

---

# FDGNN: FULLY DYNAMIC GRAPH NEURAL NETWORK

---

A PREPRINT

Alice Moallem-Oureh, Silvia Beddar-Wiesing, Rüdiger Nather, Josephine M. Thomas

Department of Computer Science

Intelligent Embedded Systems

University of Kassel

34121 Kassel

{amoallem, s.beddarwiesing, rnather, jthomas}@uni-kassel.de

<https://gain-group.de/>

June 8, 2022

## ABSTRACT

Dynamic Graph Neural Networks recently became more and more important as graphs from many scientific fields, ranging from mathematics, biology, social sciences, and physics to computer science, are dynamic by nature. While temporal changes (dynamics) play an essential role in many real-world applications, most of the models in the literature on Graph Neural Networks (GNN) process static graphs. The few GNN models on dynamic graphs only consider exceptional cases of dynamics, e.g., node attribute-dynamic graphs or structure-dynamic graphs limited to additions or changes to the graph’s edges, etc. Therefore, we present a novel Fully Dynamic Graph Neural Network (FDGNN) that can handle fully-dynamic graphs in continuous time. The proposed method provides a node and an edge embedding that includes their activity to address added and deleted nodes or edges, and possible attributes. Furthermore, the embeddings specify Temporal Point Processes for each event to encode the distributions of the structure- and attribute-related incoming graph events. In addition, our model can be updated efficiently by considering single events for local retraining.

**Keywords** Dynamic Graph Neural Network · Dynamic Graphs · Temporal Point Process · Graph Attention · Dynamic Graph Embedding

## 1 Introduction

Graph Neural Networks (GNNs) have become baseline models for learning on structured data in the past years. Many models have been developed for static graph data that can include structural information into the learning process for different problems. However, they cannot capture the evolution of a dynamic graph [1], although many real-world problems require the representations of *dynamic* changes in the graph data. Furthermore, the GNNs that can handle dynamics are restricted to specific graph types. More specifically, they can either handle node attribute changes [2], are restricted to additions of nodes or edges (growing graphs) [3, 4], or can only handle edge-structure dynamics and node attribute changes [5]. Further examples are provided in [6, 7].

To address the missing possibility of handling deletions of nodes and edges and changing attributes, we introduce a Fully Dynamic Graph Neural Network (FDGNN) that can handle node and edge additions and deletions and node and edge attribute changes (dynamic events). The approach of FDGNN is based on the model DyREP [4], which only handles edge-growing graphs. Thus, FDGNN is an extension of DyREP, including the possibility to process deletions of nodes and edges and node/edge attribute changes. FDGNN uses a node and edge activation function that captures the existence of all the nodes/edges in each timestamp. Furthermore, the node/edge embeddings use two attention mechanisms (self- and neighborhood attention) to gather historical and structural information and application-specific node/edge attribute embedding functions. A Temporal Point Process (TPP) [8] is used for every dynamic event type

that encodes the graph’s dynamics to enable a continuous-time representation. With every incoming event, the model is updateable by local retraining, which makes the training procedure efficient.

In this paper, the theoretical development of the FDGNN is proposed. Since the implementation is still in progress, the experimental results and further optimizations will be available in a second updated version of the paper. The paper is structured as follows: In §2 the preliminaries are introduced that cover the main definitions used in the paper and specify which modules are used to incorporate the structure dynamics with activity functions (§2.2), the attribute dynamics (§2.3) and the continuous-time representation via Temporal Point Processes (§2.4). For better readability we have included our conventions of notation in §2.5. In §3 the architecture of the proposed FDGNN is introduced which includes the attention mechanisms, the intensity functions for each event type, and the final embedding functions. The learning procedure for the model is based on [4] and adapted to the model in §4. The future work and further discussions can be found in §5. Section §6 concludes the paper.

## 2 Preliminaries

The proposed FDGNN is constructed to process dynamic graphs in continuous-time representation. For this purpose, it is required to define graph streams, where dynamic events operate on static graph snapshots. The advantage of such continuous-time representation is the ability to update the GNN model, i.e., the model gets updated with every new incoming event, and, therefore, no retraining from scratch is required.

### 2.1 Graph Stream and Graph Snapshot

Given that a static attributed graph is represented as a tuple  $(\mathcal{V}, \mathcal{E}, \alpha, \beta)$  where  $\mathcal{V}$  is a set of nodes and  $\mathcal{E}$  is a set of edges with attribute functions  $\alpha : \mathcal{V} \rightarrow \mathcal{A}$ ,  $\beta : \mathcal{E} \rightarrow \mathcal{B}$ , the following terms are defined accordingly.

Let  $G = (\mathcal{G}_0 := (\mathcal{V}_0, \mathcal{E}_0, \alpha_0, \beta_0), \mathcal{O})$  be a dynamic graph in continuous-time representation (**graph stream**), with static start graph  $\mathcal{G}_0$  and  $o_k \in \mathcal{O}$  being events/observations of the following form:

Event Type	Node Event	Edge Event
Addition	$o_0 := (v, \alpha_t^v, t)_{\text{add}}$	$o_2 := (e, \beta_t^e, t)_{\text{add}}$
Deletion	$o_1 := (v, \alpha_t^v, t)_{\text{del}}$	$o_3 := (e, \beta_t^e, t)_{\text{del}}$
Attribute Change	$o_4 := (v, \alpha_t^v, t)_{\text{atr}}$	$o_5 := (e, \beta_t^e, t)_{\text{atr}}$

Table 1: Events in a continuous-time representation of a graph.

where  $v \in \mathcal{V}_{\leq t} := \bigcup_{\hat{t} \in [t_0, t]} \mathcal{V}_{\hat{t}}$  are nodes appearing in  $G$  up to timestamp  $t$ , and  $e \in \mathcal{E}_{\leq t} := \bigcup_{\hat{t} \in [t_0, t]} \mathcal{E}_{\hat{t}}$  are edges appearing in  $G$  up to timestamp  $t$ . Furthermore,  $\alpha_t : \mathcal{V}_t \rightarrow \mathcal{A}$ ,  $\beta_t : \mathcal{E}_t \rightarrow \mathcal{B}$  are node and edge attribute functions where  $\mathcal{A}, \mathcal{B}$  are arbitrary attribute sets, and  $t \in \mathcal{T} \subseteq \mathbb{R}_{\geq 0}$  are timestamps (abbreviated as  $\alpha_t^v := \alpha_t(v)$  in the event tuples).

A **graph snapshot** or graph slice  $\mathcal{G}_t = (\mathcal{V}_t, \mathcal{E}_t, \alpha_t, \beta_t)$  is defined to be the graph that results from  $\mathcal{G}_0$  incorporating all the observations within  $[t_0, t]$ .

### 2.2 Node and Edge Activities

Given that FDGNN should deal with structure-dynamics, the main idea for the node and edge activity function is to capture which nodes and edges exist at any given time. To not lose the information of deleted nodes or edges, the output of the activity function for a freshly deleted object is 0. The deletion history is important for the model to learn the deletion behavior of the graph.

Let for any timestamp  $t \in \mathcal{T}$  and any node or edge  $x \in \mathcal{V} \cup \mathcal{E}$  the timestamp  $\bar{t}$  refer to the timestamp of an event affecting  $x$  right before  $t$ . Furthermore, let  $o_{\bar{t}}(x)$  be an event at time  $\bar{t}$  that affects the node or edge  $x \in \mathcal{V}_{\bar{t}} \cup \mathcal{E}_{\bar{t}}$ , and  $\xi : (\mathcal{V} \cup \mathcal{E}) \times \mathcal{T} \rightarrow \{0, 1\}$  be the **activity function** of nodes and edges at time  $t > \bar{t}$  defined as

$$\xi_t^x = \begin{cases} 1, & \text{if } o_{\bar{t}}(x) \in \{o_0, o_2, o_4, o_5\} \quad \text{addition and attribute change events} \\ 0, & \text{if } o_{\bar{t}}(x) \in \{o_1, o_3\} \quad \text{deletion events (inactivity)} \end{cases}$$

for all  $x \in \mathcal{V}_{\bar{t}} \cup \mathcal{E}_{\bar{t}}$ .

The stati of the nodes or edges indicated by the activity are used in the proposed model to determine in which situations a specific observation can be made, e.g., just existing nodes can be deleted.

### 2.3 Attribute Embeddings

To incorporate attribute dynamics the idea is to first consider a "node/edge-independent" attribute embedding as formalized in Eq. (1). Among other modules introduced later in this paper, the attribute embeddings are used to compute the final node and edge embeddings separately (cf. Eq. (7), Eq. (8)).

The **attribute embedding** encodes the dynamic behaviour of the node or edge attribute respectively by considering the historical transformed embedding vector  $\mathbf{u}_{\bar{t}}(\gamma) \in \mathbb{R}^{\bar{u}}$  and the current embedding vector  $\mathbf{Z}_t^\gamma \in \mathbb{R}^{\bar{z}}$ . For attribute  $\gamma \in \mathcal{A} \cup \mathcal{B}$  at timestamp  $t \in \mathcal{T}$ , the embedding  $\mathbf{u}_t(\gamma) \in \mathbb{R}^{\bar{u}}$  in a  $\bar{u}$ -dimensional space is obtained recursively w.r.t. time by

$$\mathbf{u}_t(\gamma) = \mathbf{M}_0 \cdot \begin{pmatrix} \mathbf{u}_{\bar{t}}(\gamma) \\ \mathbf{Z}_t^\gamma \end{pmatrix} + \mathbf{b}_u. \quad (1)$$

The parameters  $\mathbf{M}_0 \in \mathbb{R}^{\bar{u} \times (\bar{u} + \bar{z})}$  and  $\mathbf{b}_0 \in \mathbb{R}^{\bar{u}}$  are the learnable weight matrix and bias in the FDGNN. The current embedding vector  $\mathbf{Z}_t^\gamma$  is the output from application-specific chosen attribute embedding functions, e.g., CNN's for images, text2vec for text [9] etc.

### 2.4 Temporal Point Process

A set of observations is given in form of a temporal point pattern. Given the set of observations we cannot tell at what time and in which amount an observation will occur. Especially when it comes to real-world applications, the occurrence of the different events can follow complex rules and are not necessarily fully random (e.g., an already deleted node cannot be deleted again, or a non-existing edge cannot change its attributes). Assuming that an observation is dependent on the historical point patterns, the underlying process can be modelled by a **temporal point process (TPP)** [8]. The term point implies that events only occur at a certain timestamp and thus pose a data point in a time series. Furthermore, there may be additional information given for the events, as for example the event types  $k \in \{0, \dots, 5\}$  in the definition above, which are called marks in TPP literature. Therefore, only marked TPPs are considered in this work and simply denoted as TPPs. In addition, we assume here that at each timestamp only one event can occur, which makes the TPP *simple*. The consideration of simultaneously occurring events is left for future work.

To model a (marked) temporal point process and to specify the dependency of a present observation at time  $t$  on the history of events  $\mathbf{H}_{t_n}$  for any  $t > t_n$ , a **conditional intensity function** can be used. Let  $f(t|\mathbf{H}_{t_n})$  be the conditional density function and  $F(t|\mathbf{H}_{t_n})$  be its corresponding cumulative distribution function for any  $t > t_n$ .

For the marked case, let  $f(k|t) = f(k|t, \mathbf{H}_{\bar{t}})$  be the conditional density function<sup>1</sup> of event type  $k$ , i.e., specifying the distribution of the event type  $k$  given  $t$  and the history  $\mathbf{H}_{\bar{t}}$ , which now includes information of both times and marks of past events. Then, the **(marked) conditional intensity function** is defined in [8] by

$$\lambda(k, t) = \lambda(t)f(k|t) = \frac{f(t, k|\mathbf{H}_{t_n})}{1 - F(t|\mathbf{H}_{t_n})}, \quad (2)$$

where  $f(t, k|\mathbf{H}_{t_n})$  is the joint density of  $t$  and the  $k$ , and  $F(t|\mathbf{H}_{t_n})$  is the conditional cumulative distribution function of  $t$ , both conditional on the past times and event types.  $\lambda(t)$  corresponds to the unmarked conditional intensity function.

Assuming that no events can occur simultaneously, the marked conditional intensity function can be heuristically interpreted as the expected number of events occurring at a timestamp given the historical information:

$$\lambda(t, k)dt = \mathbb{E}[N([t, t + \delta t] \times k)|\mathbf{H}_t], \quad (3)$$

with a function  $N([t, t + \delta t] \times k)$  giving the number of events of type  $k$  occurring in the infinitesimal interval  $[t, t + \delta t]$  with a time difference  $\delta t > 0$ . Here, the events are discrete and contain events such as addition/deletion of nodes/edges in 1.

Note that, for events where node or edge attributes are changing ( $k = 4, 5$ ), two tasks are possible. On the one hand, considering Eq. (3) the events cover if an attribute is changing, or not. If the attribute sets  $\mathcal{A}$  and  $\mathcal{B}$  are metric spaces the following approach allows modelling the change of an attribute as mark  $m$  during an event  $o_4, o_5$  by

$$\lambda(t, m)dt = \mathbb{E}[N([t, t + \delta t] \times B_{\delta m}(m))|\mathbf{H}_t],$$

where  $[t, t + \delta t]$  is as above and  $B_{\delta m}(m)$  is an open ball of radius  $\delta m$  around  $m$  in the respective attribute spaces.

<sup>1</sup>In general, if  $k$  is a continuous random variable, this is the usual (conditional) density function, but if it is a discrete random variable, this is its (conditional) probability function.

## 2.5 Notation

The notation used throughout this work is listed in the following.

$\mathbb{N}_0$	natural numbers starting at 0
$\mathbb{R}_{\geq 0}$	non-negative real numbers
$\mathbb{R}^k$	$\mathbb{R}$ vector space of dimension $k$
$ M $	number of elements of a set $M$
$M^\top$	matrix $M$ transposed
$[i, j]$	closed interval from $i$ to $j$
$\subseteq$	subset
$\times$	factor set of two sets
$\sigma$	sigmoid activation function

In general, we use bold capital letters for matrices, such as  $\mathbf{A}$ , bold small type letters for vectors such as  $\mathbf{a}$  and small type normal letters for scalars or nodes or vectors such as  $a$ .

## 3 Model

In this section the individual parts of the model are introduced, i.e., the self attention and neighborhood attention modules (§3.1), the embedding function for attributes, nodes and edges (§3.2) and the intensity functions for each event-type (add/delete node/edge and change node/edge attribute, cf. §3.3).

In what follows, let for any timestamp  $t \in \mathcal{T}$  and any node or edge  $x \in \mathcal{V} \cup \mathcal{E}$  the timestamp  $\bar{t}$  refer to the timestamp of an event affecting  $x$  right before  $t$ .

### 3.1 Attention Mechanisms

There are two attention mechanisms defined in the following: self-attention and neighborhood attention. These attentions are used for the embedding computation in Section 3.2 via weighting the historical information of a node/edge and the neighborhood information in the update, and therefore to improve the explainability of the model.

**Self-Attention** The self-attention takes for each node/edge its historical behavior into account by using the activity status of the current timestamp  $t$  and the one right beforehand  $\bar{t}$ . Furthermore, a similarity between the attributes of the two timestamps is considered. Moreover, the attention is normalized by the attentions of all known timestamps before.

The self-attention is defined for all  $x \in (\mathcal{V}_t \cap \mathcal{V}_{\bar{t}}) \cup (\mathcal{E}_t \cap \mathcal{E}_{\bar{t}})$  as

$$p_t^x = \frac{\exp(\xi_t^x \cdot \xi_{\bar{t}}^x \cdot \bar{\tau}_t^x)}{\sum_{t_i \in \mathcal{T}} \exp(\xi_{t_i}^x \cdot \xi_{\bar{t}_i}^x \cdot \bar{\tau}_{t_i}^x)} \quad (4)$$

using a similarity function between the two attribute embeddings  $\mathbf{u}$  from Eq. (1) given by

$$\bar{\tau}_t^x = \begin{cases} \text{sim}(\mathbf{u}_t(\alpha_t^x), \mathbf{u}_{\bar{t}}(\alpha_{\bar{t}}^x)), & \text{if } x \in \mathcal{V}_t \cap \mathcal{V}_{\bar{t}} \\ \text{sim}(\mathbf{u}_t(\beta_t^x), \mathbf{u}_{\bar{t}}(\beta_{\bar{t}}^x)), & \text{if } x \in \mathcal{E}_t \cap \mathcal{E}_{\bar{t}}. \end{cases} \quad (5)$$

for  $x$  being a node or an edge respectively. The similarity function  $\text{sim}$  can be an arbitrary similarity measure, as, e.g., the cosine similarity used in [10].

**Neighborhood Attention** The neighborhood attention takes for each node/edge the activity of its neighborhood from the previous timestamp into account. It respects the active nodes and edges within its direct neighborhood and the corresponding edge attributes. With a similarity measure used on edge attributes, the focus is set on similar attributes.

Let  $u_i \in \mathcal{N}(v) := \{u \mid \exists e \in \mathcal{E} : u, v \in e\}$  for  $i = 0, \dots, |\mathcal{N}(v)|$  be the neighbors of node  $v$ . The attention for a node  $v$  and a neighbor  $u_i$  is determined by their activities and the activity of the connecting edge, together with a similarity

measure between the two node attribute embeddings<sup>2</sup>, as in (5). Then, the attention is normalized by the attentions of all neighbors of  $v$  and is determined by

$$q_t^{\{v, u_i\}} = \frac{\exp(\xi_t^v \cdot \xi_t^{u_i} \cdot \xi_t^{\{u_i, v\}} \cdot \tau_t^{\{u_i, v\}})}{\sum_{m \in \mathcal{N}(v)} \exp(\xi_t^m \cdot \xi_t^v \cdot \xi_t^{\{m, v\}} \cdot \tau_t^{\{m, v\}})}. \quad (6)$$

### 3.2 Embeddings

To incorporate attribute dynamics the idea is to first consider a "node/edge-independent" attribute embedding (cf. Eq. (1)) and compute the node/edge embeddings (cf. Eq. (7), Eq. (8)) using the attribute embeddings, the local neighborhood information and attentions according to Eq. (6) and Eq. (4). In Section §3.3 the intensity functions of the node and edge embeddings are processed to obtain models for the temporal behavior of the graph events.

**Node Embeddings** The node embedding consists of four additive components. The first is the local embedding propagation scheme that cumulates the neighborhood node embeddings and attributes embeddings. The second summand describes a self-propagation component that propagates the node's embedding from the timestamp before being multiplied by a self-attention factor. In the third component, the exogenous drive includes the time difference between the timestamps of the current event and the one beforehand that takes temporal dependencies into account. Moreover, the last component includes the attribute embedding of the node. Summarized, this gives the node embedding determined by

$$\mathbf{Z}^v(t) = \sigma \left( \underbrace{\mathbf{M}_1 \mathbf{h}_{loc}(v, \bar{t})}_{\text{loc. embd. prop.}} + \underbrace{\mathbf{M}_2 \mathbf{Z}_t^v \cdot p_t^v}_{\text{self-prop.}} + \underbrace{\mathbf{m}_3 \cdot (t - \bar{t})}_{\text{exogenous drive}} + \underbrace{\mathbf{M}_4 \cdot \mathbf{u}_t(\alpha_t^v)}_{\text{attribute prop.}} \right). \quad (7)$$

The used local embedding propagation scheme is given by

$$\mathbf{h}_{loc}(v, \bar{t}) = \sum_{u \in \mathcal{N}(v)} \left( \underbrace{m_5^u \mathbf{Z}_t^u}_{\text{neighborhood cum.}} + \underbrace{m_6^u \mathbf{u}_t(\alpha_t^u)}_{\text{neighborhood attr. cum.}} \right) \cdot q_t^{\{v, u\}},$$

where  $m_5^u, m_6^u \in \mathbb{R}$  are learnable weights and  $p_t^v, q_t^{\{v, u\}}$  are the attentions according to Eq. (4) and (6).

**Edge Embeddings** The edge embedding is similar to the node embedding. The only substantial difference is the choice of the local embedding propagation scheme, which includes the embeddings of the incident nodes, their attributes, and the embedding of the edge attribute. In total, the edge embedding is defined by

$$\mathbf{Z}^{u,v}(t) = \sigma \left( \underbrace{\mathbf{M}_7 \mathbf{h}_{loc}(u, v, \bar{t})}_{\text{loc. embd. prop.}} + \underbrace{\mathbf{M}_8 \mathbf{Z}_t^{u,v} \cdot p_t^{\{u, v\}}}_{\text{self-prop.}} + \underbrace{\mathbf{m}_9 \cdot (t - \bar{t})}_{\text{exogenous drive}} + \underbrace{\mathbf{M}_{10} \cdot \mathbf{u}_t(\beta_t^{u,v})}_{\text{attribute prop.}} \right) \quad (8)$$

with the local embedding propagation scheme

$$\mathbf{h}_{loc}(u, v, \bar{t}) = \underbrace{m_{11} \mathbf{Z}_t^u p_t^u + m_{12} \mathbf{Z}_t^v p_t^v}_{\text{incident nodes cum.}} + \underbrace{m_{13} \mathbf{u}_t(\alpha_t^u) p_t^u + m_{14} \mathbf{u}_t(\alpha_t^v) p_t^v}_{\text{inc. node attr. cum.}},$$

where  $m_{11}, m_{12}, m_{13}, m_{14} \in \mathbb{R}$  are learnable weights and  $p_t^u, p_t^v$  are the attentions according to Eq. (4).

**Remark** In case that the activity function of a node or edge  $x$  is undefined up to a timestamp  $t$ , the initial embedding only includes the attribute propagation and the local embedding propagation at time  $t$ , i.e.,

$$\mathbf{Z}^x(t) = \sigma \left( \underbrace{\mathbf{M}' \mathbf{h}_{loc}(x, t)}_{\text{loc. embd. prop.}} + \underbrace{\mathbf{M}'' \cdot \mathbf{u}_t(\gamma_t^x)}_{\text{attribute prop.}} \right).$$

<sup>2</sup>In case of knowledge about the edge attributes an additional similarity factor can be added that respects the edge attribute similarity between the incident edges of  $v$ .

### 3.3 Intensity Functions/Point Processes of the Events

This section defines the intensity functions that model the temporal behavior of the dynamic events occurring on the graph. For this purpose, consider the event types  $o_k \in \mathcal{O}$  defined in Tab. (1). Such events can be interpreted as Temporal Point Processes (TPP) that are formalized through an intensity function in the following. For each event type  $o_k$  there is a separate intensity function  $\lambda_k^x(t)$  in use for  $x \in \mathcal{V}_t \cup \mathcal{E}_t$  dependent on a timestamp  $t$ . The intensity functions are determined for active nodes and edges, while they are constant 0 if the event can not occur. Thus, each intensity function models the number of events of a certain type occurring at time  $t$  as a function of the graph status at the timestamp  $\bar{t}$  of the event right before an event at time  $t$ .

Let  $f_k$  be an activation function,  $\mathbf{Z}_t^x$  an attribute embedding of a node or edge  $x$  at time  $t$  and the matrices  $\mathbf{W}_\star$  be learnable weights. Then, the intensity functions are defined as follows.

**Add Node** Additions of nodes at timestamp  $t$  are possible if and only if the nodes did not exist right before at timestamp  $\bar{t}$ . Therefore, the activity of an existing node is 1 and thus the intensity function of node additions is 0, i.e.,  $\xi_t^v = 1 \iff \lambda_0^v(t) = 0$ . Otherwise, the intensity function uses a function of the form  $\mathbf{W}_0^\top \mathbf{Z}_t^v$  that assigns a score to a node attribute considering the node attribute embedding. Hence,  $\lambda_0^v(t)$  is defined as

$$\lambda_0^v(t) = \begin{cases} f_0(\underbrace{\mathbf{W}_0^\top \mathbf{Z}_t^v}_{\text{scores attribute}}), & \text{if } (\xi_t^v = 0 \vee \xi_t^v = \perp) \\ 0, & \text{if } (\xi_t^v = 1). \end{cases}$$

**Delete Node** For the opposite event, delete node, it is important to switch the cases, i.e., a node can just be deleted at timestamp  $t$  if and only if it existed before in timestamp  $\bar{t}$ . The corresponding intensity function is given as

$$\lambda_1^v(t) = \begin{cases} f_1(\underbrace{\mathbf{W}_1^\top \mathbf{Z}_t^v}_{\text{scores attribute}}), & \text{if } (\xi_t^v = 1) \\ 0, & \text{if } (\xi_t^v = 0). \end{cases}$$

**Add Edge** Analogously to the addition of nodes, an edge  $e$  can only be added when it did not exist the timestamp before and, additionally, both nodes  $u, v \in e$  exist. The inner function of the intensity scores the node and edge attributes, respectively.

$$\lambda_2^e(t) = \begin{cases} f_2\left(\underbrace{\mathbf{W}_2^\top \begin{pmatrix} \mathbf{Z}_t^u \\ \mathbf{Z}_t^v \\ \mathbf{Z}_t^e \end{pmatrix}}_{\text{scores node \& edge attributes}}\right), & \text{if } (\xi_t^e = 0 \wedge \xi_t^u = 1 \wedge \xi_t^v = 1 \vee \xi_t^e = \perp) \\ 0, & \text{if } (\xi_t^e = 1 \vee \xi_t^u = 0 \vee \xi_t^v = 0) \end{cases}$$

**Delete Edge** The deletion of an edge  $e$  with nodes  $u, v \in e$  will definitely take place, when the edge existed at time  $\bar{t}$ , and at least one of the incident nodes has been deleted. Thus, when the edge does not exist at time  $\bar{t}$ , it cannot be deleted, i.e.,  $\lambda_3^e(t) = 1 \iff \xi_t^{(u,v)} = 1 \wedge (\xi_t^u = 0 \vee \xi_t^v = 0)$ . Otherwise, the inner function of the intensity again scores the node and edge attributes. All in all, the intensity function to delete edges  $\lambda_3$  is defined by

$$\lambda_3^e(t) = \begin{cases} f_3\left(\underbrace{\mathbf{W}_3^\top \begin{pmatrix} \mathbf{Z}_t^u \\ \mathbf{Z}_t^v \\ \mathbf{Z}_t^e \end{pmatrix}}_{\text{scores node \& edge attributes}}\right), & \text{if } (\xi_t^e = 1 \wedge \xi_t^u = 1 \wedge \xi_t^v = 1) \\ 1, & \text{if } (\xi_t^e = 1 \wedge (\xi_t^u = 0 \vee \xi_t^v = 0)) \\ 0, & \text{if } (\xi_t^e = 0). \end{cases}$$

**Change Node Attributes** The attributes of nodes can only change if and only if the node whose attributes are in consideration exists. In this case a scoring  $\mathbf{W}_4^\top \begin{pmatrix} \mathbf{Z}_t^v \\ \mathbf{Z}_t^v \end{pmatrix}$  is needed that takes the node’s attribute change into account.

$$\lambda_4^v(t) = \begin{cases} \underbrace{f_4 \left( \mathbf{W}_4^\top \begin{pmatrix} \mathbf{Z}_t^v \\ \mathbf{Z}_t^v \end{pmatrix} \right)}_{\text{scores attribute change}}, & \text{if } (\xi_t^v = 1) \\ 0, & \text{if } (\xi_t^v = 0). \end{cases}$$

**Change Edge Attributes** For a node change event, the edge attribute change requires the existence of an edge at  $\bar{t}$  is not sufficient. Given that an edge  $e \in \mathcal{E}$  cannot exist without the incident nodes  $u, v \in e$ , the activity status of both the edge and its incident nodes is required to be 1. First, it is essential to score the events between nodes to ensure that something is happening on the edge between them. Second, if there is an event happening on the edge, the edge attributes change is scored. Altogether, the intensity function is defined by

$$\lambda_5^{u,v}(t) = \begin{cases} f_5 \left( \underbrace{\mathbf{W}_5^\top \begin{pmatrix} \mathbf{Z}_t^u \\ \mathbf{Z}_t^v \\ \mathbf{Z}_t^e \end{pmatrix}}_{\text{scores events between u,v}} + \underbrace{\hat{\mathbf{W}}_5^\top \begin{pmatrix} \mathbf{Z}_t^e \\ \mathbf{Z}_t^e \end{pmatrix}}_{\text{scores edge attribute change}} \right), & \text{if } (\xi_t^e = 1 \wedge \xi_t^u = 1 \wedge \xi_t^v = 1) \\ 0, & \text{if } (\xi_t^e = 0 \vee \xi_t^u = 0 \vee \xi_t^v = 0). \end{cases}$$

## 4 Learning Procedure

Depending on the application and the specific dynamic graph dataset, there are different suitable learning procedures, and only the necessary intensity functions have to be included in the learning. However, if every dynamic event is present in the observation set  $\mathcal{O}$ , all the intensity functions have to be used for training and updating the model.

**Training** The training procedure is done analogously to the training described in [4]. Let  $\lambda_k$  be the intensity function corresponding to the event  $o_k \in \mathcal{O}$  and  $x$  corresponds to the node or edge referred to in  $o_k$ . Then the model parameters including the learnable weights  $\mathbf{M}_\star, m_\star$ , and  $\mathbf{W}_\star$  from Sec. 3 are learnt by minimizing the negative log likelihood

$$\mathcal{L} := - \sum_{o_k \in \mathcal{O}} \log(\lambda_k^x(t)) + \int_0^T \Lambda(\delta) d\delta, \quad (9)$$

where

$$\Lambda(\delta) = \sum_{x \in \mathcal{V}_\delta \cup \mathcal{E}_\delta} \sum_{k=0}^5 \lambda_k^x(\delta)$$

represents the total probability for all events including also those that did not occur in the current timestamp (which is referred to as *survival probability*).

To keep the training effort low and not to optimize  $\mathcal{L}$  w.r.t. all events at once, mini-batch stochastic gradient descent can be utilized together with sampling techniques to approximate the integral<sup>3</sup>. Furthermore, the training is conducted in the Back Propagation Through Time (BPTT) manner. Assuming that several intensity functions share the weights independently from each other and can thus be computed in parallel, the training can be further accelerated.

**Model Update** The model can be updated with the aid of an event by applying the mini-batch stochastic gradient descent on the negative log likelihood (9) by only considering the event and the affected nodes or edges. This local retraining approach minimizes the update effort as only a small subset of the parameters from the intensity functions have to be touched.

<sup>3</sup>For example, Monte Carlo approximation.

## 5 Discussion and Future Work

The proposed FDGNN model provides opportunities for further development. In the following, the model’s generativity, explainability, and generalization are briefly discussed.

**Generative Attribute Embedding** When it comes to the forecast of new nodes or edges and their attributes, the attribute embedding in Eq. (1) has to be generative. Therefore, a generative model has to be selected concerning the application carefully.

**Explainability** GNNs are called explainable if the model outputs an explanation for the predicted result or reasoning that can be inferred based on the model architecture. In the future, the goal is to make the fully dynamic model more explainable to simplify the application of the model for real-world problems. For this purpose, it is necessary to include more components in the architecture that improve explanations for the model output considering the input graphs. For the beginning, the attention mechanisms (cf. §3.1) allow for a better understanding of the model’s behavior.

**Definite Deletions vs. Inactivity** Furthermore, the activity function from Section 2.2 will be extended to the consideration of definite deletions, in addition to activities and inactivities of nodes and edges. Thereby, nodes and edges lose the possibility to occur again in the graph.

**More complex graph types** The structural graph properties of the dynamic graphs in this paper are rather elementary so that the FDGNN can be extended to more complex dynamical graph structures in the future. Given the results in [11], every graph type can be transformed into another graph type without any loss of information. The fact that the FDGNN model can handle undirected attributed dynamic graphs, together with the transformation costs of graph types [11] advocates for the transformation of any more complex graph types into an undirected attributed graph, due to memory savings, with subsequent application of FDGNN. However, the question remains plausible if a proper extension, e.g., to dynamic hypergraphs, would be more efficient than applying FDGNN after a graph transformation. A realization of such could be done by using the incidence matrix (instead of the adjacency matrix) and extending the intensity functions in §3.3 to more than two involved nodes.

## 6 Conclusion

This paper proposes a Fully Dynamic Graph Neural Network (FDGNN) that can handle structure- and attribute dynamics. I.e., node and edge attributes can change over time, and node and edge additions and deletions are supported. The approach of FDGNN is based on the model DyRep [4], which only handles edge-growing graphs. The proposed model can thus be seen as an extension of DyRep to other dynamic graph properties.

FDGNN models the temporal behavior of the structural dynamics of graphs with the aid of an activity function and temporal point processes formalized by intensity functions. To incorporate the attribute dynamics, different intensities are defined that catch the dynamics of the node and edge attributes. The constructed node and edge embeddings used for the intensities include the local neighborhood and historical information, the exogenous drive, and the corresponding attributes. Furthermore, the utilization of attention mechanisms for the neighborhood and the historical information provides the first step towards an explainable model.

Assuming that several intensities share their weights and others are computationally unrelated, the computation can be done in parallel, and the learning phase can be reduced. Furthermore, the update of the model is efficient due to using only one event and the corresponding nodes or edges for local retraining.

## Contributions

- Conceptualization: SBW, AMO
- Methodology: SBW, AMO, RN
- Resources: AMO
- Writing (Original Draft): SBW, AMO
- Writing (Review & Editing): SBW, AMO, RN
- Supervision: RN, JT
- Project administration: SBW, AMO
- Funding acquisition: JT

Alice Moallem-Oureh is mainly responsible for the incorporation of attribute-dynamics, while Silvia Beddar-Wiesing is mainly responsible for modeling the structure-dynamics.



## Acknowledgement

The project is funded by the German Ministry of Education and Research (BMBF) with the grant number 01IS20047A.

## References

- [1] William L. Hamilton. *Graph Representation Learning*. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2020.
- [2] Youngjoo Seo, Michaël Defferrard, Pierre Vandergheynst, and Xavier Bresson. Structured sequence modeling with graph convolutional recurrent networks. In Long Cheng, Andrew Chi-Sing Leung, and Seiichi Ozawa, editors, *Neural Information Processing - 25th International Conference, ICONIP 2018, Siem Reap, Cambodia, December 13-16, 2018, Proceedings, Part I*, volume 11301 of *Lecture Notes in Computer Science*, pages 362–373. Springer, 2018.
- [3] Rakshit Trivedi, Hanjun Dai, Yichen Wang, and Le Song. Know-evolve: Deep temporal reasoning for dynamic knowledge graphs. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3462–3471. PMLR, 2017.
- [4] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. Dyrep: Learning representations over dynamic graphs. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.
- [5] Daheng Wang, Zhihan Zhang, Yihong Ma, Tong Zhao, Tianwen Jiang, Nitesh V. Chawla, and Meng Jiang. Learning attribute-structure co-evolutions in dynamic graphs. *CoRR*, abs/2007.13004, 2020.
- [6] Josephine M. Thomas, Alice Moallem-Oureh, Silvia Beddar-Wiesing, and Clara Holzhüter. Graph neural networks designed for different graph types: A survey. *CoRR*, abs/2204.03080, 2022.
- [7] Seyed Mehran Kazemi, Rishab Goel, Kshitij Jain, Ivan Kobyzev, Akshay Sethi, Peter Forsyth, and Pascal Poupart. Representation Learning for Dynamic Graphs: A Survey. *Journal of Machine Learning Research*, 21(70):1–73, 2020.
- [8] Jakob Gulddahl Rasmussen. Lecture Notes: Temporal Point Processes and the Conditional Intensity Function. *CoRR*, abs/1806.00221, 2018.
- [9] Selivanov, Dmitriy and Bickel, Manuel and Wang, Qing. Package ‘text2vec’, 2020.
- [10] Huifa Liao, Jie Hu, Tianrui Li, Shengdong Du, and Bo Peng. Deep linear graph attention model for attributed graph clustering. *Knowl. Based Syst.*, 246:108665, 2022.
- [11] Josephine M. Thomas, Silvia Beddar-Wiesing, Alice Moallem-Oureh, and Rüdiger Nather. Graph type expressivity and transformations. *CoRR*, abs/2109.10708, 2021.