# EGAD: Evolving Graph Representation Learning with Self-Attention and Knowledge Distillation for Live Video Streaming Events

Stefanos Antaris
KTH Royal Institute of Technology
HiveStreaming AB
Sweden
antaris@kth.se

Dimitrios Rafailidis
Maastricht University
Netherlands
dimitrios.rafailidis@maastrichtuniversity.nl

Sarunas Girdzijauskas
KTH Royal Institute of Technology
Sweden
sarunasg@kth.se

*Abstract*—In this study, we present a dynamic graph representation learning model on weighted graphs to accurately predict the network capacity of connections between viewers in a live video streaming event. We propose EGAD, a neural network architecture to capture the graph evolution by introducing a self-attention mechanism on the weights between consecutive graph convolutional networks. In addition, we account for the fact that neural architectures require a huge amount of parameters to train, thus increasing the online inference latency and negatively influencing the user experience in a live video streaming event. To address the problem of the high online inference of a vast number of parameters, we propose a knowledge distillation strategy. In particular, we design a distillation loss function, aiming to first pretrain a teacher model on offline data, and then transfer the knowledge from the teacher to a smaller student model with less parameters. We evaluate our proposed model on the link prediction task on three real-world datasets, generated by live video streaming events. The events lasted 80 minutes and each viewer exploited the distribution solution provided by the company Hive Streaming AB. The experiments demonstrate the effectiveness of the proposed model in terms of link prediction accuracy and number of required parameters, when evaluated against state-of-the-art approaches. In addition, we study the distillation performance of the proposed model in terms of compression ratio for different distillation strategies, where we show that the proposed model can achieve a compression ratio up to 15:100, preserving high link prediction accuracy. For reproduction purposes, our evaluation datasets and implementation are publicly available at https://stefanosantaris.github.io/EGAD.

*Index Terms*—Graph representation learning, live video streaming, evolving graphs, knowledge distillation

## I. INTRODUCTION

Nowadays, live video streaming has emerged as a prominent communication solution for several companies worldwide. For example, live video streaming is employed for corporate internal communications, marketing announcements, and so on [1], [2]. Delivering a high quality video to enterprise offices is a challenging task, which stems from the bandwidth requirement, increasing along with the number of viewers in each office. To overcome this challenge, distributed live video
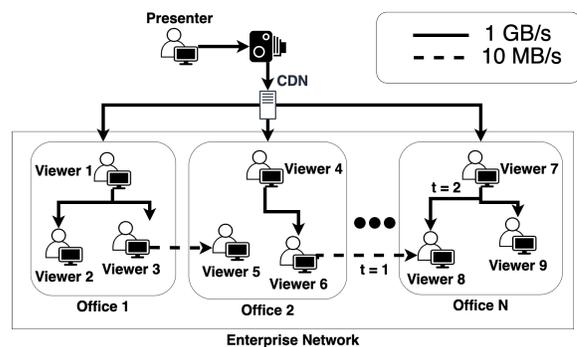
Fig. 1. A distributed live video streaming process in enterprise networks.

streaming solutions were proposed (e.g. by Hive Streaming AB) to deliver high quality video content to several enterprise offices [3], [4]. As shown in Figure 1, Viewers 1, 4 and 7 download the video content of the presenter directly from the Content Delivery Network (CDN) server. Thereafter, Viewers 1, 4 and 7 have to distribute the video content to the rest of the viewers, that is Viewers 2, 3, 5, 6, 8, and 9. To efficiently distribute the video content, each viewer should establish connections with other viewers of the same office and exploit the internal high-bandwidth network of the office (1 GB/s). However, to efficiently establish connections between viewers, Viewer 1 requires the information that Viewers 2 and 3 share the same office. Without this information, Viewer 3 might erroneously establish a connection to Viewer 5 of a different office through a low bandwidth network (10 MB/s). This will negatively impact the video distribution process of Viewer 5, as the only established connection of Viewer 5 will not satisfy the bandwidth requirements of a high quality live video streaming event [4]. Nonetheless, this requires the information of the customers' network topology during the live video streaming event, for instance, Viewers 1, 2 and 3 are in Office 1. However, it is not always feasible to acquire this information, for example, large enterprises provide limited information about their network topologies for security reasons,

or enterprises constantly adapt their networks to assure the desired business outcomes and improve the user experience [5], [6]. In addition, complying with the recent data protection regulations (GDPR) [7], live video streaming providers, such as Hive Streaming AB, are prohibited to retrieve certain network characteristics, such as private and public internet protocol (IP) addresses. Therefore, it is important to predict the network capacity of each connection - bandwidth during a live video streaming event, based on the limited information provided by the already established connections. In doing so, we can infer if the viewers are located in the same office so as to establish connection through the internal high bandwidth network.

During a live video streaming event, each viewer has a limited number of connections. In addition, the viewers adapt their connections in real-time so as to improve the distribution of the video content [4]. For example, in Figure 1, Viewer 6 and 8 are connected with a low bandwidth network at time step $t = 1$. As a consequence, Viewer 8 drops the connection with Viewer 6 at time step $t = 2$ to establish a connection with Viewer 7 via a high bandwidth network connection. The effectiveness of a distributed live video streaming solution depends on the accuracy of each viewers' predictions, that is to predict the connections between viewers in the same office. Moreover, the predictions of the viewers' connections have to be performed in a nearly real-time computational time, otherwise it will negatively impact the user experience during the live video streaming event. In this study, we model an enterprise live video streaming event as a dynamic undirected and weighted graph, where the edges weight correspond to the throughput of the connection between two nodes/viewers. The graph nodes/viewers emerge and leave at unexpected rate and each node/viewer adapts their edges/connections, so as to identify the nodes/viewers that are located in the same office and efficiently distribute the video content. Provided that an enterprise live video streaming event has thousands of viewers, such graphs are highly dimensional and sparse.

**Graph representation learning.** Recently, graph representation learning approaches emerged that compute compact latent node representations to solve the graph dimensionality problem [8]–[10]. Calculating the latent node/viewer representations has proven a successful means to address the link prediction problem on graphs [8], [11]–[15]. Baseline graph representation learning approaches exploit random walks to learn the latent node/viewer representations [8], [10]. More recently, several studies design different neural network architectures to calculate complex patterns in graph structures [13], [15]–[19]. However, these neural network architectures work on static graphs. To capture the graph evolution, recent approaches employ Recurrent Neural Networks (RNN) [12], [20] and self-attention mechanisms [11] between consecutive graph snapshots. Although dynamic graph representation learning approaches achieve high accuracy in link prediction, the underlying neural networks require to train a large amount of parameters. Therefore, these approaches incur high latency during the online inference of the node/viewer representations

due to the large model sizes of the underlying neural network architectures [21]–[26]. As a consequence, state-of-the-art approaches are not applicable to real-world live video streaming solutions, as the high online latency inference increases the computational time of link prediction during a live video streaming event, resulting in high complexity when adapting the viewers' connections.

**Knowledge distillation.** Alternatively, to reduce the high online latency inference, graph representation learning approaches could employ neural networks of smaller sizes with less parameters. However, such models might fail to accurately capture the structure of an evolving graph, resulting in low link prediction accuracy. *Knowledge distillation* has been recently introduced as a model-independent strategy to generate a small model that exhibits low online latency inference, while preserving high accuracy [21], [23], [27]. The main idea of knowledge distillation is to train a large model, namely *teacher*, as an offline training process. The teacher model is a neural network architecture that requires to train a large number of parameters, so as to learn the structure of offline data. Having pretrained the teacher model, the knowledge distillation strategies compute a smaller *student* model with less parameters, that is more suitable for deployment in production. In particular, the student model is trained on online data, and distills the knowledge of the pretrained teacher model. This means that the student model mimics the teacher model and preserves the high prediction accuracy, while at the same time reduces the online inference of the model parameters due to its small size [27], [28]. A few attempts have been made on graph representation learning with knowledge distillation strategies to reduce the model sizes of the underlying neural network architectures [29]–[31]. As we will show in Section IV-D, such approaches fail to achieve a high compression ratio on the student model, that is the size of the student model remains high when compared with the size of the teacher model. This occurs because these approaches learn low dimensional representations on static graphs, which do not correspond to the dynamic case of live video streaming events.

**Contribution.** To overcome the limitations of existing models, in this work we present a knowledge distillation strategy for dynamic graph representation learning, namely EGAD, for the link prediction task during live video streaming events. Our main contributions are summarized as follows:

- EGAD employs a self-attention mechanism on the weights of consecutive Graph Convolutional Networks (GCNs), to capture the graph evolution and learn accurate latent node/viewer representations, during a live video streaming event.
- To the best of our knowledge we are the first to study knowledge distillation for dynamic graph representation learning. We train the EGAD teacher model in an offline process and formulate a distillation loss function to transfer the pretrained knowledge to a smaller student model on online data. In doing so, we significantly reduce the number of parameters when training the student model on online data, and achieve high link prediction accuracy.
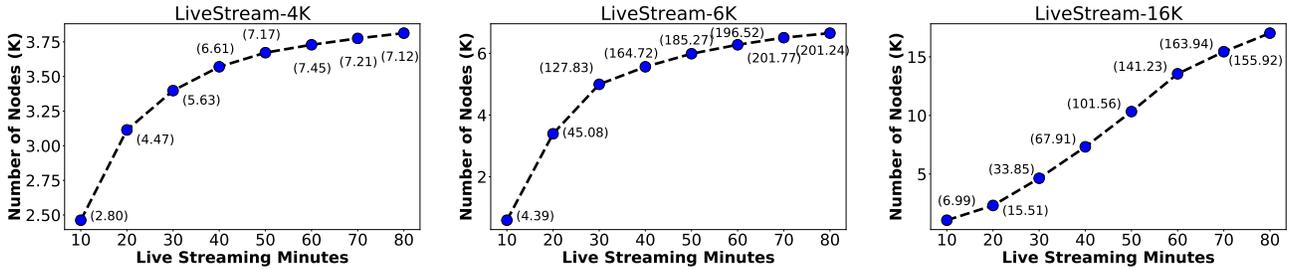
Fig. 2. Number of nodes/viewers during the three live video streaming events. In each parenthesis, we denote the respective number of connections/edges (K) among the viewers at a certain snapshot.

Our experiments on real-world datasets of live video streaming events demonstrate the superiority of the proposed model to accurately capture the evolution of the graph and reduce the online latency inference of the model parameters, when compared with other state-of-the-art methods.

The remainder of the paper is organized as follows: in Section II we present the collected live video streaming data in Hive Streaming AB, and in Section III we detail the proposed model. Our experimental evaluation is presented in Section IV, and we conclude the study in Section V.

## II. LIVE VIDEO STREAMING DATA

During a live video streaming event in Hive Streaming AB, various data are collected such as connections per viewer, throughput per connection, and so on, to provide valuable insights to customers. Each viewer periodically reports the data to centralized servers. To evaluate the performance of the proposed model, we collected real-world datasets based on the reports of three live video streaming events, that is LiveStream-4K, LiveStream-6K and LiveStream-16K. All datasets are anonymized and publicly available. The duration of each live video streaming event is 80 minutes. Each generated dataset consists of 8 weighted undirected graph/viewing snapshots, corresponding to the viewers' connections every 10 minutes. A weight of a graph/viewing edge corresponds to the throughput of the connection among two viewers at each snapshot. The LiveStream-4K dataset has $3,813$ viewers, distributed to $15$ different offices, and $11,066$ connections. In the LiveStream-6K dataset, $6,655$ viewers attended the live video streaming event from $29$ different offices. The viewers established $787,291$ connections. The LiveStream-16K dataset consists of $17,026$ viewers and $482,185$ connections in total. The viewers participated in the live video streaming event from $46$ different offices.

Figure 2 illustrates the different patterns of how viewers emerge during the three live video streaming events. LiveStream-4K has more viewers than LiveStream-6K and LiveStream-16K, during the first $10$ minutes of the live video streaming event. This indicates that in LiveStream-4K the majority of the viewers started to attend the live video streaming event from the beginning. In LiveStream-6K, the first 2 graph/viewing snapshots significantly change in terms of number of viewers, for $0 - 20$ minutes 2.8K new

viewers emerged, whereas in LiveStream-4K and LiveStream-16K 0.5K and 1K viewers emerged, respectively. LiveStream-4K is less informative as the viewers establish the lowest number of connections. Finally, we can observe that viewers in LiveStream-16K emerge at the lowest pace during the live video streaming event. As we will demonstrate in Section IV-C, the effectiveness of the proposed knowledge distillation strategy and baseline approaches not only depends on the graph sizes but also on different patterns that viewers emerge during the live video streaming events.

## III. PROPOSED METHOD

A live video streaming event is represented as a sequence of $K$ graph/viewing snapshots $\mathcal{G} = \{\mathcal{G}_1, \ldots, \mathcal{G}_K\}$. $\forall \ k = 1, \ldots, K$ snapshot we consider the graph $\mathcal{G}_k = (\mathcal{V}_k, \mathcal{E}_k, \mathbf{X}_k)$, where $\mathcal{V}_k$ corresponds to the set of $n_k = |\mathcal{V}_k|$ viewers, $\mathcal{E}_k$ is the set of connections, and $\mathbf{X}_k \in \mathbb{R}^{n_k \times m}$ is the matrix of the $m$ features of each viewer. For each graph $\mathcal{G}_k$, we consider a weighted adjacency matrix $\mathbf{A}_k \in \mathbb{R}^{n_k \times n_k}$, where $A(u, v) > 0$ for the viewers $u \in \mathcal{V}_k$ and $v \in \mathcal{V}_k$, if $e_k(u, v) \in \mathcal{E}_k$. The weight $A(u, v)$ corresponds to the bandwidth measured between viewers $u \in \mathcal{V}_k$ and $v \in \mathcal{V}_k$ at the $k$-th snapshot. Given a sequence of $l$ graph/viewing snapshots[1] $\{\mathcal{G}_{k-l}, \ldots, \mathcal{G}_k\}$, the goal of the proposed model is to compute $d$-dimensional latent representations $\mathbf{Z}_k \in \mathbb{R}^{n_k \times d}$, with $d \ll m$ [11], [12], [14]. The constructed latent representations should capture both the structure of the graph at the graph/viewing snapshot $k$ and the evolutionary behavior of the viewers up to the $k$-th minute.

Dynamic graph representation learning models employ deep neural network architectures, requiring to train a large amount of parameters [12], [32], [33]. Such models are computationally expensive to deploy to a large number of viewers in live video streaming events as they incur significant online latency to calculate the viewers' representations [22], [23], [29], [30]. The problem of knowledge distillation is to generate a smaller online *student* model $\mathcal{S}$ than a pretrained offline large teacher model $\mathcal{T}$. The goal is to reduce the number of trainable parameters of the student model $\mathcal{S}$ to minimize the online

---

[1]The reason for not accounting for all the previous snapshots from the beginning of the live video streaming event, and consider only a certain time window $l$ is because we observed in our experiments that large values of $l$ do not necessarily increase the prediction accuracy, while at the same time significantly increase the number of the model parameters. The influence of $l$ on the performance of the proposed model and the baseline approaches is studied in Table IV.

latency inference [27], [34]. In practice, the teacher model $\mathcal{T}$ is pretrained using a computationally expensive deep neural network architecture to calculate the latent representations $\mathbf{Z}_k^{\mathcal{T}}$ of the offline data. Having trained the teacher model offline, the student model $\mathcal{S}$ learns the latent representations $\mathbf{Z}_k^{\mathcal{S}}$ by minimizing a distillation loss function $L^{\mathcal{D}}$. The distillation loss function $L^{\mathcal{D}}$ calculates the prediction error of the student model $\mathcal{S}$ and the deviation from the latent representations $\mathbf{Z}_k^{\mathcal{T}}$ generated by the teacher model $\mathcal{T}$. This means that the student model $\mathcal{S}$ is able to mimic the already pretrained teacher model $\mathcal{T}$ with fewer parameters [27], [28]. In Section III-A we present the offline teacher model EGAD-$\mathcal{T}$, and then in Section III-B we describe the distillation process of the online student model EGAD-$\mathcal{S}$.

## A. EGAD-$\mathcal{T}$ Teacher Model

The teacher model EGAD-$\mathcal{T}$ learns the viewer representations $\mathbf{Z}_k^{\mathcal{T}}$ at the $k$-th graph/viewing snapshot using $l$ consecutive Graph Convolutional Network (GCN) models [12], [35], [36], with EGAD-$\mathcal{T} = \{GCN_{k-l}, \ldots GCN_k\}$, and $l$ being the number of previous graph/viewing snapshots. The input of each $GCN_k$ model is the normalized adjacency matrix $\hat{\mathbf{A}}_k \in \mathbb{R}^{n_k \times n_k}$ and the viewers' features $\mathbf{X}_k$. Provided that the graphs during the live video streaming events have nodes with no features, the node feature matrix $\mathbf{X}_k$ is replaced by the identity matrix $\mathbf{I} \in \mathbb{R}^{n \times n}$, with $m = n$. Each $GCN_k$ model calculates the viewers representations $\mathbf{Z}_k^{\mathcal{T}}$ by applying two convolution layers to $\hat{\mathbf{A}}_k$ and $\mathbf{X}_k$, as follows:

$$\mathbf{Z}_k^{\mathcal{T}} = \hat{\mathbf{A}}_k ReLU(\hat{\mathbf{A}}_k \mathbf{X}_k \mathbf{W}_k^1) \mathbf{W}_k^2 \tag{1}$$

where $\mathbf{W}_k^i \in \mathbb{R}^{d_{i-1} \times d_i}$ is the weight parameter matrix of the $i$-th convolutional layer, with $d_i < d_{i-1} < m$. Following [13], [35] we employ two convolutional layers ($i = 1, 2$), to learn the weight parameter matrices $\mathbf{W}_k^1 \in \mathbb{R}^{m \times d_1}$ and $\mathbf{W}_k^2 \in \mathbb{R}^{d_1 \times d_2}$, with $d_2 = d$, so as to compute the $d$-dimensional representations $\mathbf{Z}_k^{\mathcal{T}}$. The symmetrically normalized adjacency matrix $\hat{\mathbf{A}}_k$ is calculated as follows:

$$\begin{aligned} \hat{\mathbf{A}}_k &= \mathbf{D}_k^{-\frac{1}{2}} \tilde{\mathbf{A}}_k \mathbf{D}_k^{-\frac{1}{2}} \\ \tilde{\mathbf{A}}_k &= \mathbf{A}_k + \mathbf{I} \\ \mathbf{D}_k &= diag(\sum_j A_k(u, v)) \end{aligned} \tag{2}$$

The $l$ consecutive GCN models are connected in a sequential manner through the weights $\mathbf{W}_{k-1}^1$ and $\mathbf{W}_k^1$ of the first convolutional layers [37]. For each node $u \in \mathcal{V}_k$ we calculate $h$ independent *self-attention heads*, that is vectors $\mathbf{z}_k^j(u) \in \mathbb{R}^{d_1}$, with $j = 1, \ldots, h$, based on the $d_1$-dimensional weights $\mathbf{W}_{k-1}^1(u) \in \mathbb{R}^{d_1}$. To compute the weights $\mathbf{W}_k^1$ of each $GCN_k$ model, we average the $h$ independent self-attention vectors $\mathbf{z}_k^j(u)$ [16], as follows:

$$\begin{aligned} \mathbf{W}_k^1(u) &= ELU\big(\tfrac{1}{h} \sum_{j=1}^h \mathbf{z}_k^j(u)\big) \\ \mathbf{z}_k^j(u) &= \sum_{v \in \mathcal{N}_u} \alpha_{u,v} \mathbf{H}_k \mathbf{W}_{k-1}^1(v) \end{aligned} \tag{3}$$

where ELU is the Exponential Linear Unit activation function [38]. Variable $\mathbf{H}_k \in \mathbb{R}^{d_1 \times d_1}$ is the shared weight transformation matrix applied to the previous weights $\mathbf{W}_{k-1}^1(u)$ of each node $u \in \mathcal{V}_k$, $\mathcal{N}_u$ is the neighborhood set of the node $u$. Variable $\alpha_{u,v}$ is the normalized attention coefficient between $u \in \mathcal{V}_k$ and $v \in \mathcal{N}_u$, which is calculated based on the softmax function [11], as follows:

$$\alpha_{u,v} = \frac{exp\big(\sigma(A_k(u,v) \cdot \mathbf{a}_k^T [\mathbf{H}_k \mathbf{W}_{k-1}^1(u) || \mathbf{H}_k \mathbf{W}_{k-1}^1(v)])\big)}{\sum\limits_{w \in \mathcal{N}_u} exp\big(\sigma(A_k(u,w) \cdot \mathbf{a}_k^T [\mathbf{H}_k \mathbf{W}_{k-1}^1(u) || \mathbf{H}_k \mathbf{W}_{k-1}^1(w)])\big)} \tag{4}$$

where $\sigma$ is the sigmoid function, $A_k(u, v)$ is the edge weight between $u$ and $v$, $\mathbf{a}_k^T \in \mathbb{R}^{2d_1}$ is a $2d_1$-dimensional weight vector which is applied to the attention process between nodes $u$ and $v$ [11], [16], and $||$ is the concatenation operation. The attention coefficient $\alpha_{u,v}$ measures the importance of the connection between nodes $u \in \mathcal{V}_k$ and $v \in \mathcal{N}_u$. A high attention coefficient value $\alpha_{u,v}$ corresponds to a connection $e_k(u, v) \in \mathcal{E}_k$ which is maintained over several consecutive graph/viewing snapshots and has high edge weight in the adjacency matrix $A_k(u, v)$. This means that the learned weights $\mathbf{W}_k^1$ reflect on the importance of the existing connection between node $u$ and $v$, processing the convolution accordingly.

To train the teacher model EGAD-$\mathcal{T}$, we initialize $l$ GCN models and connect the consecutive GCN models using the self-attention mechanism in Equation 3. As aforementioned, each of the $k$-th $GCN$ models takes as an input the normalized adjacency matrix $\tilde{\mathbf{A}}_k$ and the feature vectors $\mathbf{X}_k$. When training EGAD-$\mathcal{T}$, each $GCN_k$ model computes the weights $\mathbf{W}_k^1$ in Equation 3, and then calculates the latent representations $\mathbf{Z}_k$ based on Equation 1. Note that the weights $\mathbf{W}_0^1$ for the first $GCN$ model are randomly initialized. To train our teacher model EGAD-$\mathcal{T}$, we adopt the Root Mean Square Error loss function with respect to the latent representations $\mathbf{Z}_k$ generated by the last GCN model [37], as follows:

$$\min_{\mathbf{Z}_k} L^{\mathcal{T}} = \sqrt{\frac{1}{n_k}\big(\sigma(\mathbf{Z}_k^{\mathcal{T}^\top} \cdot \mathbf{Z}_k^{\mathcal{T}}) - \mathbf{A}_k\big)^2} \tag{5}$$

where $\cdot$ represents the inner product operation between all the possible pairs of latent representations, and the term $\sigma(\mathbf{Z}_k^{\mathcal{T}^\top} \cdot \mathbf{Z}_k^{\mathcal{T}}) - \mathbf{A}_k$ calculates the error of the latent representations $\mathbf{Z}_k^{\mathcal{T}}$ to capture the structure of the graph snapshot $\mathcal{G}_k$. In our implementation, we optimize the parameters $\mathbf{H}_k$ and $\mathbf{a}_k$ between consecutive GCN models, based on the loss function in Equation 5 and the backpropagation algorithm.

## B. EGAD-$\mathcal{S}$ Student Model

We train the student model EGAD-$\mathcal{S}$ to compute the online latent representations $\mathbf{Z}_k^{\mathcal{S}}$, by exploiting the knowledge of the pretrained teacher model EGAD-$\mathcal{T}$. As we train the student model only on online data, the student model EGAD-$\mathcal{S}$ requires significantly less number of trainable parameter weights, compared with the teacher model EGAD-$\mathcal{T}$. The student model EGAD-$\mathcal{S}$ consists of $l$ consecutive GCN models, with EGAD-$\mathcal{S} = \{GCN_{k-l}, \ldots, GCN_k\}$. We calculate the

weights $\mathbf{W}_k^1$ (Equation 4), and compute the latent representations based on Equation 1.

The knowledge acquired by the teacher model EGAD-$\mathcal{T}$ is transferred to the student model EGAD-$\mathcal{S}$ via the distillation loss function $L^{\mathcal{D}}$, adopted by the student model during the online training process. We formulate the distillation loss function as a minimization problem for the student model EGAD-$\mathcal{S}$ as follows:

$$\min_{\mathbf{z}_k^{\mathcal{S}}} L^{\mathcal{D}} = (1 - \gamma)L^{\mathcal{T}} + \gamma L^{\mathcal{S}} \qquad (6)$$

where $L^{\mathcal{T}}$ is the inference error of the teacher model in Equation 5, and $L^{\mathcal{S}}$ is the root mean squared error with the latent representations $\mathbf{Z}_k^{\mathcal{S}}$ generated by the student model. Hyper-parameter $\gamma \in [0, 1]$ balances the training of the student model EGAD-$\mathcal{S}$ when inferring the knowledge of the teacher model EGAD-$\mathcal{T}$. A higher value of $\gamma$ emphasizes more on the student model EGAD-$\mathcal{S}$ and distillates less knowledge from the teacher model EGAD-$\mathcal{T}$. The distillation loss function $L^{\mathcal{D}}$ in Equation 6 allows the student model EGAD-$\mathcal{S}$ to overcome any bias introduced by the teacher model EGAD-$\mathcal{T}$ [21], [23], [27], [28]. This means that EGAD-$\mathcal{S}$ can achieve similar or better accuracy than the teacher model EGAD-$\mathcal{T}$. As we will show later in Section IV-D, the student model EGAD-$\mathcal{S}$ consistently outperforms the teacher model EGAD-$\mathcal{T}$ in terms of accuracy, by significantly downsizing the number of parameters.

## IV. EXPERIMENTAL EVALUATION

### A. Evaluation Setup

In our experiments we evaluate the performance of the proposed model on the link prediction task. To examine the two different components of our model, we train the teacher and student models EGAD-$\mathcal{T}$ and EGAD-$\mathcal{S}$, using $l$ consecutive graph/viewing snapshots up to the $k$-th graph $\mathcal{G}_k$. The task of link prediction is to forecast the unobserved connections, denoted by $\mathcal{O}_{k+1} = \mathcal{E}_{k+1} \backslash \{\mathcal{E}_{k-l}, \ldots, \mathcal{E}_k\}$, that will occur in the next graph/viewing snapshot $\mathcal{G}_{k+1}$. Following the evaluation protocol of [11], [12], [14], we concatenate the latent representations $\mathbf{Z}_k(u)$ and $\mathbf{Z}_k(v)$ based on the Hadamard operator, for the unobserved connection $o(u, v) \in \mathcal{O}_{k+1}$ of the viewers $u$ and $v \in \mathcal{V}_k$. The concatenated latent representations are then applied to a Multi-Layer Perceptron (MLP), to calculate the weight of the connection. To measure the online inference efficiency, we report the number of parameters that each model requires to train. Moreover, regarding the prediction accuracy we evaluate the examined models based on the metrics Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE):

$$MAE = \frac{\sum\limits_{o(u,v) \in \mathcal{O}_{k+1}} \left| A_{k+1}(u,v) - \mathbf{Z}_k^T(u)\mathbf{Z}_k(v) \right|}{|\mathcal{O}_{k+1}|}$$

$$RMSE = \sqrt{(\frac{1}{|\mathcal{O}_{k+1}|}) \sum\limits_{o(u,v) \in \mathcal{O}_{k+1}} (\mathbf{Z}_k^T(u)\mathbf{Z}_k(v) - \mathbf{A}_{k+1}(u,v))^2}$$

$$\qquad (7)$$

Following [11], [12], [37], for each snapshot $k$ we train each examined model on $l$ previous graph/viewing snapshots $\mathcal{G}_{k-l}, \ldots, \mathcal{G}_k$, which are considered the offline data for each time step. We randomly select $20\%$ of the unobserved links $\mathcal{O}_{k+1}$ for validation set to tune the model hyper-parameters. The remaining $80\%$ of the unobserved links are considered as the test set, which are the online data for each time step. We repeated our experiments five times, and we report the average RMSE and MAE over the five trials.

### B. Examined Models

We compare the performance of the proposed EGAD-$\mathcal{T}$ and EGAD-$\mathcal{S}$ models with the following baseline strategies:

- **DynVGAE** [14] is a dynamic joint learning model that shares the trainable parameters between consecutive variational graph auto-encoders [35]. We implemented DynVGAE from scratch and publish our code[2], as there is no publicly available implementation.
- **EvolveGCN**[3] [12] is a dynamic graph representation learning model with Gated Recurrent Units (GRUs) between the convolutional weights of consecutive GCNs.
- **DySAT**[4] [11] is a dynamic self-attention model that captures the evolution of the graph using multi-head self-attention between consecutive graph snapshots.
- **DMTKG-$\mathcal{T}$** [29] is the teacher model of the DMTKG knowledge distillation strategy. DMTKG-$\mathcal{T}$ employs Heat Kernel Signature (HKS) on static graph/viewing snapshots and uses Convolutional Neural Network layers to calculate the latent representations based on Deep-Graph [39]. To ensure fair comparison, we train DMTKG-$\mathcal{T}$ per snapshot, with each snapshot containing aggregated graph history up to $k$-th snapshot. As the source code of the DMKTG distillation strategy is not available, we made our implementation publicly available[5].
- **DMTKG-$\mathcal{S}$** [29] is the student model of the DMTKG strategy, where the goal is to minimize a distillation loss function based on the weighted cross entropy.

**Settings.** In Tables III-V we report the performance of each examined model in terms of RMSE when calibrating the hyper-parameters of the examined models, following a cross-validation strategy. For each model, we tuned the hyper-parameters based on a grid selection strategy and select the best configuration. In particular, in DynVGAE we set the
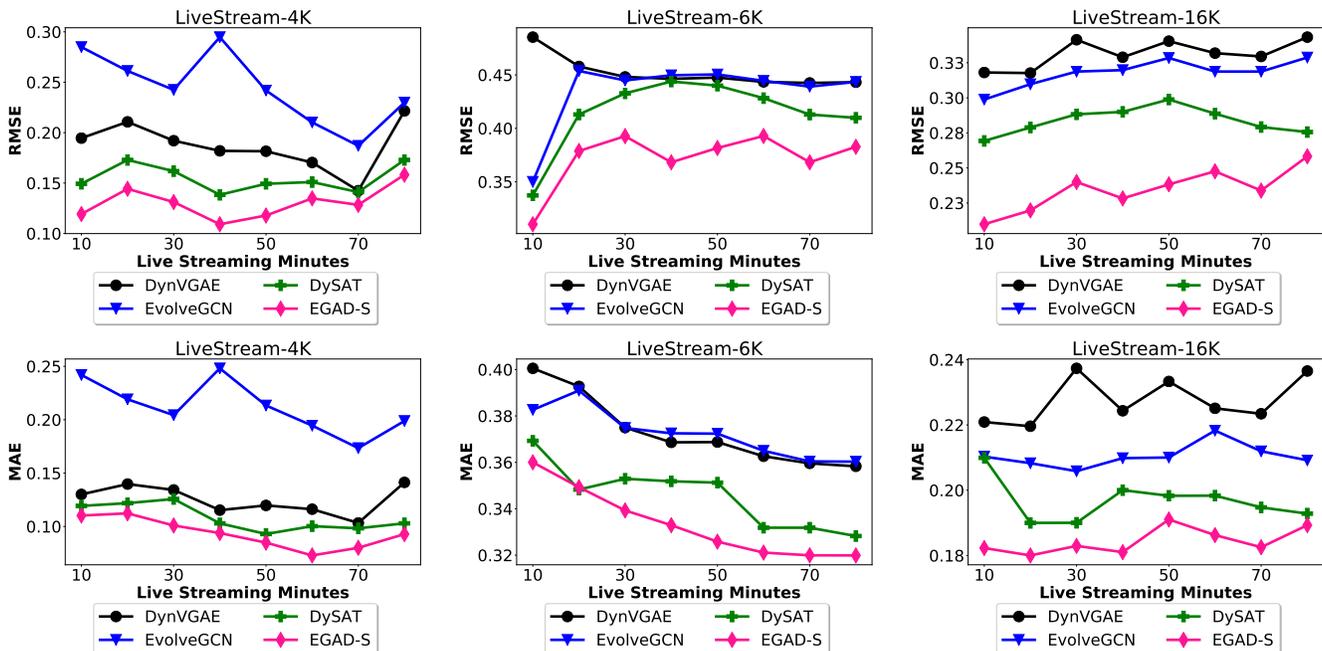
---

[2]https://github.com/stefanosantaris/DynVGAE
[3]https://github.com/IBM/EvolveGCN
[4]https://github.com/aravindsankar28/DySAT
[5]https://github.com/stefanosantaris/DMTKG

Fig. 3. Performance evaluation of EGAD-$\mathcal{S}$ against the non-distillation strategies in terms of RMSE and MAE in LiveStream-4K, LiveStream-6K and LiveStream-16K.

TABLE I
MODEL PARAMETERS IN MILLIONS FOR THE LIVE STREAMING MINUTES, WHEN COMPARING THE PROPOSED EGAD-$\mathcal{S}$ MODEL WITH THE NON-DISTILLATION STRATEGIES.

| | LiveStream-4K | | | | LiveStream-6K | | | | LiveStream-16K | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Stream. Min. | DynVGAE | EvolveGCN | DySAT | EGAD-$\mathcal{S}$ | DynVGAE | EvolveGCN | DySAT | EGAD-$\mathcal{S}$ | DynVGAE | EvolveGCN | DySAT | EGAD-$\mathcal{S}$ |
| 10 | 0.332 | 37.034 | 0.845 | **0.081** | 0.092 | 2.288 | 0.367 | **0.022** | 0.152 | 6.994 | 0.485 | **0.036** |
| 20 | 0.830 | 117.829 | 1.227 | **0.105** | 0.902 | 70.021 | 1.299 | **0.114** | 0.626 | 65.472 | 1.023 | **0.079** |
| 30 | 0.903 | 140.150 | 1.515 | **0.117** | 1.312 | 151.137 | 1.924 | **0.168** | 1.220 | 260.283 | 1.832 | **0.156** |
| 40 | 0.974 | 154.609 | 1.774 | **0.125** | 1.458 | 187.345 | 2.285 | **0.189** | 1.907 | 646.523 | 2.734 | **0.245** |
| 50 | 0.973 | 163.519 | 1.800 | **0.128** | 1.567 | 216.853 | 2.394 | **0.202** | 2.678 | 1285.467 | 3.510 | **0.341** |
| 60 | 0.988 | 168.607 | 1.815 | **0.130** | 1.641 | 238.277 | 2.468 | **0.212** | 3.503 | 2211.020 | 4.330 | **0.445** |
| 70 | 0.100 | 172.770 | 1.827 | **0.132** | 1.670 | 255.987 | 2.527 | **0.219** | 3.986 | 2867.735 | 3.813 | **0.505** |
| 80 | 1.001 | 176.247 | 1.836 | **0.133** | 1.737 | 267.638 | 2.564 | **0.224** | 4.392 | 3486.341 | 5.219 | **0.556** |

representation size to $d = 64$ and the window size $l = 2$ for all datasets. In EvolveGCN, the representation size is set to $d = 32$, with $l = 2$ previous graph/viewing snapshots. DySAT uses $l = 2$ consecutive graph/viewing snapshots and employs $h = 3$ attention heads for the LiveStream-4K and LiveStream-6K datasets. For the LiveStream-16K dataset, we use $h = 4$ attention heads. The representation size is fixed to $d = 64$ in all datasets. In DMTKG-$\mathcal{T}$, the representation size is fixed to $d = 64$ for LiveStream-4K and LiveStream-6K, while in LiveStream-16K we use $d = 128$. In DMTKG-$\mathcal{S}$, we reduce the model size by setting the size of latent representations to $d = 32$ for LiveStream-4K and LiveStream-6K, and $d = 64$ for LiveStream-16K. Regarding the proposed model, we train EGAD-$\mathcal{T}$ on $l = 3$ consecutive graph/viewing snapshots and set the number of head attentions to $h = 3$ in Equation 3, with $d = 64$-dimensional latent representations. In the student model EGAD-$\mathcal{S}$, we reduce the number of heads $h = 1$ and fix the latent representation size to $d = 16$ for all datasets. The influence of $\gamma$ on the distillation loss function (Equation 6) is further studied in Section IV-D. We initialize the learning rate to $1e-03$ based on the Adam Optimizer with 200 epochs. All experiments were performed on an Intel(R) Xeon(R) Bronze 3106 CPU 1.70GHz machine and GPU accelerated with the GEFORCE RTX 2080 Ti graph card.

### C. Performance Evaluation

In Figure 3, we evaluate the performance of the student model EGAD-$\mathcal{S}$ against the non-distillation strategies, that is DynVGAE, EvolveGCN and DySAT, in terms of RMSE and MAE. We observe that all models have a higher prediction error in terms of RMSE and MAE in LiveStream-6K than the other datasets. This occurs because the viewers in the LiveStream-6K dataset attended the live video streaming event in a completely different pattern (Section II) than LiveStream-4K and LiveStream-16K. More precisely, in LiveStream-6K the number of viewers that emerge in 0-20 minutes is significantly higher than the other events, which negatively impacts the prediction accuracy of the examined models.

The student model EGAD-$\mathcal{S}$ significantly outperforms the baseline approaches in all datasets. This suggests that the proposed student model EGAD-$\mathcal{S}$ can efficiently capture the evolution of the graph in the learned latent representations

$\mathbf{Z}_k^{\mathcal{S}}$. The second best approach is DySAT, demonstrating the ability of self-attention mechanisms to generate accurate latent representations. DySAT calculates the latent representations $\mathbf{Z}_k$ by applying self-attentional aggregations to the local node neighborhoods. Instead, the proposed EGAD-$\mathcal{S}$ model performs self-attention to the convolutional weights between consecutive GCNs. Thus, our model is able to efficiently capture the different graph evolution patterns of the live video streaming events. Compared to the second best method DySAT, the proposed EGAD-$\mathcal{S}$ model achieves relative drops $9.8$ and $13.5\%$ in terms of RMSE and MAE in the LiveStream-4K dataset. Similarly, EGAD-$\mathcal{S}$ achieves relative drops $10.2$ and $3.5\%$ in LiveStream-6K, and $17.3$ and $6.2\%$ relative drops in LiveStream-16K.

In Table I, we present the numbers of parameters in millions that are required to train the examined models. As aforementioned in Section II, the majority of the viewers in the LiveStream-4K dataset started to attend the live video streaming event from the first 10 minutes. Therefore, all models have fewer trainable parameters on the first graph snapshots $k = 0 - 30$ minutes in LiveStream-6K and LiveStream-16K than in the LiveStream-4K dataset. We observe that EGAD-$\mathcal{S}$ clearly outperforms the baseline approaches in terms of the required parameters. Evaluated against DynVGAE, EvolveGCN and DySAT, the average compress ratios of the student model EGAD-$\mathcal{S}$ are 12:100, 7:1000, and 7:100, respectively. Provided that EGAD-$\mathcal{S}$ constantly outperforms all the baseline approaches in terms of RMSE and MAE, the high compression ratios demonstrate the ability of the proposed knowledge distillation strategy to significantly reduce the model size in terms of required parameters. Moreover, it is clear that EvolveGCN model requires a significant amount of trainable parameters to generate the latent representations. This means that EvolveGCN does not scale well when increasing the number of viewers in live video streaming events. As DySAT employs multi-head attention on consecutive graph/viewing snapshots, and not on consecutive GCNs as the proposed EGAD- model does, DySAT requires a much a higher number of parameters by following a non-distillation strategy.

### D. Distillation Evaluation

In Figure 4, we study the impact of the proposed knowledge distillation strategy on the student model EGAD-$\mathcal{S}$ in terms of RMSE, when compared with the teacher model EGAD-$\mathcal{T}$. In addition, in this set of experiments we evaluate our model against DMTKG [29], a baseline graph representation approach with knowledge distillation, comparing with both the teacher model DMTKG-$\mathcal{T}$ and student model DMTKG-$\mathcal{S}$.

On inspection of Figure 4, we observe that the EGAD-$\mathcal{T}$ and EGAD-$\mathcal{S}$ models outperform DMTKG-$\mathcal{T}$ and DMTKG-$\mathcal{S}$ in all datasets. This occurs because DMTKG applies knowledge distillation on top of DeepGraph [39], which is a static graph representation learning approach. Therefore, DMTKG ignores the graphs' evolution when learning the latent representations. An interesting observation is that the student models EGAD-$\mathcal{S}$ and DMTKG-$\mathcal{S}$ achieve higher performance

than the respective teacher models EGAD-$\mathcal{T}$ and DMTKG-$\mathcal{T}$. This indicates the effectiveness of the examined distillation strategies to correctly transfer the knowledge of the teacher models to the respective student models. This occurs because the student models remove the bias of the teacher models to the offline data, and achieve high prediction accuracy, complying with similar observations that have been made in relevant studies [21], [40]. Compared to the EGAD-$\mathcal{T}$ model, EGAD-$\mathcal{S}$ achieves $6.5$, $3.6$ and $5.7\%$ relative drops in terms of RMSE for LiveStream-4K, LiveStream-6K and LiveStream-16K, respectively.

In Table II, we present the maximum number of parameters in millions that are required to train the examined models during the live video streaming events. EGAD-$\mathcal{S}$ significantly reduces the number of required parameters, achieving compression ratios 15:100, 17:100 and 21:100, on average, in LiveStream-4K, LiveStream-6K and LiveStream-16K, respectively. This occurs because the student model EGAD-$\mathcal{S}$ uses a lower number of attention heads $h$ and representation size $d$ than the teacher model EGAD-$\mathcal{T}$ (Section IV-C). Therefore, EGAD-$\mathcal{S}$ has lower online inference latency, compared with the teacher model EGAD-$\mathcal{T}$. Instead, the DMTKG distillation strategy achieves an average 1:2 compression ratio for the student model DMTKG-$\mathcal{S}$. The DMTKG distillation strategy is not able to further reduce the student model size, because DMTKG is designed for static graphs. This indicates that the DMTKG-$\mathcal{S}$ model requires more trainable parameters to learn accurate latent representations than the proposed EGAD-$\mathcal{S}$ model.

TABLE II
THE MAXIMUM NUMBERS OF REQUIRED PARAMETERS IN MILLIONS OF
THE EXAMINED KNOWLEDGE DISTILLATION STRATEGIES DURING THE
LIVE VIDEO STREAMING EVENTS. IN THE PARENTHESES, WE DENOTE THE
AVERAGE COMPRESSION RATIOS OF THE STUDENT MODELS, WHEN
COMPARED WITH THE RESPECTIVE TEACHER MODELS.

| Model | LiveStream-4K | LiveStream-6K | LiveStream-16K |
|---|---|---|---|
| **DMTKG**-$\mathcal{T}$ | 3.673 | 5.129 | 10.437 |
| **DMTKG**-$\mathcal{S}$ | 1.836 (1:2) | 2.564 (1:2) | 5.219 (1:2) |
| **EGAD**-$\mathcal{T}$ | 0.918 | 1.282 | 2.609 |
| **EGAD**-$\mathcal{S}$ | 0.133 (15:100) | 0.224 (17:100) | 0.556 (21:100) |

In Figure 5, we evaluate the influence on the hyper-parameter $\gamma$ of Equation 6 on the student model EGAD-$\mathcal{S}$. We vary the hyper-parameter $\gamma$ from $0.1$ to $0.9$ by a step of $0.1$, to balance the impact of the student $L^{\mathcal{S}}$ and teacher $L^{\mathcal{T}}$ losses on the distillation loss function $L^{\mathcal{D}}$. For each parameter $\gamma$, we report the averaged RMSE over all the graph snapshots of the live video streaming events. In all datasets, the student model EGAD-$\mathcal{S}$ achieves the highest performance when we equally balance the influence of the student and teacher models ($\gamma = 0.5$). For larger values of parameter $\gamma$, the student model EGAD-$\mathcal{S}$ emphasizes more on the loss $L^{\mathcal{S}}$ than the loss $L^{\mathcal{T}}$. As a consequence, the student model EGAD-$\mathcal{S}$ distills less knowledge from the teacher model EGAD-$\mathcal{T}$, which negatively impacts the performance of the EGAD-$\mathcal{S}$ model in terms of RMSE. Instead, decreasing the hyper-parameter $\gamma$ prevents the student model EGAD-$\mathcal{S}$ from training on the
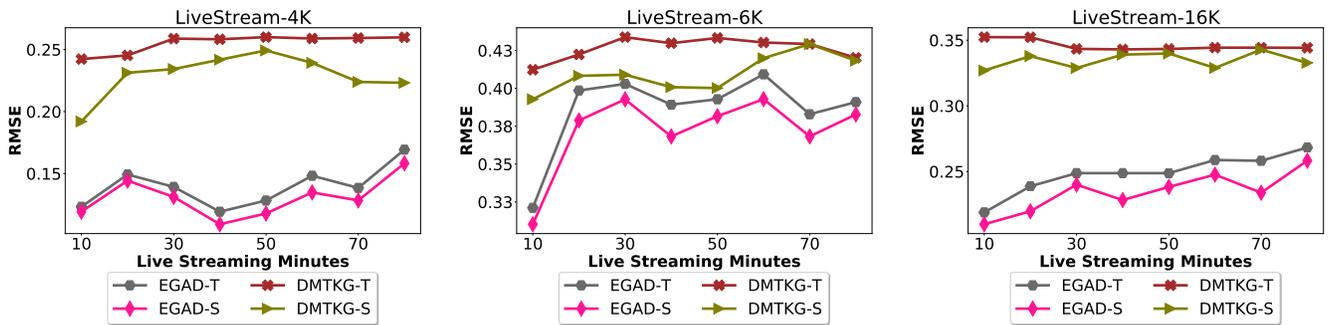
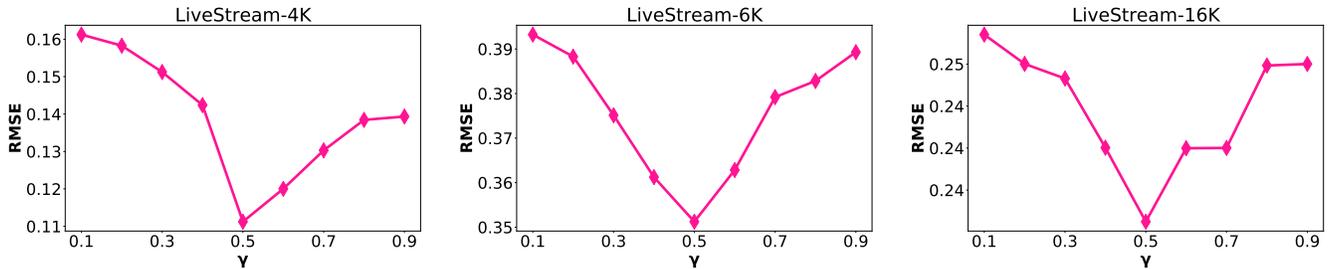Fig. 4. Comparison of student and teacher models for the examined knowledge distillation strategies.



Fig. 5. Impact of $\gamma$ on the prediction accuracy of student model EGAD-$\mathcal{S}$.

online graph data and at the same time introduces the bias to the offline data of EGAD-$\mathcal{T}$. This means that for small values of $\gamma$ EGAD-$\mathcal{S}$ mainly distills the knowledge of the teacher model EGAD-$\mathcal{T}$, resulting in limited prediction accuracy.

## V. CONCLUSION

In this paper, we presented a knowledge distillation strategy, to overcome the problem of online latency inference of dynamic graph representation learning approaches in live video streaming events. Evaluated against several baseline approaches on three real-world live video streaming events, the proposed model achieves 7:100 compression ratio on average. Moreover, the proposed student model preserves high prediction accuracy, achieving average relative drops $12.4$ and $7.7\%$ in terms of RMSE and MAE in all events, when compared with the second best approach. Distributed live video streaming providers, such as Hive Streaming AB, can significantly benefit from our model by significantly reducing the required parameters/computational time in the link prediction task. In doing so, viewers can exploit the offices' internal high bandwidth network from the beginning of the live video streaming event, by avoiding to establish low bandwidth connections. Provided that several offices have limited network capacity, our model can significantly reduce the generated network traffic. Therefore, enterprises can distribute high quality video content to their offices without any network limitations, improving user experience. Moreover, the proposed model allows enterprises to distribute video content of high resolution, such as 4K.

There are several interesting future directions to graph representation learning for live video streaming events.

- For instance, as future work we plan to evaluate the performance of the proposed model on evolving graphs of social networks. In particular, provided the limited duration of live video stream events the main challenge resides on identifying the differences of how viewers emerge during live video streaming events and at what pace users establish connections in social networks over time.
- Another interesting future direction is to study the performance of our model on graph snapshots over time steps with different duration. For example, in a live video streaming event the duration of time steps between two consecutive snapshots might vary, depending on the network demand. This means that different time steps might require an adaptive learning strategy of the time window $w$ when training our model.
- In our model, training is performed on the graph data of a single live video streaming event. In practice though there are several live video streaming events that take place on a daily basis. The question that we have to answer is how to exploit the knowledge acquired from different live video streaming events, when training our model on a new event. More precisely, we plan to study various transfer learning strategies to exploit the knowledge from different events, when training our model. This is a challenging task, because not only the internal network topologies of several companies vary, but also viewers emerge at various paces during different live video streaming events.

## REFERENCES

[1] Q. Fan, H. Yin, G. Min, P. Yang, Y. Luo, Y. Lyu, H. Huang, and L. Jiao, "Video delivery networks: Challenges, solutions and future directions," *Comput. Electr. Eng.*, vol. 66, pp. 332–341, 2018.

[2] M. F. Majeed, S. H. Ahmed, S. Muhammad, H. Song, and D. B. Rawat, "Multimedia streaming in information-centric networking: A survey and future perspectives," *Comput. Networks*, vol. 125, pp. 103–121, 2017.

TABLE III

EFFECT ON RMSE WHEN VARYING THE REPRESENTATION SIZE $d$ OF THE EXAMINED MODELS. WE REPORT AVERAGE RMSE OVER THE GRAPH SNAPSHOTS DURING THE LIVE VIDEO STREAMING EVENT. BOLD VALUES DENOTE THE BEST CONFIGURATION FOR EACH MODEL.

| Representation size $d$ | DynVGAE | EvolveGCN | DySAT | DMTKG-$\mathcal{T}$ | DMTKG-$\mathcal{S}$ | EGAD-$\mathcal{T}$ | EGAD-$\mathcal{S}$ |
|---|---|---|---|---|---|---|---|
| **LiveStream-4K** | | | | | | | |
| 16 | $0.23 \pm 0.14$ | $0.27 \pm 0.09$ | $0.18 \pm 0.07$ | $0.28 \pm 0.14$ | $0.27 \pm 0.15$ | $0.17 \pm 0.09$ | $\mathbf{0.13 \pm 0.09}$ |
| 32 | $0.21 \pm 0.12$ | $\mathbf{0.25 \pm 0.12}$ | $0.16 \pm 0.10$ | $0.26 \pm 0.13$ | $\mathbf{0.23 \pm 0.17}$ | $0.16 \pm 0.05$ | $0.15 \pm 0.10$ |
| 64 | $\mathbf{0.19 \pm 0.13}$ | $0.26 \pm 0.18$ | $\mathbf{0.15 \pm 0.09}$ | $\mathbf{0.25 \pm 0.16}$ | $0.24 \pm 0.12$ | $\mathbf{0.14 \pm 0.08}$ | $0.16 \pm 0.06$ |
| 128 | $0.20 \pm 0.17$ | $0.26 \pm 0.12$ | $0.17 \pm 0.08$ | $0.26 \pm 0.17$ | $0.25 \pm 0.16$ | $0.15 \pm 0.06$ | $0.16 \pm 0.08$ |
| 256 | $0.21 \pm 0.15$ | $0.27 \pm 0.15$ | $0.18 \pm 0.05$ | $0.27 \pm 0.19$ | $0.26 \pm 0.14$ | $0.16 \pm 0.07$ | $0.17 \pm 0.06$ |
| **LiveStream-6K** | | | | | | | |
| 16 | $0.48 \pm 0.16$ | $0.47 \pm 0.12$ | $0.41 \pm 0.15$ | $0.48 \pm 0.18$ | $0.48 \pm 0.14$ | $0.41 \pm 0.12$ | $\mathbf{0.36 \pm 0.06}$ |
| 32 | $0.46 \pm 0.16$ | $\mathbf{0.44 \pm 0.11}$ | $0.40 \pm 0.17$ | $0.47 \pm 0.18$ | $\mathbf{0.41 \pm 0.12}$ | $0.39 \pm 0.09$ | $0.37 \pm 0.09$ |
| 64 | $\mathbf{0.45 \pm 0.18}$ | $0.45 \pm 0.13$ | $\mathbf{0.39 \pm 0.18}$ | $0.43 \pm 0.16$ | $0.43 \pm 0.15$ | $\mathbf{0.37 \pm 0.10}$ | $0.39 \pm 0.10$ |
| 128 | $0.46 \pm 0.17$ | $0.45 \pm 0.14$ | $0.41 \pm 0.16$ | $0.45 \pm 0.14$ | $0.44 \pm 0.16$ | $0.38 \pm 0.11$ | $0.40 \pm 0.08$ |
| 256 | $0.46 \pm 0.18$ | $0.46 \pm 0.18$ | $0.42 \pm 0.14$ | $0.46 \pm 0.12$ | $0.46 \pm 0.18$ | $0.38 \pm 0.11$ | $0.42 \pm 0.07$ |
| **LiveStream-16K** | | | | | | | |
| 16 | $0.36 \pm 0.13$ | $0.35 \pm 0.18$ | $0.29 \pm 0.16$ | $0.39 \pm 0.15$ | $0.39 \pm 0.12$ | $0.27 \pm 0.10$ | $\mathbf{0.23 \pm 0.07}$ |
| 32 | $0.35 \pm 0.14$ | $\mathbf{0.33 \pm 0.19}$ | $0.29 \pm 0.14$ | $0.37 \pm 0.14$ | $0.36 \pm 0.11$ | $0.26 \pm 0.09$ | $0.25 \pm 0.09$ |
| 64 | $\mathbf{0.33 \pm 0.12}$ | $0.34 \pm 0.11$ | $\mathbf{0.27 \pm 0.14}$ | $0.36 \pm 0.15$ | $\mathbf{0.33 \pm 0.11}$ | $\mathbf{0.24 \pm 0.09}$ | $0.26 \pm 0.06$ |
| 128 | $0.34 \pm 0.11$ | $0.35 \pm 0.13$ | $0.28 \pm 0.13$ | $\mathbf{0.35 \pm 0.12}$ | $0.35 \pm 0.15$ | $0.25 \pm 0.10$ | $0.26 \pm 0.08$ |
| 256 | $0.34 \pm 0.17$ | $0.36 \pm 0.11$ | $0.29 \pm 0.12$ | $0.38 \pm 0.13$ | $0.36 \pm 0.18$ | $0.26 \pm 0.10$ | $0.27 \pm 0.10$ |

TABLE IV

IMPACT OF THE WINDOW SIZE $l$ ON THE PERFORMANCE OF EACH EXAMINED MODEL IN TERMS OF RMSE.

| Window size $l$ | DynVGAE | EvolveGCN | DySAT | DMTKG-$\mathcal{T}$ | DMTKG-$\mathcal{S}$ | EGAD-$\mathcal{T}$ | EGAD-$\mathcal{S}$ |
|---|---|---|---|---|---|---|---|
| **LiveStream-4K** | | | | | | | |
| 1 | $0.32 \pm 0.16$ | $0.35 \pm 0.19$ | $0.21 \pm 0.12$ | N/A | N/A | $0.18 \pm 0.06$ | $0.16 \pm 0.05$ |
| 2 | $\mathbf{0.19 \pm 0.13}$ | $\mathbf{0.25 \pm 0.12}$ | $\mathbf{0.15 \pm 0.09}$ | N/A | N/A | $0.16 \pm 0.09$ | $0.14 \pm 0.06$ |
| 3 | $0.24 \pm 0.18$ | $0.28 \pm 0.14$ | $0.19 \pm 0.14$ | N/A | N/A | $\mathbf{0.14 \pm 0.08}$ | $\mathbf{0.13 \pm 0.09}$ |
| 4 | $0.29 \pm 0.12$ | $0.32 \pm 0.16$ | $0.20 \pm 0.12$ | N/A | N/A | $0.17 \pm 0.10$ | $0.16 \pm 0.04$ |
| 5 | $0.36 \pm 0.17$ | $0.42 \pm 0.19$ | $0.24 \pm 0.14$ | N/A | N/A | $0.20 \pm 0.09$ | $0.19 \pm 0.06$ |
| **LiveStream-6K** | | | | | | | |
| 1 | $0.48 \pm 0.16$ | $0.56 \pm 0.11$ | $0.42 \pm 0.14$ | N/A | N/A | $0.41 \pm 0.11$ | $0.39 \pm 0.07$ |
| 2 | $\mathbf{0.45 \pm 0.18}$ | $\mathbf{0.44 \pm 0.11}$ | $\mathbf{0.39 \pm 0.18}$ | N/A | N/A | $0.38 \pm 0.08$ | $0.37 \pm 0.09$ |
| 3 | $0.46 \pm 0.14$ | $0.49 \pm 0.13$ | $0.41 \pm 0.16$ | N/A | N/A | $\mathbf{0.37 \pm 0.10}$ | $\mathbf{0.36 \pm 0.06}$ |
| 4 | $0.52 \pm 0.19$ | $0.51 \pm 0.12$ | $0.43 \pm 0.12$ | N/A | N/A | $0.40 \pm 0.09$ | $0.39 \pm 0.04$ |
| 5 | $0.54 \pm 0.12$ | $0.54 \pm 0.18$ | $0.49 \pm 0.15$ | N/A | N/A | $0.43 \pm 0.12$ | $0.40 \pm 0.05$ |
| **LiveStream-16K** | | | | | | | |
| 1 | $0.42 \pm 0.12$ | $0.38 \pm 0.16$ | $0.33 \pm 0.16$ | N/A | N/A | $0.31 \pm 0.10$ | $0.29 \pm 0.06$ |
| 2 | $\mathbf{0.33 \pm 0.12}$ | $\mathbf{0.33 \pm 0.19}$ | $\mathbf{0.27 \pm 0.14}$ | N/A | N/A | $0.29 \pm 0.11$ | $0.27 \pm 0.08$ |
| 3 | $0.37 \pm 0.11$ | $0.36 \pm 0.17$ | $0.29 \pm 0.12$ | N/A | N/A | $\mathbf{0.24 \pm 0.09}$ | $\mathbf{0.23 \pm 0.07}$ |
| 4 | $0.39 \pm 0.18$ | $0.39 \pm 0.19$ | $0.31 \pm 0.16$ | N/A | N/A | $0.30 \pm 0.07$ | $0.29 \pm 0.09$ |
| 5 | $0.46 \pm 0.14$ | $0.41 \pm 0.15$ | $0.39 \pm 0.14$ | N/A | N/A | $0.38 \pm 0.10$ | $0.24 \pm 0.10$ |

[3] R. Roverso, R. Reale, S. El-Ansary, and S. Haridi, "Smoothcache 2.0: Cdn-quality adaptive http live streaming on peer-to-peer overlays," in *MMSys*, 2015, p. 61–72.

[4] R. Roverso, S. El-Ansary, and M. Högqvist, "On http live streaming in large enterprises," in *SIGCOMM*, 2013, p. 489–490.

[5] T. V. Phan, M. Hajizadeh, N. T. Khai, and T. Bauschert, "Destination-aware adaptive traffic flow rule aggregation in software-defined networks," in *NetSys*, 2019, pp. 1–6.

[6] P. S. Rivera, J. Griffioen, Z. Fei, and J. H. Hayes, "Expressing and managing network policies for emerging HPC systems," in *PEARC*, 2019, pp. 36:1–36:7.

[7] "GDPR Regulation Europe," https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:32016R0679, 2016, [Online; accessed 01-April-2020].

[8] A. Grover and J. Leskovec, "node2vec: Scalable feature learning for networks," in *KDD*, 2016, pp. 855–864.

[9] W. L. Hamilton, R. Ying, and J. Leskovec, "Representation learning on graphs: Methods and applications," *IEEE Data Eng. Bull.*, vol. 40, no. 3, pp. 52–74, 2017.

[10] B. Perozzi, R. Al-Rfou, and S. Skiena, "Deepwalk: Online learning of social representations," in *KDD*, 2014, pp. 701–710.

[11] A. Sankar, Y. Wu, L. Gou, W. Zhang, and H. Yang, "Dysat: Deep neural representation learning on dynamic graphs via self-attention networks," in *WSDM*, 2020, pp. 519–527.

[12] A. Pareja, G. Domeniconi, J. Chen, T. Ma, T. Suzumura, H. Kanezashi, T. Kaler, T. B. Schardl, and C. E. Leiserson, "EvolveGCN: Evolving graph convolutional networks for dynamic graphs," in *AAAI*, 2020.

[13] A. Hasanzadeh, E. Hajiramezanali, K. R. Narayanan, N. Duffield, M. Zhou, and X. Qian, "Semi-implicit graph variational auto-encoders," in *NeurIPS*, 2019, pp. 10 711–10 722.

[14] S. Mahdavi, S. Khoshraftar, and A. An, "Dynamic joint variational graph autoencoders," in *ECML*, 2019, pp. 385–401.

[15] W. L. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *NIPS*, 2017, pp. 1024–1034.

[16] P. Velickovic, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, "Graph attention networks," in *ICLR*, 2018.

[17] S. Cao, W. Lu, and Q. Xu, "Grarep: Learning graph representations with global structural information," in *CIKM*, 2015, p. 891–900.

TABLE V

EFFECT ON RMSE WHEN VARYING THE NUMBER OF HEADS $h$ OF THE SELF-ATTENTION MECHANISMS OF DYSAT, EGAD-$\mathcal{T}$ AND EGAD-$\mathcal{S}$.

| Number of heads $h$ | DynVGAE | EvolveGCN | DySAT | DMTKG-$\mathcal{T}$ | DMTKG-$\mathcal{S}$ | EGAD-$\mathcal{T}$ | EGAD-$\mathcal{S}$ |
|---|---|---|---|---|---|---|---|
| **LiveStream-4K** | | | | | | | |
| 1 | N/A | N/A | $0.19 \pm 0.06$ | N/A | N/A | $0.16 \pm 0.05$ | $\mathbf{0.13 \pm 0.09}$ |
| 2 | N/A | N/A | $0.17 \pm 0.07$ | N/A | N/A | $0.15 \pm 0.04$ | $0.14 \pm 0.10$ |
| 3 | N/A | N/A | $\mathbf{0.15 \pm 0.09}$ | N/A | N/A | $\mathbf{0.14 \pm 0.08}$ | $0.16 \pm 0.08$ |
| 4 | N/A | N/A | $0.16 \pm 0.08$ | N/A | N/A | $0.16 \pm 0.07$ | $0.17 \pm 0.09$ |
| 5 | N/A | N/A | $0.18 \pm 0.08$ | N/A | N/A | $0.17 \pm 0.09$ | $0.20 \pm 0.03$ |
| **LiveStream-6K** | | | | | | | |
| 1 | N/A | N/A | $0.47 \pm 0.12$ | N/A | N/A | $0.42 \pm 0.07$ | $\mathbf{0.36 \pm 0.06}$ |
| 2 | N/A | N/A | $0.45 \pm 0.15$ | N/A | N/A | $0.40 \pm 0.12$ | $0.38 \pm 0.06$ |
| 3 | N/A | N/A | $\mathbf{0.39 \pm 0.18}$ | N/A | N/A | $\mathbf{0.37 \pm 0.10}$ | $0.41 \pm 0.04$ |
| 4 | N/A | N/A | $0.41 \pm 0.14$ | N/A | N/A | $0.38 \pm 0.10$ | $0.43 \pm 0.09$ |
| 5 | N/A | N/A | $0.46 \pm 0.17$ | N/A | N/A | $0.39 \pm 0.11$ | $0.43 \pm 0.08$ |
| **LiveStream-16K** | | | | | | | |
| 1 | N/A | N/A | $0.32 \pm 0.17$ | N/A | N/A | $0.28 \pm 0.08$ | $\mathbf{0.23 \pm 0.07}$ |
| 2 | N/A | N/A | $0.29 \pm 0.18$ | N/A | N/A | $0.26 \pm 0.10$ | $0.25 \pm 0.05$ |
| 3 | N/A | N/A | $0.28 \pm 0.15$ | N/A | N/A | $\mathbf{0.24 \pm 0.09}$ | $0.26 \pm 0.09$ |
| 4 | N/A | N/A | $\mathbf{0.27 \pm 0.14}$ | N/A | N/A | $0.25 \pm 0.06$ | $0.26 \pm 0.10$ |
| 5 | N/A | N/A | $0.29 \pm 0.12$ | N/A | N/A | $0.27 \pm 0.07$ | $0.27 \pm 0.08$ |

[18] D. Wang, P. Cui, and W. Zhu, "Structural deep network embedding," in *KDD*, 2016, pp. 1225–1234.

[19] Z. T. Kefato and S. Girdzijauskas, "Gossip and attend: Context-sensitive graph representation learning," in *ICWSM*, 2020.

[20] P. Goyal, S. R. Chhetri, and A. Canedo, "dyngraph2vec: Capturing network dynamics using dynamic graph representation learning," *Knowl. Based Syst.*, vol. 187, 2020.

[21] J. Tang and K. Wang, "Ranking distillation: Learning compact ranking models with high performance for recommender system," in *KDD*, 2018, p. 2289–2298.

[22] J. Ba and R. Caruana, "Do deep nets really need to be deep?" in *NIPS*, 2014, pp. 2654–2662.

[23] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," in *NIPS*, 2015.

[24] Y. Liu, J. Cao, B. Li, C. Yuan, W. Hu, Y. Li, and Y. Duan, "Knowledge distillation via instance relationship graph," in *CVPR*, 2019, pp. 7096–7104.

[25] J. Wang, L. Gou, W. Zhang, H. Yang, and H. Shen, "Deepvid: Deep visual interpretation and diagnosis for image classifiers via knowledge distillation," *IEEE Trans. Vis. Comput. Graph.*, vol. 25, no. 6, pp. 2168–2180, 2019.

[26] J. Yim, D. Joo, J. Bae, and J. Kim, "A gift from knowledge distillation: Fast optimization, network minimization and transfer learning," in *CVPR*, 2017, pp. 7130–7138.

[27] C. Bucila, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *KDD*, 2006, pp. 535–541.

[28] M. Phuong and C. Lampert, "Towards understanding knowledge distillation," in *ICML*, 2019, pp. 5142–5151.

[29] J. Ma and Q. Mei, "Graph representation learning via multi-task knowledge distillation," in *NeurIPS*, 2019.

[30] C. Lassance, M. Bontonou, G. B. Hacene, V. Gripon, J. Tang, and A. Ortega, "Deep geometric knowledge distillation with graphs," in *ICASSP*, 2020, pp. 8484–8488.

[31] S. Lee and B. C. Song, "Graph-based knowledge distillation by multi-head attention network," in *BMVC*, 2019, p. 141.

[32] L. Zhou, Y. Yang, X. Ren, F. Wu, and Y. Zhuang, "Dynamic network embedding by modeling triadic closure process," in *AAAI*, 2018, pp. 571–578.

[33] P. Goyal, N. Kamra, X. He, and Y. Liu, "Dyngem: Deep embedding method for dynamic graphs," vol. abs/1805.11273, 2018.

[34] R. Anil, G. Pereyra, A. Passos, R. Ormándi, G. E. Dahl, and G. E. Hinton, "Large scale distributed neural network training through online distillation," in *ICLR*, 2018.

[35] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *ICLR*, 2017.

[36] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson, "Structured sequence modeling with graph convolutional recurrent networks," in *ICONIP*, 2018, pp. 362–373.

[37] S. Antaris and D. Rafailidis, "VStreamDRLS: Dynamic graph representationlearning with self-attention for enterprisedistributed video streaming solutions," in *ASONAM*, 2020.

[38] D. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," in *ICLR*, 2016.

[39] C. Li, X. Guo, and Q. Mei, "Deepgraph: Graph structure predicts network growth," 2016.

[40] Y. Kim and A. M. Rush, "Sequence-level knowledge distillation," in *EMNLP*, 2016, pp. 1317–1327.