

# SPEED: Streaming Partition and Parallel Acceleration for Temporal Interaction Graph Embedding

Xi Chen<sup>†\*</sup>, Yongxiang Liao<sup>†\*</sup>, Yun Xiong<sup>†§</sup>, Yao Zhang<sup>†</sup>, Siwei Zhang<sup>†</sup>, Jiawei Zhang<sup>‡</sup>, Yiheng Sun<sup>||</sup>

<sup>†</sup>Shanghai Key Laboratory of Data Science, School of Computer Science, Fudan University, Shanghai, China

<sup>‡</sup>IFM Lab, University of California, Davis, CA, United States

<sup>||</sup>Tencent Weixin Group, Shenzheng, China

{x\_chen21, liaoyx21}@m.fudan.edu.cn, {yunx, yaozhang}@fudan.edu.cn,  
swzhang22@m.fudan.edu.cn, jiawei@ifmlab.org, sunyihengcn@gmail.com

**Abstract**—Temporal Interaction Graphs (TIGs) are widely employed to model intricate real-world systems such as financial systems and social networks. To capture the dynamism and interdependencies of nodes, existing TIG embedding models need to process edges sequentially and chronologically. However, this requirement prevents it from being processed in parallel and struggle to accommodate burgeoning data volumes to GPU. Consequently, many large-scale temporal interaction graphs are confined to CPU processing. Furthermore, a generalized GPU scaling and acceleration approach remains unavailable. To facilitate large-scale TIGs’ implementation on GPUs for acceleration, we introduce a novel training approach namely Streaming Edge Partitioning and Parallel Acceleration for Temporal Interaction Graph Embedding (SPEED). The SPEED is comprised of a Streaming Edge Partitioning Component (SEP) which addresses space overhead issue by assigning fewer nodes to each GPU, and a Parallel Acceleration Component (PAC) which enables simultaneous training of different sub-graphs, addressing time overhead issue. Our method can achieve a good balance in computing resources, computing time, and downstream task performance. Empirical validation across 7 real-world datasets demonstrates the potential to expedite training speeds by a factor of up to 19.29x. Simultaneously, resource consumption of a single-GPU can be diminished by up to 69%, thus enabling the multiple GPU-based training and acceleration encompassing millions of nodes and billions of edges. Furthermore, our approach also maintains its competitiveness in downstream tasks.

**Index Terms**—Temporal Interaction Graph, Graph Partitioning, Graph Embedding, Data Mining

## I. INTRODUCTION

Real-world systems featuring sequences of interaction behavior with timestamps—such as social networks, financial trades, and recommendation systems—can all be conceptualized as Temporal Interaction Graphs (TIGs) [1]–[5]. Given a series of timestamped interactions, existing TIG embedding models [1]–[5] represent the objects as nodes and the interaction behaviors among objects with time information as edges, an example is demonstrated in Fig.1. These models encode historical interaction, i.e., event, information in nodes message and memory modules [1]–[5]. The nodes’ embedding can be generated and applied to various downstream tasks by consuming nodes’ history information.

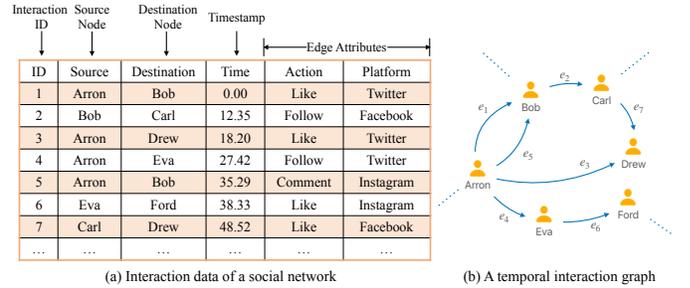


Fig. 1. Interaction data of a social network and its corresponding temporal interaction graphs (TIG).  $e_i$  in (b) refers to the edge which contains the information of time and edge attributes.

As these real-world systems evolve, the associated data scale will also expand correspondingly. Previous research on temporal interaction graph embedding [1]–[5] has primarily focused on enhancing downstream task performance on small datasets, often neglecting the training of large-scale data and efficiency considerations. Based on the updating mechanism of the memory module of existing models, as the number of nodes increases, a larger amount of memory information will need to be stored, leading to higher computing resource requirements. Simultaneously, an upsurge in interaction behaviors leads to considerable overhead in both computational memory resources and time costs.

Traditional single-GPU or CPU training methods [1]–[5] will encounter significant challenges when handling large-scale temporal interaction graphs due to substantial time and space overhead. The time overhead is considerable, especially when dealing with graphs characterized by a large number of edges. From a spatial perspective, a single-GPU will struggle to accommodate to and store memory information for an extensive number of nodes. Consequently, we aspire to construct a training approach capable of addressing both time and space overhead problems.

A straightforward solution to mitigate time overhead is through parallel acceleration using multiple GPUs. However, conventional TIG models [1]–[5] pose significant challenges for parallel training of TIGs. Unlike in static graphs, the message-passing operations in TIGs must adhere to time-based restrictions, preventing a node from gathering information

\*Xi Chen and Yongxiang Liao contributed equally to this paper.

§Yun Xiong is the corresponding author.

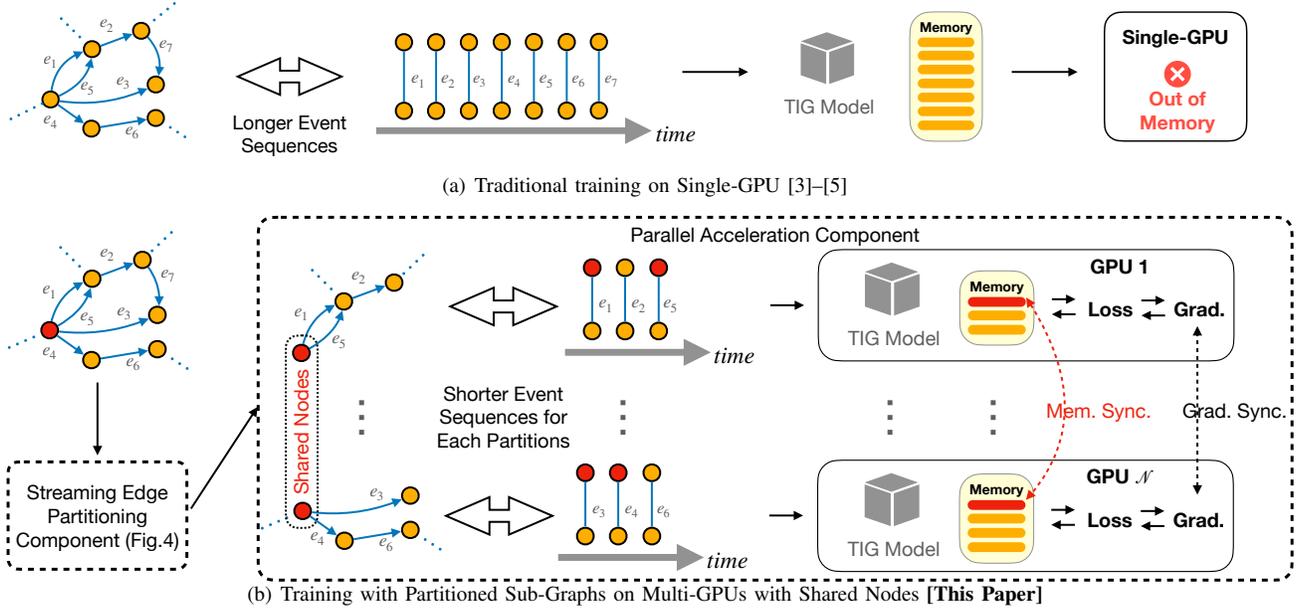


Fig. 2. A comparison between our training approach and conventional single-GPU training. Our approach has much lower training memory and time costs by deploying fewer edge data on each GPU and executing the training in parallel. Benefiting from Streaming Edge Partitioning Component (SEP) (illustrated in Fig.4), we can also accommodate to larger datasets, since the memory required for nodes memory module in our approach per GPU is smaller than in traditional single-GPU training. Note that the outputs of SEP are the partitioned sub-graphs which are the inputs of Parallel Acceleration Component.

from future neighbors. This leads to **Challenge 1**: Temporal and structural dependencies—how to overcome time-based message-passing restrictions in TIGs’ implementation and train TIG models in parallel while preserving the complex interplay between temporal and structural dependencies. Additionally, in TIGs, all events of nodes are interconnected [5], [6]. This necessitates the model to traverse past edges sequentially and chronologically to keep nodes memory up-to-date. It will result in **Challenge 2**: Training interaction sequences in parallel—how to handle multiple temporal interaction sequences with interconnected events while preventing information loss among them and updating nodes’ memory in a distributed parallel training manner. Existing TIG models are struggle to handle the large-scaled TIG data studied in this paper, due to the bottleneck of training in parallel [6].

The space overhead issue primarily originates from the memory module. In most existing baseline methods, a memory slot is maintained for each node to update its representation. While storing memory module in GPU can accelerate computation, as the number of nodes increases, the storage for these nodes’ memory may also lead to excessive GPU memory consumption. This poses **Challenge 3**: Space overhead caused by memory module—how to accommodate and manage nodes’ memory storage for large-scale TIGs on GPUs with a limited memory size.

In response to all above challenges, we introduce an effective and efficient approach, namely Streaming Edge Partitioning and Parallel Acceleration for Temporal Interaction Graph Embedding (SPEED). The SPEED consists of two functional components, i.e., the Streaming Edge Partitioning Component (SEP) and the Parallel Acceleration Component

(PAC). The TIG partitioning strategy entails assigning fewer graph nodes to each GPU, thereby regulating the nodes memory module’s GPU memory consumption. This helps to address the issue of space overhead. Another advantage of this lies in the reduction of corresponding edges per GPU. Utilizing the multi-GPU parallel component enables simultaneous training of different edge sequence data, thereby accelerating the training process.

The graph partitioning is crucial for handling large-scale TIGs. However, as shown in Tab.I, current graph partitioning algorithms fail to simultaneously satisfy the following requirements: i) temporal information consideration: they neglect the temporal aspect of TIG, disregarding the diverse messages brought by edges at varying timestamps; ii) low replication factor: replicated nodes are added to decrease information loss, however, these algorithms lack control over the number of these nodes, leading to space overhead issues; iii) load balancing: ensuring an even distribution of edges and nodes to balance the training time and resource usage across different GPUs; iv) good scalability: a requirement for low algorithm overhead and the ability to scale efficiently to large-scale temporal interaction graphs.

Hence, we propose the Streaming Edge Partitioning Component (SEP). Firstly, we introduce an exponential time decay strategy to incorporate temporal information. It aims to capture the recent trend of interactions between nodes. Then, we estimate the centrality of each node by aggregating the weights of its historically connected edges. Moreover, to minimize the replication factor, we designate nodes with the top- $k$  centrality values as *hubs*. We treat the input graph as stream of edges, sequentially assigning edge to a specific partition. During this

TABLE I  
COMPARISON OF DIFFERENT GRAPH PARTITION ALGORITHMS

Partition Algorithms	Support Time Information	Low Replication Factor	Load Balance	Scalable
METIS [7] & KL [8]	✗	✓	✗	✗
Random [9] & LGD [10]	✗	✓	✗	✓
ROC [11]	✗	✗	✗	✓
Libra [12] & Greedy [13] & HDRF [14]	✗	✗	✓	✓
Ours	✓	✓	✓	✓

phase, only nodes classified as *hubs* can be duplicated across different partitions as “shared nodes”. It ensures that the vital information will be maintained across all partitions without unnecessary replication of data. Furthermore, we apply a greedy heuristic to maintain the load balance among partitions. The combination of the above strategies enables SEP to satisfy all the requirements previously outlined in Tab.I.

Given the time-sensitive nature of TIGs, the training data is fed into the model in a chronological order, aligned with the timestamp of each edge. Therefore, if the graph is merely divided and allocated to different GPUs for independent computation, the GPU processing for edges with later timestamps will need to wait for those with earlier timestamps. This creates a bottleneck for the multi-GPU training for TIGs.

To mitigate GPU waiting and adjust the SEP component, we propose the Parallel Acceleration Component (PAC). Initially, each sub-graph is assigned to a distinct GPU, ignoring inter-node edges across different GPUs. This means some edges get deleted, causing potential information loss and model effectiveness reduction. To counterbalance this, we take advantage of the shared nodes. The memory of shared nodes is synchronized after each training epoch as a compensatory measure. We further introduce a random shuffling approach. Specifically, we partition the graph into smaller sub-graphs and then shuffle and amalgamate them to fit the number of GPUs. As this precedes every epoch, various “deleted” edges between small partitions can be recovered by merging them into larger partitions and trained across epochs. The PAC allows us to train different sub-graphs in parallel on multi-GPUs, and alleviate information loss.

Our SPEED overcomes the bottlenecks of the existing TIG embedding models for large-scale graphs in terms of computing efficiency (time) and computing resources (device memory). To further provide readers with the outline of the SPEED, Fig.2 offers a comparative illustration of our training approach and the conventional single-GPU approach, highlighting our approach’s enhanced capacity to handle large datasets. Additionally, as demonstrated in Fig.3 and supported by extensive experimental studies, our approach significantly accelerates the training process while maintaining competitive

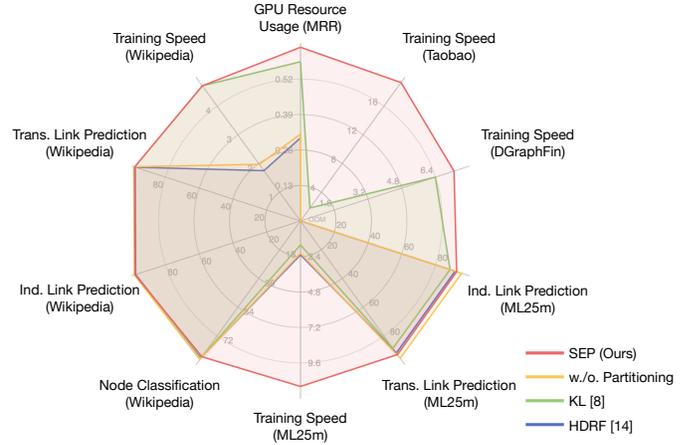


Fig. 3. Comparison between graph partitioning methods. The axes are arranged counterclockwise, with graph sizes increasing. Our method outperforms others in terms of training speed (especially as the graph size increases) and GPU resource utilization (on one GPU), while achieving comparable performance on downstream tasks (link prediction in both transductive and inductive styles, as well as node classification on representative datasets. For full experimental results, please refer to Sec.III). The data are based on averaged experimental results with the TIGE [5]. The data refers to mean reciprocal ranking; Training Speed refers to training speed-up(x) compared with CPU training; OOM refers to out-of-memory. Since the maximum number for each axis differs, the ticks may also vary between axes.

performance on downstream tasks.

The contributions of this paper are as follows:

- 1) We propose a novel large-scale Temporal Interaction Graph embedding approach which achieves a balance among computational resources, time costs and downstream task performance.
- 2) We design a steaming partitioning strategy, specifically tailored for parallel training. This strategy not only captures time information, but also maintains load balance and low replication factor.
- 3) We present a parallelable acceleration method for training graphs with billions of temporal interactions using multi-GPUs. With the partitions shuffling and memory synchronizing across sub-graphs, our approach alleviates the information loss raised by graph partitioning.
- 4) Extensive experimental results demonstrate that the proposed approach significantly expedites training speeds and reduces resource consumption while maintaining its competitiveness in downstream tasks.

## II. PROPOSED METHODS

### A. Notations

Given a set of nodes  $\mathcal{V} = \{1, 2, \dots, N\}$ , nodes features are noted as  $\mathbf{v}_i \in \mathbb{R}^d$ .  $\mathcal{E} = \{e_{ij}(t) = (i, j, t) | i, j \in \mathcal{V}\}$  is a series of interaction behaviours (as shown in Fig.1), where interaction event  $e_{ij}(t) = (i, j, t)$  happens between nodes  $i, j$  at time  $t \in [0, t_{\max}]$ . The interaction feature is noted as  $\mathbf{e}_{ij}(t) \in \mathbb{R}^d$ . Then, we have a temporal interaction graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , which stores the timestamps in edge features.  $\mathcal{E}_t = \{(i, j, \tau) \in \mathcal{E} | \tau < t\}$  contains interaction events record

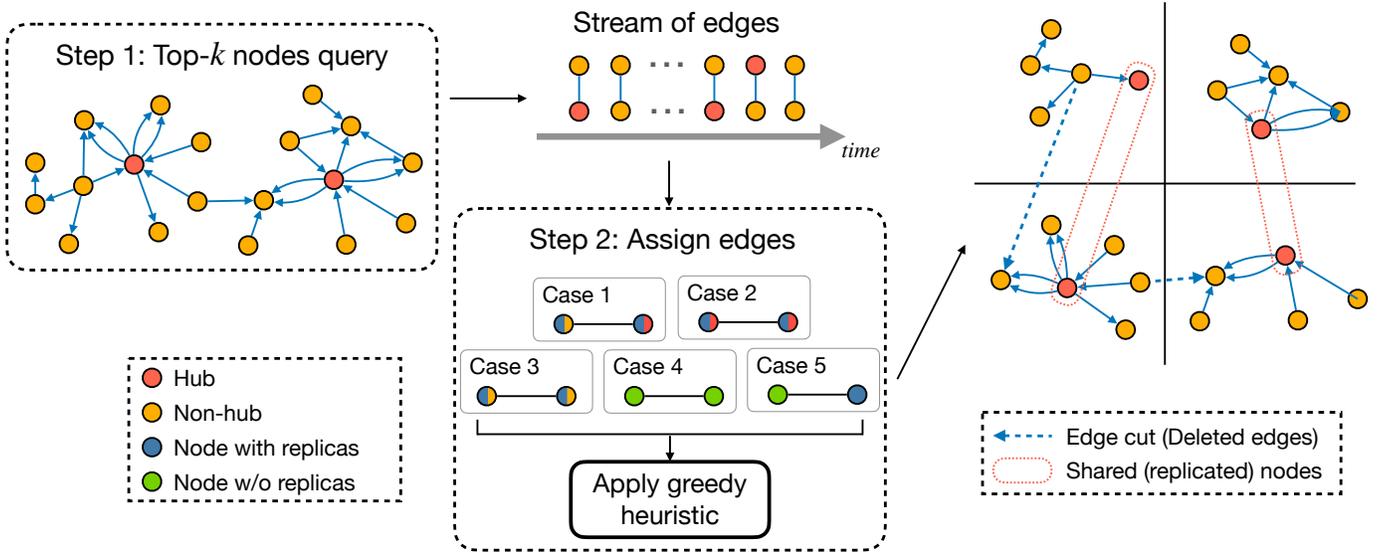


Fig. 4. Overview of the workflow of Streaming Edge Partitioning Component. Please refer to Alg.1 for details of the different cases in step 2.

before time  $t$ . In the case of a non-attributed graph, we make the assumption that node and edge features are zero vectors. For the sake of simplicity, we maintain a consistent dimensionality of  $d$  for node and edge features, as well as for node representations.

### B. Streaming Edge Partitioning Component (SEP)

As graph partitioning is a preprocessing step in distributed training, the quality of the partitioning may have great impact on the quality of the distributed training. To attain the objectives of the four dimensions as outlined in Tab. I, we employ the node-cut based streaming partitioning algorithm, complemented by two key innovations: First, in order to introduce temporal information, we employ an *exponential time decay* strategy to define the centrality of nodes. Second, we control the number of shared nodes in the process of edge assignment to avoid high replication factor while minimizing the information loss. An illustration can be found in Fig.4.

**Calculation of node’s centrality.** In TIG model training, whenever an event occurs (i.e., an edge is added), the information from that edge serves as an input, updating the representations of the involved nodes. Temporal neighbor sampling is essential to ensure that a node can acquire ample information about its neighbors, thus mitigating the “staleness problem” [4]. A common method for temporal neighbor sampling is sampling only the most recent neighbors. Our intuition is that the more recent events often have a greater impact on node’s future actions [3]. Therefore, before the graph partitioning phase, we will assign a greater weight to more recently appearing edges. Based on the observation above, we propose an *exponential time decay* strategy for edge weight computation, which is commonly applied in temporal neighbor sampling [15]–[17]. Concurrently, the node’s centrality is determined by aggregating the weights of all its historically connected edges. Let  $\mathcal{T}(i)$  denotes the set of timestamps for all historical edges

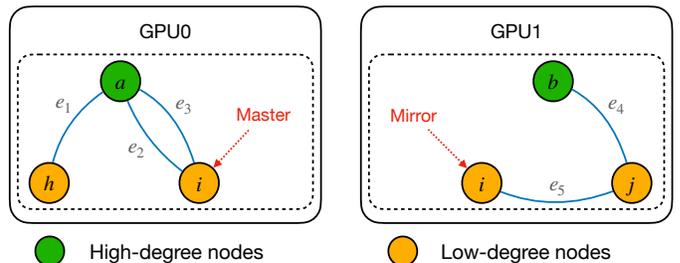


Fig. 5. An example of a low-degree node  $i$  being replicated in HDRF [14]. The order of occurrence of the edges is  $(e_1, e_2, e_3, e_4, e_5)$ . The mirror node  $i$  is replicated in GPU1 when edge  $e_5 = (i, j)$  appears because HDRF tends to replicate node with higher degree.

of node  $i$  and  $t_{max}$  as the last timestamp. For node  $i$ , its centrality is defined as:

$$Cent(i) = \sum_{t \in \mathcal{T}(i)} \exp(\beta(t - t_{max})), \quad (1)$$

where  $\beta \in (0, 1)$  is a scalar hyper-parameter.

**Streaming node-cut graph partitioning algorithm.** The standard greedy heuristic introduced in [13] may yield highly imbalanced partition sizes and a high replication factor [14], as it treats nodes with varying degrees as equals. In the HDRF algorithm [14], the degree of nodes is factored into the greedy heuristic. But the HDRF only considers the static graph and takes it as stream of edges that are input in some order like DFS or BFS. In real-world TIG partitioning, however, the appearance of some unpredictable edges can cause a large number of low-degree nodes to be replicated. As shown in Fig.5, low-degree node  $i$  is replicated in GPU1 when edge  $e_5 = (i, j)$  appears.

To avoid such situations, we first query the top- $k$  most important nodes (a hyper-parameter of our method, noted as

---

**Algorithm 1: Streaming Edge Partitioning**

---

**Input:** Edge set  $\mathcal{E}$ , Node set  $\mathcal{V}$ , Proportion of replicable nodes  $k$

**Output:** Partition of nodes  $\mathcal{P} = (p_1, \dots, p_n)$  and a shared nodes list  $\mathcal{S}$

```
1 Scan the input edge set  $\mathcal{E}$  to calculate the centrality
  value  $Cent(i)$  of each node  $i$  in  $\mathcal{V}$ . Return the set of
  nodes  $\mathcal{T}$  with the first  $k * |\mathcal{V}|$  largest centrality value;
2 for each edge  $e = (i, j, t)$  in  $\mathcal{E}$  do
3    $A(i)$  is the set of partitions to which node  $i$  has
   been assigned;
4   if  $A(i) \neq \emptyset$  and  $A(j) \neq \emptyset$  then
5     Case 1:  $i(j) \in \mathcal{T}$  and  $j(i) \notin \mathcal{T}$ ;
6     Assign  $e$  to the partition  $\hat{p}$  where  $j(i)$  resides;
7     Case 2:  $i \in \mathcal{T}$  and  $j \in \mathcal{T}$ ;
8     Assign  $e$  to the partition  $\hat{p}$  with maximum
     score  $C(i, j, p)$ ;
9     Case 3:  $i \notin \mathcal{T}$  and  $j \notin \mathcal{T}$ ;
10    if  $A(i) == A(j) == \hat{p}$  then
11      Assign  $e$  to the partition  $\hat{p}$ ;
12    else
13      Discard the edge  $e$ ;
14  else
15    Case 4 & 5: either both nodes are unassigned
    or one of the nodes is assigned;
16    Assign  $e$  to the partition  $\hat{p}$  with maximum
    score  $C(i, j, p)$ ;
17 for each node  $i \in \mathcal{V}$  do
18   if  $|A(i)| > 1$  then
19     Add  $i$  to the shared nodes list  $\mathcal{S}$ ;
20     Add  $i$  to all the partitions in  $\mathcal{P}$ ;
21   else
22     Add  $i$  to the  $p \in A(i)$ 
```

---

$top_k$ ) as *hubs*, based on the node centrality calculation in the previous step. Then we take the input TIG as stream of edges and sequentially assign edge to a specific partition. During the edge assigning phase, we restrict replication to nodes in *hubs*, thereby reducing the replication factor and preserving edge information as much as possible. Moreover, to maintain load balance, we employ a greedy heuristic at lines 8 and 16 in Alg.1. Specifically, when an input edge  $e = (i, j, t)$  is being processed, the normalized centrality value of nodes  $i, j$  are defined as:

$$\theta(i) = \frac{Cent(i)}{Cent(i) + Cent(j)} = 1 - \theta(j). \quad (2)$$

Then we compute a score  $C(i, j, p)$  for each partition  $p$  and greedily assign the edge to the partition with the maximum score  $C$  defined as follows:

$$C(i, j, p) = C_{REP}(i, j, p) + C_{BAL}(p). \quad (3)$$

The first term  $C_{REP}$  is to ensure that edge is assigned to the partition where the lower centrality node resides (line 8 in Alg.1) or to the partition where the node already be assigned (line 16 in Alg.1). The second term  $C_{BAL}$  is to maintain load balancing by assigning edge to the smallest partition. More specifically,  $C_{REP}$  is defined as:

$$C_{REP}(i, j, p) = h(i, p) + h(j, p). \quad (4)$$

For node  $i$  and partition  $p$ ,  $h(i, p)$  is defined as:

$$h(i, p) = \begin{cases} 1 + (1 - \theta(i)), & \text{if } p \in A(i) \\ 0, & \text{otherwise} \end{cases}, \quad (5)$$

where  $A(i)$  denotes to the set of partitions to which node  $i$  has been assigned. The second term  $C_{BAL}$  is defined as:

$$C_{BAL}(p) = \lambda \cdot \frac{maxsize - |p|}{\epsilon + maxsize - minsize}. \quad (6)$$

The parameter  $\lambda$  ( $\lambda > 0$ ) manages the impact of load balancing in greedy heuristic. Meanwhile,  $\epsilon$  is a small constant added to avoid zero denominators in the calculations, with *maxsize* and *minsize* defining the upper and lower limits of the partition size, respectively.

At lines 17-22 in Alg.1, after edge assigning, we add the nodes with replicas in more than one partition to the shared nodes list, which is used as input for subsequent distributed training.

**Theoretical Analysis.** In this section, we perform a theoretical analysis of SEP, particularly examining the worst-case scenarios for the replication factor (*RF*) and edge cuts (*EC*).

**Metrics.** We use the metrics as follows:

$$RF = \frac{Total\ node\ replicas}{Total\ number\ of\ nodes}, \quad (7)$$

$$EC = \frac{Total\ edge\ cuts\ between\ partitions}{Total\ number\ of\ edges}. \quad (8)$$

**Theorem 1.** When partitioning a graph with  $|\mathcal{V}|$  nodes on  $|\mathcal{P}|$  partitions, our algorithm achieves a upper bound of *RF* as:

$$RF \leq k|\mathcal{P}| + (1 - k). \quad (9)$$

*Proof.* The first term signifies the replicas generated by the fraction  $k$  of nodes (referred as *hubs*) with the highest centrality in the graph. In the worst-case scenario, all nodes within the *hubs* have duplicated copies across all partitions. The second term consider the replicas created by non-*hubs*. In our algorithm, they each can create at most one replica.  $\square$

Cohen et al. [18] demonstrated that when a fraction  $k$  of nodes with the highest degrees is removed from a power-law graph (along with their edges), the maximum node degree  $M$  in the remaining graph can be approximated as:

$$M = mk^{\frac{1}{1-\alpha}}, \quad (10)$$

where  $m$  is the minimum node degree in the graph and  $\alpha$  is a parameter that indicates the skewness of the graph.

To provide a theoretical upper bound on *EC*, we directly employ the degree of a node as its centrality value.

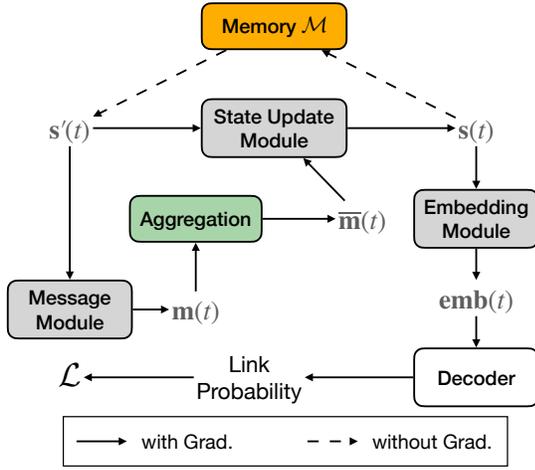


Fig. 6. Illustration of data flows of TIG models [5]. Note that the memory module is constantly being updated.

**Theorem 2.** *When partitioning a power-law graph with  $|\mathcal{V}|$  nodes and  $|\mathcal{E}|$  edges on  $|\mathcal{P}|$  partitions, our algorithm achieves an upper bound of EC as:*

$$EC \leq \frac{1}{|\mathcal{E}|} \sum_{q=0}^{|\mathcal{V}|(1-k)-1} m(k + \frac{q}{|\mathcal{V}|})^{\frac{1}{1-\alpha}}. \quad (11)$$

*Proof.* According to Alg.1, edge dropping only occurs during the execution of *Case 3*. It happens when both nodes are non-hubs and have replicas in different partitions. Therefore, edges connected to hubs are preserved. The worst-case scenario is that all edges connecting two non-hubs are dropped. This scenario occurs when all non-hubs are connected to hubs upon their first appearance, and all edges between two non-hubs cross partitions. Given that the largest degree of the graph excluding hubs is  $mk^{\frac{1}{1-\alpha}}$ , when we remove a non-hub node with the highest degree, the highest degree in the remaining graph is  $m(k + \frac{1}{|\mathcal{V}|})^{\frac{1}{1-\alpha}}$ . Therefore, we can bound EC by counting all edges connecting two non-hubs.  $\square$

### C. Parallel Acceleration Component (PAC)

**Training approach for TIGs.** The data flows of most models are shown in Fig.6. A TIG model typically adopts an Encoder-Decoder structure, with the Encoder composed of four key modules. To avoid repeat calculation, the *Memory Module*  $\mathcal{M} \in \mathbb{R}^{N \times d}$  is employed to capture the historical interaction information for each node  $i$ , represented as  $\mathcal{M}_i$ . The *Message Module* is used to compute the current state, i.e., message  $\mathbf{m}_i(t)$  of nodes. For an interaction event  $e_{ij}(t)$ , messages are computed by the previous states  $\mathbf{s}'_i(t)$ ,  $\mathbf{s}'_j(t)$ , event feature vector  $\mathbf{e}_{ij}(t)$ , and time encoding computed by  $\Delta t$ . The message computing functions (*MSG*) are learnable and can be chosen from different neural networks or just simply concatenate the inputs. We use node  $i$  as an example, where  $\Phi$  is the time encoding function [3]:  $\mathbf{m}_i(t) = MSG(\mathbf{s}'_i(t), \mathbf{s}'_j(t), \Phi(\Delta t), \mathbf{e}_{ij}(t))$ . Previous states are fetched from node memory, i.e.,  $\mathbf{s}'_i(t) \leftarrow \mathcal{M}_i$  and  $\mathbf{s}'_j(t) \leftarrow \mathcal{M}_j$ .

All messages in the same batch of a node can be aggregated by an aggregation function, which can be simply mean message or other neural networks, e.g., RNN, and output the aggregated message  $\bar{\mathbf{m}}_i(t)$ . After an event which involves node  $i$  happened, the new node representation can be updated by the node state before the event  $\mathbf{s}'_i(t)$  and the aggregated message  $\bar{\mathbf{m}}_i(t)$  by applying the *State Update Module*:  $\mathbf{s}_i(t) = UPD(\mathbf{s}'_i(t), \bar{\mathbf{m}}_i(t))$ , where  $\mathbf{s}'_i(t) \leftarrow \mathcal{M}_i$ . The state update function can also be chosen from different learnable neural networks, e.g., GRU, RNN. Upon computing the new state, the memory of node  $i$  is updated by overwriting the corresponding value  $\mathcal{M} \leftarrow \mathbf{s}_i(t)$ . Finally, the *Embedding Module* is used to compute the node embedding  $\mathbf{emb}_i(t)$  at a specific time  $t$ :  $\mathbf{emb}_i(t) = \sum_{j \in neighbor_i^k([0, t])} f(\mathbf{s}_i(t), \mathbf{s}_j(t), \mathbf{e}_{ij}(t))$ , where  $f$  is a learnable function, e.g., identity, time projection or attention. The Decoder  $g$  can calculate the probability of the edge existence between two nodes:  $p_{ij}(t) = g(\mathbf{emb}_i(t), \mathbf{emb}_j(t))$ , which then provides the self-supervised signals for training.

In our approach to training various TIG models, we first establish a general architecture for most TIG models. This is done by integrating different modules into a single architecture. This means all implemented models are specific instances of our approach. Moreover, our approach allows these models to be easily extended to accommodate other new models.

**Distributed Parallel Training.** Our approach primarily employs multi-GPU computation acceleration to facilitate parallel training of TIG models. To make this possible, the original large graph is divided into several partitions using our SEP component. An illustration is in Fig.2(b).

The outputs of SEP component are nodes lists  $\{\mathcal{V}_1, \dots, \mathcal{V}_p\}$  from which we construct new sub-graphs  $\{\mathcal{G}_1, \dots, \mathcal{G}_p\}$  by identifying edges  $\mathcal{E}_k = \{(i, j, t) \in \mathcal{E} | i, j \in \mathcal{V}_k\}$  in the original dataset, and we need  $\mathcal{N}$  partitions for training. We can choose to divide the original graph to  $\mathcal{N}$  partitions directly ( $|\mathcal{P}| = \mathcal{N}$ ). However, after the graph is partitioned, some edges are inevitably deleted. Note that we have  $\mathcal{V}_a \cup \mathcal{V}_b$  with edge sets being  $\mathcal{E}_a \cup \mathcal{E}_b \cup \mathcal{D}\mathcal{E}_{ab}$ , where  $\mathcal{D}\mathcal{E}_{ab} = \{e_{ij}(t) | i \in \mathcal{V}_a, j \in \mathcal{V}_b\}$  refers to deleted edges between sub-graph  $\mathcal{G}_a$  and  $\mathcal{G}_b$ . We proposed two strategies to relieve the information loss caused by the edge deletion. As outlined in Sec.II-B, shared nodes are added to reduce information loss. We can also initially divide the graph into more parts  $\{\mathcal{V}_1, \dots, \mathcal{V}_{|\mathcal{P}|}\}$ ,  $|\mathcal{P}| > \mathcal{N}$ . Prior to each training epoch, we randomly shuffle all parts and combine them to form  $\mathcal{N}$  partitions. When two small partitions are combined, the combined partition will contain the “deleted” edges between the two small partitions *combined*( $\mathcal{V}_a, \mathcal{V}_b$ ) with edge sets being  $\mathcal{E}_a \cup \mathcal{E}_b \cup \mathcal{D}\mathcal{E}_{ab}$ . Through random shuffling, the “deleted” edges between different small partitions can be restored when they are combined, allowing them to be trained across different epochs.

For distributed parallel training based on graph partitioning, we have  $\mathcal{N}$  GPUs, and the model will be duplicated into  $\mathcal{N}$  copies and deployed on each GPU. The graph data used for training on different GPUs are different sub-graphs which is one of the partitions of the original training graph, i.e., the original graph is partitioned into  $\mathcal{N}$  parts. Assuming

---

**Algorithm 2: Training Loop**

---

**Data:** Test Dataloader

```
1  $i\_batch = 0$ ;  
2 while  $True$  do  
3   for  $i$ ,  $data$  in  $enumerate(Test\ Dataloader)$  do  
4      $loop\_start = i == 0$ ;  
5      $loop\_end = i == len(dataloader) - 1$ ;  
6     if  $loop\_start$  then  
7        $\lfloor$  reset nodes memory;  
8     Training and compute loss;  
9     Backward and Optimize;  
10    if  $loop\_end$  then  
11       $\lfloor$  backup current nodes memory state;  
12     $i\_batch += 1$ ;  
13  if All GPUs have finished at least one complete  
    traversal for all mini batches then  
14     $\lfloor$  break
```

---

we have a partitioning of nodes represented as  $\{\mathcal{V}_1, \dots, \mathcal{V}_N\}$ , the corresponding interactions, i.e., sub-graphs can be written as  $\mathcal{E}_k = \{(i, j, t) \in \mathcal{E} | i, j \in \mathcal{V}_k\}$ . Thus we can train sub-graphs parallel at the same time on different GPUs. While only  $\mathcal{M}^{(k)} \in \mathbb{R}^{N^{(k)} \times d}$  memory module is needed for a single-GPU.

In order to balance GPU computational resource utilization and allow for training on larger graphs, our TIG partitioning algorithms ensure that the node counts in each partition are balanced. Based on this setting, we can initialize a memory store module for each GPU with only maximization of all GPUs nodes count. This help us to put graph with very large number of nodes on GPUs.

However, the interaction events, i.e. edges, assigned to different sub-graphs are not exactly the same. Therefore, in order to synchronize the training and the backward of gradients between different GPUs, we design a new training approach. In each epoch, all edges on each GPU are traversed at least once. On GPUs with fewer edges, a loop is made within the epoch. Since the end states of different GPUs may not all be at the end state of a complete data cycle when the entire epoch ends training, this will result in incomplete memory updates of parts of nodes. Thus, we create a backup of the node memory lists at the end of each GPU’s data cycle. At the end of the entire epoch, we then restore the node memory across all GPUs to match these latest memory backups. The pseudocode of the training phase is shown in Alg.2.

If there are shared nodes, for them, we ensure node memory synchronization across all GPUs. There are two ways of node memory synchronization. The first approach sets the memory value of all shared nodes on every GPU to match memory with the largest timestamp recorded across all GPUs. The second approach resets the memory of all shared nodes by taking an average across all GPUs. After experimental testing, we found that the two synchronization methods have little impact on the

TABLE II

DATASET STATISTIC.  $d_n$  AND  $d_e$  INDICATE THE DIM OF NODES AND EDGES, RESPECTIVELY. CLASSES MEANS THE NUMBERS OF LABELS.

	# Nodes	# Edges	$d_n$	$d_e$	Classes
Wikipedia	9,227	157,474	172	172	2
Reddit	10,984	672,447	172	172	2
MOOC	7,144	411,749	172	172	2
LastFM	1,980	1,293,103	172	172	-
ML25m	221,588	25,000,095	100	1	-
DGraphFin	4,889,537	4,300,999	100	11	4
Taobao	5,149,747	100,135,088	100	4	9,439

performance of downstream tasks, and we adopted the former in our experiments.

Our distributed parallel acceleration approach allows existing TIG models to be adapted and accelerated. This not only reduces the consumption of computational resources but also significantly decreases computation time. Models which applied our approach also exhibits competitive performance in downstream tasks.

### III. EXPERIMENTS

#### A. Datasets, Models, Basic Settings

We applied SPEED on 7 temporal interaction graph datasets - ML25m [19], DGraphFin [20], Taobao [21], Wikipedia, Reddit, MOOC and LastFM [1]. Among them, ML25m, DGraphFin and Taobao are three large datasets. Some datasets are lack of node features or edge attributes, for which we follow the processing method in previous works [1], [4] by using zero vectors to represent them. State changes indicators of nodes, i.e., dynamic labels, are included in Wikipedia, Reddit and MOOC, which makes them support dynamic node classification task. The statistics of all datasets can be found in Tab.II. To avoid information leakage, we chronologically divided the edges into 70% for training, 15% for validation, and 15% for testing before implementing our SEP. To prove the efficiency of our methods, we conduct a series of experiments using on 4 different models, i.e., Jodie [1], DyRep [2], TGN [4], TIGE [5].

Given the variations between datasets, we apply distinct experimental settings to the large and the small datasets. To optimize training time, we use larger batch sizes for the larger datasets. Specifically, we apply a batch size of 200 to the 4 small datasets, and batch sizes of 2,000, 2,000, and 1,000 are employed when training on ML25m, DGraphFin, and Taobao, respectively. Given that the number of edges in the large datasets exceeds that of the small datasets by over 10 times, we apply smaller maximum training epochs and patience values for the 3 large datasets.

All experiments are conducted on a single server with 72-core CPU, 128GB of memory, and 4 Nvidia Tesla V100 GPUs.

#### B. Main Experiments

In our experiments, we manipulate the number of shared nodes by adjusting the parameter  $top_k$ . The choice of  $top_k$  is

TABLE III

RUNNING (TRAINING) TIMES (PER EPOCH IN SECONDS) AND GPU MEMORY RESOURCES RESERVED (PER GPU IN GB) FOR 3 BIG DATASETS. USES TRAINING ON CPU AS SPEED COMPARISON BASELINE. (RESULTS ARE BASED ON SELF-SUPERVISED TRAINING OF LINK PREDICTION. ALL OUR METHODS ARE TRAINED ON A 4X NVIDIA TESLA V100 MACHINE AND GRAPHS ARE DIVIDED INTO 4 PARTITIONS.)

		ML25m			DGraphFin			Taobao		
		Training Time	Speed -up	GPU Mem. Reserved	Training Time	Speed -up	GPU Mem. Reserved	Training Time	Speed -up	GPU Mem. Reserved
Jodie	$top_k = 0$	226.50	<b>5.90x</b>	<b>0.88</b>	102.70	<b>5.00x</b>	<b>10.19</b>	917.29	<b>8.88x</b>	<b>12.25</b>
	$top_k = 1$	414.73	3.22x	0.91	105.02	4.89x	13.44	2096.44	3.89x	16.59
	$top_k = 5$	773.50	1.73x	1.02	117.47	4.37x	14.72	4866.02	1.67x	19.05
	$top_k = 10$	1045.46	1.28x	1.17	134.26	3.82x	15.98	7326.25	1.11x	24.62
	HDRF	1177.96	1.13x	1.90	OOM	OOM	OOM	OOM	OOM	OOM
	Single-GPU	1313.66	1.02x	2.82	OOM	OOM	OOM	OOM	OOM	OOM
	CPU	1335.61	1x	-	513.30	1x	-	8147.59	1x	-
Dyrep	$top_k = 0$	236.99	<b>7.12x</b>	<b>1.25</b>	105.18	<b>7.69x</b>	<b>10.20</b>	961.75	<b>15.64x</b>	<b>12.26</b>
	$top_k = 1$	431.77	3.91x	1.28	107.37	7.54x	12.92	2219.06	6.78x	16.13
	$top_k = 5$	812.03	2.08x	1.45	121.85	6.64x	14.57	5186.08	2.90x	20.79
	$top_k = 10$	1095.72	1.54x	1.61	136.58	5.92x	15.90	7964.41	1.89x	24.62
	HDRF	1238.93	1.36x	2.30	OOM	OOM	OOM	OOM	OOM	OOM
	Single-GPU	1403.57	1.20x	3.42	OOM	OOM	OOM	OOM	OOM	OOM
	CPU	1687.00	1x	-	809.19	1.00	-	15045.97	1x	-
TGN	$top_k = 0$	238.35	<b>11.13x</b>	<b>1.26</b>	107.49	<b>8.56x</b>	<b>11.89</b>	976.55	<b>18.83x</b>	<b>14.08</b>
	$top_k = 1$	441.81	6.00x	1.27	111.95	8.22x	14.98	2266.76	8.11x	18.04
	$top_k = 5$	821.24	3.23x	1.46	126.33	7.29x	16.35	5312.83	3.46x	21.69
	$top_k = 10$	1109.66	2.39x	1.64	140.46	6.55x	17.84	8058.26	2.28x	25.50
	HDRF	1255.80	2.11x	2.33	OOM	OOM	OOM	OOM	OOM	OOM
	Single-GPU	1358.06	1.95x	3.92	OOM	OOM	OOM	OOM	OOM	OOM
	CPU	2652.76	1x	-	920.37	1x	-	18389.03	1x	-
TIGE	$top_k = 0$	251.35	<b>11.20x</b>	<b>1.26</b>	108.10	<b>8.00x</b>	<b>11.89</b>	984.37	<b>19.27x</b>	<b>14.08</b>
	$top_k = 1$	443.86	6.34x	1.27	112.56	7.69x	14.98	2273.37	8.34x	18.04
	$top_k = 5$	825.69	3.41x	1.44	125.76	6.88x	16.35	5357.19	3.54x	21.69
	$top_k = 10$	1109.07	2.54x	1.63	142.37	6.08x	17.85	8110.54	2.34x	25.50
	HDRF	1231.67	2.28x	2.31	OOM	OOM	OOM	OOM	OOM	OOM
	Single-GPU	1284.87	2.19x	3.92	OOM	OOM	OOM	OOM	OOM	OOM
	CPU	2813.96	1x	-	865.23	1x	-	18967.77	1x	-

crucial since the proportion of “important nodes” varies across real-world datasets. A suitable value of  $top_k$  is necessary to strike an optimal balance between edge-cuts across partitions and the workload of the machines in the cluster.

1) *Efficiency and computing resources*: To illustrate our advantages in training efficiency and computing resource utilization, we conduct a group of experiments on the 3 big datasets, wherein both training times and the GPU memory allocated by PyTorch modules are recorded. The results are presented in Tab.III.

Across all experiments, our methods proved to be the fastest and consumed the fewest GPU resources (on each single GPU). Additionally, the performance of the downstream tasks is highly competitive. By employing our method, training times can be accelerated by up to 8.56x on ML25m, 11.20x on DGraphFin, and 19.27x on Taobao. In terms of computing resources, our approach is essential for managing large datasets such as DGraphFin and Taobao. Their extensive number of nodes would otherwise lead to out-of-memory (OOM) issues during direct training, particularly due to the storage requirements for extensive node memory (detailed analysis is in III-B3). With the assistance of our method, we distribute the large graph across different GPUs for training to address the OOM problem. Simultaneously, our method alleviates

the computational burden on a single GPU and accelerates training. The efficiency advantages of our partitioning method, compared to the static graph partitioning method, are discussed in Sec.III-D2.

2) *Performance on down-stream tasks*: **Temporal link prediction**. Following the experimental settings previously outlined in [5], we observe performance on temporal link prediction tasks executed in both transductive and inductive manners. Average precision (AP) score is being used as the evaluation metric. In the case of transductive tasks, the nodes of the predicted edges have appeared during training. Conversely, inductive tasks focus on predicting temporal links between nodes that have not previously been encountered. The results for different tasks are shown in Tab.IV. It is evident that our approach manages to deliver competitive results in downstream tasks with a faster training speed. However, due to OOM issues, some results for larger datasets, that do not applied our method, are unavailable. The best performance on link prediction for both transductive and inductive settings tends to be more focused when  $top_k$  is higher. The algorithm HDRF [14] cannot control the number of shared nodes. This may also lead to OOM issues, as excessive node replication and distribution across GPUs occur. With smaller  $top_k$ , there are also part of best performance. Our method achieves equi-

TABLE IV

AVERAGE PRECISION (%) FOR FUTURE EDGE PREDICTION TASK IN TRANSDUCTIVE AND INDUCTIVE SETTINGS.  $top_k$  REFERS TO THE HYPER-PARAMETER THAT CONTROL THE PERCENTAGE OF SHARED NODES. WHEN THERE IS NO RESTRICTION FOR  $top_k$  THE ALGORITHM DEGENERATES TO HDRF. THE BEST AP OF DIFFERENT  $top_k$  ON DIFFERENT DATASETS AND BACKBONES ARE IN BOLD. PARTS OF THE RESULTS ARE UNAVAILABLE FOR ML25M, DUE TO THE ORIGINAL METHODS CANNOT EFFICIENTLY PROCESS THE DATASET DIRECTLY.

		Wikepeida	Reddit	MOOC	LastFM	ML25m	DGraphFin	Taobao	
Transductive	Jodie	$top_k = 0$	94.87±0.7	94.13±0.6	61.54±0.8	65.12±5.1	92.22±2.1	74.75±1.7	91.43±1.1
		$top_k = 1$	91.12±0.3	90.29±1.7	61.87±1.2	67.69±0.9	93.70±0.5	<b>74.77±0.6</b>	91.28±0.6
		$top_k = 5$	94.48±1.2	94.07±1.5	<b>63.49±0.4</b>	66.89±2.8	93.73±0.9	72.69±1.9	<b>91.92±0.6</b>
		$top_k = 10$	<b>95.21±0.5</b>	<b>94.56±1.1</b>	63.33±0.4	<b>69.91±1.1</b>	94.67±0.1	71.84±1.5	90.05±1.2
		HDRF	94.02±0.9	94.91±0.1	64.78±1.6	70.65±2.6	95.01±0.1	OOM	OOM
		w/o Partitioning	94.62±0.5	97.11±0.3	76.50±1.8	68.77±3.0	N/A	OOM	OOM
	DyRep	$top_k = 0$	90.68±0.9	93.07±0.4	59.69±3.2	58.66±1.9	<b>93.99±1.4</b>	<b>78.70±0.2</b>	68.71±12.2
		$top_k = 1$	81.82±1.7	87.14±3.8	60.15±1.3	55.77±3.8	87.38±4.6	77.86±0.1	87.07±1.9
		$top_k = 5$	90.29±1.3	<b>93.47±1.2</b>	<b>62.76±6.6</b>	53.53±2.1	92.74±0.7	77.73±0.1	<b>89.45±0.5</b>
		$top_k = 10$	<b>91.58±1.5</b>	92.24±0.3	60.27±8.4	<b>61.56±2.1</b>	92.01±1.4	78.15±0.1	83.31±3.8
		HDRF	91.21±0.7	<b>93.61±1.0</b>	59.16±1.5	71.43±2.2	93.76±0.5	OOM	OOM
		w/o Partitioning	94.59±0.2	97.98±0.1	75.37±1.7	68.77±2.1	N/A	OOM	OOM
	TGN	$top_k = 0$	97.50±0.2	<b>96.69±0.3</b>	61.54±0.8	55.97±4.3	93.39±1.0	82.16±0.3	85.75±1.3
		$top_k = 1$	95.90±0.1	92.50±1.1	77.34±5.4	<b>62.79±6.0</b>	91.32±5.5	<b>82.27±0.2</b>	84.22±2.9
		$top_k = 5$	97.35±0.1	95.84±1.2	81.44±2.7	58.98±3.7	94.20±0.7	82.08±0.2	81.93±8.6
		$top_k = 10$	<b>97.61±0.1</b>	95.53±1.4	<b>85.66±0.9</b>	53.64±5.4	<b>94.51±0.5</b>	81.98±0.1	<b>85.93±2.0</b>
		HDRF	96.36±1.9	94.59±2.0	78.37±8.4	53.47±3.6	94.53±0.2	OOM	OOM
		w/o Partitioning	98.46±0.1	98.70±0.1	85.88±3.0	71.76±5.3	N/A	OOM	OOM
	TIGE	$top_k = 0$	98.34±0.0	98.19±0.1	75.32±1.1	<b>83.05±0.7</b>	93.09±2.4	82.11±0.1	88.03±7.4
		$top_k = 1$	98.18±0.1	97.81±0.2	84.80±1.0	82.20±0.7	92.32±2.3	82.41±0.1	88.81±3.0
$top_k = 5$		98.22±0.1	<b>98.29±0.4</b>	86.80±1.9	82.43±0.6	92.73±2.8	<b>82.75±0.1</b>	<b>89.88±4.2</b>	
$top_k = 10$		<b>98.50±0.0</b>	98.07±0.1	<b>88.28±1.6</b>	82.59±0.3	<b>94.12±0.5</b>	82.57±0.2	89.04±2.9	
HDRF		97.98±0.4	97.89±0.3	86.17±0.6	82.78±0.8	91.98±0.0	OOM	OOM	
w/o Partitioning		98.83±0.1	99.04±0.0	89.64±0.9	87.85±0.9	N/A	OOM	OOM	
Inductive	Jodie	$top_k = 0$	93.59±0.7	92.86±0.2	61.71±0.3	70.43±5.6	91.75±2.2	68.76±1.6	79.70±2.1
		$top_k = 1$	88.95±0.7	91.39±1.2	60.37±1.3	74.58±0.7	93.21±0.6	68.91±0.1	79.71±0.7
		$top_k = 5$	92.99±1.4	93.18±1.7	61.82±0.5	74.53±3.3	93.31±1.0	69.21±0.5	<b>81.18±1.0</b>
		$top_k = 10$	<b>93.96±0.5</b>	<b>93.46±1.1</b>	<b>62.98±1.1</b>	<b>79.68±1.2</b>	<b>94.27±0.1</b>	<b>69.59±0.2</b>	78.25±2.7
		HDRF	92.29±0.6	93.05±0.6	64.90±0.9	82.67±0.9	94.66±0.1	OOM	OOM
		w/o Partitioning	93.11±0.4	94.36±1.1	77.83±2.1	82.55±1.9	N/A	OOM	OOM
	DyRep	$top_k = 0$	89.40±0.5	93.11±0.4	60.05±4.1	69.82±1.6	<b>93.68±1.3</b>	<b>64.28±0.1</b>	59.07±6.8
		$top_k = 1$	81.33±1.4	90.00±2.9	58.68±1.5	63.70±6.7	86.72±4.9	64.18±0.3	75.25±2.8
		$top_k = 5$	89.12±1.4	<b>93.78±0.8</b>	<b>62.82±6.1</b>	61.03±3.1	92.43±0.7	64.20±0.3	<b>78.35±0.8</b>
		$top_k = 10$	<b>90.14±1.2</b>	92.11±0.5	59.57±7.4	<b>74.42±2.9</b>	91.63±1.5	64.12±0.2	71.27±4.7
		HDRF	89.70±0.2	92.62±0.9	59.45±0.5	82.67±0.9	93.66±0.4	OOM	OOM
		w/o Partitioning	92.05±0.3	95.68±0.2	78.55±1.1	81.33±2.1	N/A	OOM	OOM
	TGN	$top_k = 0$	96.86±0.2	<b>94.84±0.5</b>	76.14±1.0	59.32±3.8	93.29±0.9	66.11±0.7	<b>74.08±0.2</b>
		$top_k = 1$	95.44±0.4	90.71±0.3	77.53±3.5	<b>66.94±8.4</b>	91.09±5.5	66.51±0.7	68.58±2.3
		$top_k = 5$	96.62±0.1	93.95±1.7	81.20±3.2	63.79±5.9	93.80±0.8	66.82±0.4	70.23±6.0
		$top_k = 10$	<b>96.90±0.1</b>	93.49±1.9	<b>86.29±0.7</b>	55.40±8.3	<b>94.10±0.5</b>	<b>67.39±0.8</b>	73.34±2.6
		HDRF	96.05±1.4	92.31±2.6	79.47±8.3	52.91±10.8	94.07±0.2	OOM	OOM
		w/o Partitioning	97.81±0.1	97.55±0.1	85.55±2.9	80.42±4.9	N/A	OOM	OOM
	TIGE	$top_k = 0$	98.05±0.1	97.40±0.3	79.30±3.6	<b>85.32±1.1</b>	92.59±2.6	<b>65.98±0.4</b>	79.43±4.9
		$top_k = 1$	97.87±0.1	96.86±0.3	84.41±1.7	85.27±0.9	92.27±2.1	65.41±0.2	76.78±4.5
$top_k = 5$		97.91±0.1	<b>97.50±0.4</b>	86.30±1.8	84.84±1.1	92.31±3.0	65.98±0.7	<b>79.46±2.6</b>	
$top_k = 10$		<b>98.17±0.0</b>	97.25±0.2	<b>87.90±1.6</b>	85.12±0.4	<b>93.79±0.5</b>	65.60±0.3	78.37±2.3	
HDRF		97.68±0.2	96.96±0.2	86.38±0.6	84.44±0.6	91.62±0.1	OOM	OOM	
w/o Partitioning		98.45±0.1	98.39±0.1	89.51±0.7	90.14±1.0	N/A	OOM	OOM	

librium by regulating the value of  $top_k$ . A larger  $top_k$  permits more shared nodes, thereby collating additional information beneficial to downstream tasks. Conversely, a smaller  $top_k$  helps filter out potential noise or irrelevant edges, further enhancing downstream task performance. Overall, managing  $top_k$  guarantees both acceleration on large datasets and advantages for downstream tasks.

**Node classification.** Considering node classification tasks

require dynamic labels, we conduct our experiments on datasets with available labels, i.e., Wikipedia, Reddit and MOOC. The performances evaluated by AUROC (area under the ROC curve) are shown in Tab.V, as the same as the previous works used [1]–[5]. As the results show, the application of our methods can yield results that surpass those achieved using the original models without partitioning, across various models and datasets. This further underscores the effectiveness

TABLE V

AUROC (%) FOR DYNAMIC/STATIC NODE CLASSIFICATION TASK. WE USE THE RESULTS REPORTED IN [5] WHICH TRAINED WITHOUT GRAPH PARTITIONING AS BASELINES AND PRESENT THE RESULTS OF HDRF ALGORITHM FOR COMPARING. THE BEST AUROC OF DIFFERENT  $top_k$  ON DIFFERENT DATASETS AND BACKBONES ARE IN BOLD.

		Wikepeida	Reddit	MOOC
Jodie	$top_k = 0$	87.54±0.7	58.98±3.6	64.27±1.2
	$top_k = 1$	<b>87.83±0.2</b>	62.04±2.9	62.28±0.9
	$top_k = 5$	87.52±1.1	<b>64.00±5.5</b>	61.11±0.7
	$top_k = 10$	86.97±1.0	61.90±0.4	<b>65.36±1.8</b>
	HDRF	87.82±0.5	63.26±5.5	65.72±1.6
	w/o Partitioning	84.84±1.2	61.83±2.7	66.87±0.4
DyRep	$top_k = 0$	85.92±0.8	<b>64.27±1.4</b>	63.58±1.6
	$top_k = 1$	86.53±0.9	62.29±5.1	62.12±3.1
	$top_k = 5$	<b>86.63±2.7</b>	62.01±1.4	64.12±1.5
	$top_k = 10$	86.26±1.2	61.88±4.4	<b>65.43±3.1</b>
	HDRF	86.95±0.4	63.59±1.3	63.90±3.8
	w/o Partitioning	84.59±2.2	62.91±2.4	67.76±0.5
TGN	$top_k = 0$	<b>86.39±1.0</b>	<b>67.77±4.5</b>	<b>73.61±1.1</b>
	$top_k = 1$	83.71±2.5	62.81±1.3	72.04±0.7
	$top_k = 5$	82.40±1.8	67.31±1.1	72.95±1.3
	$top_k = 10$	82.64±2.0	66.88±3.2	71.81±0.8
	HDRF	83.93±1.2	64.41±2.0	72.40±1.1
	w/o Partitioning	87.81±0.3	67.06±0.9	59.54±1.0
TIGE	$top_k = 0$	85.59±3.2	<b>64.16±0.5</b>	73.08±1.8
	$top_k = 1$	<b>86.89±1.4</b>	63.59±2.4	72.50±0.5
	$top_k = 5$	85.32±0.8	63.44±3.6	<b>73.63±0.7</b>
	$top_k = 10$	85.01±0.4	62.23±0.6	72.86±1.1
	HDRF	85.91±2.5	63.74±1.1	72.69±0.9
	w/o Partitioning	86.92±0.7	69.41±1.3	72.35±2.3

TABLE VI

EDGE CUT%, AVERAGE NODES PORTION AND STANDARD DEVIATION OF EDGES AND NODES ON EACH GPU, I.E., PARTITIONS, OF DATASET TAobao. NOTE THAT THE TAobao HAS 1E8 EDGES AND 5E7 NODES.

Dataset: Taobao	Edge Statistics		Node Statistics		
	Total Cut	Std.	Avg. Portion	Std.	
KL	20.5%	3.2e7	25.0%	0.5	
Ours	$top_k = 0$	69.5%	3.1e3	25.0%	5.9e3
	$top_k = 1$	40.1%	3.1e3	25.5%	1.1e4
	$top_k = 5$	21.4%	4.1e2	28.7%	4.9e3
	$top_k = 10$	8.5%	1.9e2	32.4%	8.8e3
	HDRF	0%	2.2	50.5%	2.1e3
	Random	75.1%	8.6e5	25.0%	1.2e3

of SPEED in enhancing the performance of downstream tasks.

3) *Information loss and load balancing analysis:* As the TIG partitioning deletes edges, and information loss occurs, the performance of downstream tasks will inevitably be affected. While adding shared nodes to different sub-graphs can alleviate some of the information loss caused by edge deletion, synchronizing node memory will also result in information loss from the shared nodes. Simultaneously, the imbalance in the distribution of edges and nodes can slow down the entire training process and impact GPU resource allocation, respectively.

Thus, we choose the largest dataset Taobao as an example for further investigation (statistics are shown in Tab.VI). we

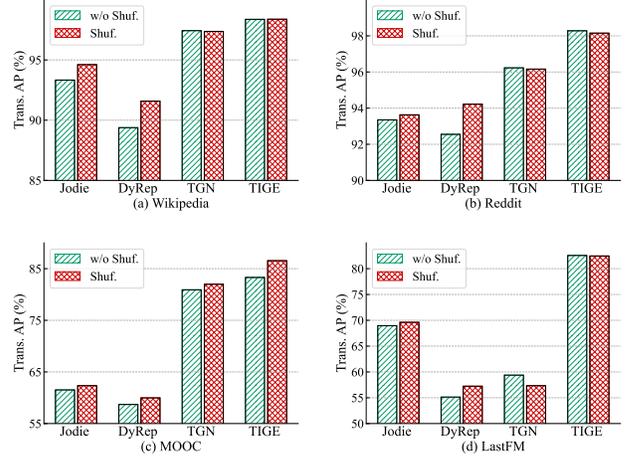
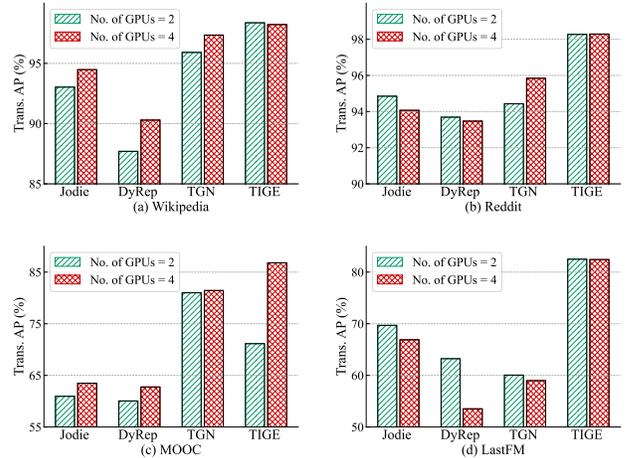


Fig. 7. Impacts of partition shuffling.

Fig. 8. Impacts of Changing the Number of GPUs ( $N$ ).

compute the total number of deleted edges (shown as “total cut” in the table), the edges on different sub-graphs, i.e, different GPUs (shown as the “edges std.” in the table) and the nodes on different sub-graphs (shown as the “nodes std.” in the table). We can find that as the number of shared nodes increases ( $top_k$ ), edge deletion decreases (edge cut). Our method is balancing the load in terms of both the number of edges and nodes. The balance of the number of edges enables our method to have a better acceleration effect and speed up the overall training time. At the same time, the balance of the number of nodes enables our method to better balance the use of different GPU resources, so that graphs with a larger number of nodes can be accelerated for training.

Compared to the KL method, due to the imbalanced distribution of edges, training speed is slower than with graphs partitioned by our SEP. Moreover, HDRF does not control the number of shared nodes, which results in a large number of nodes being distributed on GPUs. This results that although the graph with large number of nodes is partitioned by HDRF,

TABLE VII

LINK PREDICTION TASK PERFORMANCE RESULTS ARE IN AVERAGE PRECISION (%). TRAINING TIMES ARE IN SECONDS PER EPOCH.  $top_k = 0$  REFERS TO OUR METHODS. NOTE THAT OUR METHODS WITH OTHER  $top_k$  VALUE MAY OUTPERFORM KL ON DOWNSTREAM TASKS.

		ML25m			DGraphFin			Taobao		
		Transductive	Inductive	Training Time (Speed-up)	Transductive	Inductive	Training Time (Speed-up)	Transductive	Inductive	Training Time (Speed-up)
KL	Jodie	<b>93.01±0.6</b>	<b>93.31±0.4</b>	1658.6	74.53±0.7	<b>68.97±0.7</b>	115.5	89.73±4.0	78.22±5.1	9602.8
	DyRep	92.10±2.5	92.49±1.8	1723.9	78.56±0.1	<b>64.98±0.2</b>	118.8	<b>77.97±6.7</b>	<b>68.30±5.8</b>	10347.0
	TGN	91.58±1.7	92.06±1.4	1721.0	<b>82.74±0.4</b>	<b>66.68±1.2</b>	122.0	83.15±3.9	71.29±1.2	10466.8
	TIGE	88.69±3.1	88.67±3.3	1732.8	<b>82.74±0.1</b>	<b>66.23±0.2</b>	123.1	<b>89.62±3.8</b>	<b>80.21±2.6</b>	10492.7
$top_k = 0$	Jodie	92.22±2.1	91.75±2.2	<b>226.5 (7.3x)</b>	<b>74.75±1.7</b>	68.76±1.6	<b>102.7 (1.1x)</b>	<b>91.43±1.1</b>	<b>79.70±2.1</b>	<b>917.3 (10.5x)</b>
	DyRep	<b>93.99±1.4</b>	<b>93.68±1.3</b>	<b>237.0 (7.3x)</b>	<b>78.70±0.2</b>	64.28±0.1	<b>105.2 (1.1x)</b>	68.71±12.2	59.07±6.8	<b>961.8 (10.7x)</b>
	TGN	<b>93.39±1.0</b>	<b>93.29±0.9</b>	<b>238.4 (7.2x)</b>	82.16±0.3	66.11±0.7	<b>107.5 (1.1x)</b>	<b>85.75±1.3</b>	<b>74.08±0.2</b>	<b>976.6 (10.7x)</b>
	TIGE	<b>93.09±2.4</b>	<b>92.59±2.6</b>	<b>251.4 (6.9x)</b>	82.11±0.1	65.98±0.4	<b>108.1 (1.1x)</b>	88.03±7.4	79.43±4.9	<b>984.4 (10.7x)</b>

TABLE VIII

PARTITIONING TIME (IN SECONDS) COMPARISON BETWEEN SEP AND KL FOR FOUR DATASETS. EXPERIMENTS ARE PERFORMED ON CPU.

	Wikipedia	DGraphFin	ML25m	Taobao
	Time (Speed-up)	Time (Speed-up)	Time (Speed-up)	Time (Speed-up)
KL	11.01	238.46	866.85	14862.33
SEP	<b>0.27 (41x)</b>	<b>3.12 (76.5x)</b>	<b>10.67 (81.24x)</b>	<b>157.18 (94.57x)</b>

as shown in Tab.III and Tab.IV, it is not feasible to train it on a multi-GPU setup.

### C. Comparative Experiments

We also conduct two sets of comparative experiments to demonstrate the effect of our components or different experimental settings on down-stream tasks.

1) *Shuffle Partitions*: This set of experiments is designed to investigate the effects of the shuffle partitions method on various datasets and models. Specifically, all graphs are initially partitioned into eight small partitions. During training, these eight partitions are shuffled and combined into 4 partitions. Alternatively, without shuffling, they are directly combined into 4 partitions. These combined 4 partitions are then trained on 4 GPUs in parallel. We present the results on the 4 datasets by setting  $top_k = 5$ , as shown in Fig.7. similar trend holds for the other datasets which cannot be shown due to space limitation. The results of the experiments illustrate that the shuffle partitions is effective in the majority of cases.

2) *Change of Number of GPUs ( $\mathcal{N}$ )*: To delve deeper into the impact of changing the number of partitions or GPUs (i.e.,  $\mathcal{N}$ ) on downstream tasks, we conduct another set of experiments. (The results are shown in Fig.8.) The original graph is partitioned into either 2 or 4 parts, with the corresponding number of GPUs (2 or 4) being used directly for training. The results illustrate the impact of changing number of partitions and number of GPU cards. Meanwhile, comparing the experimental results in shuffle partitions, the impact of partitioning graph into more parts and then randomly splicing them or directly partition the graph into number of parts equals to  $\mathcal{N}$ .

An increase in the  $\mathcal{N}$  leads to an increase in edge deletions, as the  $\mathcal{N}$  should correspond to the number of sub-graphs. Thus, an increase in the number of sub-graphs results in an increased number of deleted edges, implying more information loss. As previously mentioned, the effectiveness of model with deleted edges might be impacted. This is because the deleted edges could potentially contain noise which the deletion of them can contribute to the performance.

### D. Compare with Static Graph Partitioning Algorithm

In the realm of static graph partitioning algorithms, both KL [8] and HDRF [14] are considered representative approaches. However, as HDRF is a special case within our approach, we mainly use KL as the representative method for comparison. Static graph partitioning algorithms generally achieve low edge cuts because they can access global graph information. However, a competent partitioning algorithm should not only deliver quality results but also be time-efficient and achieve load balancing. KL [8] is a representative algorithm for static graph partitioning. Tab.VI demonstrates that the KL performs well on edge cuts without shared nodes, but performs worst on load balancing, i.e., standard deviation of edges. Furthermore, we train models using our proposed PAC to compare the speed-up of training time and performance in downstream tasks. We also evaluate partitioning time across different datasets. The detailed analysis is as follows.

1) *Performance and Training Time Speed-up on Downstream Tasks*: As presented in Tab.VII, training times using the KL algorithm are comparatively longer than those of our method (for which we present the results for  $top_k = 0$ , since KL also does not use shared nodes). This discrepancy is more pronounced in datasets with a larger number of edges. The reason for this is that KL ignores edge balancing, resulting in an uneven distribution of edges across different GPUs. Our training approach will loop over epochs for GPUs with fewer edge data, resulting in these data portions being trained more times compared to other GPUs with more edge data. While, GPUs with more edges may trained only one cycle and took longer time. This leads to two problems, the first one is the training time increases, and the second is the unbalanced data traversing. The second problem will also case the uncertainty

of downstream task performances. Our methods accelerate training up to 7.2x on ML25m, 1.1x on DGraphFin, and 10.7x on Taobao compared to KL. Our methods with  $top_k = 0$  outperform other approaches on most datasets and models. Note that increasing  $top_k$  from 0 to a higher value might enhance performance due to information loss reduced.

2) *Efficiency on Graph partitioning*: As is evident from the Tab.VIII, the efficiency advantage of our SEP over the KL becomes more pronounced as the size of the dataset increases. SEP can enhance the graph partitioning speed by up to a factor of 94.57x compared to the KL algorithm. In scenarios where real-time performance is required, especially when the graph is dynamically changing, the additional overhead associated with the re-partitioning of the KL algorithm is not feasible.

#### IV. RELATED WORK

**TIG Embedding.** TIG models capture the dynamic nature of graphs, thereby enabling superior modelling of TIGs. Jodie [1] employs two Recurrent Neural Networks (RNNs) to dynamically update node representations. DyRep [2] proposes a deep temporal point process model that utilizes a two-time scale approach to capture both association and communication information. TGN [4] introduces a memory-based approach for TIG embedding. TIGE [5] puts forward a model that incorporates a dual-memory module for more effective aggregation of neighbour information. Given that most existing models are constrained to single-GPU training, there exists a compelling motivation for the proposal of a distributed training approach for TIG models.

**GNN Training Acceleration.** For static Graph Neural Networks (GNNs), numerous studies [22]–[25] have attempted to implement large graph sampling for training. However, as the graph size and the number of model layers expand, they invariably encounter the “neighborhood explosion” problem.

Efforts have been made to achieve distributed full batch training [11], [26]–[28], but these often compromise model convergence and accuracy. Distributed GNN mini-batch training represents an alternative platforms like AliGraph [29] and AGL [30], though industrial-scale, do not utilize GPU acceleration. DistDGL [31] employs synchronized stochastic gradient descent for distributed training and maintains a Key-Value store to efficiently acquire graph information. BGL [32] introduces a dynamic caching mechanism to minimize CPU-GPU communication overhead. ByteGNN [33] enables the mini-batch sampling phase to be parallelizable by viewing it as a series of Directed Acyclic Graphs with small tasks.

A handful of studies have concentrated on accelerating temporal GNNs training. EDGE [6] improves computational parallelism by selecting and duplicating specific  $d$ -nodes,

thereby eliminating certain computational dependencies. However, its applicability is confined to Jodie [1], limiting generalizability. TGL [34] introduces a Temporal-CSR data structure, coupled with a parallel sampler, to sample neighboring nodes efficiently for mini-batch training. However, it is not tailored for distributed training and thus orthogonal to our work.

**Graph Partitioning in GNNs.** METIS [7] is a multi-stage static partitioning method designed to minimize edge cuts. It is used by [25] to construct a batch during training and by [27], [29], [31] to partition large graphs for distributed training. NeuGraph [26] utilizes KL [8] to maximize the assignment of edges connected to the same node into the same partition. However, such static graph partitioning methods have high time complexity and require re-partitioning when the graph changes. Euler [9] and Roc [11] apply methods such as random partitioning and linear regression-based techniques. They ignore the graph structural information, resulting in lower quality of partitioning as well as unbalanced computational load.

Streaming graph partitioning methods aim to perceive the graph as an edge-stream or node-stream input. AliGraph [29] incorporates Linear Deterministic Greedy (LGD) [10], an edge-cut based method suited for partitioning dynamically evolving graphs. DistGNN [28] uses a node-cut based method Libra [12]. However, it relies on a hash function to randomly assign edges, thereby ignoring the structural information of the graph and resulting in a high edge-cut ratio. Greedy [13] and HDRF [14] have been shown to have better partitioning quality [35]. However, they either only suitable for static graphs or regard edges at different timestamps equivalently, failing to utilize the characteristics of temporal interaction graphs. Also, they face an excessive number of replica nodes when partitioning real-world graph data. This insight drives us to propose a novel partitioning method tailored for TIGs.

#### V. CONCLUSION

In this paper, we propose a novel Temporal Interaction Graph embedding approach consisting of a streaming edge partitioning method, accompanied by a corresponding distributed parallel training component. By applying our approach, we can efficiently train very large-scale temporal interaction graphs on GPUs. Moreover, our approach can be accelerated using distributed parallel training with multiple GPUs. Our experiments demonstrate that our methods can handle TIGs with millions of nodes and billions of edges. In contrast, previous methods are unable to directly train such large graphs due to computing resource limitations. In future work, we intend to further investigate the impact of edge deletion and strive to provide more interpretability to the information loss issue, concentrating on eliminating noisy or unimportant edges while retaining valid ones.

## REFERENCES

- [1] Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1269–1278, 2019.
- [2] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. Dyrep: Learning representations over dynamic graphs. In *International conference on learning representations*, 2019.
- [3] Da Xu, Chuanwei Ruan, Evren Korpeoglu, Sushant Kumar, and Kannan Achan. Inductive representation learning on temporal graphs. *arXiv preprint arXiv:2002.07962*, 2020.
- [4] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. Temporal graph networks for deep learning on dynamic graphs. *arXiv preprint arXiv:2006.10637*, 2020.
- [5] Yao Zhang, Yun Xiong, Yongxiang Liao, Yiheng Sun, Yucheng Jin, Xuehao Zheng, and Yangyong Zhu. Tiger: Temporal interaction graph embedding with restarts. In *Proceedings of the ACM Web Conference 2023*, pages 478–488, 2023.
- [6] Xinshi Chen, Yan Zhu, Haowen Xu, Mengyang Liu, Liang Xiong, Muhan Zhang, and Le Song. Efficient dynamic graph representation learning at scale. *arXiv preprint arXiv:2112.07768*, 2021.
- [7] George Karypis and Vipin Kumar. A fast and high quality multilevel scheme for partitioning irregular graphs. *SIAM Journal on scientific Computing*, 20(1):359–392, 1998.
- [8] Brian W Kernighan and Shen Lin. An efficient heuristic procedure for partitioning graphs. *The Bell system technical journal*, 49(2):291–307, 1970.
- [9] Euler github. <https://github.com/alibaba/euler>.
- [10] Isabelle Stanton and Gabriel Kliot. Streaming graph partitioning for large distributed graphs. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1222–1230, 2012.
- [11] Zhihao Jia, Sina Lin, Mingyu Gao, Matei Zaharia, and Alex Aiken. Improving the accuracy, scalability, and performance of graph neural networks with roc. *Proceedings of Machine Learning and Systems*, 2:187–198, 2020.
- [12] Cong Xie, Ling Yan, Wu-Jun Li, and Zhihua Zhang. Distributed power-law graph computing: Theoretical and empirical analysis. *Advances in neural information processing systems*, 27, 2014.
- [13] Joseph E Gonzalez, Yucheng Low, Haijie Gu, Danny Bickson, and Carlos Guestrin. Powergraph: Distributed graph-parallel computation on natural graphs. In *Presented as part of the 10th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 12)*, pages 17–30, 2012.
- [14] Fabio Petroni, Leonardo Querzoni, Khuzaima Daudjee, Shahin Kamali, and Giorgio Iacoboni. Hdrf: Stream-based partitioning for power-law graphs. In *Proceedings of the 24th ACM international on conference on information and knowledge management*, pages 243–252, 2015.
- [15] Amauri Souza, Diego Mesquita, Samuel Kaski, and Vikas Garg. Provably expressive temporal graph networks. *Advances in Neural Information Processing Systems*, 35:32257–32269, 2022.
- [16] Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. Inductive representation learning in temporal networks via causal anonymous walks. *arXiv preprint arXiv:2101.05974*, 2021.
- [17] Yiming Li, Yanyan Shen, Lei Chen, and Mingxuan Yuan. Zebra: When temporal graph neural networks meet temporal personalized pagerank. *Proceedings of the VLDB Endowment*, 16(6):1332–1345, 2023.
- [18] Reuven Cohen, Keren Erez, Daniel Ben-Avraham, and Shlomo Havlin. Breakdown of the internet under intentional attack. *Physical review letters*, 86(16):3682, 2001.
- [19] F Maxwell Harper and Joseph A Konstan. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- [20] Xuanwen Huang, Yang Yang, Yang Wang, Chunping Wang, Zhisheng Zhang, Jiarong Xu, and Lei Chen. DGraph: A large-scale financial dataset for graph anomaly detection. In *Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2022.
- [21] Jingwei Zhuo, Ziru Xu, Wei Dai, Han Zhu, Han Li, Jian Xu, and Kun Gai. Learning optimal tree models under beam search. In *International Conference on Machine Learning*, pages 11650–11659. PMLR, 2020.
- [22] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. *Advances in neural information processing systems*, 30, 2017.
- [23] Jianfei Chen, Jun Zhu, and Le Song. Stochastic training of graph convolutional networks with variance reduction. *arXiv preprint arXiv:1710.10568*, 2017.
- [24] Jie Chen, Tengfei Ma, and Cao Xiao. Fastgcn: fast learning with graph convolutional networks via importance sampling. *arXiv preprint arXiv:1801.10247*, 2018.
- [25] Wei-Lin Chiang, Xuanqing Liu, Si Si, Yang Li, Samy Bengio, and Cho-Jui Hsieh. 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*, pages 257–266, 2019.
- [26] Lingxiao Ma, Zhi Yang, Youshan Miao, Jilong Xue, Ming Wu, Lidong Zhou, and Yafei Dai. Neugraph: Parallel deep neural network computation on large graphs. In *USENIX Annual Technical Conference*, pages 443–458, 2019.
- [27] Zhenkun Cai, Xiao Yan, Yidi Wu, Kaihao Ma, James Cheng, and Fan Yu. Dgcl: an efficient communication library for distributed gnn training. In *Proceedings of the Sixteenth European Conference on Computer Systems*, pages 130–144, 2021.
- [28] Vasimuddin Md, Sanchit Misra, Guixiang Ma, Ramanarayan Mohanty, Evangelos Georganas, Alexander Heinecke, Dhiraj Kalamkar, Nesreen K Ahmed, and Sasikanth Avancha. Distgcn: Scalable distributed training for large-scale graph neural networks. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–14, 2021.
- [29] Rong Zhu, Kun Zhao, Hongxia Yang, Wei Lin, Chang Zhou, Baole Ai, Yong Li, and Jingren Zhou. Aligraph: A comprehensive graph neural network platform. *arXiv preprint arXiv:1902.08730*, 2019.
- [30] Dalong Zhang, Xin Huang, Ziqi Liu, Zhiyang Hu, Xianzheng Song, Zhibang Ge, Zhiqiang Zhang, Lin Wang, Jun Zhou, Yang Shuang, et al. Agl: a scalable system for industrial-purpose graph machine learning. *arXiv preprint arXiv:2003.02454*, 2020.
- [31] Da Zheng, Chao Ma, Minjie Wang, Jinjing Zhou, Qidong Su, Xiang Song, Quan Gan, Zheng Zhang, and George Karypis. Distdgl: distributed graph neural network training for billion-scale graphs. In *2020 IEEE/ACM 10th Workshop on Irregular Applications: Architectures and Algorithms (IA3)*, pages 36–44. IEEE, 2020.
- [32] Tianfeng Liu, Yangrui Chen, Dan Li, Chuan Wu, Yibo Zhu, Jun He, Yanghua Peng, Hongzheng Chen, Hongzhi Chen, and Chuanxiong Guo. {BGL}::{GPU-Efficient}{GNN} training by optimizing graph data {I/O} and preprocessing. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*, pages 103–118, 2023.
- [33] Chenguang Zheng, Hongzhi Chen, Yuxuan Cheng, Zhezheng Song, Yifan Wu, Changji Li, James Cheng, Hao Yang, and Shuai Zhang. Bytegcn: efficient graph neural network training at large scale. *Proceedings of the VLDB Endowment*, 15(6):1228–1242, 2022.
- [34] Hongkuan Zhou, Da Zheng, Israt Nisa, Vasileios Ioannidis, Xiang Song, and George Karypis. Tgl: A general framework for temporal gnn training on billion-scale graphs. *arXiv preprint arXiv:2203.14883*, 2022.
- [35] Zainab Abbas, Vasiliki Kalavri, Paris Carbone, and Vladimir Vlassov. Streaming graph partitioning: an experimental study. *Proceedings of the VLDB Endowment*, 11(11):1590–1603, 2018.