

# DyVGRNN: DYnamic mixture Variational Graph Recurrent Neural Networks

Ghazaleh Niknam<sup>a,\*</sup>, Soheila Molaei<sup>b,\*</sup>, Hadi Zare<sup>a,\*\*</sup>, Shirui Pan<sup>c</sup>, Mahdi Jalili<sup>d</sup>, Tingting Zhu<sup>b</sup>, David Clifton<sup>b,e</sup>

<sup>a</sup>*Department of Data Science and Technology, University of Tehran*

<sup>b</sup>*Department of Engineering Science, University of Oxford*

<sup>c</sup>*School of Information and Communication Technology, Griffith University*

<sup>d</sup>*School of Engineering, RMIT University*

<sup>e</sup>*Oxford-Suzhou Institute of Advanced Research (OSCAR), Suzhou, China*

---

## Abstract

Although graph representation learning has been studied extensively in static graph settings, dynamic graphs are less investigated in this context. This paper proposes a novel integrated variational framework called DYnamic mixture Variational Graph Recurrent Neural Networks (DyVGRNN), which consists of extra latent random variables in structural and temporal modelling. Our proposed framework comprises an integration of Variational Graph Auto-Encoder (VGAE) and Graph Recurrent Neural Network (GRNN) by exploiting a novel attention mechanism. The Gaussian Mixture Model (GMM) and the VGAE framework are combined in DyVGRNN to model the multimodal nature of data, which enhances performance. To consider the significance of time steps, our proposed method incorporates an attention-based module. The experimental results demonstrate that our method greatly outperforms state-of-the-art dynamic graph representation learning methods in terms of link prediction and clustering.<sup>1</sup>

**Keywords:** Dynamic Graph Representation Learning, Dynamic Node Embedding, Variational Graph Auto-Encoder, Graph Recurrent Neural Network, Attention Mechanism

---

\*Equal Contribution

\*\*Corresponding author

<sup>1</sup>The source code is available at <https://github.com/GhazalehNiknam/DyVGRNN>.

---

## 1. Introduction

Many real and man-made systems can be represented as graph structures where individual entities are connected through links. Graph structures play a key role in many real-world applications. The recommendation in social networks [1], traffic forecasting in transportation networks [2], and pattern recognition in biological networks [3] are some of these applications. Due to their complexity and high dimensions, these structures are difficult to study. To deal with this problem, representation learning approaches are used [4]. These methods aim to map high-dimensional vectors to low ones in latent space so that these latent vectors capture the structural information of the graph as well as each node’s features.

Downstream machine learning tasks can then use these latent vectors as feature inputs [5, 6]. For example, COOL [7], and GHNN [8] employ graph representations in their classification task, Modularity-aware VGAE [9], and GCN-LP [10] in their link prediction tasks, and SOLI [11] in its clustering task. Although many real-world graphs, known as dynamic graphs, evolve over time, the bulk of existing graph representation learning algorithms concentrates on static graphs, in which the set of nodes and edges does not change over time. This work aims to capture the underlying dynamics of the network.

Our proposed method, “DYnamic mixture Variational Graph Recurrent Neural Networks (DyVGRNN)”, integrates a variational framework with a Graph Recurrent Neural Network (GRNN) to simultaneously capture the evolution of the dynamic graph topology and node attributes. The DyVGRNN can model the addition/removal of nodes and edges in dynamic graphs and can be applied to simple or attributed networks. While conventional variational frameworks can capture hidden and hierarchical dependencies, they are tussling with multimodal data.

Multimodality arises when in a dataset with an overall population and various subpopulations, we are unable to dedicate each subpopulation to an indi-

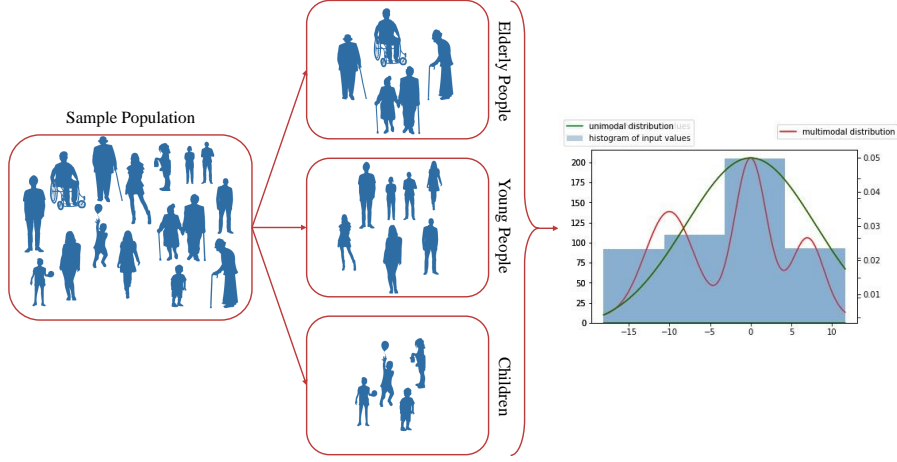


Figure 1: Examining the effect of considering unknown subpopulations on modelling. Here, if the green curve is utilised for modelling and the age of the population under investigation is not considered, some specific information could be lost. On the other hand, a thorough knowledge of the input data is given if the red curve is employed for modelling.

30 visual observation. Mixture models such as Gaussian Mixture Models (GMM)  
 31 are an absolute solution for these kinds of datasets. These models describe the  
 32 probability distribution of observations in the whole population [12, 13, 14].  
 33 Technically, mixture models are a principled modelling approach to handle such  
 34 complex data and are a universal approximator of densities [15, 16].

35 For more clarification, consider a study that examines how an advertisement  
 36 impacts a sample group of people. Some important data, like the effect of age,  
 37 may be lost if the study employs the population while omitting subpopulations  
 38 and models the data using a unimodal distribution. More flexibility and a  
 39 more in-depth understanding of the input data can be obtained by employing a  
 40 mixture model. Figure 1 shows this affection on a synthetic dataset.

41 In this paper, we employ GMM to model the prior and posterior distribution  
 42 in the Graph Variational Auto-Encoder (GVAE). With this combination, it is  
 43 possible to capture the distribution of the input data more effectively and to get  
 44 a deeper knowledge of it. Furthermore, a module based on the attention mecha-

nism on graph snapshots is introduced in our proposed method to demonstrate the significance of time steps. Our experiments show DyVGRNN’s superior performance in dynamic link prediction tasks in several real-world dynamic graphs compared to the state-of-the-art methods. Our contributions to this work are as follows:

- We propose a novel integrated variational framework consisting of extra latent random variables in structural and temporal modelling.
- We combine variational inference based on GMM with the proposed framework to infer the multimodal nature of data and improve the comprehension of the model.
- We introduce a module according to the attention mechanism of graph snapshots to consider the importance of time steps.
- Our experiments show the superior performance of the proposed DyVGRNN in several real-world dynamic graphs compared to the state-of-the-art methods.

## 2. Related Work

To build a solid understanding of dynamic graph representation learning methods, it is important to first delve into the foundational concepts of static methods. Therefore, we will begin by exploring static methods before progressing to dynamic methods.

### 2.1. Static Graph Representation Learning

Shallow embedding methods, which are based on matrix factorisation and random walks, were the first attempts to learn graph representation on static graphs. Matrix factorisation methods such as Graph Factorisation (GF) [17], GraRep [18], and HOPE [19] are inspired by dimensionality reduction techniques. The key distinction among these three methods is the measure used to determine node similarity. On the other hand, in random walk methods [20, 21],

72 nodes have similar representations when they tend to occur together in short  
73 random walks on the graph. In contrast to matrix factorisation approaches that  
74 use deterministic node similarity measures, random walk methods use flexible  
75 stochastic node similarity measures.

76 DeepWalk [20], and node2vec [21] fall into the random walks category, which  
77 optimise embeddings to encode random walk statistics instead of decoding de-  
78 terministic measures of node similarity. While shallow embedding methods have  
79 been quite popular in the last decade, they have significant drawbacks, including  
80 the inability to handle parameter sharing, difficulty with node attributes, and  
81 transductive behaviour [22]. To overcome the limitations of shallow embedding  
82 methods, Graph Neural Networks (GNNs) have been proposed as powerful deep  
83 embedding approaches [22, 6].

84 GNNs are categorised into three types: those based on Graph Recurrent Neu-  
85 ral Networks (GRNN), those based on Graph Convolutional Networks (GCN),  
86 and those based on Graph Auto Encoders (GAE) [23]. The first structure  
87 presented in the context of GNNs is the GRNN. This structure received little  
88 attention prior to the advent of dynamic graphs. The primary assumption in  
89 the GRNN is that messages are exchanged between nodes and their neighbours  
90 until a stable equilibrium is reached.

91 GCNs generalise the convolutions to graph-structured data [24], which plays  
92 a leading role in the construction of many other GNNs. The GCN-based ap-  
93 proaches extract high-level node representations by stacking multiple graph  
94 convolution layers [25]. Following GCNs, GAE-based methods are presented  
95 which include an encoder (mainly based on GCN) to learn representations and  
96 a decoder to reconstruct input data [24, 26]. Variational Graph Auto-Encoder  
97 (VGAE) is a variant of GAE comprising a probabilistic encoder and a proba-  
98 bilistic decoder to model the uncertainty of node representation for more gen-  
99 eralisation of inference [27].

## 100 2.2. Dynamic Graph Representation Learning

101 Dynamic graphs can be represented in two different ways: discretely and  
102 continuously. A discrete dynamic graph is represented as a set of static graphs  
103 taken at predetermined intervals, referred to as snapshots. Continuous graphs  
104 contain no summarisation and provide whole temporal information. Continu-  
105 ous methods cannot be utilised on discrete networks, whereas discrete methods  
106 can be applied on continuous networks. Therefore, discrete techniques are more  
107 flexible than continuous ones [28]. While the discrete representation learning  
108 approach is our focus, we also briefly touch on the continuous representation  
109 learning approaches.

110  
111 **Continuous Methods.** Continuous dynamic graph representation learning  
112 approaches are categorised into two groups: RNN-based and temporal point-  
113 based approaches. RNNs are used in the first category to continually maintain  
114 node embeddings. Every time an event or network change occurs, RNN-based  
115 approaches all update the embeddings of the interacting nodes. DyGNN [29]  
116 falls into this category, which consists of two components: an update component  
117 that updates the states of the nodes involved in an interaction and a propagation  
118 component that propagates the update to those nodes’ neighbours. JODIE [30]  
119 is another RNN-based approach designed for user-item interaction networks in  
120 recommender systems. This method uses one RNN for users and the other for  
121 items. JODIE updates the embeddings when an interaction happens between a  
122 user and an item.

123 The utilisation of the Temporal Point Process (TPP), parametrised by neu-  
124 ral networks, is a recurring feature of temporal point-based techniques. For  
125 example, DyREP [28] uses a two-time scale TPP, which is parametrised by an  
126 RNN. This two-time scale TPP expresses the dynamics of the network (realised  
127 as topological evolution) as well as dynamics on the network (realised as node  
128 communication). Utilising temporal information, the attention coefficient for  
129 a structural edge between nodes is computed. Using these coefficients, the ag-  
130 gregate quantity required for embedding propagation is then determined. In

131 addition, the Latent Dynamic Graph (LDG) [31] extends DyREP using the  
 132 Neural Relational Inference (NRI) [32] model.

133

134 **Discrete Methods.** The most straightforward way for modelling discrete dy-  
 135 namic graphs began with a single GNN in each snapshot [23]. The output of  
 136 each GNN is subsequently sent into the time-series modelling module as input.  
 137 For example, GCRNM1 [33] modelled structural features using the GCN varia-  
 138 tion described in [34] and graph evolution using the peephole LSTM introduced  
 139 in [35]. RgCNN [36] used PATCHY-SAN, a GCN-based approach for modelling  
 140 structural properties, and stacked this with a standard LSTM for modelling  
 141 temporal properties.

142 DyGGNN [37] leveraged a Gated Graph Neural Network (GGNN) and a long  
 143 short-term memory network (LSTM) in its framework to model the topology  
 144 of dynamic graphs and temporal information among them. Waterfall Dynamic-  
 145 GCN and Concatenated Dynamic-GCN [38] are two architectures exploiting a  
 146 GCN and an LSTM in the stacked form by applying them to each node sepa-  
 147 rately. The extra skip connection of the GCN in the Concatenated Dynamic-  
 148 GCN distinguishes these designs. Also, DySAT [39] is another stacked architec-  
 149 ture that uses self-attention blocks to capture structural and temporal proper-  
 150 ties.

151 The techniques mentioned earlier all offer a stacked architecture with a sep-  
 152 arate GNN for processing each snapshot of the dynamic graph and a time series  
 153 module for processing the outputs of these GNNs. By integrating structural  
 154 and temporal modelling into a single layer and capturing both concurrently, dy-  
 155 namic graphs can better capture growing relationships [23]. EvolveGCN [40] is  
 156 an integrated framework consisting of a GCN and an RNN that GCN’s weights  
 157 are updated with the RNN.

158 Another integrated framework is GC-LSTM [41], which combines an LSTM  
 159 with a GCN. The graph snapshots are fed into LSTM in this framework, and  
 160 then a spectral graph convolution is performed on the hidden layer of LSTM.  
 161 LRGCN [42] leverages an R-GCN to jointly address intra-time and inter-time

relationships and an LSTM to capture the time dependency between graph snapshots. Recurrent Event Network (RE-NET) [43] is an auto-regressive architecture for modelling dynamic knowledge graphs and integrating an R-GCN in several RNNs.

Inspired by the success of the static GAE framework, dynamic GAE-based methods have emerged. The Dynamic Graph Embedding model (DynGEM) [44] modifies the static GAE to initialise it with the weights of the previous snapshot, and substantial modifications are not permitted from one snapshot to the next. Based on DynGEM, Dyngraph2vec [45] is introduced. This framework employs the  $l$  time window that defines the  $l$  most recent snapshots for encoding. Chen et al. [46] proposed Encoder-LSTM-Decoder (E-LSTM-D), which combines an LSTM with an encoder-decoder architecture. They stacked LSTM on GAE to learn graph evolution patterns.

All the above dynamic graph representation learning techniques employ deterministic vectors to represent each node in a low-dimensional space. These deterministic representations cannot reflect the uncertainty of the node representation. Although GAE-based methods perform effectively, they disregard data distribution and may lead to overfitting and poor representations [47, 48]. The combination of the GAE framework and deep generative models has been introduced for this purpose. Deep generative models have the ability to represent complex dependencies and interactions between input and output data by considering the distribution of data [49].

GCN-GAN [50] is a generative adversarial-based method for applying GCN to examine the topological properties of each snapshot and an LSTM to characterise the evolution of the dynamic graph. This component is a generator, while a dense feed-forward network is a discriminator. SI-VGRNN [51] is a generative approach that uses a VGAE in each snapshot. They consider a GRNN to model the temporal evolution of the graph. Our proposed framework contains an integration of VGAE and GRNN by exploiting a novel attention mechanism. Moreover, a natural assumption of multimodality of observed data is applied in our modelling [15, 16].



Earlier efforts modelled the uncertainty of the observed data using a unimodal Gaussian distribution. Under this assumption, modelling complex data with properties like multimodality is inefficient. Although SI-VGRNN develops semi-implicit variational inference for greater modelling flexibility, they only regard this assumption on their posterior modelling and not their prior. Hence, the improvement in their results is marginal [51]. To capture multimodality in the input data, our proposed DyVGRNN leverages GMM to model prior and posterior.

Furthermore, most previous works treat timed snapshots equally, despite the fact that assessing differences in snapshot significance may lead to more accurate results. SI-VGRNN assigns a fixed priority to different time series modelling snapshots, even though these snapshots may affect them differently. Here, we propose an attention-based module for examining the importance of snapshots. Unlike the traditional application of the attention mechanism in static graph representation learning, where the input is a matrix of nodes and the attention mechanism examines the importance of each node’s neighbouring nodes, the input in our module is a matrix of information for each time step, and the importance of time steps is examined.

### 3. The Proposed Model

#### 3.1. Notation and Problem definition

Let’s represent a dynamic graph  $G$  as  $G = \{G^{(1)}, G^{(2)}, \dots, G^{(T)}\}$ , where  $G^{(t)} = (V^{(t)}, E^{(t)})$  denotes a graph at time step  $t$ . Here  $V^{(t)}$  and  $E^{(t)}$  represent sets of nodes and edges, and  $T$  denotes the number of time steps. Since we intend to model a possible node or edge set change, the number of nodes and/or edges can change over time. Thus,  $(V^{(t)}, E^{(t)})$  and  $(V^{(t+1)}, E^{(t+1)})$  can be completely different. The input of the proposed method is a sequence of variable-length adjacency matrices in the form of  $\mathbf{A} = \{\mathbf{A}^{(1)}, \mathbf{A}^{(2)}, \dots, \mathbf{A}^{(T)}\}$  where  $\mathbf{A}^{(t)} \in R^{N_t \times N_t}$  and  $N_t$  denotes the number of nodes in this snapshot. Furthermore, there is a sequence of variable-length feature matrices in the form

of  $\mathbf{X} = \{\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(T)}\}$  as input, if the nodes have features. Here, each  $\mathbf{X}^{(t)}$  is a  $N_t \times F$  matrix, where  $F$  denotes the number of features. We assume  $F$  is constant over time. Table 1 summarises the notations used in this paper.

Table 1: The notation summary. This table summarises the notations used in this paper and provides a brief explanation for each.

Symbols	Meaning
$G$	Dynamic graph
$T$	Total number of snapshots
$G^{(t)}$	A snapshot of $G$ at time step $t$
$V^{(t)}$	Set of nodes in $G^{(t)}$
$E^{(t)}$	Set of edges in $G^{(t)}$
$\mathbf{A}^{(t)}$	The adjacency matrix of $G^{(t)}$
$N_t$	Number of nodes in $G^{(t)}$
$\mathbf{X}^{(t)}$	The features matrix of $G^{(t)}$
$F$	Number of features in $\mathbf{X}^{(t)}$
$\mathbf{Z}, \mathbf{W}, \mathbf{C}$	The latent variables in GMM
$\phi$	The parameters of encoder neural networks
$\theta$	The parameters of decoder neural networks
$\beta$	The parameters of the GNN related to each GMM component
$\phi_Z$	The parameters of the GNN related to $\mathbf{Z}$
$\phi_W$	The parameters of the GNN related to $\mathbf{W}$
$H$	The dimension of the representation embedding size

224

### 225 3.2. DyVGRNN

Figure 2 shows a high-level overview of our proposed method, DyVGRNN. The proposed method consists of three main modules described in this section. First, integrating GMM and VGAE used to model each graph snapshot is examined. Following, the process of modelling the evolution is described. Finally, we discuss the attention-based module for considering the importance of each

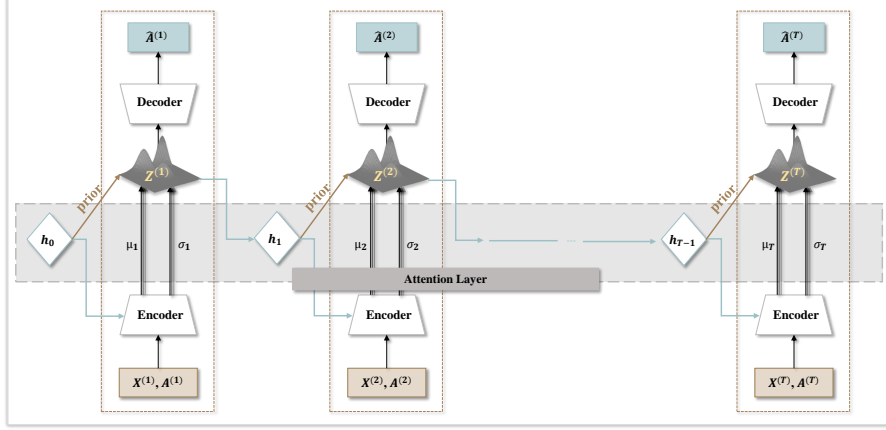


Figure 2: A high-level overview of our method. A VGAE integrated with GMM performs on each time step. The prior distribution of the VGAE is a function of the previous time step and a GRNN structure with extra hidden variables of the prior time step acts as a backbone of the entire framework. GRNN captures the dynamics of both graph topology and the node features jointly. The hidden state of GRNN is also added to latent random variables of GM-VGAE, making it capable of modelling variations in the topology or graph properties over time. Moreover, an attention-based module measures the importance of each graph snapshot in modelling evolution over time.

graph snapshot in modelling evolution over time.

### 3.2.1. Integration of GMM and VGAE

Our model defines three hidden variables  $\mathbf{Z}$ ,  $\mathbf{W}$ , and  $\mathbf{C}$  for integrating GMM and VGAE into a framework called Gaussian Mixture Variational Graph Auto Encoder (GM-VGAE). In this case, the inference model of standard VGAE for snapshot  $t$ , generalises and follows the process shown in the Equation (1)

$$\begin{aligned}
 \mathbf{W}^{(t)} &\sim \mathcal{N}(0, \mathbf{I}) \\
 \mathbf{C}^{(t)} &\sim \text{Cat}(\pi) \\
 \mathbf{Z}^{(t)} | \mathbf{C}^{(t)}, \mathbf{W}^{(t)} &\sim \prod_{k=1}^K \mathcal{N}(\mu_{c_k^{(t)}}(\mathbf{W}^{(t)}; \beta), \Sigma_{c_k^{(t)}}(\mathbf{W}^{(t)}; \beta))^{c_k^{(t)}}
 \end{aligned} \tag{1}$$

Here,  $K$  is a hyperparameter of the model, which denotes the number of components in the mixture model.  $\mathbf{W}^{(t)}$  is one of the latent variables of snapshot

239  $t$  that follows a Gaussian distribution with mean zero and covariance matrix  
 240  $\mathbf{I}$ .  $\mathbf{C}^{(t)}$  is a one-hot vector denoting the mixing coefficients of the Gaussian  
 241 mixture components of snapshot  $t$ . This vector is sampled from  $\pi$  (the mixing  
 242 probability), which indicates one of the Gaussian mixture components.

243  $\mathbf{W}^{(t)}$  is fed to a GNN parametrised by  $\beta$ . The output of this neural network  
 244 is a set of  $K$  ( $\boldsymbol{\mu}_{c_k}^{(t)}$ ) and  $K$  ( $\boldsymbol{\Sigma}_{c_k}^{(t)}$ ). Each  $\boldsymbol{\mu}_{c_k}^{(t)}$  and  $\boldsymbol{\Sigma}_{c_k}^{(t)}$  in these sets are calculated  
 245 by a GNN. An inner product between latent variables is used for reconstructing  
 246 the adjacency matrix, as shown in Equation (2).

$$\begin{aligned} p(\mathbf{A}^{(t)}|\mathbf{Z}^{(t)}) &= \prod_{i=1}^N \prod_{j=1}^N p(A_{ij}^{(t)}|\mathbf{z}_i^{(t)}, \mathbf{z}_j^{(t)}) \\ p(A_{ij}^{(t)}|\mathbf{z}_i^{(t)}, \mathbf{z}_j^{(t)}) &= \text{Sigmoid}(\mathbf{z}_i^{(t)T} \mathbf{z}_j^{(t)}) \end{aligned} \quad (2)$$

247 Based on the mean-field variational family, the general form of posterior can be  
 248 factorised as Equation (3).

$$\begin{aligned} q(\mathbf{Z}^{(t)}, \mathbf{W}^{(t)}, \mathbf{C}^{(t)}|\mathbf{A}^{(t)}) &= \\ \prod_{i=1}^{N_t} q_{\phi_Z}(\mathbf{z}_i^{(t)}|\mathbf{A}_i^{(t)}) q_{\phi_W}(\mathbf{w}_i^{(t)}|\mathbf{A}_i^{(t)}) q_{\beta}(\mathbf{z}_i^{(t)}|\mathbf{c}_i^{(t)}, \mathbf{w}_i^{(t)}) \end{aligned} \quad (3)$$

249 In this equation,  $\phi_Z$ ,  $\phi_W$ , and  $\beta$  are the parameters of neural networks, and the  
 250 output of these networks is the parameters of the variational distributions. The  
 251  $\mathbf{C}$ -posterior is as follows,

$$\begin{aligned} p_{\beta}(\mathbf{c}_j = 1|\mathbf{Z}, \mathbf{W}) &= \frac{p(\mathbf{c}_j = 1)p(\mathbf{Z}|\mathbf{c}_j = 1, \mathbf{W})}{\sum_{k=1}^K p(\mathbf{c}_k = 1)p(\mathbf{Z}|\mathbf{c}_k = 1, \mathbf{W})} \\ &= \frac{\pi_j \mathcal{N}(\mathbf{Z}|\mu_j(\mathbf{W}; \beta), \sigma_j(\mathbf{W}; \beta))}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{Z}|\mu_k(\mathbf{W}; \beta), \sigma_k(\mathbf{W}; \beta))} \end{aligned} \quad (4)$$

### 252 3.2.2. *Modelling the Evolution*

253 In contrast to standard VGAE that samples prior from a standard Gaussian  
 254 distribution ( $\mathcal{N}(0, I)$ ), the proposed VGAE (GM-VGAE) has a new prior extrac-  
 255 tion process that allows the parameter of the prior distribution to be modelled  
 256 by a function of the previous time step. In other words, the prior distribution  
 257 parameters are based on the information of the previous hidden state rather

258 than deterministic parameters. The construction of the prior distribution can  
 259 be written as shown in Equation (5).

$$\begin{aligned}
 \{\boldsymbol{\mu}_{prior}^{(t)}, \boldsymbol{\Sigma}_{prior}^{(t)}\} &= F^{prior}(\mathbf{h}_{t-1}) \\
 \mathbf{W}^{(t)} &\sim \mathcal{N}(\boldsymbol{\mu}_{prior}^{(t)}, \boldsymbol{\Sigma}_{prior}^{(t)}) \\
 \mathbf{C}^{(t)} &\sim Cat(\pi) \\
 \mathbf{Z}^{(t)} | \mathbf{C}^{(t)}, \mathbf{W}^{(t)} &\sim \prod_{k=1}^K \mathcal{N}(\boldsymbol{\mu}_{\mathbf{c}_k^{(t)}}(\mathbf{W}^{(t)}; \beta), \boldsymbol{\Sigma}_{\mathbf{c}_k^{(t)}}(\mathbf{W}^{(t)}; \beta))^{\mathbf{c}_k^{(t)}}
 \end{aligned} \tag{5}$$

260 Here  $\boldsymbol{\mu}_{prior}^{(t)}$  and  $\boldsymbol{\Sigma}_{prior}^{(t)}$  represent the parameters of the prior distribution.  $F^{prior}$   
 261 is a function that produces the parameters of prior distribution based on the  
 262 previous hidden state. This function can be a neural network. The prior dis-  
 263 tribution of the first step is assumed to be a standard multivariate Gaussian  
 264 distribution as  $\mathcal{N}(0, \mathbf{I})$ . If node addition occurs at each snapshot, the prior  
 265 distribution of the added node is defined as  $\mathcal{N}(0, \mathbf{I})$ . Eliminating a node can  
 266 be conceived as removing all edges connected to the node. In this way, prior  
 267 probabilities are unaffected.

268 The GRNN structure acts as a chain in the whole framework to capture the  
 269 dynamics of graph topology and features of the nodes. The GRNN update rule  
 270 is defined as shown in Equation (6).

$$\mathbf{h}_t = f(\mathbf{A}^{(t)}, \mathbf{X}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{h}_{t-1}) \tag{6}$$

271 Here  $f$  can be one of the Recurrent Neural Network (RNN) frameworks,  
 272 such as long short-term memory (LSTM) or gated recurrent units (GRU). In  
 273 this paper, we use LSTM-Attention for this purpose. If node addition occurs at  
 274 snapshot  $t$ , the hidden state of the node at snapshot  $t - 1$  is considered being  
 275 zero. The Z-posterior of the model is shown in Equation (7).

$$\begin{aligned}
 q(\mathbf{Z}^{(t)} | \mathbf{A}^{(t)}, \mathbf{X}^{(t)}, \mathbf{h}_{t-1}) &\sim \prod_{k=1}^K N(\boldsymbol{\mu}_{\mathbf{c}_k^{(t)}, enc}^{(t)}, \boldsymbol{\Sigma}_{\mathbf{c}_k^{(t)}, enc}^{(t)})^{\mathbf{c}_k^{(t)}} \\
 \boldsymbol{\mu}_{enc}^{(t)} &= GNN_{\mu}(A^{(t)}, CONCAT(\mathbf{X}^{(t)}, \mathbf{h}_{t-1})) \\
 \boldsymbol{\Sigma}_{enc}^{(t)} &= GNN_{\Sigma}(A^{(t)}, CONCAT(\mathbf{X}^{(t)}, \mathbf{h}_{t-1}))
 \end{aligned} \tag{7}$$

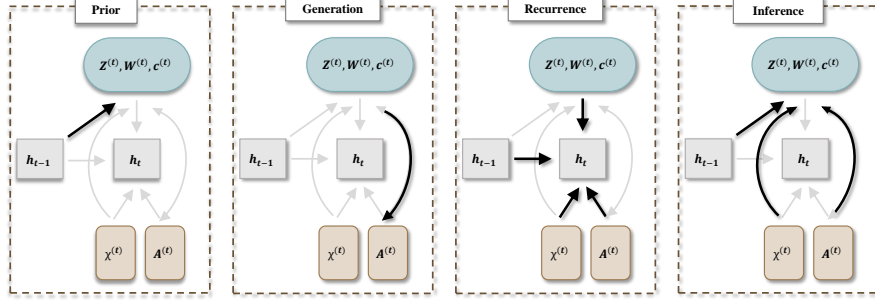


Figure 3: Graphical illustrations for Prior, Inference, Recurrence, and Generation of DyVGRNN. Arrows indicate the dependency of each component on the other component. The drawn arrow for Prior suggests the source of prior parameters, which is the previous hidden state of the model. The arrows of Inference and Recurrence indicate the resources needed to infer the latent variables and update the hidden state, respectively. The arrow of generation shows that the adjacency matrix can be reconstructed by having latent variables.

Here  $\mu_{enc}^{(t)}$  and  $\Sigma_{enc}^{(t)}$  represent the parameters of the posterior distribution, respectively.  $GNN_{\mu}(\cdot)$  and  $GNN_{\Sigma}(\cdot)$  can be any kind of GNN. We use a two-layer GCN for this purpose. The graphical illustrations for Prior, Inference, Recurrence, and Generation of DyVGRNN are shown in Figure 3. To carry out the learning process, the standard ELBO formulation is generalised as Equation (8) [13].

$$L_{ELBO} = \mathbb{E}_q \left[ \frac{p(\mathbf{A}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{W}^{(t)}, \mathbf{C}^{(t)})}{q(\mathbf{Z}^{(t)}, \mathbf{W}^{(t)}, \mathbf{C}^{(t)} | \mathbf{A}^{(t)})} \right] \quad (8)$$

in which,

$$p(\mathbf{A}^{(t)}, \mathbf{Z}^{(t)}, \mathbf{W}^{(t)}, \mathbf{C}^{(t)}) = p(\mathbf{W}^{(t)})p(\mathbf{C}^{(t)})p(\mathbf{Z}^{(t)} | \mathbf{W}^{(t)}, \mathbf{C}^{(t)})p(\mathbf{A}^{(t)} | \mathbf{Z}^{(t)}) \quad (9)$$

Based on the mean-field variational family, shown in Equation (3) and Equation (9), the lower bound for each snapshot can be written as Equation (10).

$$\begin{aligned} L_{ELBO}^{(t)} = & \mathbb{E}_{q(\mathbf{Z} | \mathbf{A}, \mathbf{X})} [\log p(\mathbf{A}^{(t)} | \mathbf{Z}^{(t)})] - \\ & \mathbb{E}_{q(\mathbf{W} | \mathbf{A}, \mathbf{X})p(\mathbf{C} | \mathbf{Z}, \mathbf{W})} [D_{KL}(q_{\phi_Z}(\mathbf{Z}^{(t)} | \mathbf{A}^{(t)}, \mathbf{X}^{(t)}) || p_{\beta}(\mathbf{Z}^{(t)} | \mathbf{W}^{(t)}, \mathbf{C}^{(t)}))] - \\ & D_{KL}(q_{\phi_W}(\mathbf{W}^{(t)} | \mathbf{A}^{(t)}, \mathbf{X}^{(t)}) || p(\mathbf{W}^{(t)})) - \\ & \mathbb{E}_{q(\mathbf{Z} | \mathbf{A}, \mathbf{X})q(\mathbf{W} | \mathbf{A}, \mathbf{X})} [D_{KL}(p_{\beta}(\mathbf{C}^{(t)} | \mathbf{Z}^{(t)}, \mathbf{W}^{(t)}) || p(\mathbf{C}^{(t)}))] \end{aligned} \quad (10)$$

285 This equation consists of four terms representing the reconstruction error term,  
 286 prior conditional term, **W**-prior term, and **C**-prior term. The total loss function  
 287 of the model is calculated as the sum of the loss functions of each snapshot.  
 288 Thus, the loss function can be written as Equation (11).

$$L_{ELBO}^{(total)} = \sum_{t=1}^T L_{ELBO}^{(t)} \quad (11)$$

### 289 3.2.3. *Attention Module*

290 The attention mechanism was first introduced by [52] in the field of Nat-  
 291 ural Language Processing (NLP). This work became the basis for [53], which  
 292 attracted much attention. Recent studies in NLP have emphasised that the use  
 293 of the attention mechanism improves the efficiency and performance of models  
 294 [54, 55, 56]. Other fields have also been positively influenced by the capabil-  
 295 ity of this mechanism [57, 39], and we endeavour to use the potential of this  
 296 mechanism.

297 We add an attention module to the proposed model, which receives as input  
 298 the hidden states and structural information of all time steps. The attention  
 299 module’s output hidden state is considered the model’s final hidden state. Struc-  
 300 tural information is then used to calculate the loss function like Equation (10)  
 301 with the parameters gained by the attention mechanism. Then, backpropaga-  
 302 tion of the gradients of the loss function leads to updating the weights. In this  
 303 way, the importance of each snapshot is taken into consideration in the learning  
 304 process.

305 Here, the mean and standard deviation matrices are remarked as the struc-  
 306 tural information of each snapshot. Thereupon, the received information is  
 307 converted into a matrix, each row showing one snapshot’s information. This  
 308 operation is fulfilled for both the mean and standard deviation matrices. The  
 309 un-normalised attention scores between two snapshots are calculated according

310 to Equation (12).

$$\begin{aligned}
e_{i,j}^{\mu} &= \text{LeakyReLU}(a(\text{CONCAT}(\mu_i, \mu_j))) \\
e_{i,j}^{\sigma} &= \text{LeakyReLU}(a(\text{CONCAT}(\sigma_i, \sigma_j))) \\
e_{i,j}^{\mathbf{h}} &= \text{LeakyReLU}(a(\text{CONCAT}(\mathbf{h}_i, \mathbf{h}_j)))
\end{aligned} \tag{12}$$

311 Here  $a$  is a learnable weight vector. The normalised attention scores calculate by  
312 applying a Softmax to un-normalised attention scores as shown in Equation (13).  
313 Eventually, these  $\alpha$  sets determine the importance of each time step.

$$\begin{aligned}
\alpha_{i,j}^{\mu} &= \frac{\exp e_{i,j}^{\mu}}{\sum_{k \in \mu} \exp e_{i,k}^{\mu}}, \quad \mu = \{\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(T)}\} \\
\alpha_{i,j}^{\sigma} &= \frac{\exp e_{i,j}^{\sigma}}{\sum_{k \in \sigma} \exp e_{i,k}^{\sigma}}, \quad \sigma = \{\sigma^{(1)}, \sigma^{(2)}, \dots, \sigma^{(T)}\} \\
\alpha_{i,j}^{\mathbf{h}} &= \frac{\exp e_{i,j}^{\mathbf{h}}}{\sum_{k \in \mathbf{h}} \exp e_{i,k}^{\mathbf{h}}}, \quad \mathbf{h} = \{\mathbf{h}^{(1)}, \mathbf{h}^{(2)}, \dots, \mathbf{h}^{(T)}\}
\end{aligned} \tag{13}$$

314 Both DyREP [28] and DySAT [39] leverage the attention mechanism as  
315 part of their method. They employ node-based attention mechanisms in their  
316 framework. DyREP computes the attention coefficient and evaluates the im-  
317 portance of each node’s neighbours using temporal information. DySAT applies  
318 one attention layer to focus on each node’s immediate neighbours, and a sec-  
319 ond attention layer to focus on each node’s temporal history in each snapshot.  
320 While our attention module is based on graphs, these two methods use a node-  
321 based attention module. In fact, in their methods, the input would be a matrix  
322 of nodes and the attention mechanism examines the importance of the neigh-  
323 bouring nodes of each node. Whereas the input of our module is a matrix of  
324 information for each time step, and the importance of time steps is examined.

#### 325 4. Experimental Details

326 In this section, the results of the experiments are presented. First, the datasets,  
327 the state-of-the-art methods, and the studied tasks and metrics are introduced.  
328 Then, the results of the experiments are described.



#### 4.1. Datasets

Our experiments are performed on five real-world graph datasets. Table 2 presents a summary of the employed datasets.

**Facebook.** This dataset contains information about Facebook posts. The Facebook dataset is collected by [58], and the procedure of cleaning and preparing the data is similar to the procedure in [59, 60]. This dataset has 663 nodes and 1068 edges but does not contain node or edge attributes.

**LFB.** This dataset is a larger-scale version of the Facebook dataset containing 45435 nodes and 180011 edges. The procedure of cleaning and preparing the data in this version is also similar to the procedure in [59, 60]. 36 snapshots of the activations throughout the last three years are included in the dataset. In the LFB dataset, there are a large number of users but not many links between them.

**Enron emails (Enron).** This dataset contains 500,000 emails exchanged between Enron employees from 1998 to 2002 [61]. The nodes represent 184 employees, and the edges represent the emails exchanged between pairs of employees in the graph created from this dataset. The steps of cleaning and producing the appropriate structure for applying the algorithm are done according to the procedure in [59, 60, 51]. This dataset has no node or edge attributes.

**Collaboration (Colab).** There is information about co-authorship relationships between 315 authors in this dataset. Each node represents an author, and each edge demonstrates co-authorship relationships between a pair of authors

Table 2: Summary of the employed datasets. “-” in “Number of Edge” column means the number changes across different snapshots.

Dataset	Number of Snapshots	Number of Nodes	Number of Edges	Number of Node Attributes
<b>Enron</b>	11	184	115-266	-
<b>Colab</b>	10	315	165-308	-
<b>Facebook</b>	6	663	844-1068	-
<b>UCI</b>	7	537-1899	59835	-
<b>Cora</b>	6	500-2708	406-5429	1433
<b>LFB</b>	36	45435	180011	-
<b>AS733</b>	30	6628	13512	-

351 from 2000 to 2009 [60]. This dataset has no node or edge attributes.

352 **UCI.** This dataset was aggregated by the University of California, Irvine [61].

353 In this dataset, message interaction information between students based on an

354 online community has been collected. Nodes represent students, and edges rep-

355 resent the sending of a message between two students. This information was

356 collected over a 7-day period. Each day denotes one snapshot of the graph. This

357 dynamic graph starts at 537 nodes, ends at 1899 nodes, has 59835 edges, and

358 has no node properties.

359 **Cora.** This dataset is a static citation graph in which the nodes represent

360 the publications, and the edges denote the citation [62]. Cora consists of 2,708

361 nodes with a 1,433-dimensional binary attribute vector. To make use of Cora

362 dynamically, we preprocess the data in the same way as described in [63, 51].

363 In the dynamic network, we added 500 nodes with their accompanying edges at

364 each temporal snapshot (208 nodes for the last snapshot), using the indexes of

365 the nodes as their arrival order, and six snapshots of the dynamic graph were

366 taken, starting with 500 nodes and ending with 2708 nodes.

367 **AS733.** This dataset is a communication network containing Autonomous Sys-

368 tems (AS) and traffic flows between them that show who communicates with

369 whom. AS733 was gathered from the Route Views Project at the University

370 of Oregon, which contains 733 daily instances spanning 785 days between 1997

371 and 2000 [64]. There are 6628 nodes and 13512 edges in this dataset.

## 372 4.2. Baselines

373 We compare DyVGRNN with the following baselines and state-of-the-art

374 methods. We use the original implementation of the methods introduced in

375 their paper. To ensure a fair comparison, the hyperparameters are adjusted

376 based on the suggestion in their papers.

## 377 4.3. Discrete Dynamic Graph Representation Learning Methods

378 **DynAE (Dynamic Auto-Encoder) [45]:** This model is an auto-encoder com-

379 posed of multiple fully connected layers as the encoder and decoder. These layers

are used to capture nonlinear interactions between nodes at each snapshot and across multiple snapshots.

**DynRNN (Dynamic Recurrent Neural Network) [45]:** This model consists of an LSTM encoder and an LSTM decoder. These encoder and decoder allow capturing the long-term dependencies in dynamic graphs.

**DynAERNN (Dynamic Auto-Encoder Recurrent Neural Network) [45]:** This model includes a fully connected layer connected to an LSTM as the encoder. The fully connected layer generates initial low-dimensional hidden representations, which are then fed to LSTM. Here, the decoder is a fully connected network.

**SI-VGRNN (Variational Graph Recurrent Neural Networks) [51]:** This method was the inspiration for this paper that is based on VGAE, which is combined with GRNN to capture topology and node feature changes in dynamic graphs. This paper suggested regarding and disregarding the semi-implicit part as an SI-VGRNN and VGRNN, respectively.

**DySAT (Dynamic Self-Attention Network) [39]:** This method computes node representations through self-attention blocks that capture structural and temporal properties.

**HTGN (Hyperbolic Temporal Graph Network) [65]:** This approach maps the dynamic graph in hyperbolic space and combines a hyperbolic GNN and a hyperbolic GRNN to capture network evolution while implicitly maintaining hierarchical information.

#### 4.4. Continuous Dynamic Graph Representation Learning Methods

**DyREP [28]:** This model uses a two-time scale Temporal Point Process (TPP) model, which is parametrised by an RNN.

**JODIE [30]:** This model uses RNNs to predict representations in the future. Since the method was originally proposed for bipartite graphs, we modified it for standard graphs in accordance with [66].

**TGAT [67]:** This model is based on the self-attention mechanism and develops a functional time encoding technique based on the classical Bochner’s theorem.

Table 3: AP scores of link prediction on dynamic graphs. The best results are highlighted.

Model	Enron	Colab	Facebook	LFB	UCI	Cora	AS733
<b>DynAE</b>	76.00	64.02	56.04	58.90	91.12	57.11	74.23
<b>DynRNN</b>	85.61	78.95	75.88	75.28	89.21	80.75	87.53
<b>DynAERNN</b>	89.37	81.84	78.55	78.27	89.92	82.93	88.77
<b>DySAT</b>	93.06	90.40	80.39	80.39	85.01	87.73	96.72
<b>HTGN</b>	94.31	91.91	83.80	83.80	86.72	90.12	98.41
<b>VGRNN</b>	93.29	87.77	89.04	81.40	91.83	93.32	96.69
<b>SI-VGRNN</b>	94.44	88.36	90.19	82.01	93.16	96.68	97.13
<b>DyVGRNN</b>	<b>97.28</b>	<b>96.77</b>	<b>92.70</b>	<b>86.22</b>	<b>95.07</b>	<b>97.48</b>	<b>99.10</b>

#### 4.5. Tasks

We perform the link prediction and clustering tasks in this study to evaluate our method. The link prediction task in dynamic graphs is defined differently than in static graphs. Given a dynamic graph  $G = \{G^{(1)}, G^{(2)}, \dots, G^{(T)}\}$ , the link prediction is divided into two categories: 1) dynamic link prediction attempts to identify the unobserved links in  $G^{(T)}$ , and 2) dynamic new link prediction tries to predict links in  $G^{(T+1)}$  which does not exist in  $G^{(T)}$ .

#### 4.6. Metrics

We use the Average Precision (AP) and the Area Under the receiver operating characteristic Curve (AUC) [27] metrics to compare our proposed method with state-of-the-art methods in link prediction and new link prediction tasks. To calculate these measures, all edges of  $G^T$  are considered as actual links (positive samples), and on the other hand, the pairs of nodes without an edge imply false links (negative samples). Furthermore, the silhouette criterion is applied for the evaluation of the clustering results to interpret and validate data consistency within clusters.

Table 4: AUC scores of link prediction on dynamic graphs. The best results are highlighted.

Model	Enron	Colab	Facebook	LFB	UCI	Cora	AS733
<b>DynAE</b>	74.22	63.14	56.06	57.18	91.89	57.13	73.84
<b>DynRNN</b>	86.41	75.7	73.18	73.98	89.27	80.10	86.11
<b>DynAERNN</b>	87.43	76.06	76.02	75.28	90.08	78.00	88.37
<b>DySAT</b>	93.06	87.25	76.88	76.88	86.73	85.3	95.06
<b>HTGN</b>	94.17	89.26	83.70	83.7	87.25	89.73	98.75
<b>VGRNN</b>	93.10	85.95	89.47	79.11	92.01	94.41	95.17
<b>SI-VGRNN</b>	93.93	85.45	90.94	80.27	93.5	97.17	96.37
<b>DyVGRNN</b>	<b>96.59</b>	<b>95.80</b>	<b>93.17</b>	<b>86.73</b>	<b>95.15</b>	<b>98.74</b>	99.19

#### 4.7. Settings

The proposed model uses the LSTM-attention with a single hidden layer of 32 units for the GRNN. The  $GNN_{\mu}$  and  $GNN_{\Sigma}$  are set to be two-layer GCN with 32 and 16 units, respectively. Our model is initialised using Glorot initialisation [68]. The learning rate for training our model is set to be 0.01. Model training is done in 1000 epochs using the Adam SGD optimiser [69]. Moreover, we use a validation set for the early stopping. Therefore, the training will terminate if the validation accuracy does not improve in 10 consecutive stages. The mean of the evaluation metrics is reported based on 10 runs of the model under different random seeds.

#### 4.8. Results Analysis

**Dynamic Link Prediction.** Tables 3 and 4 represent the comparison results in terms of AP and AUC on the link prediction task. The results of the dominant algorithm are highlighted. DyVGRNN shows significant improvement in results compared to the other methods. The enhancement of our method using the AP criterion compared to the first method is 21.28% in the Enron dataset, 32.75% in the Colab dataset, 36.66% in the Facebook dataset, and 40.37% in the Cora dataset. Large datasets like LFB and AS733 show improvements of 27.32% and 24.87%, respectively. Likewise, in the UCI dataset, where the first method

Table 5: AUC scores of new link prediction on dynamic graphs. The best results are highlighted.

Model	Enron	Colab	Facebook	LFB	UCI	Cora	AS733
DynAE	66.10	58.14	54.62	56.34	89.94	56.27	68.93
DynRNN	83.20	71.71	73.32	74.15	87.27	79.94	74.72
DynAERNN	83.77	71.99	76.35	76.55	88.29	77.36	76.63
DySAT	87.94	79.74	74.97	74.97	84.2	86.11	82.84
HTGN	91.26	81.74	82.21	82.21	84.98	87.85	96.62
VGRNN	88.43	77.09	87.20	76.33	89.93	94.94	81.86
SI-VGRNN	88.60	77.95	87.74	77.42	90.45	96.36	83.27
DyVGRNN	<b>94.26</b>	<b>92.71</b>	<b>92.51</b>	<b>85.26</b>	<b>94.17</b>	<b>97.16</b>	<b>97.89</b>

performed well, our proposed method boosts the result by 3.95%. If we compare the AUC criteria, the results are also significantly improved. For example, the results of a comparison with SI-VGRNN, which on average provided the best results among the previous methods, show that the proposed method leads to 2.66% improvement in the Enron dataset, 10.35% in the Colab dataset, 2.23% in the Facebook dataset, 1.65% in the UCI dataset, and eventually 1.57% in the Cora dataset. Large datasets LFB and AS733 have improvements of 6.46.21% and 2.82%, respectively.

**Dynamic New Link Prediction.** Tables 5 and 6 represent the results of comparisons regarding AUC and AP on the new link prediction task. The proposed method has achieved significant results in all datasets. A similar analysis for the link prediction task can be provided for the new link prediction task. In general, it can be noted that the proposed method can have a high potential for predicting the overall structure of the graph in the new snapshot.

To point out some significant improvements, we can mention the progress of more than 40% in the Cora dataset or the increase of over 37% in the Facebook dataset in both criteria compared to DynAE. In addition, our method performed superior to SI-VGRNN, which indicates a positive effect of the assumption of GMM and the proposed attention module. A comparison of the proposed DyV-

Table 6: AP scores of new link prediction on dynamic graphs. The best results are highlighted.

Model	Enron	Colab	Facebook	LFB	UCI	Cora	AS733
<b>DynAE</b>	66.50	58.82	54.57	54.91	89.65	56.65	69.12
<b>DynRNN</b>	80.96	75.34	75.52	76.01	86.86	80.01	75.12
<b>DynAERNN</b>	85.16	77.68	78.70	78.27	88.15	82.34	76.87
<b>DySAT</b>	86.83	83.47	78.34	78.34	83.94	87.15	89.07
<b>HTGN</b>	90.62	84.06	81.70	81.7	84.26	89.83	95.52
<b>VGRNN</b>	87.57	79.63	86.30	79.61	89.48	93.21	88.59
<b>SI-VGRNN</b>	87.88	81.26	86.72	80.12	90.07	95.32	89.49
<b>DyVGRNN</b>	<b>94.44</b>	<b>93.65</b>	<b>91.81</b>	<b>85.00</b>	<b>94.11</b>	<b>96.82</b>	96.83

GRNN and VGRNN is presented in Appendix A in order to further analyse the effectiveness of the methods.

**Clustering.** For further investigation, we provide a clustering comparison as well. The proposed approach is compared against SI-VGRNN, which achieves the highest result among various methods, and DySAT, which performs best among deterministic ones. To this end, the silhouette criterion is utilised for clustering the Cora dataset. This criterion is 0.32 for DySAT, 0.36 for SI-VGRNN, and 0.43 for our approach. Demonstrating a transparent view, we visualise the representations of these three methods in a two-dimensional space as shown in Figure 4. Compared to the raw features, the trained representations in two-dimensional space for our method indicate well-separated clustering compared to SI-VGRNN. In addition, modelling uncertainty in SI-VGRNN and DyVGRNN yields superior clustering outcomes compared to DySAT, which is a deterministic-based method. We also provide a classification comparison in Appendix C.

**Comparison with Continuous Methods.** We compare our model to state-of-the-art methods in the category of continuous dynamic graph representation learning in terms of dynamic link prediction. The results of this comparison are shown in Figure 5. As demonstrated in the Figure 5, our proposed method out-

483 performs other continuous methods. DyRep has the best performance among  
 484 existing continuous approaches, which our method enhances.

#### 485 4.9. Complexity and Running Time

486 To compute the time complexity of our method, the analysis of [70] is fol-  
 487 lowed. For this purpose, the proposed DyVGRNN can be divided into three  
 488 main parts. 1) modelling node temporal attributes by LSTM which the time  
 489 complexity is  $O(T|V|H^2)$ . 2) modelling node structural properties by VGAE,  
 490 which consists of GCN structure in its encoder and an inner product decoder.  
 491 Time complexity of GCN is  $O(|V|H^2 + (|V| + |E|)H)$ . Since  $H$  and  $|V|$  are

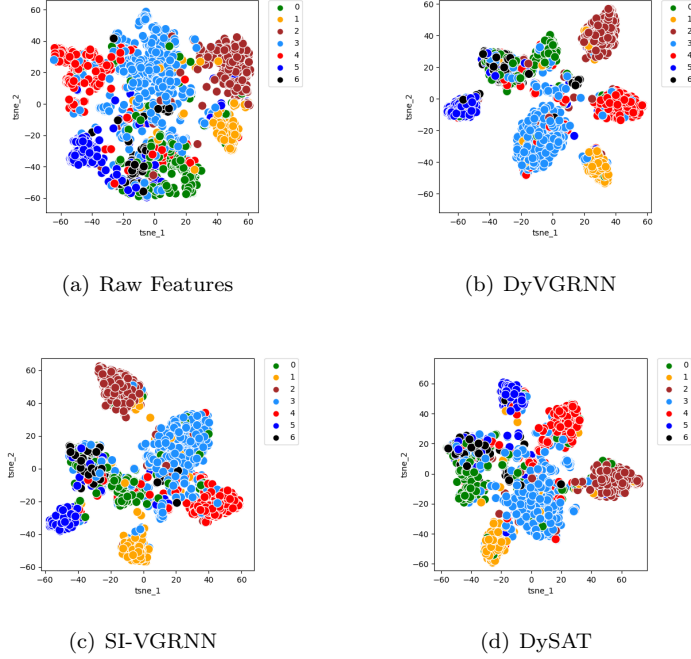


Figure 4: Cluster visualisation for embeddings of Cora dataset in 2D space. a) Raw feature cluster visualisation demonstrates the inability to differentiate between clusters. b) Cluster visualisation of DyVGRNN embeddings showing distinct clusters. c) Cluster visualisation of SI-VGRNN embeddings indicates more indiscernible clusters compared to DyVGRNN. d) Cluster visualisation of DySAT embedding also reveals more undetectable clusters compared to the two other methods.



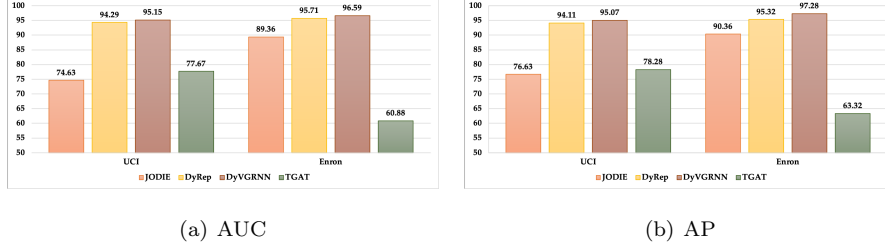


Figure 5: The comparison of the proposed DyVGRNN and continuous methods for the task of dynamic link prediction. As shown in the charts below, different methods are represented by different colours. These results are performed on UCI and Enron. Each chart shows the results for UCI and Enron in the left and right groups, respectively. a) The results of comparing in terms of AUC score. b) The results of comparing in terms of AP score.

492 relatively small w.r.t. to  $|E|$ , the time cost is indeed  $O(|E|)$ .

493 Moreover, the time complexity of the inner product decoder is  $O(|E|)$ . As  
 494 a result, the time complexity of VGAE is  $O(|E|)$ . 3) Considering the attention  
 495 mechanism which has an order of  $O(EH^2)$ . Eventually, the time complexity of  
 496 our proposed method is  $O(T|V|H^2) + O(EH^2)$ . Table 7 lists the time complexity  
 497 of some methods evaluated in our work on LFB dataset. In addition, Figure 6  
 498 contrasts the running times of SI-VGRNN, DySAT, and DyVGRNN. As seen,  
 499 our approach runs faster than DySAT but lower than SI-VGRNN. Although  
 500 compared to VGRNN, this is seen as a shortcoming for our model, accuracy at

Table 7: Time Complexity of different methods.

Method	Time Complexity
DynAE	$O(T( E  +  V ))$
DynRNN	$O(T V H^2)$
DynAERNN	$O(T V H^2 + T( E  +  V ))$
HTGN	$O(T V H^2 +  E H^2)$
DySAT	$O(T V H^2 +  E H^2)$
VGRNN	$O(T V H^2) + O( E )$
DyVGRNN	$O(T V H^2) + O( E H^2)$

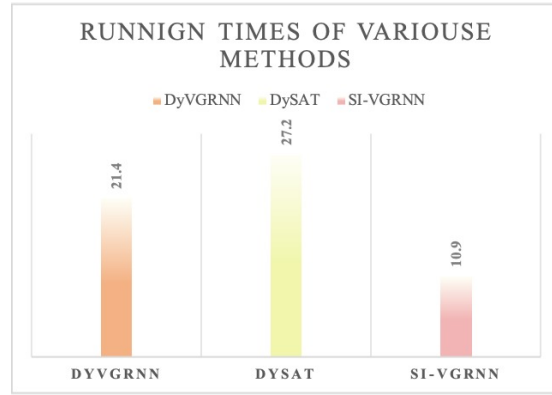


Figure 6: Comparison of running times of different methods on the LFB dataset. The colours represent various methods in the colour scheme.

501 inference time is more important in many cases.

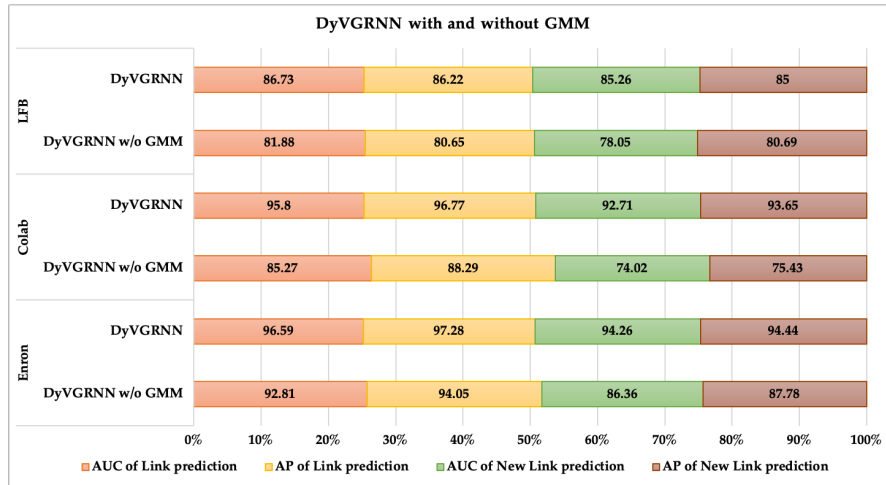


Figure 7: The effect of GMM on proposed DyVGRNN. The outcomes of running our model in two modes with and without GMM are shown in this figure. Two tasks, dynamic link prediction, and new dynamic link prediction are performed on Enron, Colab, and Facebook with the results. The various criteria for these two tasks are represented by the colours in accordance with the colour scheme.

Table 8: Effect of parameter  $K$  on the DyVGRNN outcome. The results of dynamic link prediction and dynamic new link prediction by adjusting  $K$  to different values are given in this table. “Mean of AUCs” in each dataset category show the mean of AUC of link prediction and AUC of link prediction for different  $K$ .

Dataset	Metric	$K = 2$	$K = 3$	$K = 4$	$K = 5$	$K = 6$	$K = 7$
Enron	AUC of link prediction	95.80	96.59	95.70	96.60	95.92	95.82
	AP of link prediction	96.77	97.28	96.34	97.10	96.73	96.64
	AUC of new link prediction	92.71	94.26	93.10	93.58	93.64	92.98
	AP of new link prediction	93.65	94.44	93.36	93.38	94.25	93.57
	Mean of AUCs	94.25	95.42	94.4	95.09	94.78	94.4
	Mean of APs	95.21	95.86	94.85	95.24	95.49	95.10
Facebook	AUC of link prediction	93.17	90.20	92.55	91.37	92.61	93.02
	AP of link prediction	92.70	88.67	92.17	90.66	92.18	92.47
	AUC of new link prediction	92.51	89.79	91.95	90.70	91.93	92.55
	AP of new link prediction	91.81	88.11	91.67	89.81	91.17	92.04
	Mean of AUCs	92.84	89.99	92.25	91.03	92.27	92.78
	Mean of APs	92.25	88.39	91.92	90.23	91.67	92.25
Colab	AUC of link prediction	95.80	90.20	92.55	91.37	92.61	93.02
	AP of link prediction	96.77	88.67	92.17	90.66	92.18	92.47
	AUC of new link prediction	92.71	89.79	91.95	90.70	91.93	92.55
	AP of new link prediction	93.65	88.11	91.67	89.81	91.17	92.04
	Mean of AUCs	92.84	89.99	92.25	91.03	92.27	92.78
	Mean of APs	92.25	88.39	91.92	90.23	91.67	92.25

#### 4.10. Ablation Study

In this section, we conduct ablation studies to verify the effectiveness of the key components of the proposed model.

##### Selection of $K$

Since each dataset has various properties, we need to select the hyperparameter  $K$  according to the unique properties of each dataset. To this end, this study compares the results by examining the various values of  $K$  and selecting the best value. Table 8 shows the results of these comparisons. The best value of  $K$  for Enron, UCI, Cora, and AS733 datasets was 3, and for Facebook, Colab, and LFB datasets were 2. The first column of these tables shows the situation where the GMM does not affect the results. As can be seen, at  $K = 2$ , i.e.



Figure 8: Impact of the attention module on DyVGRNN. In this figure, the results of running our model in two different modes with and without the attention module are depicted. To achieve this, two tasks—dynamic link prediction and new dynamic link prediction—are carried out. The colours, in accordance with the colour scheme, represent the various criteria for these two tasks. The results of the comparison on the a) Colab, b) Facebook, c) Enron datasets.

513 applying the GMM, a significant improvement in the results is achieved. This  
514 improvement demonstrates the validity of our claim that the use of GMM pos-  
515 itively affects outcomes.

#### 516 **Impact of the GMM**

517 The effects of utilising a GMM to handle multimodality are examined in this  
518 section. This is accomplished by considering the proposed DyVGRNN in two  
519 different scenarios: first, without using GMM, and second, using GMM. Fig-  
520 ure 7 shows the result of comparisons in these two modes. As seen in this figure,  
521 GMM leads to improving the results.

#### 522 **Impact of the Attention Module**

523 To assess the effectiveness of the attention module, we have divided the proposed  
524 model into two modes: with and without using it. To emphasise the attention  
525 module, we investigated the proposed method without considering GMM. Fig-

ure 8 shows the result of comparisons in these two modes.

### Impact of Features

A noteworthy point in examining the results is the effect of the node features on the results. Figure 9 shows the performance of DyVGRNN in the Cora in two modes: with and without features. The performance is significantly improved with the presence of node features, which indicates our proposed method can capture long-term dependencies in both the topological evolution and dynamics of node features.



Figure 9: The results of comparing the proposed method on Cora with and without using features. Dynamic link prediction is used to accomplish this. a) The result of comparison in terms of AUC. Results are enhanced by the presence of node features. b) The result of comparison in terms of AP. Results are improved when node features are present.

## 5. Conclusion and Future Works

We proposed DyVGRNN, an integrated variational GRNN for learning node representations of dynamic graphs. DyVGRNN has additional random latent variables in the GRNN framework for capturing the evolution of graph structures and node attributes. We have shown that the combination of variational inference based on GMM and the proposed framework leads to a high level of validity and knowledge of the model. We also introduced an attention module to consider each snapshot’s importance, leading to improved performance. The experiments’ results showed our model’s superiority over baseline and state-of-the-art methods. In the future, we are looking to apply a probabilistic decoder

to the VGAE structure than a simple inner product decoder. In our proposed method, VGAEs reconstruct the adjacency matrix, but not the features matrix. Therefore, considering the reconstruction of the feature matrix and adjacency matrix would lead to a raise in accuracy. We believe it is a worthwhile area to explore more. In addition, it would intrigue to study the impact of other GNN frameworks, such as GAT, GraphSAGE, and GIN, with different layer numbers for the encoder and perhaps the decoder.

## 6. Acknowledgements

DAC was supported by an NIHR Research Professorship, an RAEng Research Chair, the InnoHK Hong Kong Centre for Cerebro-cardiovascular Health Engineering (COCHE), the NIHR Oxford Biomedical Research Centre (BRC), and the Pandemic Sciences Institute at the University of Oxford.

## References

- [1] Y. Liu, X. Shi, L. Pierce, X. Ren, Characterizing and forecasting user engagement with in-app action graph: A case study of snapchat, in: Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, 2019, pp. 2023–2031.
- [2] L. Zhao, Y. Song, C. Zhang, Y. Liu, P. Wang, T. Lin, M. Deng, H. Li, T-gcn: A temporal graph convolutional network for traffic prediction, IEEE Transactions on Intelligent Transportation Systems 21 (9) (2019) 3848–3858.
- [3] A. Fout, J. Byrd, B. Shariat, A. Ben-Hur, Protein interface prediction using graph convolutional networks, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 6533–6542.
- [4] R. Angles, C. Gutierrez, Survey of graph database models, ACM Computing Surveys (CSUR) 40 (1) (2008) 1–39.

- [5] D. Bacciu, F. Errica, A. Micheli, M. Podda, A gentle introduction to deep learning for graphs, *Neural Networks* 129 (2020) 203–221.
- [6] W. L. Hamilton, Graph representation learning, *Synthesis Lectures on Artificial Intelligence and Machine Learning* 14 (3) (2020) 1–159.
- [7] S. Molaei, N. G. Bousejin, H. Zare, M. Jalili, S. Pan, Learning graph representations with maximal cliques, *IEEE Transactions on Neural Networks and Learning Systems* (2021) 1–8.
- [8] W. Ju, X. Luo, Z. Ma, J. Yang, M. Deng, M. Zhang, Ghnn: Graph harmonic neural networks for semi-supervised graph-level classification, *Neural Networks* 151 (2022) 70–79.
- [9] G. Salha-Galvan, J. F. Lutzeyer, G. Dasoulas, R. Hennequin, M. Vazirgianis, Modularity-aware graph autoencoders for joint community detection and link prediction, *Neural Networks* (2022) 474–495.
- [10] T. B. Mudiyansele, X. Lei, N. Senanayake, Y. Zhang, Y. Pan, Predicting circrna disease associations using novel node classification and link prediction models on graph convolutional networks, *Methods* 198 (2022) 32–44.
- [11] S. Molaei, N. G. Bousejin, H. Zare, M. Jalili, Deep node clustering based on mutual information maximization, *Neurocomputing* 455 (2021) 274–282.
- [12] J. Chen, A. Zhang, Hgmf: heterogeneous graph-based fusion for multimodal data with incompleteness, in: *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2020, pp. 1295–1305.
- [13] N. Dilokthanakul, P. A. Mediano, M. Garnelo, M. C. Lee, H. Salimbeni, K. Arulkumaran, M. Shanahan, Deep unsupervised clustering with gaussian mixture variational autoencoders, *ICLR* (2017) 1–12.
- [14] G. Niknam, S. Molaei, H. Zare, D. Clifton, S. Pan, Graph representation learning based on deep generative gaussian mixture models, *Neurocomputing* (2022) 157–169.

- 599 [15] N. Kostantinos, Gaussian mixtures and their applications to signal process-  
600 ing, *Advanced signal processing handbook: theory and implementation for*  
601 *radar, sonar, and medical imaging real time systems* (2000) 3–1.
- 602 [16] I. Goodfellow, Y. Bengio, A. Courville, *Deep learning*, MIT press, 2016.
- 603 [17] A. Ahmed, N. Shervashidze, S. Narayanamurthy, V. Josifovski, A. J. Smola,  
604 Distributed large-scale natural graph factorization, in: *Proceedings of the*  
605 *22nd international conference on World Wide Web*, 2013, pp. 37–48.
- 606 [18] S. Cao, W. Lu, Q. Xu, Grarep: Learning graph representations with global  
607 structural information, in: *Proceedings of the 24th ACM international on*  
608 *conference on information and knowledge management*, 2015, pp. 891–900.
- 609 [19] M. Ou, P. Cui, J. Pei, Z. Zhang, W. Zhu, Asymmetric transitivity pre-  
610 serving graph embedding, in: *Proceedings of the 22nd ACM SIGKDD in-*  
611 *ternational conference on Knowledge discovery and data mining*, 2016, pp.  
612 1105–1114.
- 613 [20] B. Perozzi, R. Al-Rfou, S. Skiena, Deepwalk: Online learning of social  
614 representations, in: *Proceedings of the 20th ACM SIGKDD international*  
615 *conference on Knowledge discovery and data mining*, 2014, pp. 701–710.
- 616 [21] A. Grover, J. Leskovec, node2vec: Scalable feature learning for networks,  
617 in: *Proceedings of the 22nd ACM SIGKDD international conference on*  
618 *Knowledge discovery and data mining*, 2016, pp. 855–864.
- 619 [22] W. L. Hamilton, R. Ying, J. Leskovec, Representation learning on graphs:  
620 *Methods and applications*, *IEEE Data(base)* (2017) 1–24.
- 621 [23] J. Skarding, B. Gabrys, K. Musial, Foundations and modeling of dynamic  
622 networks using dynamic graph neural networks: A survey, *IEEE Access* 9  
623 (2021) 79143–79168.
- 624 [24] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolu-  
625 tional networks, *ICLR* (2017) 1–14.



- [25] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, S. Y. Philip, A comprehensive survey on graph neural networks, *IEEE transactions on neural networks and learning systems* 32 (1) (2020) 4–24.
- [26] D. P. Kingma, M. Welling, Auto-encoding variational bayes, *ICLR* (2014) 14–16.
- [27] T. N. Kipf, M. Welling, Variational graph auto-encoders, *NIPS Workshop on Bayesian Deep Learning* (2016) 1–12.
- [28] J. Skarding, B. Gabrys, K. Musial, Foundations and modeling of dynamic networks using dynamic graph neural networks: A survey, *IEEE Access* 9 (2021) 79143–79168.
- [29] Y. Ma, Z. Guo, Z. Ren, J. Tang, D. Yin, Streaming graph neural networks, in: *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2020, pp. 719–728.
- [30] S. Kumar, X. Zhang, J. Leskovec, Predicting dynamic embedding trajectory in temporal interaction networks, in: *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 1269–1278.
- [31] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, R. Zemel, Neural relational inference for interacting systems, *Proceedings of Machine Learning Research* 80 (2018) 2688–2697.
- [32] Z. Han, J. Jiang, Y. Wang, Y. Ma, V. Tresp, The graph hawkes network for reasoning on temporal knowledge graphs, in: *Learning with Temporal Point Processes Workshop at the 33rd Conference on Neural Information Processing Systems (NeurIPS 2019)* NeurIPS 2019, 2019.
- [33] Y. Seo, M. Defferrard, P. Vandergheynst, X. Bresson, Structured sequence modeling with graph convolutional recurrent networks, in: *International Conference on Neural Information Processing*, Springer, 2018, pp. 362–373.

- [34] M. Defferrard, X. Bresson, P. Vandergheynst, Convolutional neural networks on graphs with fast localized spectral filtering, in: Proceedings of the 30th International Conference on Neural Information Processing Systems, 2016, pp. 3844–3852.
- [35] F. A. Gers, N. N. Schraudolph, J. Schmidhuber, Learning precise timing with lstm recurrent networks, *Journal of machine learning research* 3 (Aug) (2002) 115–143.
- [36] A. Narayan, P. H. Roe, Learning graph dynamics using deep neural networks, *IFAC-PapersOnLine* 51 (2) (2018) 433–438.
- [37] A. Taheri, K. Gimpel, T. Berger-Wolf, Learning to represent the evolution of dynamic graphs with recurrent models, in: Companion Proceedings of The 2019 World Wide Web Conference, 2019, pp. 301–307.
- [38] F. Manessi, A. Rozza, M. Manzo, Dynamic graph convolutional networks, *Pattern Recognition* 97 (2020) 107000.
- [39] A. Sankar, Y. Wu, L. Gou, W. Zhang, H. Yang, Dysat: Deep neural representation learning on dynamic graphs via self-attention networks, in: Proceedings of the 13th International Conference on Web Search and Data Mining, 2020, pp. 519–527.
- [40] A. Pareja, G. Domeniconi, J. Chen, T. Ma, T. Suzumura, H. Kanezashi, T. Kaler, T. Schardl, C. Leiserson, Evolvegcnn: Evolving graph convolutional networks for dynamic graphs, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 34, 2020, pp. 5363–5370.
- [41] J. Chen, X. Wang, X. Xu, Gc-lstm: Graph convolution embedded lstm for dynamic network link prediction, *Applied Intelligence* 52 (7) (2022) 7513–7528.
- [42] J. Li, Z. Han, H. Cheng, J. Su, P. Wang, J. Zhang, L. Pan, Predicting path failure in time-evolving graphs, in: Proceedings of the 25th ACM SIGKDD

- 680 International Conference on Knowledge Discovery & Data Mining, 2019,  
681 pp. 1279–1289.
- 682 [43] W. Jin, M. Qu, X. Jin, X. Ren, Recurrent event network: Autoregressive  
683 structure inference over temporal knowledge graphs, in: Proceedings of the  
684 2020 Conference on Empirical Methods in Natural Language Processing  
685 (EMNLP), 2020, pp. 6669–6683.
- 686 [44] P. Goyal, N. Kamra, X. He, Y. Liu, Dyngem: Deep embedding method for  
687 dynamic graphs, CoRR abs/1805.11273.  
688 URL <http://arxiv.org/abs/1805.11273>
- 689 [45] P. Goyal, S. R. Chhetri, A. Canedo, dyngraph2vec: Capturing network  
690 dynamics using dynamic graph representation learning, Knowledge-Based  
691 Systems 187 (2020) 104816.
- 692 [46] J. Chen, J. Zhang, X. Xu, C. Fu, D. Zhang, Q. Zhang, Q. Xuan, E-lstm-  
693 d: A deep learning framework for dynamic network link prediction, IEEE  
694 Transactions on Systems, Man, and Cybernetics: Systems 51 (6) (2019)  
695 3699–3712.
- 696 [47] S. Pan, R. Hu, G. Long, J. Jiang, L. Yao, C. Zhang, Adversarially regular-  
697 ized graph autoencoder for graph embedding, in: Proceedings of the 27th  
698 International Joint Conference on Artificial Intelligence, 2018, pp. 2609–  
699 2615.
- 700 [48] D. Charte, F. Charte, M. J. del Jesus, F. Herrera, An analysis on the use of  
701 autoencoders for representation learning: Fundamentals, learning task case  
702 studies, explainability and challenges, Neurocomputing 404 (2020) 93–107.
- 703 [49] D. J. Rezende, S. Mohamed, D. Wierstra, Stochastic backpropagation and  
704 approximate inference in deep generative models, in: International confer-  
705 ence on machine learning, PMLR, 2014, pp. 1278–1286.
- 706 [50] K. Lei, M. Qin, B. Bai, G. Zhang, M. Yang, Gcn-gan: A non-linear tempo-  
707 ral link prediction model for weighted dynamic networks, in: IEEE INFO-

- 708 COM 2019-IEEE Conference on Computer Communications, IEEE, 2019,  
709 pp. 388–396.
- 710 [51] E. Hajiramezanali, A. Hasanzadeh, N. Duffield, K. Narayanan, M. Zhou,  
711 X. Qian, Variational graph recurrent neural networks, in: Proceedings of  
712 the 33rd International Conference on Neural Information Processing Sys-  
713 tems, 2019, pp. 10701–10711.
- 714 [52] D. Bahdanau, K. Cho, Y. Bengio, Neural machine translation by jointly  
715 learning to align and translate, ICLR (2015) 1–15.
- 716 [53] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez,  
717 L. Kaiser, I. Polosukhin, Attention is all you need, in: Proceedings of the  
718 31st International Conference on Neural Information Processing Systems,  
719 2017, pp. 6000–6010.
- 720 [54] T. Shen, J. Jiang, T. Zhou, S. Pan, G. Long, C. Zhang, Disan: directional  
721 self-attention network for rnn/cnn-free language understanding, in: Pro-  
722 ceedings of the Thirty-Second AAAI Conference on Artificial Intelligence  
723 and Thirtieth Innovative Applications of Artificial Intelligence Conference  
724 and Eighth AAAI Symposium on Educational Advances in Artificial Intel-  
725 ligence, 2018, pp. 5446–5455.
- 726 [55] Z. Tan, M. Wang, J. Xie, Y. Chen, X. Shi, Deep semantic role labeling with  
727 self-attention, in: Proceedings of the Thirty-Second AAAI Conference on  
728 Artificial Intelligence and Thirtieth Innovative Applications of Artificial  
729 Intelligence Conference and Eighth AAAI Symposium on Educational Ad-  
730 vances in Artificial Intelligence, 2018, pp. 4929–4936.
- 731 [56] J. B. Tenenbaum, V. De Silva, J. C. Langford, A global geometric frame-  
732 work for nonlinear dimensionality reduction, *science* 290 (5500) (2000)  
733 2319–2323.
- 734 [57] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, Y. Bengio,

- 735 Graph attention networks, in: International Conference on Learning Rep-  
736 resentations, 2018, pp. 1–18.
- 737 [58] B. Viswanath, A. Mislove, M. Cha, K. P. Gummadi, On the evolution of  
738 user interaction in facebook, in: Proceedings of the 2nd ACM workshop on  
739 Online social networks, 2009, pp. 37–42.
- 740 [59] K. S. Xu, A. O. Hero, Dynamic stochastic blockmodels for time-evolving  
741 social networks, IEEE Journal of Selected Topics in Signal Processing 8 (4)  
742 (2014) 552–562.
- 743 [60] M. Rahman, M. Al Hasan, Link prediction in dynamic networks using  
744 graphlet, in: Joint European Conference on Machine Learning and Knowl-  
745 edge Discovery in Databases, Springer, 2016, pp. 394–409.
- 746 [61] C. E. Priebe, J. M. Conroy, D. J. Marchette, Y. Park, Scan statistics on  
747 enron graphs, Computational & Mathematical Organization Theory 11 (3)  
748 (2005) 229–247.
- 749 [62] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, T. Eliassi-Rad,  
750 Collective classification in network data, AI magazine 29 (3) (2008) 93–93.
- 751 [63] X. Liu, P.-C. Hsieh, N. Duffield, R. Chen, M. Xie, X. Wen, Real-time  
752 streaming graph embedding through local actions, in: Companion Pro-  
753 ceedings of The 2019 World Wide Web Conference, 2019, pp. 285–293.
- 754 [64] J. Leskovec, J. Kleinberg, C. Faloutsos, Graphs over time: densification  
755 laws, shrinking diameters and possible explanations, in: Proceedings of the  
756 eleventh ACM SIGKDD international conference on Knowledge discovery  
757 in data mining, 2005, pp. 177–187.
- 758 [65] M. Yang, M. Zhou, M. Kalander, Z. Huang, I. King, Discrete-time tem-  
759 poral network embedding via implicit hierarchical learning in hyperbolic  
760 space, in: Proceedings of the 27th ACM SIGKDD Conference on Knowl-  
761 edge Discovery & Data Mining, 2021, pp. 1975–1985.

- [66] Y. Wang, Y.-Y. Chang, Y. Liu, J. Leskovec, P. Li, Inductive representation learning in temporal networks via causal anonymous walks, in: International Conference on Learning Representations, 2020.
- [67] D. Xu, C. Ruan, E. Korpeoglu, S. Kumar, K. Achan, Inductive representation learning on temporal graphs, ICLR (2020) 1–12.
- [68] X. Glorot, Y. Bengio, Understanding the difficulty of training deep feed-forward neural networks, in: Proceedings of the thirteenth international conference on artificial intelligence and statistics, 2010, pp. 249–256.
- [69] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, in: ICLR (Poster), 2015, pp. 1–15.
- [70] J. Gao, B. Ribeiro, On the equivalence between temporal and static equivariant graph representations, in: International Conference on Machine Learning, PMLR, 2022, pp. 7052–7076.
- [71] W. L. Hamilton, R. Ying, J. Leskovec, Inductive representation learning on large graphs, in: Proceedings of the 31st International Conference on Neural Information Processing Systems, 2017, pp. 1025–1035.
- [72] Y. Yao, C. Joe-Wong, Interpretable clustering on dynamic graphs with recurrent graph neural networks, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 35, 2021, pp. 4608–4616.
- [73] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word representations in vector space, in: Proceedings of Workshop at ICLR, 2013.

## 784 Appendix A. Visualisation the Results of Comparison

785 In order to more thoroughly assess the performance of the proposed method,  
 786 DyVGRNN and VGRNN are compared in Figure A.10. It is evident that  
 787 DyVRNN performs better over time in practically all epochs. Despite hav-  
 788 ing close competition in the early epochs, DyVGRNN quickly passes VGRNN  
 789 and establishes its superiority.

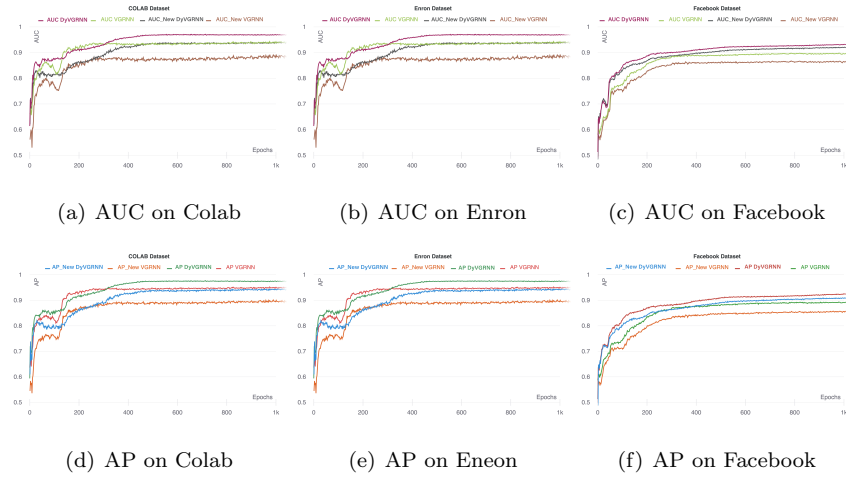


Figure A.10: Comparing the proposed method with VGRNN on different datasets in terms of AUC and AP. The colours reflect the various criteria for dynamic link prediction and dynamic new link prediction under the colour scheme. a) The comparison of two methods in terms of AUC on Colab. The early epochs are closely contested, but after epoch 300, DyVGRNN soon overtakes VGRNN. b) The comparison of two methods in terms of AP on Colab. The superiority of DyVGRNN is significant after epoch 300. c) The comparison of two methods in terms of AUC on Enron. After epoch 300, DyVGRNN's dominance becomes considerable. d) The comparison of two methods in terms of AP on Enron. Again, in the 300th period, DyVGRNN's advantage becomes substantial. e) The comparison of two methods in terms of AUC on Facebook. Even in the early epochs, DyVGRNN's supremacy was noticeable. f) The comparison of two methods in terms of AP on Facebook. From the very beginning, DyVGRNN's dominance is significant.

## 790 Appendix B. Qualitative Analysis

791 Figure B.11 displays a visualisation of the learnt embeddings over time to  
 792 show how effectively the embeddings are encoded. To do so, we use the clustering  
 793 task and the silhouette metric on synthetic data and visualised the learnt embed-  
 794 dings in a two-dimensional space throughout our training. The clusters become  
 795 more well-separated with time, as can be observed.

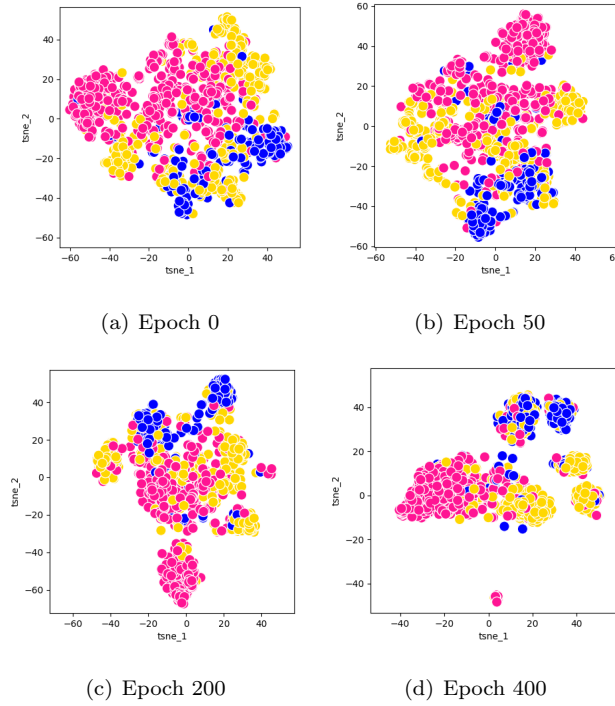


Figure B.11: Visualisation of the learnt embeddings of DyVGRNN over time. In this figure, each colour corresponds to a cluster. a) Visualisation of embedding on epoch 0 of the running. The clusters are confused. b) Visualisation of embedding on epoch 50 of the running. Clusters are hardly distinguishable. c) Visualisation of embedding on epoch 200 of the running. Clusters show themselves, but they are still intertwined. d) Visualisation of embedding on epoch 400 of the running. Clusters are almost easily distinguishable.



## 796 Appendix C. Node Classification Task

797 We compare our model to three baseline methods in order to assess its perfor-  
 798 mance on the classification task. Two of these methods, GCN [24] and Graph-  
 799 SAGE [71], are supervised techniques that relied solely on static graph struc-  
 800 tures and node attributes, ignoring temporal information. Another method,  
 801 RNNGCN [72], utilised a two-layer GCN with a decay weight as a learnable  
 802 parameter. This decay weight has applied to information from each timestep,  
 803 gradually decreasing over time. The resulting linear combination of information  
 over time is then used for classification purposes.

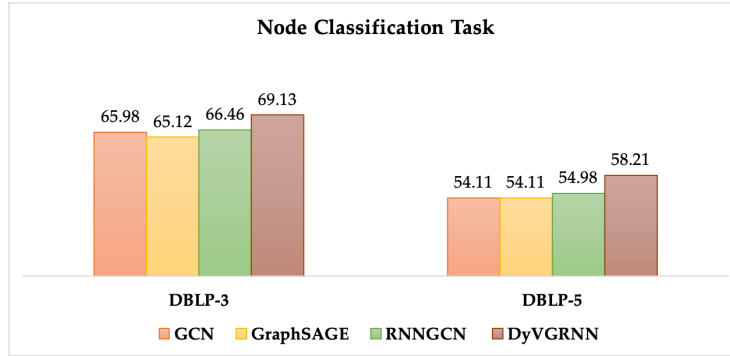


Figure C.12: The results of comparing the classification performance of the proposed method on DBLP-3 and DBLP-5 datasets with other baselines in terms of AUC. The colours, in accordance with the colour scheme, represent the various methods.

804

805 The datasets used in this task has obtained from DBLP<sup>2</sup>, a comprehen-  
 806 sive database of academic papers in various subfields of computer science. The  
 807 authors of these papers are represented as nodes in a graph, with connections be-  
 808 tween nodes indicating co-authorship. Analysing the authorship of papers pub-  
 809 lished between 2005 and 2018 resulted in the dynamic graph in these datasets,  
 810 treating each year as a snapshot. DBLP-5 has 6606 nodes, 42815 edges, and 10  
 811 snapshots, while DBLP-3 has 4257 nodes, 23540 edges, and 10 snapshots. These  
 812 datasets included node attributes extracted by word2vec [73] from authors' pa-

---

<sup>2</sup><https://dblp.org/>

813 per titles and abstracts. They both have 100 attributes. These datasets are  
814 further clustered into three and five classes, respectively, based on the research  
815 area of the authors. These classes remained static over time. Figure C.12 shows  
816 the results of our comparison in terms of the AUC. As seen, our proposed DyV-  
817 GRNN outperforms other methods in both datasets.