

# NeutronCloud: Resource-Aware Distributed GNN Training in Fluctuating Cloud Environments

Mingyi Cao<sup>1</sup>, Chunyu Cao<sup>1</sup>, Yanfeng Zhang<sup>1</sup>, Zhenbo Fu<sup>1</sup>, Xin Ai<sup>1</sup>, Qiange Wang<sup>2</sup>, Yu Gu<sup>1</sup>, Ge Yu<sup>1</sup>

<sup>1</sup>Northeastern University, China; <sup>2</sup>National University of Singapore, Singapore;

2201763@stu.neu.edu.cn, {caochunyu, fuzhenbo, aixin0}@stumail.neu.edu.cn, wang.qg@nus.edu.sg

{zhangyf, yuge, guyu}@mail.neu.edu.cn

## ABSTRACT

Graph Neural Networks (GNNs) are widely employed for their ability to efficiently extract insights from graph-structured data. To support large-scale graph training, researchers have employed distributed training techniques, partitioning the entire graph across multiple computing nodes and performing parallel training by exchanging dependency vertex through cross-node communication. However, existing GNN training systems operate on statically partitioned subgraphs, making it difficult to adapt to dynamic resource variations. In practice, dynamic resource variations in cloud environments often cause unpredictable fluctuations in computation and communication resources, posing significant challenges for aligning each worker’s workload with its available resources during GNN training.

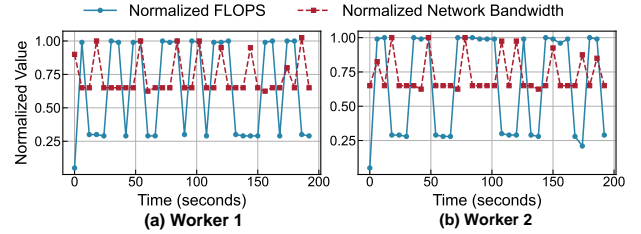
In this paper, we propose NeutronCloud, a system designed for efficient GNN training in cloud environments. First, we adopt a resource-aware workload adjustment strategy. It builds on hybrid dependency handling by obtaining dependency information through both local computation and remote communication. During training, it dynamically adjusts the ratio between locally computed and remotely fetched dependencies based on each worker’s available resources, ensuring workload-resource alignment. Second, we employ a dependency-aware partial-reduce approach reusing historical vertex embeddings and skipping the stragglers during gradient aggregation to address extreme resource fluctuations that cause some workers to lag significantly behind others in the cluster. Experimental results on the resource-fluctuating environment demonstrate that NeutronCloud achieves  $1.83\times$ – $4.43\times$  speedup compared to state-of-the-art distributed GNN systems including NeutronStar and Sancus.

## PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://anonymous.4open.science/r/NeutronCloud>.

## 1 INTRODUCTION

Graph Neural Networks (GNNs) have emerged as a fundamental approach for graph-based tasks such as social network analysis [3, 9], link prediction [21, 30, 50], and recommendation systems [47]. The GNN models iteratively perform neighbor aggregation and representation updates for each vertex to capture complex topological information. Distributed training methods are adopted to partition graph data across multiple workers to handle large-scale graphs, running GNN models in parallel. Most existing distributed GNN training systems [4, 8, 15, 23, 33, 36, 37, 41, 49, 52, 53, 55] assume a stable resource environment, where computational and network



**Figure 1: Normalized computational throughput (FLOPS) and network bandwidth measured over 200 seconds on two Alibaba Cloud shared GPU instances (each with a shared NVIDIA A10 GPU and up to 10Gbps bandwidth).**

resources remain consistent during training. However, in reality, resource fluctuations are ubiquitous in real-world applications.

Resource fluctuations can be categorized into computational resource fluctuations and network resource fluctuations. In cloud environments, computational resources fluctuate due to contention on shared physical hosts [34]. For example, burstable instances (e.g., AWS T-series or Alibaba Cloud “shared-core” VMs) [1, 35] may decrease CPU/GPU performance when neighboring VMs experience high demand. Likewise, network performance varies due to congestion and bandwidth contention. Multi-tenant clusters often face communication bottlenecks caused by overloaded data center switches and cross-worker network contention [27, 45, 51]. To illustrate these fluctuations, we evaluate computational power (FLOPS, floating point operations per second) and network bandwidth over time in a two-worker cluster on Alibaba Cloud’s shared GPU instances. As shown in Figure 1, FLOPS and network bandwidth experience unpredictable reductions of 70% and 35%, respectively.

Such fluctuations significantly impact the performance of distributed GNN training. When performing distributed GNN training in a resource-fluctuating environment, the initially assigned workload becomes mismatched with the fluctuating resources, causing computational and communication bottlenecks, and leading to the emergence of slower workers (stragglers). In addition, GNN data samples exhibit complex interdependencies, further exacerbating this issue. In distributed GNN training, graph partitioning results in local vertices needing to aggregate features from remote neighbors (called remote dependence vertices), and all worker parameters need to be synchronized at the end of each epoch. Thus, the delay caused by stragglers propagates throughout the entire cluster during training, slowing down the overall training performance.

We summarize that the key challenge in distributed GNN training under resource fluctuations is how to quickly adjust each worker’s computational and communication workload to match the available resources. For Deep Neural Network (DNN) systems, there are no dependencies among input samples. They mitigate resource-load

mismatches through real-time load migration to address resource fluctuations [12, 54]. However, for GNN systems, the data dependencies make workload migration difficult. Workload migration requires not only transferring vertex features but also adjusting and maintaining the dependencies between workers to ensure graph consistency. This leads to substantial additional overhead and reduces the efficiency of workload adjustment. The overhead of migration can even offset the benefits of workload adjustment.

In this paper, we propose NeutronCloud, a distributed GNN training system capable of resource-aware workload adjustment, addressing the above challenge through two critical strategies.

First, we propose a resource-aware workload strategy that adapts dependency handling based on real-time resource conditions. The fetching of remote dependencies takes up most of the time in the entire distributed GNN training process [2, 10, 31, 36, 43]. We propose a lightweight resource-aware workload adjustment strategy based on the hybrid dependence handling method. This strategy dynamically adjusts the processing of remote dependence vertices in the cluster to adjust computational and communication workloads. Specifically, we cache remote dependence vertices and their multi-hop neighbors, obtaining embeddings through local computation to address the reduction in communication resources. When computational resources are constrained, more remote dependence vertex embeddings are fetched via cross-worker communication. To enable flexible runtime adjustment, we trade additional storage for adaptability by pre-caching remote dependence vertices and their multi-hop neighbors during the preprocessing phase. This design eliminates the need for workload migration during adjustment.

Second, when some workers face extreme resource degradation (e.g., computation and communication resources decrease simultaneously), adjusting vertex dependencies processing is not enough to reduce the serious delay caused by the synchronization of the severe stragglers. To address this problem, we introduce a dependency-aware partial-reduce strategy, allowing local computation using cached historical embeddings when embeddings from severe stragglers cannot be received in time. Despite being slightly outdated, these embeddings still provide useful information for computation without stalling the process. During gradient aggregation, we synchronize gradients only from faster workers, skipping the severe stragglers. To ensure unbiased gradients and model convergence, we adopt a dependency-aware weighted gradient aggregation strategy and set a bound on severe stragglers that are skipped.

In summary, we make the following contributions.

- We propose a resource-aware workload adjustment strategy, which dynamically adjusts the number of remote dependencies for each worker by quantifying variations in computation and communication resources, ensuring a better match between resource and workload.
- We propose a dependency-aware partial-reduce approach to reduce synchronization overhead. By using history embeddings, we update only the faster worker’s parameters while setting a bound to ensure convergence.
- We develop NeutronCloud, an efficient GNN training in resource fluctuation environments. The experimental results show that NeutronCloud achieves  $1.83\times - 4.43\times$  speedup compared to state-of-the-art GNN systems.

## 2 BACKGROUND

### 2.1 Resource Fluctuations

Resource fluctuations are common in real-world applications and are often caused by contention on shared physical hosts or network contention. These resource fluctuations frequently impact computing power and network bandwidth, significantly affecting the performance of distributed GNN training.

To analyze the impact of resource fluctuations on the performance of distributed GNN systems, we evaluate two GNN training systems, NeutronStar [43] and Sancus [31], on resource-fluctuating and resource-stable clusters using a two-layer GCN model. The cluster consists of four workers, each equipped with NVIDIA T4 GPUs, interconnected via a 10 Gbps network.

Motivated by existing methods for simulating resource fluctuations [13, 24–26], we emulate runtime fluctuations by injecting sleep commands into workers. Specifically, in each iteration, each worker incurs additional overhead with a 10% probability, equivalent to twice the average epoch runtime. To emulate sustained fluctuation patterns (as shown in Figure 1), each injected sleep command lasts for 5 consecutive epochs.

As shown in Figure 2, the per-epoch runtime in a resource-fluctuating environment is  $3.5\times$  slower than that in a resource-stable environment. Distributed GNN training needs to synchronize data across all workers in each epoch. Resource fluctuations on any worker can impact the efficiency of the entire cluster. This indicates that resource fluctuations can significantly degrade the efficiency of distributed GNN training.

### 2.2 GNN Training

A graph neural network processes a graph as input, where each vertex and edge is associated with high-dimensional features. Through multiple layers, GNNs iteratively propagate and aggregate information from neighboring vertices, updating their representations to capture the graph structure information.

In the  $l$ -th layer, the aggregation result  $a_v^l$  of vertex  $v$  is obtained by collecting the embedding  $h_u^{l-1}$  of its neighboring vertex  $u$  in the  $(l-1)$ -th layer, as shown in the following formula:

$$a_v^l = \text{Aggregate}(h_u^{l-1}, u \in N(v)) \quad (1)$$

Subsequently, the update function uses the result  $a_v^l$  and  $h_v^{l-1}$ , which is the embedding of the vertex  $v$  in the  $(l-1)$ -th layer, to compute the representation  $h_v^l$  of the vertex  $v$  in the  $l$ -th layer, as shown in the following formula.

$$h_v^l = \text{Update}(a_v^l, h_v^{l-1}) \quad (2)$$

The final-layer vertex embeddings are fed into downstream tasks, where a task-specific loss is computed against ground truth labels. This loss triggers backward propagation to calculate gradients via automatic differentiation. These gradients drive parameter updates through gradient-based optimizers like SGD or its adaptive variants (e.g., Adam). The standard SGD update rule is:

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta} \mathcal{L}(\theta_t) \quad (3)$$

where the  $\eta$  denotes the learning rate controlling the step size, the  $\eta \nabla_{\theta} \mathcal{L}$  denotes the gradient of loss function w.r.t parameters.

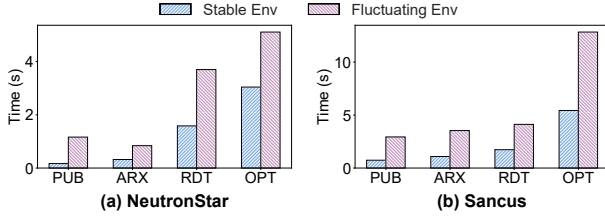


Figure 2: Per-epoch runtime comparison of NeutronStar [43] and Sancus [31] under stable environment (Stable Env) and resource-fluctuating environment (Fluctuating Env) on four datasets PUB (Pubmed), ARX (Ogbn-Arxiv), RDT (Reddit), and OPT (Ogbn-Products).

### 2.3 Distributed GNN Training Approach

Distributed GNN training partitions the input graph across multiple workers. Dependencies arise when vertices need to aggregate features from remote neighbors. The critical aspect of distributed GNN systems is efficiently handling remote dependencies.

We categorize existing dependency handling strategies into three approaches. The **Dependencies Communicated** (DepComm) approach requires each vertex to gather its neighbors' representations from remote workers via cross-worker communication [16, 37], as shown in Figure 3 (a). This method reduces storage consumption but leads to significant communication overhead. The **Dependencies Cached** (DepCache) approach caches the features of multi-hop neighbors of remote dependence vertices on the local worker in advance for multi-layer computing [55], as shown in Figure 3 (b). This eliminates inter-worker communication but results in significant redundant computation and storage overhead.

In the hybrid dependency handling approach, where some remote dependence vertices are handled using the DepCache approach, such as the vertex 0 of worker 1, while others are handled using the DepComm approach, such as the vertex 1 of worker 1, as shown in Figure 3(c). The approach can balance the use of computational and communication resources. Based on this hybrid design, we propose a dynamic adjustment mechanism that adaptively adjusts the ratio between DepCache and DepComm during training. This enables each worker to adapt its computation and communication workload to fluctuating resources.

The hybrid dependency handling approach is shown in Algorithm 1. During forward propagation, the embedding of each vertex is computed by aggregating its neighbor embeddings layer by layer based on its dependency strategy. For the DepCache approach, neighbor embeddings are computed using cached data within the worker. For the DepComm approach, neighbor embeddings are fetched from remote workers via communication (Lines 5-9). During backward propagation, the gradient for each vertex is also computed based on its dependency approach. For vertices with local gradient dependencies, the gradient is computed within the worker. For vertices with remote gradient dependencies, gradients are sent to remote workers via communication (Lines 20-23). Finally, through the All-Reduce, local gradients are exchanged between workers and aggregated to obtain global gradients for parameter updates (lines 24), ensuring consistency for the model parameter.

Therefore, the exchange of vertex embeddings and gradients at each layer, combined with the synchronization of parameter updates, creates layer-wise communication barriers, introducing significant synchronization overhead to distributed GNN training.

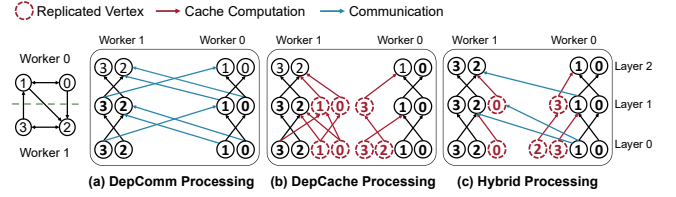


Figure 3: hybrid dependence handling method of distributed GNN training.

**Algorithm 1** DepComm-DepCache Hybrid Training for a Two-layer GCN

**Input:**  $G(V, E)$ ,  $L$ ,  $\{h_v^{(0)} \mid v \in V\}$ ,  $\{L_v \mid v \in V_L\}$ , initial model parameters  $\{W_1^{(l)}, W_2^{(l)}\}$  for each layer  $l$ .

**Output:** Updated model parameters  $\{W_1^{(l)}, W_2^{(l)}\}$ : Updated parameters for each layer.

- 1: partition  $V$  into  $m$  disjoint subsets  $\{V_1, \dots, V_m\}$  and assign  $V_i$ ,  $E_i = \{e_{u,v} \mid v \in V_i, e_{u,v} \in E\}$ ,  $\{h_v^{(0)} \mid v \in V_i\}$ , and  $\{L_v \mid v \in V_i \cap V_L\}$  to each worker  $i$
- Worker  $i = 1, 2, \dots, m$  do in parallel**
- 2: **for**  $l = 1$  to  $L$  **do**
- 3:   **for** each vertex  $v \in V$  **do**
- 4:     **for** each  $e_{u,v} \in E_i$  and  $u \notin V_i$  **do**
- 5:       **if** DepCache( $u, v$ ) **then**
- 6:           $h_v^{(l)} \leftarrow \text{Aggregate}(h_u^{l-1}, u \in N(v))$
- 7:       **else if** DepComm( $u, v$ ) **then**
- 8:          Fetch  $h_u^{(l-1)}$  from remote worker
- 9:           $h_v^{(l)} \leftarrow \text{Aggregate}(h_u^{l-1}, u \in N(v))$
- 10:        $h_v^{(l)} \leftarrow \text{Update}(h_v^l, h_v^{l-1})$
- 11:   **for** each vertex  $v \in V_L$  **do**
- 12:      $\hat{L}_v \leftarrow P(h_v^{(L)})$
- 13:    $\text{loss} \leftarrow \text{loss\_func}(\{\hat{L}_v \mid v \in V_L\}, \{L_v \mid v \in V_L\})$
- 14:   **for** each vertex  $v \in V_L$  **do**
- 15:      $\nabla h_v^{(L)} \leftarrow \frac{\partial \text{loss}}{\partial h_v^{(L)}}$
- 16:   **for**  $l = L$  to 1 **do**
- 17:     **for** each vertex  $v \in V$  **do**
- 18:        $\{\nabla h_{u,v}^{(l)} \mid e(u, v) \in E\} \leftarrow \frac{\partial \mathcal{L}}{\partial h_v^{(l)}}$
- 19:       **for** each edge  $e(u, v) \in E$  **do**
- 20:          **if** DepCache( $u, v$ ) **then**
- 21:            $\{\nabla h_u^{(l-1)}, \nabla h_v^{(l-1)}\} \leftarrow \frac{\partial \mathcal{L}}{\partial h_{u,v}^{(l)}}$
- 22:          **else if** DepComm( $u, v$ ) **then**
- 23:           Send  $\nabla h_u^{(l)}$  to remote worker
- Synchronize between workers**
- 24:   update  $\{W_1^{(l)}, W_2^{(l)}\}$  for each layer  $l$

### 2.4 Existing GNN Systems under Fluctuating Environment

Most existing distributed GNN systems are built on the assumption of stable resource availability. Systems such as DGL [42], AliGraph [55], MGG [44], ROC [16], and DGCL [2] rely on pre-defined task allocation, communication-computation pipeline, and optimizations for computation, memory, or communication efficiency typically performed during initialization. However, these designs

assume static resource conditions, which prevents them from adapting to runtime resource fluctuations. As a result, static execution plans quickly become suboptimal, and delays caused by stragglers can cascade across synchronization between workers, leading to resource underutilization.

Some recent systems attempt to mitigate the impact of stragglers caused by resource fluctuations. PipeGCN [40] adopts pipelined execution with historical embeddings to overlap computation and communication, while Sancus [31] introduces historical embeddings to reduce communication with stragglers. These approaches can partially hide the latency introduced by stragglers. However, they still rely on synchronous parameter updates (e.g., all-reduce), which are bottlenecked by the slowest worker in each iteration and limit adaptability under resource fluctuations.

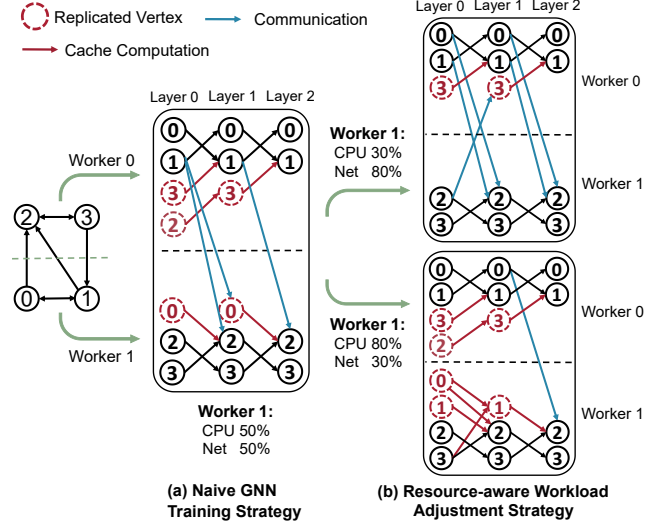
While NeutronStar [43] also adopts a hybrid dependency handling strategy, it lacks runtime resource awareness and cannot adjust dependency handling based on dynamic execution conditions. Moreover, its hybrid dependency handling method requires evaluating the benefit of each dependency and performing global sorting, which incurs high overhead and makes it unsuitable for online adjustment.

### 3 SYSTEM OVERVIEW

We introduce NeutronCloud, a distributed GNN system capable of handling resource fluctuations through two critical strategies. First, NeutronCloud provides a resource-aware workload adjustment that adjusts the computational and communication workloads of each worker to match the real-time resource conditions. Second, NeutronCloud introduces a dependency-aware partial-reduce approach to reduce synchronization overhead between fast and slow nodes.

**Resource-aware Workload Adjustment Strategy.** NeutronCloud employs a lightweight resource-aware workload adjustment algorithm based on the hybrid dependence handling method. The strategy consists of two phases. In the preprocessing phase, remote dependence vertices and their multi-hop neighbors are cached in local CPU memory. Additionally, the cost of different handling methods for remote dependence vertices is computed, and the priority of vertex adjustment is determined. In the online adjustment phase, the algorithm continuously records the computation and communication time of each worker per epoch and quantifies resource fluctuations. Based on the adjustment priorities obtained in the preprocessing phase, a binary search is performed to determine the handling method for remote dependence vertices in the next epoch, speeding up the adjustment process. Moreover, no workload migration is required during the adjustment process.

**Dependency-aware Partial-reduce Strategy.** We propose a dependency-aware partial-reduce method to address the significant synchronization time caused by severe stragglers. This method consists of partial computation and partial update. Partial computation is applied to graph propagation (including both forward and backward propagation), when computing remote dependence vertices handled by the DepCache approach, if the severe stragglers fail to provide timely embedding updates, we use the locally cached historical embeddings of dependence vertices to continue forward propagation. Similarly, historical gradients are used for backward



**Figure 4: An illustrative example of resource-aware workload adjustment.**

propagation. Additionally, a lightweight controller monitors embeddings and gradients staleness, ensuring they are used only if their divergence from the latest embeddings is below a certain threshold.

Partial update limits gradient aggregation to faster workers and skips the stragglers. A runtime controller monitors worker progress to ensure that each update involves the majority of workers. To maintain fairness and convergence, we ensure that all workers participate in parameter updates and periodically update their cached embeddings within a bounded number of epochs. And apply dependency-aware weighted parameter synchronization to ensure unbiased aggregation.

## 4 RESOURCE-AWARE WORKLOAD ADJUSTMENT

In this section, we first analyze the impact of resource fluctuation on GNN training. Then, we define an optimization objective that minimizes the runtime of each worker under resource-fluctuating environments by adjusting how dependencies are handled. Finally, we propose a lightweight algorithm to approximate this objective during training.

### 4.1 Impact Study of Resource Fluctuation

To design our workload adjustment strategy, we analyze how resource fluctuations affect each worker's total runtime, including subgraph computation, DepCache, and DepComm. Based on the impact of the processing time of the two dependency handling approaches, we formulate an objective to minimize the total runtime by adjusting their proportion.

**Training Time under Resource Fluctuation.** We decompose the per-worker training time into three components under resource fluctuation:  $T'_{\text{local},i}$  for local subgraph computation,  $T'_{\text{cache},i}$  for DepCache computation, and  $T'_{\text{comm},i}$  for DepComm communication. This decomposition is expressed as:

$$T'_i = T'_{\text{local},i} + T'_{\text{cache},i} + T'_{\text{comm},i} \quad (4)$$

*Impact of Resource Fluctuation on DepCache.* DepCache improves training efficiency by locally caching the features of multi-hop dependent neighbors to reduce communication. For each vertex, this results in a layer-wise recursive structure: at layer  $l$ , the embedding computation depends on its in-neighbors from layer  $l-1$ , and so on. As shown in Figure 4(a), computing the embedding of Vertex 3 at Layer 1 in Worker 0 requires the Layer-0 embedding of Vertex 2.

Under resource fluctuation, the cost of processing cached dependencies becomes sensitive to system performance. In particular, contention on compute and memory resources affects vertex and edge operations, which is reflected in the runtime-aware per-dimensional costs  $T'_v$  and  $T'_e$  of processing each vertex and edge. Let  $V_{\text{cache},i}^{(l)}$  and  $E_{\text{cache},i}^{(l)}$  denote the cached vertices and edges at layer  $l$  on worker  $i$ . Since multi-hop dependencies are cached across layers, the overall DepCache cost for worker  $i$  is computed by accumulating costs across multiple layers:

$$T'_{\text{cache},i} = \sum_{l=1}^{L-1} \left( |V_{\text{cache},i}^{(l)}| \cdot T'_v + |E_{\text{cache},i}^{(l)}| \cdot T'_e \right) \cdot d^{(l)} \quad (5)$$

where  $d^{(l)}$  is the feature dimension of layer  $l$ .

*Impact of Resource Fluctuation on Local Subgraph Computation.* Similar to DepCache computation, the local subgraph computation time under resource fluctuation is also composed of the multi-layer vertex and edge processing costs, as shown in the following equation:

$$T'_{\text{local},i} = \sum_{l=1}^L \left( |V_{\text{local},i}^{(l)}| \cdot T'_v + |E_{\text{local},i}^{(l)}| \cdot T'_e \right) \cdot d^{(l)} \quad (6)$$

where  $|V_{\text{local},i}^{(l)}|$  and  $|E_{\text{local},i}^{(l)}|$  represent the number of vertices and edges in the local partition of worker  $i$ .

*Impact of Resource Fluctuation on DepComm.* The DepComm cost arises from cross-partition vertex dependencies, when a vertex must fetch feature vectors from remote neighbors of other workers. As shown in Figure 4(a), computing the embedding of Vertex 2 on Worker 1 of Layer 1 requires accessing the embedding of its remote neighbor Vertex 1 located on Worker 0 of Layer 0. Such communication recurs across layers, as each GNN layer aggregates information from the remote neighbors. Under bandwidth fluctuations, the per-dimension communication cost between Worker  $i$  and Worker  $j$  is denoted by  $T'_{c,ij}$ . The total communication cost is given by:

$$T'_{\text{comm},i} = \sum_{l=1}^L \left( d^{(l-1)} \cdot \sum_{\substack{j=0 \\ j \neq i}}^{m-1} T'_{c,ij} \cdot |V_{\text{comm},ij}^{(l-1)}| \right) \quad (7)$$

where  $|V_{\text{comm},ij}^{(l-1)}|$  denotes the number of dependency vertices transferred from worker  $j$  to worker  $i$  at layer  $l-1$ , and  $m$  denotes the total number of workers.

**Optimization Objective.** To support decentralized optimization in a resource-fluctuating environment, each worker independently minimizes its own runtime-aware training time  $T'_i$ , which reflects the current execution delay under resource fluctuations. This is achieved by dynamically adjusting the assignment of dependency vertices between DepCache and DepComm, based on locally aware

resource. Meanwhile, the additional memory consumption for DepCache must respect the worker's memory constraint  $C_i$ . The optimization problem for worker  $i$  is formulated as:

$$\min T'_i = T'_{\text{comm},i} + T'_{\text{cache},i} + T'_{\text{local},i} \quad (8)$$

$$V_{\text{cache},i} \cup V_{\text{comm},i} = V_{\text{depend},i} \quad (9)$$

$$V_{\text{cache},i} \cap V_{\text{comm},i} = \emptyset \quad (10)$$

$$\sum_{v \in V_{\text{cache},i}} s(v) \leq C_i \quad (11)$$

where  $V_{\text{depend},i}$  denotes the complete set of dependency vertices required by worker  $i$ , and  $s(v)$  denotes the storage cost of vertex  $v$ .

The storage cost  $s(v)$  of a dependency vertex  $v$  is determined by the number of neighbor vertices and edges it introduces across layers when cached locally. Formally, we define:

$$s(v) = \sum_{l=1}^{L-1} \left( |V_i^{(l)}(v)| + |E_i^{(l)}(v)| \right) \cdot d^{(l)},$$

where  $|V_i^{(l)}(v)|$  and  $|E_i^{(l)}(v)|$  denote the number of vertices and edges introduced by  $v$  at layer  $l$ .

## 4.2 Lightweight Resource-aware Workload Adjustment

**NP-hardness.** We show that the dependency adjustment problem in our formulation is NP-hard by polynomially reducing the Cardinality-Constrained Knapsack Problem (CCKP) to it. Given any CCKP instance, we map each item to a dependency: item size corresponds to storage cost, and profit corresponds to the reduction in execution time variation. The cardinality and capacity constraints are preserved. Under this mapping, selecting items in CCKP to maximize total profit is equivalent to selecting dependencies handling method to minimize execution time variation without exceeding memory constraint. Since this reduction is in polynomial time and CCKP is NP-hard, our problem is also NP-hard.

**Lightweight Adaptive Adjustment Algorithm.** We propose a lightweight workload adjustment algorithm to address this NP-hard problem. The proposed algorithm performs a dynamic adjustment of dependencies to align each worker's computation and communication load with its available resources. As shown in Figure 4, when a worker faces limited CPU resources, more dependencies are retrieved using the DepComm approach (upper part of Figure 4(b)). Conversely, when a worker suffers from limited communication bandwidth, the DepCache approach is preferred (lower part of Figure 4(b)). The algorithm initially sets all dependency vertices to use DepComm, and then selects the vertices with the highest benefit to convert them to DepCache. The algorithm consists of two stages: a preprocessing stage and an online adjustment stage. The preprocessing stage is executed before training starts to cache essential data and prioritize dependent vertices through sorting by their computational volume. The online adjustment stage dynamically adjusts dependency handling approaches based on real-time resource feedback. By offloading operations with significant computational overhead to the preprocessing stage, we ensure that online adjustment remains lightweight and efficient during training.



**Preprocessing Phase.** In this phase, each worker  $i$  pre-constructs two ordered arrays:  $\mathcal{N}_i$ , which stores the number of multi-hop neighbor edges and vertices for each remote dependence vertex and their storage overhead; and  $\mathcal{S}_i$ , the prefix sum array of the storage cost, as shown in Algorithm 2. These arrays are precomputed to support fast adjustment in the online phase. Remote dependence vertices are grouped by their source remote worker  $j$  into sets  $\mathcal{V}_{ij}$ , as vertices from the same worker  $j$  share a unit transmission time  $T_{\text{comm},ij}$  (Line 1-2). Then, using the measured  $T_e$  and  $T_v$ , we compute the computation cost  $T_{\text{cache},ij}(v)$  for  $\mathcal{V}_{ij}$  (Line 4). Based on this metric, we sort  $\mathcal{V}_{ij}$  and  $\mathcal{N}_{ij}$  in ascending order (Line 5-6) and generate the prefix sum array  $\mathcal{S}_{ij}$  (Line 7-9). These arrays allow efficient lookups during the online adjustment phase. Although  $T_e$  and  $T_v$  may fluctuate, the computational cost of remote dependence vertices remains constant, ensuring that the relative order of  $T_{\text{cache},ij}(v)$  remains consistent in each epoch. We verify the assumption that  $T_v$  and  $T_e$  degrade proportionally under GPU contention. In GNN training,  $T_v$  reflects vertex-wise neural network computation, which is compute-bound and typically implemented as dense matrix multiplications (GEMM), while  $T_e$  corresponds to edge-wise aggregation, which is memory-bound and realized via sparse matrix multiplications (SpMM). We measure the performance of GEMM and SpMM kernels under both compute-intensive and memory-intensive background workloads. As shown in Table 1, both operations exhibit similar degradation across different contention types, confirming that compute and memory resources degrade proportionally and thus preserve the relative order of  $T_{\text{cache}}(v)$ . Therefore, these arrays only need to be generated once before training, significantly reducing the adjustment overhead.

**Table 1: Performance degradation under GPU contention. Each entry shows the degradation (compare to the no-contention performance) for a foreground operator when co-running with a background load.**

Background Load	SpMM (degradation)	GEMM (degradation)
SpMM	↓92%	↓85%
GEMM	↓50%	↓48%

**Online Adjustment Phase.** During training, worker  $i$  iteratively adjusts dependencies based on resource fluctuations while ensuring they remain within the storage budget  $C_i$ , as shown in Algorithm 3. First, we only need to apply a binary search on  $\mathcal{V}_{ij}$  to identify vertices with positive benefits  $\gamma_i(v) > 0$  using updated resource values  $T'_e$ ,  $T'_v$ , and  $T'_{\text{comm},ij}$  (line 2-6). And we obtain the storage overhead sum of all positive benefits vertices (line 7). The total storage requirement of these vertices is calculated as  $S_i = \sum_{j=1}^M \mathcal{S}_{ij}$  (line 9). If  $S_i$  below  $C_i$ , convert all positive benefits vertices into cache dependencies vertices, while the rest remain communication dependencies vertices (lines 12-13). If  $S_i$  exceeds  $C_i$ , the storage budget is proportionally allocated to each remote worker  $j$  as  $C_{ij} = C_i \cdot \mathcal{S}_{ij}/S_i$  (line 15-16). Then, for each  $j$ , positive benefits vertices from  $\mathcal{V}_{ij}$  that satisfy the storage budget  $C_{ij}$  are selected and converted into cache dependencies vertices (line 18-20). This two-phase approach guarantees that the most beneficial conversions are prioritized while satisfying the memory constraints.

#### Algorithm 2 Preprocessing Phase

**Input:** vertices subset  $V_i$ , edges subset  $E_i$ ; set of remote dependence vertices  $D$ ; number of remote workers  $m$ ;  $T_v$ ,  $T_e$ .  
**Output:**  $\{\mathcal{N}_{ij}\}$ ,  $\mathcal{N}_{ij}(v)$ .  $V$  is the number of multi-hop neighbors vertices of  $v$ ,  $\mathcal{N}_{ij}(v)$ .  $E$  is the number of multi-hop neighbors edges of  $v$ ,  $\mathcal{N}_{ij}(v)$ .  $S$  is the storage overhead of multi-hop neighbors,  $\{\mathcal{V}_{ij}\}$  is an array of vertices sorted by  $T_{\text{comp},ij}$ ,  $\{\mathcal{S}_{ij}\}$  is a prefix sum array of the storage cost for  $\mathcal{V}_{ij}$ 's neighbors, for  $j = 1, \dots, m$ .

```

1: for  $j = 1$  to  $m$  do
2:    $\mathcal{V}_{ij} \leftarrow \{v \in D, v \in V_j\}$ 
3:   for each  $v \in \mathcal{V}_{ij}$  do
4:      $T_{\text{cache},ij}(v) \leftarrow \mathcal{N}_{ij}(v) \cdot V \cdot T_v + \mathcal{N}_{ij}(v) \cdot E \cdot T_e$ 
5:   sort  $(\mathcal{V}_{ij})$  by  $T_{\text{cache},ij}(v)$  (ASC)
6:   Update  $\mathcal{N}_{ij}$  accordingly
7:    $\mathcal{S}_{ij}[1] \leftarrow \mathcal{N}_{ij}(\mathcal{V}_{ij}[1]) \cdot S$ 
8:   for  $k = 2$  to  $|\mathcal{V}_{ij}|$  do
9:      $\mathcal{S}_{ij}[k] \leftarrow \mathcal{S}_{ij}[k-1] + \mathcal{N}_{ij}(\mathcal{V}_{ij}[k]) \cdot S$ 
10: return  $\{\mathcal{N}_{ij}\}$ ,  $\{\mathcal{S}_{ij}\}$  ( $j = 1, \dots, m$ )

```

#### Algorithm 3 Online Adjustment Phase

**Input:**  $\{\mathcal{V}_{ij}\}$ ,  $\{\mathcal{N}_{ij}\}$ ,  $\{\mathcal{S}_{ij}\}$ ; updated times  $T'_v$ ,  $T'_e$ ,  $T'_{\text{comm},ij}$ ; total storage constraint  $C_i$ ; number of remote workers  $m$ .  
**Output:**  $\mathcal{V}_j^{\text{cache}}$  and  $\mathcal{V}_j^{\text{comm}}$  for each  $j = 1, \dots, m$ .

```

1: Initialize:  $\mathcal{V}_j^{\text{cache}} \leftarrow \emptyset$ ,  $\mathcal{V}_j^{\text{comm}} \leftarrow \emptyset \forall j = 1, \dots, m$ .
2: for  $j = 1$  to  $m$  in parallel do
3:   for each  $v \in \mathcal{V}_{ij}$  do
4:      $T'_{\text{cache},ij}(v) \leftarrow \mathcal{N}_{ij}(v) \cdot V \cdot T'_v + \mathcal{N}_{ij}(v) \cdot E \cdot T'_e$ 
5:      $\gamma_{ij}(v) \leftarrow T'_{\text{comm},ij}(v) - T'_{\text{cache},ij}(v)$ 
6:   BinarySearch in  $\mathcal{V}_{ij}$  to find all  $v$  with  $\gamma_{ij}(v) > 0$ 
7:    $\mathcal{S}_{ij} \leftarrow \mathcal{S}_{ij}[\text{vlast}]$ 
8:   //  $\text{vlast}$  is the last node that makes  $\gamma_{ij}(v) > 0$ .
9:    $S_i \leftarrow \sum_{j=1}^m \mathcal{S}_{ij}$ 
10:  if  $S_i \leq C_i$  then
11:    for  $j = 1$  to  $m$  in parallel do
12:       $\mathcal{V}_j^{\text{cache}} \leftarrow$  all  $v$  with  $\gamma_{ij}(v) > 0$ 
13:       $\mathcal{V}_j^{\text{comm}} \leftarrow \mathcal{V}_{ij} \setminus \mathcal{V}_j^{\text{cache}}$ 
14:  else
15:    for  $j = 1$  to  $m$  do
16:       $C_{ij} \leftarrow (S_{ij}/S_i) \cdot C_i$ 
17:    for  $j = 1$  to  $m$  in parallel do
18:       $\mathcal{V}_j^{\text{cache}} \leftarrow$  largest prefix of  $\mathcal{V}_j^{\text{cache}}$  fitting  $C_{ij}$  (via  $\mathcal{S}_{ij}$ )
19:       $\mathcal{V}_j^{\text{cache}} \leftarrow \mathcal{V}_j^{\text{cache}}$ 
20:       $\mathcal{V}_j^{\text{comm}} \leftarrow \mathcal{V}_{ij} \setminus \mathcal{V}_j^{\text{cache}}$ 
21: return  $\{\mathcal{V}_j^{\text{cache}}\}$ ,  $\{\mathcal{V}_j^{\text{comm}}\}$  ( $j = 1, \dots, m$ )

```

**Complexity Analysis.** In our implementation, we integrate benefit computation into the binary search process during online adjustment to eliminate redundant passes. As a result, each adjustment involves at most  $M - 1$  binary searches over sorted subsets of size  $O(N/M)$ , where  $M$  is the number of workers, and  $N$  is the number of dependency vertices. This leads to an overall time complexity of  $O(M \log(N/M))$  for the online stage. In contrast, the preprocessing phase is executed once before training. It performs sorting and array construction over all dependency vertices, incurring a one-time cost of  $O(N \log N)$ .

**Approximation guarantee.** Our method achieves a  $\frac{1}{2}$ -approximation guarantee as it is equivalent to that of the default greedy algorithm of CCKP, which ranks items by their benefit-to-cost ratio  $\gamma(v)/s(v)$  and selects the largest feasible prefix, achieving a worst-case  $\frac{1}{2}$ -approximation [18]. Our algorithm adopts a simplified strategy that ranks vertices by  $\gamma(v)$  and selects the Top-K vertices under the memory budget. In our design, ranking by  $\gamma(v)$  is effectively equivalent to ranking by  $\gamma(v)/s(v)$ , as vertices with higher  $\gamma(v)$  typically induce lower neighborhood memory costs  $s(v)$ . This allows us to omit  $s(v)$  from the ranking criterion.

## 5 DEPENDENCY-AWARE PARTIAL-REDUCE

### 5.1 Overall Workflow

**Motivation.** The previously proposed resource-aware workload adjustment strategy generally mitigates the impact of resource fluctuations. However, under some extreme resource degradation (e.g., computing and communication resources decrease simultaneously). Some workers may lag significantly behind the rest of the cluster. When such degradation persists, resource-aware workload adjustment alone may fail to mitigate the impact of **severe stragglers**.

**Challenges of Partial-reduce in GNNs.** In distributed DNNs, partial-reduce (a variant of all-reduce) reduces synchronous overhead caused by severe stragglers by synchronizing gradients among only a subset of faster workers in each epoch. However, due to complex dependencies between training samples, partial-reduce faces the following challenges in distributed GNN training.

First, partial-reduce only works on the parameter synchronization stage. However, since workers need to communicate with each other to obtain the embedding and gradients of dependency vertices, implicit synchronization points occur at each layer. This means that multiple synchronization points exist in each epoch, leading to additional layer-wise synchronization except for parameter synchronization.

Second, the number of workers involved in parameter updates may vary across iterations in partial-reduce, leading to unstable gradient scaling and slower convergence [6, 22, 32]. In distributed DNN training, such variations are typically addressed by normalizing parameter updates through weighted gradient averaging, where each worker's contribution is proportional to the number of local training samples [57]. However, the number of training vertices for different workers is different, and the aggregation operation collects non-training vertices embeddings, which will also generate gradients for these non-training vertices during backward computation, as shown in Figure 6(a). Simply performing a weighted averaging of parameters across workers based on the number of training vertices may result in significant convergence bias.

**A Two-stage Solution: Partial Computation and Partial Update.** We propose a dependency-aware partial-reduce method for GNN, consisting of partial computation and partial update, as shown in Figure 5. Partial computation reusing historical vertex embeddings and gradients for forward/backward computation. Figure 6(b) illustrates that during the backward computation of the first epoch, Worker 0 and Worker 1 proceed without waiting for the gradients from the straggling Worker 2. Partial update collects the gradients from each worker based on corresponding weight

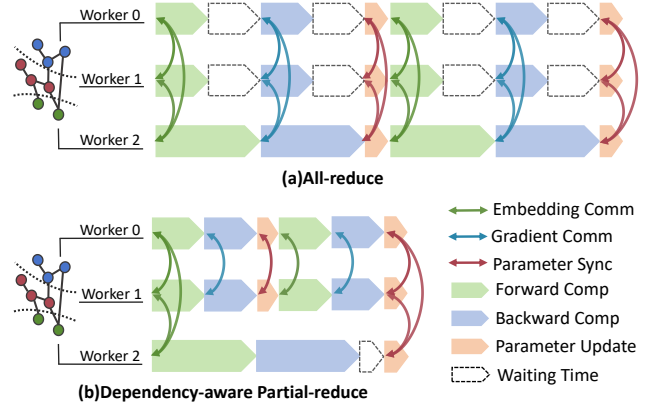


Figure 5: Partial-reduce for GNN. "Comm" indicates communication, "Comp" indicates computation, and "Sync" indicates

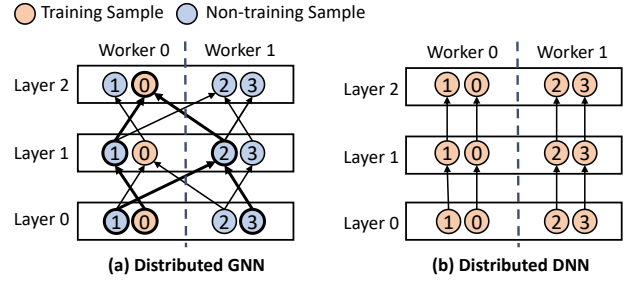


Figure 6: BFS-based partial update.

values, which are computed from the number of training vertices participating in each layer. Specifically, a Breadth-First Search (BFS) traversal is initiated from the training vertices to collect the number of neighbor vertices per worker at each layer. These vertex counts serve as normalization weights during gradient aggregation.

### 5.2 Partial Computation for Forward/Backward

We adopt a partial computation mechanism, where workers can use locally cached historical vertex embeddings instead of waiting for remote vertex embeddings when delays occur. In addition, we introduce a layer-wise timing constraint to limit the maximum training time of each layer.

Concretely, we pre-measure a baseline computation time  $T_0^{(l)}$  for each GNN layer  $l$ , and define an adaptive threshold  $\Delta T_{\max}^{(l)} \approx \frac{T_0^{(l)}}{2}$  to determine the maximum acceptable waiting time for remote vertex embeddings. Whenever a worker's computation for layer  $l$  exceeds  $T_0^{(l)} + \Delta T_{\max}^{(l)}$ , it no longer waits for remote vertex embeddings but instead utilizes locally cached vertex embeddings for computation.

To ensure embedding version consistency, we design a dynamic staleness control mechanism to track the staleness of embeddings on each worker. The controller dynamically controls the maximum allowable delay for remote embedding synchronization. Specifically, when a worker's cached embeddings fall behind by more than  $K$  iterations, the system forces the worker to wait for the latest remote embeddings. This approach ensures that the embeddings used in the computation remain up-to-date, preventing approximation errors due to outdated embeddings.

### 5.3 Partial Update for Parameter Synchronization

**Workflow of Partial Update.** Partial update employs a dual mechanism of dynamic grouping and bounded staleness constraints, as shown in Fig. 5(b). A lightweight global controller real-time monitors the runtime of each iteration in the cluster. When the current iteration time reaches  $T_0 + \Delta T$ , if the majority of workers ( $\geq \frac{M}{2} + 1$ ) are ready, the controller informs them to execute a weighted parameter synchronization in this temporary worker group (e.g., Worker 0-1 bypassing Worker 2 for direct updates). If the number of workers that completed the computation does not reach the threshold ( $\frac{M}{2} + 1$ ), continue to wait until the threshold is met.

To prevent parameter deviation caused by some workers not participating in grouping for a long time, we ensure all workers are involved in parameter synchronization every  $K$  iteration. The controller checks participation every  $K$  iteration. If any worker is missing, the next iteration of parameter synchronization will be performed across all workers in the cluster.

**BFS-based Gradient Weighting Strategy.** To address the convergence degradation caused by inconsistent workers in each iteration and uneven training vertex distribution in partial update, we propose a BFS-based gradient weighting strategy. The key idea is to normalize each worker’s gradient contribution according to the number of training vertices.

Let  $n_k^{(l)}$  denote the number of vertices contributing to gradients at layer  $l$  on worker  $k$ . These vertices are determined through a backward BFS traversal starting from the training vertices. Each worker computes its local normalized gradient as:

$$g_k^{(l)} = \frac{1}{n_k^{(l)}} \sum_{v \in \mathcal{V}_k^{(l)}} \nabla \mathcal{L}_v,$$

where  $\mathcal{V}_k^{(l)}$  is the set of training-related vertices on worker  $k$  at layer  $l$ , and  $\nabla \mathcal{L}_v$  is the gradient of the loss with respect to vertex  $v$ . The parameter update rule then becomes:

$$\theta_{t+1}^{(l)} = \theta_t^{(l)} - \eta \cdot \frac{\sum_{k \in \mathcal{S}_t} n_k^{(l)} \cdot g_k^{(l)}}{\sum_{k \in \mathcal{S}_t} n_k^{(l)}}, \quad (12)$$

where  $\eta$  is the global learning rate and  $\mathcal{S}_t$  is the set of workers active at iteration  $t$ . The dependency-aware normalization adjusts the contribution of each worker’s gradients based on the number of vertices involved in the backward computation at each layer, including both training vertices and their dependencies. This helps stabilize gradient scaling across iterations and improves convergence speed under dynamic worker participation.

**Bias Elimination under Partial Reduce.** We demonstrate that the proposed vertex-weighted gradient aggregation yields an unbiased estimate of the global gradient. We take expectations on both sides and assume the vertex-level gradients are unbiased, i.e.,

$\mathbb{E}[g_k^{(l)}] = \nabla \mathcal{L}(\theta_t)$ , we have:

$$\begin{aligned} \mathbb{E} \left[ \frac{\sum_{k \in \mathcal{S}_t} n_k^{(l)} \cdot g_k^{(l)}}{\sum_{k \in \mathcal{S}_t} n_k^{(l)}} \right] &= \frac{\sum_{k \in \mathcal{S}_t} n_k^{(l)} \cdot \mathbb{E}[g_k^{(l)}]}{\sum_{k \in \mathcal{S}_t} n_k^{(l)}} \\ &= \frac{\sum_{k \in \mathcal{S}_t} n_k^{(l)} \cdot \nabla \mathcal{L}(\theta_t)}{\sum_{k \in \mathcal{S}_t} n_k^{(l)}} = \nabla \mathcal{L}(\theta_t). \end{aligned} \quad (13)$$

In contrast, traditional GNN training directly sums local gradients across active workers without normalization, which leads to biased gradient estimates when only a subset of training vertices is covered at each iteration. The expected aggregated gradient becomes:

$$\mathbb{E} \left[ \sum_{k \in \mathcal{S}_t} g_k^{(l)} \right] = \frac{N_t^{(l)}}{N^{(l)}} \nabla \mathcal{L}(\theta_t), \quad (14)$$

where  $N^{(l)} = \sum_{k=1}^K n_k^{(l)}$  is the total number of training vertices at layer  $l$ , and  $N_t^{(l)} = \sum_{k \in \mathcal{S}_t} n_k^{(l)}$  is the number of vertices contributing to gradients at layer  $l$  among all workers involved in gradient aggregation. The factor  $N_t^{(l)}/N^{(l)}$  reflects incomplete training signal coverage and causes systematic underestimation of the true gradient.

Therefore, by incorporating dependency-aware normalization, our method eliminates this bias and ensures unbiased gradient estimation under resource fluctuation and imbalanced participation.

### 5.4 Convergence Analysis

**Proof Sketch.** We prove that training converges in expectation, even under partial computation and partial update. Specifically, the expected norm of the global gradient diminishes to zero as the number of iterations increases.

We establish convergence under the following assumptions.

**(1) Lipschitz Smoothness.** The objective function  $\mathcal{L}(\theta)$  has  $L$ -Lipschitz continuous gradients; that is,  $\|\nabla \mathcal{L}(\theta_1) - \nabla \mathcal{L}(\theta_2)\| \leq L\|\theta_1 - \theta_2\|$  for any  $\theta_1, \theta_2$ . **(2) Unbiasedness with Bounded Variance.** The stochastic gradient, even when computed with stale embeddings, is an unbiased estimator of the true gradient, with bounded variance. **(3) Diminishing Step Size.** The learning rate  $\eta_t$  satisfies standard diminishing conditions:  $\sum_t \eta_t = \infty$ ,  $\sum_t \eta_t^2 < \infty$ , ensuring convergence.

Based on these assumptions, we apply a standard telescoping-sum argument over the objective function values. The resulting convergence bound shows that the cumulative impact of stale embeddings introduces only small residual errors, which are dominated by the diminishing step size. As the number of iterations increases, these residuals vanish, and the expected norm of the global gradient converges to zero. This confirms that, despite partial reductions and stale data usage, our distributed GNN training achieves convergence to a stationary point.

**Assumptions.** We make the following assumptions:

**(1) Smoothness:** The objective function  $f(\theta) = \frac{1}{|V|} \sum_{v \in V} \ell(h_v(\theta), y_v)$  has  $L$ -Lipschitz continuous gradients. Specifically, for any two parameter vectors  $\theta_1, \theta_2 \in \mathbb{R}^d$ , we assume that

$$\|\nabla f(\theta_1) - \nabla f(\theta_2)\| \leq L\|\theta_1 - \theta_2\|$$



This condition ensures that the function does not change too abruptly and plays a crucial role in deriving descent inequalities for optimization.

**(2) Unbiasedness and Bounded Variance:** Let  $\hat{g}_t^{(l)}$  denote the weighted gradient in layer  $l$ , defined as

$$\hat{g}_t^{(l)} \triangleq \sum_{k=1}^{K'} w_k^{(l)} g_{k,t}^{(l)}, \quad \text{where } w_k^{(l)} = \frac{n_k^{(l)}}{\sum n_k^{(l)}}.$$

Due to communication delays, each worker computes gradients based on stale parameters  $\theta_{t-\tau_k}^{(l)}$  with  $\tau_k \leq \tau$ . We assume the following conditions hold:

$$\mathbb{E}[\hat{g}_t^{(l)} \mid \theta_{t-\tau_k}^{(l)}] = \nabla f(\theta_{t-\tau_k}^{(l)})$$

and that the variance of the gradient is bounded, i.e.,

$$\mathbb{E} \left\| \hat{g}_t^{(l)} - \nabla f(\theta_{t-\tau_k}^{(l)}) \right\|^2 \leq \frac{\sigma^2}{\sum_{k=1}^{K'} n_k^{(l)}}$$

Additionally, due to the smoothness assumption, we assume that the change in the parameter  $\theta$  is not too large, which implies that

$$\left\| \nabla f(\theta_t^{(l)}) - \nabla f(\theta_{t-\tau_k}^{(l)}) \right\| \leq L \left\| \theta_t^{(l)} - \theta_{t-\tau_k}^{(l)} \right\|$$

**(3) Learning Rate:** The learning rates  $\{\eta_t\}$  are chosen to satisfy  $\eta_t < \frac{2}{L}$  and

$$\sum_{t=0}^{\infty} \eta_t = \infty, \quad \sum_{t=0}^{\infty} \eta_t^2 < \infty,$$

in practice, we use  $\eta_t = O(1/\sqrt{T})$ .

**Theorem.** For a distributed training strategy, each worker  $k$  in layer  $l$  updates its local parameter via  $\theta_t^{(l),k} = \theta_{t-1}^{(l),k} - \eta_t g_t^{(l),k}$ . Then, a selected subset  $S^{(l)}(t)$  (with  $|S^{(l)}(t)| = K'$ ) performs a weighted aggregation to update the global parameter. We assume an unbiased gradient estimator  $\hat{g}_t^{(l)}$ , i.e.,  $\mathbb{E}[\hat{g}_t^{(l)} \mid \theta_t] \approx \nabla f(\theta_t)$ , with bounded variance  $\mathbb{E}[\|\hat{g}_t^{(l)} - \nabla f(\theta_t)\|^2] \leq \sigma^2 / \sum_{k=1}^{K'} n_k^{(l)}$ . Under these assumptions (plus appropriate smoothness or learning-rate conditions if needed), we aim to show that the parameter sequence  $\{\theta_t\}$  converges to a stationary point. In particular,

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla f(\theta_t)\|^2] = 0.$$

**Proof.** By the  $L$ -smoothness of  $f$ , we first establish an upper bound on the function value after an update. Applying the smoothness condition, we obtain:

$$\begin{aligned} f(\theta_{t+1}) &\leq f(\theta_t) + \langle \nabla f(\theta_t), \theta_{t+1} - \theta_t \rangle + \frac{L}{2} \|\theta_{t+1} - \theta_t\|^2 \\ &= f(\theta_t) - \eta_t \langle \nabla f(\theta_t), \hat{g}_t \rangle + \frac{L\eta_t^2}{2} \|\hat{g}_t\|^2. \end{aligned} \quad (15)$$

**Table 2: Dataset description.**

Dataset	V	E	#F	#L	#H
Yelp (YP)	716,874	13,954,819	300	100	128
Reddit (RDT)	232,965	114,615,892	602	41	128
Ogbn-products (OPT)	2,449,029	61,859,140	100	47	64
Amazon (AMZ)	1,598,960	132,169,734	200	107	128

Taking expectation conditioned on  $\theta_t$  and using the unbiasedness property  $\mathbb{E}[\hat{g}_t \mid \theta_t] \approx \nabla f(\theta_t)$ , we obtain:

$$\begin{aligned} \mathbb{E}[f(\theta_{t+1}) \mid \theta_t] &\leq f(\theta_t) - \eta_t \|\nabla f(\theta_t)\|^2 + \frac{L\eta_t^2}{2} \mathbb{E}\|\hat{g}_t\|^2 \\ &= f(\theta_t) - \eta_t \left(1 - \frac{L\eta_t}{2}\right) \|\nabla f(\theta_t)\|^2 + \frac{L\eta_t^2 \sigma^2}{2 \sum_{k=1}^{K'} n_k}. \end{aligned} \quad (16)$$

Next, summing over  $t = 0$  to  $T - 1$  and applying the telescoping sum, we obtain the following bound:

$$\sum_{t=0}^{T-1} \eta_t \left(1 - \frac{L\eta_t}{2}\right) \mathbb{E}\|\nabla f(\theta_t)\|^2 \leq f(\theta_0) - f^* + \frac{L\sigma^2}{2} \sum_{t=0}^{T-1} \frac{\eta_t^2}{\sum_{k=1}^{K'} n_k}. \quad (17)$$

To derive the final result, we divide both sides by  $\sum_{t=0}^{T-1} \eta_t$  and take the limit as  $\sum_{t=0}^{T-1} \eta_t \rightarrow \infty$ , yielding:

$$\lim_{T \rightarrow \infty} \frac{1}{\sum_{t=0}^{T-1} \eta_t} \sum_{t=0}^{T-1} \eta_t \mathbb{E}\|\nabla f(\theta_t)\|^2 = 0. \quad (18)$$

For the special case of a constant learning rate  $\eta_t = \eta$ , the above result simplifies to:

$$\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}\|\nabla f(\theta_t)\|^2 = 0. \quad (19)$$

which completes the proof, which concludes that the convergence is guaranteed.

## 6 EXPERIMENTS

### 6.1 Experimental Setup

**Environments.** Our experiments are conducted on Aliyun ECS cluster with 16 GPU nodes. Each node has 16 vCPUs, 155GB DRAM, and 1 NVIDIA Tesla T4 GPU, running Ubuntu 20.04 LTS OS. The network bandwidth is 10 Gbps.

**Datasets and GNN Algorithms.** Table 1 lists the four graph datasets used in our evaluation: Reddit [11] is based on user interactions in a social network. Ogbn-products [14] originate from similar relationships between products in an e-commerce platform. The Yelp [48] dataset is constructed from user reviews of local businesses. The Amazon [4] dataset is derived from product co-purchasing graphs, with nodes representing products and edges capturing frequently co-reviewed or co-purchased items. The vertex feature dimensions, the number of labels of datasets, and hidden layer dimensions are listed in Table 1. We use four popular GNN models, including Graph Convolutional Network (GCN) [20], Graph Attention Network (GAT) [38], GraphSage [11], and TAGCN [7] to evaluate the performance, all of them are in a 2-layer structure.

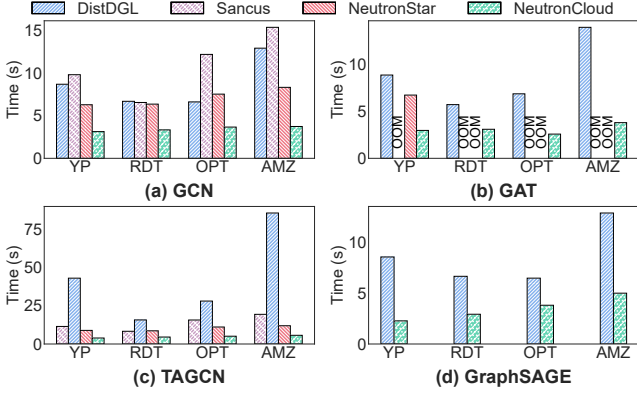


Figure 7: Per-epoch runtime of different systems on different models. "OOM" indicates the out-of-memory error.

**The Systems for Comparisons.** In our experiments, we compare NeutronCloud with two types of systems: mini-batch systems and full-graph systems. For the mini-batch system, we choose DistDGL [53] as the baseline model for comparison. DistDGL reduces computational and memory overhead through sampling. In our experiment, DistDGL selects at most 10 neighbors for the first hop of a vertex and then selects up to 15 neighbors for each of these 10 neighbors. For the full-graph system, we compare NeutronCloud with NeutronStar [43] and Sancus [31]. NeutronStar adopts a hybrid dependency-handling approach, and effectively balances the computation and communication workload, enabling high-performance GNN training. Sancus uses historical embeddings to reduce cross-worker communication. In NeutronCloud, we follow the graph partitioning strategy used in NeutronStar, adopting a chunk-based [56] approach that divides the vertex ID space into contiguous ranges. Vertex features and labels are colocated with their corresponding vertices, while edges are assigned to partitions based on their destination vertex. By default, we set the staleness bound  $K = 3$  in both accuracy and runtime performance evaluations, meaning that each worker is allowed to compute with cached embeddings and delay gradient synchronization for up to 3 epochs. All experimental results are presented in terms of the runtime per epoch, which refers to the time for forward and backward propagations, of all vertices in the graph. The results are obtained by averaging over 100 epochs. A shorter runtime per epoch indicates that the model takes less time to achieve the same accuracy.

**Heterogeneity Simulation.** Inspired by existing heterogeneous training methods, we simulated a resource-fluctuating environment in the experiment. Specifically, for each worker, we independently added a certain probability of sleep time in each epoch to reflect resource dynamic. Specifically, each worker has a 10% probability of adding 5 seconds of sleep time within an epoch, which is partially added to both forward and backward propagations.

## 6.2 Overall Comparison

We compare the performance of NeutronCloud by running GCN, GAT, GraphSage, and TAGCN on a 16-node cluster. Recording the average runtime and the average increased training time caused by resource fluctuations (called **fluctuation-induced delay time**) per epoch. The experimental results are summarized in Figure 7 and Figure 8 respectively. By meticulously logging these times,

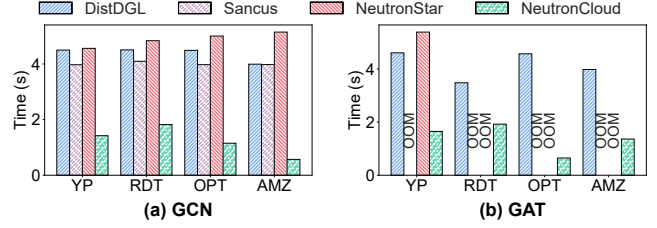


Figure 8: Fluctuation-induced delay time of different systems on GCN and GAT models. "OOM" indicates the out-of-memory error.

we evaluate the efficiency of the different systems under different algorithms.

**Per Epoch Time Comparison.** Compared to NeutronStar, DistDGL, and Sancus, NeutronCloud demonstrates superior performance across all datasets. For the average runtime per epoch, it also achieved speedups of up to 2.10 $\times$ , 2.97 $\times$  and 4.15 $\times$ , as shown in Figure 7.

Sancus sequentially triggers each worker to broadcast embeddings, sending all local embeddings to all workers, regardless of whether other partitions contain these vertices during training. This leads to considerable redundant communication and longer waiting times, further degrading performance. Therefore, communication delays caused by fluctuations in communication resources can impact training efficiency. The METIS partitioning used by DistDGL may result in certain workers having more vertices. When dynamic resources, this load imbalance can exacerbate the impact of straggling workers, further slowing down the training process. NeutronStar adopts a pipeline parallelism strategy. In each worker, the workload is partitioned into multiple chunks, and chunk-level scheduling is applied to overlap the communication and compute tasks. Exhibits better adaptability and performance compared to Sancus and DistDGL in most datasets under resource-fluctuating environment.

The static workload allocation method used by DistDGL, NeutronStar and Sancus cannot match fluctuating available resources, leading to mismatches between workload and available resources. The advantages of NeutronCloud stem from two main factors: (1) the resource-aware workload adjustment strategy enables dynamic changes in workload to match the continuously changing resources; (2) the proposed dependency-aware partial-reduce method, which effectively alleviates the negative impact of extreme resource fluctuations on training efficiency.

**Fluctuation-induced Delay Time Comparison.** The fluctuation-induced delay time can more intuitively show the impact of resource fluctuations on distributed GNN training in different systems. For the average fluctuation-induced delay time per epoch, compared to DistDGL, NeutronStar, and Sancus, NeutronCloud achieved average speedup ratios of 3.89 $\times$ , 4.81 $\times$  and 3.87 $\times$ , respectively, as shown in Figure 8. Sancus reduces communication overhead by reusing historical embeddings, making it less affected by communication resource fluctuations compared with DistDGL and NeutronStar.

## 6.3 Preprocessing Overhead Analysis

We compare the preprocessing time, including prioritizing dependent vertices through sorting and statistics the number of training vertices, to the execution time of running GCN for 100 epochs. Table 3 shows the results. The preprocessing phase adds an average

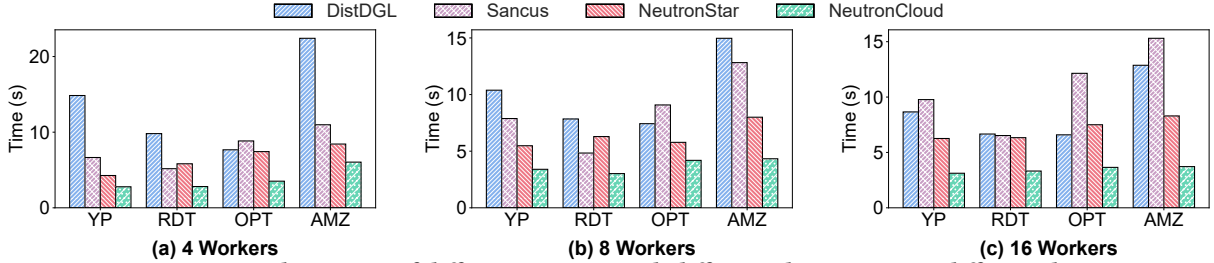


Figure 9: Per-epoch runtime of different systems with different cluster sizes on different datasets.

overhead of 3.91% to NeutronCloud. This overhead gradually decreases as the graph size increases (e.g., 5.66% on Yelp vs. 2.84% on Amazon). And This is typically amortized during training, while the techniques applied in this phase contribute to an average speedup of 2.08 $\times$  for NeutronCloud.

Table 3: Pre-processing vs. training time breakdown (seconds and portion of total) for 100 epochs. Baseline denotes the training time before applying our approach.

Dataset	Training(Baseline)	Preprocessing	Training	Total
Yelp	675.00	18.74 / 5.66%	312.00 / 94.34%	330.74
Reddit	683.00	16.62 / 4.77%	332.00 / 95.23%	348.62
Products	750.00	8.94 / 2.39%	365.00 / 97.61%	373.94
Amazon	879.00	10.86 / 2.84%	372.00 / 97.16%	382.86

## 6.4 Scalability Analysis

**Performance with varying cluster sizes.** In this experiment, We compared NeutronCloud with other systems when training GCN on four datasets with different cluster sizes. Figure 9 shows the runtime of the system on different cluster sizes per epoch, where NeutronCloud consistently outperforms other systems. Specifically, compared to DistDGL, NeutronCloud, and Sancus, NeutronCloud achieved average speedup ratios of 2.53 $\times$ -2.65 $\times$  when the cluster size increased from 4 to 16.

As the number of workers increases, the overall training time of distributed GNNs is expected to reduce. However, for full-graph systems, Sancus and NeutronStar, the training time not only fails to reduce but instead increases due to resource fluctuations. In contrast, NeutronCloud experiences a slight reduction in training time because dynamic workload allocation methods can better match fluctuating available resources. Sancus exhibits poor scalability, possibly due to its heavy reliance on serial global broadcasting and unnecessary full-data transfers. In large-scale distributed environments, the communication overhead increases as the number of workers grows, and communication overhead and network bottlenecks significantly limit its performance scalability. In contrast, NeutronCloud employs a dependency-aware partial-reduce method, which better handles severe straggler issues caused by resource fluctuations in large-scale node clusters.

**Performance with varying model layers.** In this experiment, we compare NeutronCloud with baselines when training GCN with different model layers over Ogbn-products in a 16-node cluster. For the 2, 3, and 4-layer models, the DistDGL sampling strategies were set to (25,10), (25,15,10), and (25,20,15,10), respectively. The results are shown in Figure 10(a). We observe that the performance advantage of NeutronCloud over other baselines gradually increased with the model depths. For the 2-layer model, NeutronCloud achieves an average speedup of 2.40 $\times$ . For the 3-layer and 4-layer models, the speedups were 3.15 $\times$  and 3.42 $\times$ , respectively.

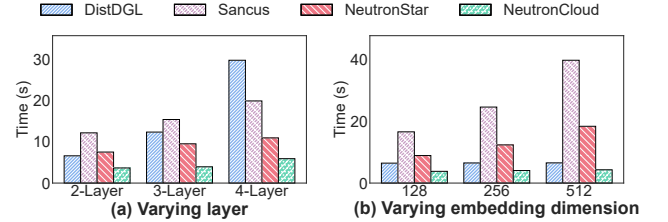


Figure 10: Per-epoch runtime of different systems on Ogbn-Products with varying model configurations.

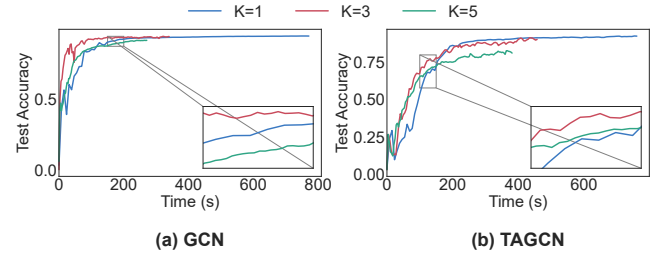


Figure 11: Accuracy for different parameters  $K$ .

**Performance with varying embedding dimensions.** In this experiment, we compare NeutronCloud with baselines when training GCN with different embedding dimensions on a 16-node cluster. As shown in Figure 10(b), NeutronCloud consistently outperforms baselines, and its advantage increases with the embedding size. Specifically, NeutronCloud achieves average speedups of 2.81 $\times$ , 3.57 $\times$ , and 5.02 $\times$  over all baselines for 128-, 256-, and 512-dimensional embeddings, respectively. Larger embedding dimensions increase communication. NeutronCloud reduces this overhead by using partial-reduce to limit embedding exchange and using resource-aware adjustment to shift dependency communication to local computations.

## 6.5 Performance Tolerance to Parameter Staleness ( $K$ ).

We empirically study the tolerance of different GNN model training for stale embeddings and gradients when using cached historical embeddings for computation and a subset of gradients for parameter synchronization. We evaluate the convergence performance of the systems on a cluster with 16 workers and show the accuracy curves for the GNN models (GCN, TAGCN) on NeutronCloud, NeutronStar, DistDGL, and Sancus under different parameters  $K$ , as shown in Figure 11. It can be seen that as the parameter  $K$  increases, the accuracy and the average training time of each epoch show a downward trend. When the value of parameter  $K$  is 3, the model achieves a relatively ideal balance between accuracy and training overhead.

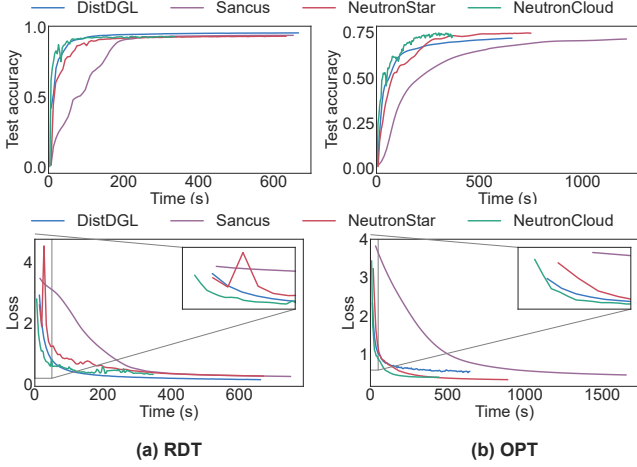


Figure 13: Time-to-loss.

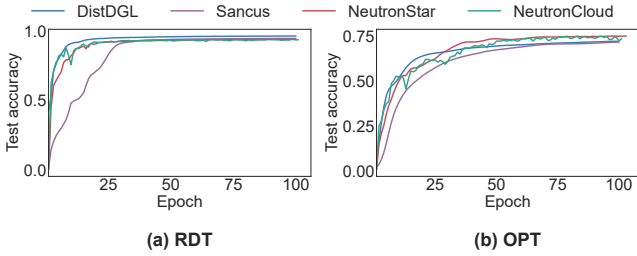


Figure 14: Epoch-to-accuracy.

## 6.6 Convergence Analysis

We evaluate the convergence of NeutronCloud, NeutronStar, DistDGL, and Sancus on a 16-worker cluster using a GCN model for node classification, and show the accuracy curves on two datasets as shown in Figure 12.

The results show that in the early stages, NeutronCloud experiences accuracy fluctuations caused by severe stragglers. After 100 epochs, the test accuracy stabilizes, and NeutronCloud achieves the same test accuracy comparable to other systems while reaching the target accuracy faster than all other systems. Similarly, Figure 13 shows the curve of loss changing over time. The convergence accuracy of NeutronCloud is comparable with that of other systems.

However, Sancus exhibits a long convergence time, which may be due to its cross-worker communication. It does not immediately transmit vertex embeddings after each update but instead transmits them at fixed epoch intervals. NeutronCloud employs a bounded control mechanism in partial computation, limiting the staleness of historical embeddings to at most  $K$  iterations, ensuring that deviation does not accumulate indefinitely. Additionally, periodic All-Reduce in partial update is performed, ensuring that all workers complete parameter synchronization within  $K$  epochs, preventing the model from degrading due to severe stragglers. As a result, NeutronCloud achieves higher accuracy.

Figure 14 shows the accuracy curves of the GCN model on NeutronStar, NeutronCloud, DistDGL, and Sancus across training epochs. We observe that under the same number of epochs, NeutronCloud achieves slightly lower model accuracy compared to other systems, but it eventually achieves convergence accuracy

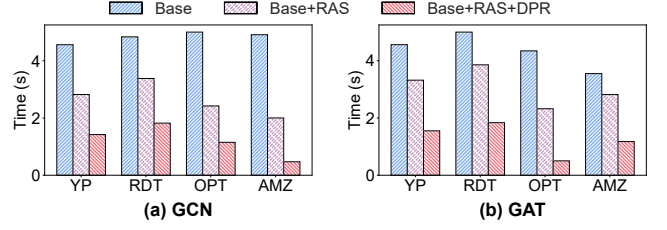


Figure 15: Performance gain analysis. "RAS" indicates the resource-aware workload adjustment strategy, "DPR" indicates the dependency-aware partial reduce strategy.

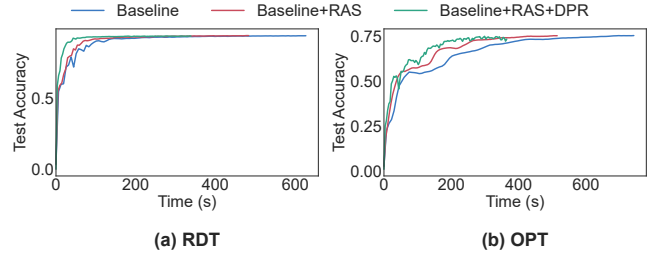


Figure 16: Test accuracy over time for the RDT and OPT datasets under three settings: Baseline, Baseline+RAS, and Baseline+RAS+DPR.

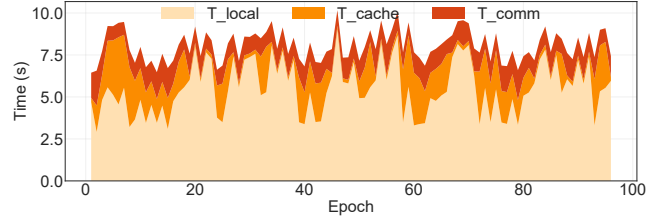


Figure 17: The distribution of overhead associated with  $T_{\text{local}}$ ,  $T_{\text{cache}}$ , and  $T_{\text{comm}}$  during a 100-epoch training.

comparable to that of other systems. This is because it uses stale embeddings for forward computation and parameter synchronization only among a subset of workers in each iteration. However, since NeutronCloud executes each epoch faster, the system still achieves faster convergence in terms of overall training time.

## 6.7 Performance Gain Analysis

**Training efficiency analysis.** To validate the effectiveness of NeutronCloud's key designs, we conduct experiments on GNN models (GCN, GAT) across four datasets, evaluating the impact of resource-aware workload adjustment strategy (RAS) and dependency-aware partial-reduce (DPR) on system efficiency. To ensure a fair comparison, we start with a foundational framework established on the NeutronCloud codebase and gradually integrate the two optimization methods.

Figure 15 shows the average fluctuation-induced delay time per epoch. Compared to the Baseline, the Baseline+RAS achieves an average speedup of 1.67 $\times$ . Figure 17 shows the training-time breakdown for the GCN model on the Amazon dataset of Baseline+RAS.



We record the sleep time used to simulation resource fluctuations as part of  $T_{\text{local}}$ . By dynamically rebalancing the workload between DEPCACHE and DEPCOMM, NeutronCloud effectively mitigates the performance degradation introduced by these fluctuations.

Compared to the Baseline+RAS, the Baseline+RAS+DPR achieves an average speedup of 2.68 $\times$ . The dependency-aware partial-reduce strategy includes two levels of mechanisms, reducing the impact of severe stragglers on overall training progress from both layer-wise synchronization and parameter synchronization dimensions under extreme conditions. It significantly reduces the prolonged synchronization overhead caused by severe stragglers.

**Accuracy evaluation and analysis** To validate the effectiveness of NeutronCloud’s designs, we evaluate the convergence performance of the Baseline and two optimization variants. Figure 16 shows the accuracy curves of the GCN model on the Reddit and Ogbn-products datasets. After 100 epochs, the test accuracy stabilizes, and both Baseline+RAS+DPR and Baseline+RAS achieve accuracy comparable to that of the Baseline.

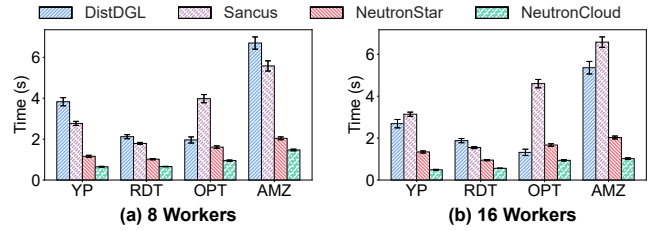
## 6.8 Performance on Real Cloud GPUs

To evaluate our method in a real heterogeneous environment, we further conduct experiments on commercial cloud platforms that support GPU sharing. With the improvement of single-GPU performance, public GPU cloud providers have started offering virtual GPU containers, such as Vultr[39] and Alibaba Cloud[5], allowing multiple users to share the same physical GPU resources. By leveraging GPU virtualization technologies such as NVIDIA MIG[28] or SR-IOV[29], these platforms allocate independent computing instances to each user. These cloud providers typically allow users to request GPU resources on demand, renting computing power at a finer granularity (e.g., 1/2, 1/4, or even 1/20 of a single GPU). This approach improves the overall utilization of GPU resources while effectively reducing user costs.

Specifically, we launch 16 GPU instances on Alibaba GPU cloud provider to validate the adaptability of our method, each instance is a virtual GPU container that contains a 1/3 NVIDIA A10 GPU and up to 20Gbps bandwidth. Since the computational performance of shared GPUs is affected by cluster scheduling policies, we observed significant dynamic heterogeneity among these instances. We run the GCN model on four datasets to compare the execution speed of NeutronCloud with other systems. To account for the variability of real cloud environments and ensure statistical reproducibility, we repeat each experiment multiple times across three independent periods and report the average per-epoch runtime. Figure 18 presents the experimental results with error bars indicating the standard deviation between numerous experiments. The result demonstrates that NeutronCloud achieves stable performance across repeated trials, confirming the reproducibility of our results. Compared to DistDGL, NeutronStar, and Sancus, NeutronCloud achieves average speedup ratios of 3.90 $\times$ , 1.83 $\times$ , and 4.43 $\times$ , respectively.

## 7 RELATED WORK

**Existing Distributed GNN Systems.** Sancus [31] reduces synchronization overhead by using historical embeddings, but adopts



**Figure 18: Per-epoch runtime comparison under real cloud heterogeneous setting.**

a broadcast-based communication scheme where each worker sequentially sends all local embeddings to all others, regardless of actual demand. This incurs redundant communication and prolonged synchronization delays. NeutronStar [43] performs full-graph training with pipeline parallelism and chunk-level scheduling to accelerate training by overlapping communication and computation. However, it cannot adapt to environments with fluctuating resources due to the parameter synchronization approach of All-Reduce and static load distribution. DistDGL [53] adopts mini-batch training and reduces overhead through neighborhood sampling. However, its reliance on static METIS [17] partitioning can cause load imbalance under resource fluctuations, leading to unstable performance in resource-dynamic environments.

**The Partial-Reduce Strategy for Distributed DNNs.** The deployment of deep learning in heterogeneous environments, such as public clouds, introduces new challenges, including communication delays and straggler effects. All-Reduce, originally designed for homogeneous clusters, suffers performance degradation under such conditions due to synchronization bottlenecks. To address this, recent studies [19, 25, 46] have proposed various strategies to mitigate the impact of stragglers by relaxing strict synchronization requirements during training. Partial-Reduce [25], for example, modifies the parameter synchronization protocol to skip delayed workers during gradient aggregation, thereby reducing both computation and communication overhead. DPAR[19] adapts gradient synchronization strategies based on the computational capacity of nodes, allowing more capable nodes to take on additional computational tasks to balance the overall training speed. RNA[46] allows nodes with varying computational capabilities to update parameters at their speeds, improving overall throughput. These methods effectively alleviate the issue of node straggling in distributed DNN training frameworks. However, due to the dependencies between samples in GNNs, it cannot be directly applied to distributed GNN training.

### Distributed DNNs under Resource-fluctuating Environments.

In practical distributed training, cluster resources are not fixed but dynamically change due to various factors. For example, workload fluctuations, communication topology adjustments, and hardware failures make traditional static task allocation and synchronization strategies insufficient for complex environments, thereby affecting training throughput and convergence speed. To address these challenges, SDpipe[26] introduces a Semi-Decentralized training framework. By leveraging heterogeneity-aware task scheduling and a dynamic gradient synchronization strategy, it enables adaptive adjustments of computation and communication loads in response



to resource dynamics within the cluster. This improves training efficiency and reduces synchronization wait time.

## 8 CONCLUSION

We present NeutronCloud, a system designed for efficient GNN training in cloud environments with dynamic and fluctuating resources. Its performance and adaptability are enabled by two key components: (1) a resource-aware workload adjustment strategy that dynamically matches computational and communication workloads to real-time resource conditions, and (2) a dependency-aware partial-reduce strategy that reuses historical vertex embeddings and skips stragglers during gradient aggregation to improve training efficiency. Compared with existing distributed GNN systems such as DistDGL and Sancus, NeutronCloud achieves end-to-end training speedups ranging from 1.83 $\times$  to 4.43 $\times$  in the real cloud environments.

## REFERENCES

- [1] Sanjith Athlur, Nitika Saran, Muthian Sivathanu, Ramachandran Ramjee, and Nipun Kwatra. 2022. Varuna: scalable, low-cost training of massive deep learning models. In *EuroSys '22: Seventeenth European Conference on Computer Systems*, Rennes, France, April 5 - 8, 2022. ACM, 472–487.
- [2] Zhenkun Cai, Xiao Yan, Yidi Wu, Kaihao Ma, James Cheng, and Fan Yu. 2021. DGCL: an efficient communication library for distributed GNN training. In *EuroSys '21: Sixteenth European Conference on Computer Systems, Online Event, United Kingdom, April 26-28, 2021*. ACM, 130–144.
- [3] Sandro Cavallari, Vincent W. Zheng, Hongyun Cai, Kevin Chen-Chuan Chang, and Erik Cambria. 2017. Learning Community Embedding with Community Detection and Node Embedding on Graphs. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management, CIKM 2017, Singapore, November 06 - 10, 2017*. ACM, 377–386.
- [4] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 2019. Cluster-GCN: An Efficient Algorithm for Training Deep and Large Graph Convolutional Networks. 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, 257–266.
- [5] Alibaba Cloud. 2025. VGPU Accelerated Instance Families - Elastic GPU Service. <https://www.alibabacloud.com/help/en/elastic-gpu-service/latest/vgpu-accelerated-instance-families>. Accessed: 2025-03-31.
- [6] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Quoc V. Le, Mark Z. Mao, Marc'Aurelio Ranzato, Andrew W. Senior, Paul A. Tucker, Ke Yang, and Andrew Y. Ng. 2012. Large Scale Distributed Deep Networks. In *Advances in Neural Information Processing Systems 25: 26th Annual Conference on Neural Information Processing Systems 2012. Proceedings of a meeting held December 3-6, 2012, Lake Tahoe, Nevada, United States*. 1232–1240.
- [7] Jian Du, Shanghang Zhang, Guanhang Wu, José M. F. Moura, and Soumya Kar. 2017. Topology adaptive graph convolutional networks. *CoRR* abs/1710.10370 (2017).
- [8] Euler 2019. Euler. <https://github.com/alibaba/euler/wiki/System-Introduction>.
- [9] Wenqi Fan, Yao Ma, Qing Li, Yuan He, Yihong Eric Zhao, Jiliang Tang, and Dawei Yin. 2019. Graph Neural Networks for Social Recommendation. In *The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019*. ACM, 417–426.
- [10] 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 2021, July 14-16, 2021*. USENIX Association, 551–568.
- [11] William L. Hamilton, Zitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*. 1024–1034.
- [12] Aaron Harlap, Henggang Cui, Wei Dai, Jinliang Wei, Gregory R. Ganger, Phillip B. Gibbons, Garth A. Gibson, and Eric P. Xing. 2016. Addressing the straggler problem for iterative convergent parallel ML. In *Proceedings of the Seventh ACM Symposium on Cloud Computing (Santa Clara, CA, USA) (SoCC '16)*. Association for Computing Machinery, New York, NY, USA, 98–111.
- [13] Aaron Harlap, Henggang Cui, Wei Dai, Jinliang Wei, Gregory R. Ganger, Phillip B. Gibbons, Garth A. Gibson, and Eric P. Xing. 2016. Addressing the straggler problem for iterative convergent parallel ML. In *Proceedings of the Seventh ACM Symposium on Cloud Computing, Santa Clara, CA, USA, October 5-7, 2016*. ACM, 98–111.
- [14] 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. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- [15] Kezhao Huang, Jidong Zhai, Zhen Zheng, Youngmin Yi, and Xipeng Shen. 2021. Understanding and bridging the gaps in current GNN performance optimizations. In *PPoPP '21: 26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming, Virtual Event, Republic of Korea, February 27- March 3, 2021*. ACM, 119–132.
- [16] Zhihao Jia, Sina Lin, Mingyu Gao, Matei Zaharia, and Alex Aiken. 2020. Improving the Accuracy, Scalability, and Performance of Graph Neural Networks with Roc. In *Proceedings of the Third Conference on Machine Learning and Systems, MLSys 2020, Austin, TX, USA, March 2-4, 2020*. mlsys.org.
- [17] George Karypis and Vipin Kumar. 1998. A Fast and High Quality Multilevel Scheme for Partitioning Irregular Graphs. *SIAM J. Sci. Comput.* 20, 1 (1998), 359–392.
- [18] Samir Khuller, Anna Moss, and Joseph Naor. 1999. The Budgeted Maximum Coverage Problem. *Inf. Process. Lett.* 70, 1 (1999), 39–45.
- [19] HyungJun Kim, Chunggeon Song, HwaMin Lee, and Heonchang Yu. 2023. Addressing Straggler Problem Through Dynamic Partial All-Reduce for Distributed Deep Learning in Heterogeneous GPU Clusters. In *IEEE International Conference on Consumer Electronics, ICCE 2023, Las Vegas, NV, USA, January 6-8, 2023*. IEEE, 1–6.
- [20] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- [21] Jérôme Kunegis and Andreas Lommatzsch. 2009. Learning spectral graph transformations for link prediction. In *Proceedings of the 26th Annual International Conference on Machine Learning, ICML 2009, Montreal, Quebec, Canada, June 14-18, 2009 (ACM International Conference Proceeding Series)*, Vol. 382. ACM, 561–568.
- [22] Mu Li, David G. Andersen, Jun Woo Park, Alexander J. Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J. Shekita, and Bor-Yiing Su. 2014. Scaling Distributed Machine Learning with the Parameter Server. In *11th USENIX Symposium on Operating Systems Design and Implementation, OSDI '14, Broomfield, CO, USA, October 6-8, 2014*. USENIX Association, 583–598.
- [23] Zhiqi Lin, Cheng Li, Youshan Miao, Yunxin Liu, and Yinlong Xu. 2020. PaGraph: Scaling GNN training on large graphs via computation-aware caching. In *SoCC '20: ACM Symposium on Cloud Computing, Virtual Event, USA, October 19-21, 2020*. ACM, 401–415.
- [24] Qinyi Luo, Jiaao He, Youwei Zhuo, and Xuehai Qian. 2020. Prague: High-Performance Heterogeneity-Aware Asynchronous Decentralized Training. In *ASPLOS '20: Architectural Support for Programming Languages and Operating Systems, Lausanne, Switzerland, March 16-20, 2020*. ACM, 401–416.
- [25] Xupeng Miao, Xiaonan Nie, Yingxia Shao, Zhi Yang, Jiawei Jiang, Lingxiao Ma, and Bin Cui. 2021. Heterogeneity-Aware Distributed Machine Learning Training via Partial Reduce. In *SIGMOD '21: International Conference on Management of Data, Virtual Event, China, June 20-25, 2021*. ACM, 2262–2270.
- [26] Xupeng Miao, Yining Shi, Zhi Yang, Bin Cui, and Zhihao Jia. 2023. SDPipe: A Semi-Decentralized Framework for Heterogeneity-aware Pipeline-parallel Training. *Proc. VLDB Endow.* 16, 9 (2023), 2354–2363.
- [27] Jayashree Mohan, Amar Phanishayee, Janardhan Kulkarni, and Vijay Chidambaram. 2021. Synergy: Resource Sensitive DNN Scheduling in Multi-Tenant Clusters. *CoRR* abs/2110.06073 (2021).
- [28] NVIDIA. 2025. Multi-Instance GPU (MIG) Technology. <https://www.nvidia.com/en-sg/technologies/multi-instance-gpu/>. Accessed: 2025-03-31.
- [29] PCI-SIG. 2010. Single Root I/O Virtualization (SR-IOV) Specification. <https://pcsig.com/specifications/iov>. Accessed: 2025-03-31.
- [30] Hao Peng, Jianxin Li, Yu He, Yaopeng Liu, Mengjiao Bao, Lihong Wang, Yangqiu Song, and Qiang Yang. 2018. Large-Scale Hierarchical Text Classification with Recursively Regularized Deep Graph-CNN. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW 2018, Lyon, France, April 23-27, 2018*. ACM, 1063–1072.
- [31] Jingshu Peng, Zhao Chen, Yingxia Shao, Yanyan Shen, Lei Chen, and Jiannong Cao. 2023. Sancus: Staleness-Aware Communication-Avoiding Full-Graph Decentralized Training in Large-Scale Graph Neural Networks (Extended Abstract). In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI 2023, 19th-25th August 2023, Macao, SAR, China*. ijcai.org, 6480–6485.
- [32] Benjamin Recht, Christopher Ré, Stephen J. Wright, and Feng Niu. 2011. Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent. In *Advances in Neural Information Processing Systems 24: 25th Annual Conference on Neural Information Processing Systems 2011. Proceedings of a meeting held 12-14 December 2011, Granada, Spain*. 693–701.
- [33] Jaeyong Song, Hongsun Jang, Hunseong Lim, Jaewon Jung, Youngsok Kim, and Jinho Lee. 2024. GraNNDIS: Fast Distributed Graph Neural Network Training Framework for Multi-Server Clusters. In *Proceedings of the 2024 International*

- Conference on Parallel Architectures and Compilation Techniques, PACT 2024, Long Beach, CA, USA, October 14-16, 2024. ACM, 91–107.
- [34] Yujia Song, Ruyue Xin, Peng Chen, Rui Zhang, Juan Chen, and Zhiming Zhao. 2023. Identifying performance anomalies in fluctuating cloud environments: A robust correlative-GNN-based explainable approach. *Future Gener. Comput. Syst.* 145 (2023), 77–86.
- [35] John Thorpe, Pengzhan Zhao, Jonathan Eyolfson, Yifan Qiao, Zhihao Jia, Minjia Zhang, Ravi Netravali, and Guoqing Harry Xu. 2023. Bamboo: Making Preemptible Instances Resilient for Affordable Training of Large DNNs. In *20th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2023, Boston, MA, April 17-19, 2023*. USENIX Association, 497–513.
- [36] Alok Tripathy, Katherine A. Yelick, and Aydin Buluç. 2020. Reducing communication in graph neural network training. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2020, Virtual Event / Atlanta, Georgia, USA, November 9-19, 2020*. IEEE/ACM, 70.
- [37] Md. Vasimuddin, Sanchit Misra, Guixiang Ma, Ramanarayan Mohanty, Evangelos Georganas, Alexander Heinecke, Dhiraj D. Kalamkar, Nesreen K. Ahmed, and Sasikanth Avancha. 2021. DistGNN: scalable distributed training for large-scale graph neural networks. In *International Conference for High Performance Computing, Networking, Storage and Analysis, SC 2021, St. Louis, Missouri, USA, November 14-19, 2021*. ACM, 76.
- [38] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- [39] Vultr. 2025. Vultr: High Performance SSD Cloud. Retrieved from <https://www.vultr.com>. Accessed: 2025-03-31.
- [40] Cheng Wan, Youjie Li, Cameron R. Wolfe, Anastasios Kyrillidis, Nam Sung Kim, and Yingyan Lin. 2022. PipeGCN: Efficient Full-Graph Training of Graph Convolutional Networks with Pipelined Feature Communication. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*.
- [41] Xinchun Wan, Kaiqiang Xu, Xudong Liao, Yilun Jin, Kai Chen, and Xin Jin. 2023. Scalable and Efficient Full-Graph GNN Training for Large Graphs. *Proc. ACM Manag. Data* 1, 2 (2023), 143:1–143:23.
- [42] Minjie Wang, Lingfan Yu, Da Zheng, Quan Gan, Yu Gai, Zihao Ye, Mufei Li, Jinjing Zhou, Qi Huang, Chao Ma, Ziyue Huang, Qipeng Guo, Hao Zhang, Haibin Lin, Junbo Zhao, Jinyang Li, Alexander J. Smola, and Zheng Zhang. 2019. Deep Graph Library: Towards Efficient and Scalable Deep Learning on Graphs. *CoRR abs/1909.01315* (2019).
- [43] Qiang Wang, Yanfeng Zhang, Hao Wang, Chaoyi Chen, Xiaodong Zhang, and Ge Yu. 2022. NeutronStar: Distributed GNN Training with Hybrid Dependency Management. In *SIGMOD '22: International Conference on Management of Data, Philadelphia, PA, USA, June 12 - 17, 2022*. ACM, 1301–1315.
- [44] Yuke Wang, Boyuan Feng, Zheng Wang, Tong Geng, Kevin J. 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 2023, Boston, MA, USA, July 10-12, 2023*. USENIX Association, 779–795.
- [45] Qizhen Weng, Wencong Xiao, Yinghao Yu, Wei Wang, Cheng Wang, Jian He, Yong Li, Liping Zhang, Wei Lin, and Yu Ding. 2022. MLaaS in the Wild: Workload Analysis and Scheduling in Large-Scale Heterogeneous GPU Clusters. In *19th USENIX Symposium on Networked Systems Design and Implementation, NSDI 2022, Renton, WA, USA, April 4-6, 2022*. USENIX Association, 945–960.
- [46] Donglin Yang, Wei Rang, and Dazhao Cheng. 2020. Mitigating Stragglers in the Decentralized Training on Heterogeneous Clusters. In *Middleware '20: 21st International Middleware Conference, Delft, The Netherlands, December 7-11, 2020*. ACM, 386–399.
- [47] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2018, London, UK, August 19-23, 2018*. ACM, 974–983.
- [48] Hanqing Zeng, Hongkuan Zhou, Ajitesh Srivastava, Rajgopal Kannan, and Viktor K. Prasanna. 2020. GraphSAINT: Graph Sampling Based Inductive Learning Method. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- [49] Dalong Zhang, Xin Huang, Ziqi Liu, Jun Zhou, Zhiyang Hu, Xianzheng Song, Zhibang Ge, Lin Wang, Zhiqiang Zhang, and Yuan Qi. 2020. AGL: A Scalable System for Industrial-purpose Graph Machine Learning. *Proc. VLDB Endow.* 13, 12 (2020), 3125–3137.
- [50] Muhan Zhang and Yixin Chen. 2018. Link Prediction Based on Graph Neural Networks. In *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*. 5171–5181.
- [51] Yihao Zhao, Yuanqiang Liu, Yanghua Peng, Yibo Zhu, Xuanzhe Liu, and Xin Jin. 2022. Multi-resource interleaving for deep learning training. In *SIGCOMM '22: ACM SIGCOMM 2022 Conference, Amsterdam, The Netherlands, August 22 - 26, 2022*. ACM, 428–440.
- [52] Chengguang Zheng, Hongzhi Chen, Yuxuan Cheng, Zhezheng Song, Yifan Wu, Changji Li, James Cheng, Hao Yang, and Shuai Zhang. 2022. ByteGNN: Efficient Graph Neural Network Training at Large Scale. *Proc. VLDB Endow.* 15, 6 (2022), 1228–1242.
- [53] Da Zheng, Xiang Song, Chengru Yang, Dominique LaSalle, and George Karypis. 2022. Distributed Hybrid CPU and GPU training for Graph Neural Networks on Billion-Scale Heterogeneous Graphs. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*. ACM, 4582–4591.
- [54] Qihua Zhou, Song Guo, Haodong Lu, Li Li, Minyi Guo, Yanfei Sun, and Kun Wang. 2021. Falcon: Addressing Stragglers in Heterogeneous Parameter Server Via Multiple Parallelism. *IEEE Trans. Comput.* 70, 1 (2021), 139–155.
- [55] 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. *Proc. VLDB Endow.* 12, 12, 2094–2105.
- [56] Xiaowei Zhu, Wenguang Chen, Weimin Zheng, and Xiaosong Ma. 2016. Gemini: A Computation-Centric Distributed Graph Processing System. In *12th USENIX Symposium on Operating Systems Design and Implementation, OSDI 2016, Savannah, GA, USA, November 2-4, 2016*. Kimberly Keeton and Timothy Roscoe (Eds.). USENIX Association, 301–316.
- [57] Martin Zinkevich, Markus Weimer, Alexander J. Smola, and Lihong Li. 2010. Parallelized Stochastic Gradient Descent. In *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010. Proceedings of a meeting held 6-9 December 2010, Vancouver, British Columbia, Canada*. Curran Associates, Inc., 2595–2603.