (Figure 24c) employs a two-stage pipeline per rank, but without inter-rank load balancing, it suffers from severe imbalance—especially on datasets with relatively small per-graph node counts (e.g., Arxiv, Reddit, Wiki). In contrast, **REAL-G** (Figure 24d) combines load balancing with triple-stream overlap on each rank, reducing inter-rank imbalance and improving overall GPU utilization.

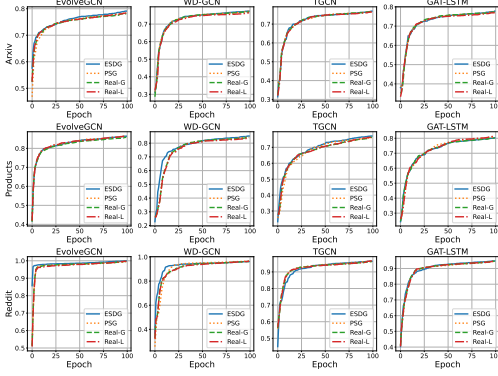


Fig. 25. Model accuracy.

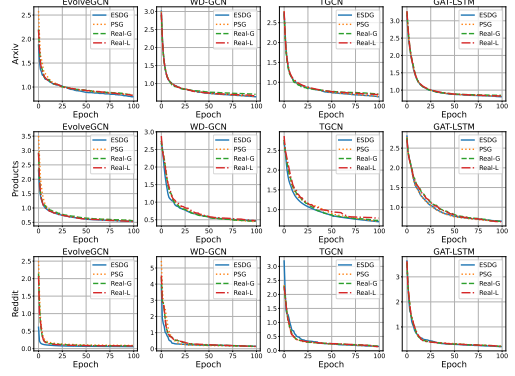


Fig. 26. Training loss.

E Model Accuracy and Loss

The impact on model accuracy between REAL-L and REAL-G versus PSG and ESDG is directly analogous to increasing the training batch size. To ensure that this does not lead to any accuracy degradation, we compared the model accuracy and loss over 100 epochs between REAL-L, REAL-G, PSG and ESDG methods on the Arxiv, Products, and Reddit datasets, as shown in Figures 25 and 26. The Stackoverflow, Stackoverflow-a2q and wiki dataset, lacking labels, is not included in the accuracy comparison. From Figure 25, it is evident that both REAL-G and REAL-L achieve accuracy levels comparable to ESDG throughout the entire training process. The accuracy curves for all three methods exhibit minimal divergence, suggesting that the scheduling techniques implemented in REAL-G and REAL-L do not degrade the model's ability to learn effective representations. Similarly, Figure 26 illustrates the training loss trajectories over the course of 100 epochs. The loss curves for REAL-G and REAL-L closely align with that of ESDG, demonstrating nearly identical convergence behavior. This alignment indicates that the optimization process remains stable and efficient under the proposed scheduling strategies. This consistency in accuracy and loss highlights the robustness of the proposed scheduling approaches, as they maintain the model's performance without introducing significant deviations from the baseline method.

F Convergence Analysis

F.1 Notation and Assumptions

We restate the relevant notation and assumptions for convenience:

- $\mathbf{x} \in \mathbb{R}^d$: model parameter vector.
- ζ : randomly sampled training data.
- Loss function $L(\mathbf{x}) = \mathbb{E}_{\zeta}[f(\mathbf{x}; \zeta)]$, where $f(\mathbf{x}; \zeta)$ is the loss on sample ζ .
- Model parameters at iteration t : \mathbf{x}_t .

- Learning rate: η .
- Number of GPUs: N .
- Number of batches processed per GPU per iteration: K (in our specific algorithm, $K = 2$).

Assumption F.1 (Independence). For all i , the groups of samples $\{\zeta_t^{i,1}, \zeta_t^{i,2}, \dots, \zeta_t^{i,K}\}$ processed by GPU i at iteration t are independent and also independent of other GPUs' samples.

Assumption F.2 (Bounded Gradients). There exists a constant $G > 0$ such that for all \mathbf{x} and ζ ,

$$\|\nabla f(\mathbf{x}; \zeta)\| \leq G.$$

Assumption F.3 (Lipschitz Continuity of Gradients). The gradient of $L(\mathbf{x})$ is L -Lipschitz continuous, i.e., for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,

$$\|\nabla L(\mathbf{x}) - \nabla L(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|.$$

Assumption F.4 (Synchronization Consistency). At each synchronization step, all GPUs have consistent model parameters, i.e., for all i , $\mathbf{x}_t^i = \mathbf{x}_t$.

F.2 Algorithm Description

Definition F.5 (REAL). At each iteration t , GPU i processes K mini-batches $\zeta_t^{i,1}, \zeta_t^{i,2}, \dots, \zeta_t^{i,K}$ at the current model parameters \mathbf{x}_t and computes:

$$\tilde{\nabla} f_i(\mathbf{x}_t) = \frac{1}{K} \sum_{k=1}^K \nabla f(\mathbf{x}_t; \zeta_t^{i,k}).$$

After each GPU i finishes, they synchronize to obtain:

$$\tilde{\nabla} f(\mathbf{x}_t) = \frac{1}{N} \sum_{i=1}^N \tilde{\nabla} f_i(\mathbf{x}_t).$$

Then the model is updated as:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \tilde{\nabla} f(\mathbf{x}_t).$$

F.3 Detailed Steps

Lemma F.6 (Unbiasedness of the Gradient Estimator). *Under the independence assumption, we have*

$$\mathbb{E}[\tilde{\nabla} f(\mathbf{x}_t)] = \nabla L(\mathbf{x}_t).$$

PROOF. By definition:

$$\tilde{\nabla} f_i(\mathbf{x}_t) = \frac{1}{K} \sum_{k=1}^K \nabla f(\mathbf{x}_t; \zeta_t^{i,k}).$$

Since each $\zeta_t^{i,k}$ is drawn independently from the same distribution of ζ , we have:

$$\mathbb{E}[\nabla f(\mathbf{x}_t; \zeta_t^{i,k})] = \nabla L(\mathbf{x}_t).$$

Thus:

$$\mathbb{E}[\tilde{\nabla} f_i(\mathbf{x}_t)] = \frac{1}{K} \sum_{k=1}^K \mathbb{E}[\nabla f(\mathbf{x}_t; \zeta_t^{i,k})] = \nabla L(\mathbf{x}_t).$$

Averaging over N GPUs:

$$\mathbb{E}[\tilde{\nabla} f(\mathbf{x}_t)] = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[\tilde{\nabla} f_i(\mathbf{x}_t)] = \nabla L(\mathbf{x}_t),$$

as required. \square

Lemma F.7 (Variance Bound). *Let σ^2 be an upper bound on the variance of a single-sample stochastic gradient, i.e.,*

$$\mathbb{E}[\|\nabla f(\mathbf{x}_t; \zeta) - \nabla L(\mathbf{x}_t)\|^2] \leq \sigma^2.$$

Then for the estimator $\tilde{\nabla}f(\mathbf{x}_t)$, we have

$$\mathbb{E}[\|\tilde{\nabla}f(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] \leq \frac{\sigma^2}{NK}.$$

PROOF. Each $\tilde{\nabla}f_i(\mathbf{x}_t)$ is the average of K independent stochastic gradients:

$$\tilde{\nabla}f_i(\mathbf{x}_t) = \frac{1}{K} \sum_{j=1}^K \nabla f(\mathbf{x}_t; \zeta_{ij}).$$

Since each $\nabla f(\mathbf{x}_t; \zeta_{ij})$ is independent and has variance bounded by σ^2 , the variance of $\tilde{\nabla}f_i(\mathbf{x}_t)$ can be computed as:

$$\mathbb{E}[\|\tilde{\nabla}f_i(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] = \mathbb{E}\left[\left\|\frac{1}{K} \sum_{j=1}^K (\nabla f(\mathbf{x}_t; \zeta_{ij}) - \nabla L(\mathbf{x}_t))\right\|^2\right].$$

Expanding the squared norm and using the independence of the gradients, we have:

$$\mathbb{E}\left[\left\|\frac{1}{K} \sum_{j=1}^K (\nabla f(\mathbf{x}_t; \zeta_{ij}) - \nabla L(\mathbf{x}_t))\right\|^2\right] = \frac{1}{K^2} \sum_{j=1}^K \mathbb{E}[\|\nabla f(\mathbf{x}_t; \zeta_{ij}) - \nabla L(\mathbf{x}_t)\|^2].$$

By the assumption that $\mathbb{E}[\|\nabla f(\mathbf{x}_t; \zeta_{ij}) - \nabla L(\mathbf{x}_t)\|^2] \leq \sigma^2$, it follows that:

$$\mathbb{E}[\|\tilde{\nabla}f_i(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] \leq \frac{\sigma^2}{K}.$$

The estimator $\tilde{\nabla}f(\mathbf{x}_t)$ is the average of N independent $\tilde{\nabla}f_i(\mathbf{x}_t)$:

$$\tilde{\nabla}f(\mathbf{x}_t) = \frac{1}{N} \sum_{i=1}^N \tilde{\nabla}f_i(\mathbf{x}_t).$$

Since each $\tilde{\nabla}f_i(\mathbf{x}_t)$ is independent and has variance bounded by $\frac{\sigma^2}{K}$, the variance of $\tilde{\nabla}f(\mathbf{x}_t)$ satisfies:

$$\mathbb{E}[\|\tilde{\nabla}f(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] = \mathbb{E}\left[\left\|\frac{1}{N} \sum_{i=1}^N (\tilde{\nabla}f_i(\mathbf{x}_t) - \nabla L(\mathbf{x}_t))\right\|^2\right].$$

Expanding the squared norm and using the independence of the $\tilde{\nabla}f_i(\mathbf{x}_t)$, we have:

$$\mathbb{E}\left[\left\|\frac{1}{N} \sum_{i=1}^N (\tilde{\nabla}f_i(\mathbf{x}_t) - \nabla L(\mathbf{x}_t))\right\|^2\right] = \frac{1}{N^2} \sum_{i=1}^N \mathbb{E}[\|\tilde{\nabla}f_i(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2].$$

Substituting the bound $\mathbb{E}[\|\tilde{\nabla}f_i(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] \leq \frac{\sigma^2}{K}$, we obtain:

$$\mathbb{E}[\|\tilde{\nabla}f(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] \leq \frac{\sigma^2}{NK}.$$

□

Theorem F.8 (Non-Convex Convergence Rate).

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] \leq O\left(\frac{1}{\sqrt{T}}\right).$$

PROOF. From L -smoothness:

$$L(\mathbf{x}_{t+1}) \leq L(\mathbf{x}_t) + \langle \nabla L(\mathbf{x}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{L}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2.$$

Since $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \tilde{\nabla} f(\mathbf{x}_t)$:

$$L(\mathbf{x}_{t+1}) \leq L(\mathbf{x}_t) - \eta \langle \nabla L(\mathbf{x}_t), \tilde{\nabla} f(\mathbf{x}_t) \rangle + \frac{L\eta^2}{2} \|\tilde{\nabla} f(\mathbf{x}_t)\|^2.$$

Taking expectation and using Lemma F.6:

$$\mathbb{E}[L(\mathbf{x}_{t+1})] \leq \mathbb{E}[L(\mathbf{x}_t)] - \eta \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] + \frac{L\eta^2}{2} \mathbb{E}[\|\tilde{\nabla} f(\mathbf{x}_t)\|^2].$$

Decompose $\|\tilde{\nabla} f(\mathbf{x}_t)\|^2$:

$$\|\tilde{\nabla} f(\mathbf{x}_t)\|^2 = \|\nabla L(\mathbf{x}_t) + (\tilde{\nabla} f(\mathbf{x}_t) - \nabla L(\mathbf{x}_t))\|^2.$$

Taking expectations:

$$\mathbb{E}[\|\tilde{\nabla} f(\mathbf{x}_t)\|^2] = \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] + \mathbb{E}[\|\tilde{\nabla} f(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2].$$

By Lemma F.7:

$$\mathbb{E}[\|\tilde{\nabla} f(\mathbf{x}_t) - \nabla L(\mathbf{x}_t)\|^2] \leq \frac{\sigma^2}{NK}.$$

Thus:

$$\mathbb{E}[L(\mathbf{x}_{t+1})] \leq \mathbb{E}[L(\mathbf{x}_t)] - \eta \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] + \frac{L\eta^2}{2} \left(\mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] + \frac{\sigma^2}{NK} \right).$$

Rearranging:

$$\eta \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] \leq \mathbb{E}[L(\mathbf{x}_t) - L(\mathbf{x}_{t+1})] + \frac{L\eta^2}{2} \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] + \frac{L\eta^2\sigma^2}{2NK}.$$

So:

$$\left(\eta - \frac{L\eta^2}{2}\right) \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] \leq \mathbb{E}[L(\mathbf{x}_t) - L(\mathbf{x}_{t+1})] + \frac{L\eta^2\sigma^2}{2NK}.$$

Summing over $t = 1$ to T :

$$\left(\eta - \frac{L\eta^2}{2}\right) \sum_{t=1}^T \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] \leq L(\mathbf{x}_1) - L(\mathbf{x}_{T+1}) + \frac{L\sigma^2\eta^2}{2NK} \cdot T.$$

Rearranging:

$$\frac{1}{T} \left(\eta - \frac{L\eta^2}{2}\right) \sum_{t=1}^T \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] \leq \frac{L(\mathbf{x}_1) - L(\mathbf{x}_{T+1})}{T} + \frac{L\sigma^2\eta^2}{2NK}.$$

Since $L(\mathbf{x})$ is bounded below, let $L(\mathbf{x}^*)$ be a lower bound:

$$L(\mathbf{x}_1) - L(\mathbf{x}_{T+1}) \leq L(\mathbf{x}_1) - L(\mathbf{x}^*).$$

Then:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] \leq \frac{2(L(\mathbf{x}_1) - L(\mathbf{x}^*))}{T\eta(2 - L\eta)} + \frac{L\sigma^2\eta}{2NK}.$$

Let $\eta = \frac{c}{\sqrt{T}}$:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\|\nabla L(\mathbf{x}_t)\|^2] &\leq \frac{2\sqrt{T}(L(\mathbf{x}_1) - L(\mathbf{x}^*))}{Tc(2 - \frac{Lc}{\sqrt{T}})} + \frac{L\sigma^2 c}{2NK\sqrt{T}} \\ &= \frac{2(L(\mathbf{x}_1) - L(\mathbf{x}^*))}{c(2\sqrt{T} - Lc)} + \frac{L\sigma^2 c}{2NK\sqrt{T}} = O\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

□

G Simulation of Scalability

In this section, we describe our comprehensive simulation setup and methodology for predicting snapshot execution times. Our experimental framework encompasses both the generation of dynamic graph snapshots and the development of an accurate time prediction model for computational resource optimization.

For the simulation of dynamic graph evolution, we followed the established methodology of DyGraph, generating 10,000 snapshots together with their corresponding diffmap representations to ensure a robust characterization of temporal network dynamics.

To develop an accurate prediction model for snapshot execution times, we conducted detailed profiling experiments on real GPUs. During profiling, we recorded both the total execution time and the potential reuse time for each snapshot, where reuse time represents computations that could be shared between snapshots. Our analysis revealed a strong linear correlation between snapshot training time and the graph's structural characteristics, specifically the number of nodes and edges. This observation led us to develop a linear regression model for predicting snapshot training times:

$$t_{snapshot_i} = \alpha_1 \cdot N_{node_i} + \alpha_2 \cdot N_{edge_i} + \alpha_3 \quad (16)$$

where $t_{snapshot_i}$ represents the predicted training time for snapshot i , N_{node_i} and N_{edge_i} denote the number of nodes and edges in snapshot i , respectively, and α_1 , α_2 , and α_3 are learned coefficients that capture the relationship between graph structure and computational requirements. For snapshots that are not the first in a group, we instead simulate their execution time using the corresponding diffmap, by replacing N_{node_i} and N_{edge_i} with the number of changed nodes and edges. For estimating the reusable computation time between snapshot groups, we developed a ratio-based model that accounts for temporal overlap:

$$\mathcal{R}_{i,j} = \sum_{k=i+1}^{\max(i, i+tw-(j-i))} t_{snapshot_k} \cdot \text{RATIO} \quad (17)$$

where tw denotes the time window size. This allows us to estimate the effective execution time when combining unprofiled snapshot groups using:

$$t_{combined} = t_{group_i} + t_{group_j} - \mathcal{R}_{i,j} \quad (18)$$

where $t_{combined}$ represents the predicted execution time for two combined snapshot groups, t_{group_i} and t_{group_j} are their individual predicted times, and $\mathcal{R}_{i,j}$ quantifies the computational overlap between groups i and j based on our ratio model. This comprehensive approach enables us to effectively estimate the computational resources needed for processing both individual snapshots and combined groups, facilitating efficient resource allocation and workload planning.

Using another set of data for extrapolation, we found that the prediction error was less than 5%, as illustrated in Figure 27. This model was then used to simulate the execution time of each snapshot in our evaluation.

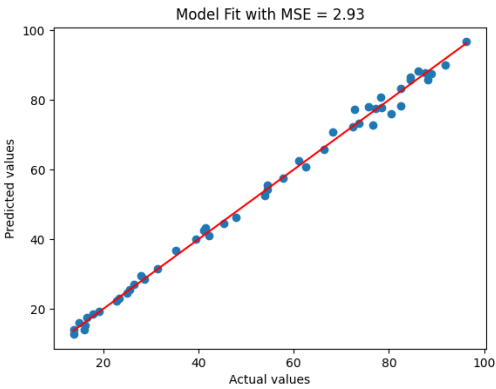


Fig. 27. Comparison between predicted and measured snapshot training time.

Table 4. Preprocessing overhead: DiffMap construction time.

Dataset	Arxiv	Products	Reddit	Stack	Stack-a2q	Wiki
DiffMap Time (s)	7.69	121.77	91.08	78.99	47.06	22.99

Table 5. GPU memory footprint comparison. “Red” values indicate sizes that typically exceed single-GPU memory capacity.

Component	Arxiv	Products	Reddit	Stack	Stack-a2q	Wiki
Full Graph	2.7 GB	55 GB	101 GB	35 GB	17 GB	9.4 GB
DiffMap	0.48 GB	12 GB	11 GB	6.3 GB	3.3 GB	1.1 GB
Node Degree	0.06 GB	1.16 GB	0.16 GB	0.60 GB	0.42 GB	0.14 GB
Total	0.54 GB	13.16 GB	11.16 GB	6.90 GB	3.72 GB	1.24 GB
Saved (%)	80.0%	76.1%	88.9%	80.3%	78.1%	86.8%

H Overhead Analysis

We analyze the system overhead introduced by REAL from two aspects: the preprocessing latency for constructing DiffMaps and the additional memory footprint for maintaining auxiliary data structures.

Preprocessing Latency. Table 4 reports the wall-clock time for DiffMap construction. Even for the largest dataset (Products), the process completes in approximately 121 seconds. Given that training typically spans several hours, this one-time offline cost is negligible and effectively amortized without impacting online throughput.

Memory Efficiency. Table 5 compares the aggregate memory footprint of the entire dataset under the Full Graph representation versus our DiffMap representation. For large-scale datasets like Reddit, the Full Graph representation requires a massive 101 GB of storage capacity. In contrast, REAL reduces this footprint by 88.9% (to 11.16 GB) by storing only topological deltas. While REAL incurs a minor overhead for storing node degrees (necessary for normalizing edge weights in dynamic GCNs), this cost is negligible compared to the substantial memory savings. Crucially, this

compact representation significantly reduces the data volume required for GPU loading, thereby minimizing the bottleneck of expensive host-to-device data transfers.

1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813