

DEEP INTERACTION PROCESSES FOR TIME-EVOLVING GRAPHS

Anonymous authors

Paper under double-blind review

ABSTRACT

Time-evolving graphs are ubiquitous such as online transactions on an e-commerce platform and user interactions on social networks. While neural approaches have been proposed for graph modeling, most of them focus on static graphs. In this paper we present a principled deep neural approach that models continuous time-evolving graphs at multiple time resolutions based on a temporal point process framework. To model the dependency between latent dynamic representations of each node, we define a mixture of temporal cascades in which a node’s neural representation depends on not only this node’s previous representations but also the previous representations of related nodes that have interacted with this node. We generalize LSTM on this temporal cascade mixture and introduce novel time gates to model time intervals between interactions. Furthermore, we introduce a selection mechanism that gives important nodes large influence in both k -hop subgraphs of nodes in an interaction. To capture temporal dependency at multiple time-resolutions, we stack our neural representations in several layers and fuse them based on attention. Based on the temporal point process framework, our approach can naturally handle growth (and shrinkage) of graph nodes and interactions, making it inductive. Experimental results on interaction prediction and classification tasks – including a real-world financial application – illustrate the effectiveness of the time gate, the selection and attention mechanisms of our approach, as well as its superior performance over the alternative approaches.

1 INTRODUCTION

Representation learning over graph data has become a core machine learning task with a wide range of applications including e-commerce, finance, social networks, and bioinformatics. Various neural graph representations such as (Perozzi et al., 2014; Grover & Leskovec, 2016; Wang et al., 2016; Kipf & Welling, 2017; Defferrard et al., 2016; Scarselli et al., 2009; Ying et al., 2018; Hamilton et al., 2017b; Monti et al., 2017; Den Berg et al., 2017) have been proposed to learn from *static* graph data and successfully used for downstream tasks (*e.g.*, classification). Graph data, however, are often dynamic in practice; nodes and interactions between them can grow and shrink. A straightforward approach to handle dynamic graphs is to compress them into one or several static graphs. The drawbacks of this approach are multifold; we not only blur temporal structural information but also miss time information that can be critical for real-world applications. An illustrative example is given in figure 1.

To handle continuous time-evolving graph, we can approximate a by a sequence of snapshot graphs, each of which includes all interactions that occur during a user-specified discrete-time interval, as shown in (Leskovec et al., 2007; Hamilton et al., 2016; Kulkarni et al., 2015; Goyal et al., 2018). This treatment reduces time resolution and it is tricky to specify the appropriate aggregation granularity. To avoid these problems, Nguyen et al. (2018) proposed continuous-time dynamic networks (CTDNE) that generalize deep walk methods to learn time-dependent network embedding. As a *transductive* method, CTDNE cannot handle the growth of new nodes. Dai et al. (2016) applied temporal point processes to model time-evolving graphs and, as a nonparametric Bayesian approach, their approach can naturally cope with the growth of new nodes and interactions. They used recurrent neural networks (RNNs) to define an intensity function in temporal point processes. These RNN models are shallow and one-step unrolled, making it easy to compute but relatively limited in modeling power. Trivedi

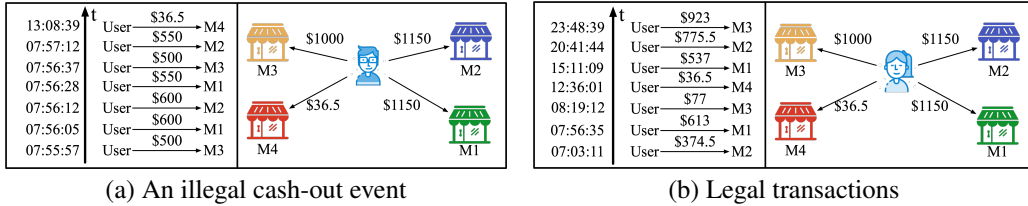


Figure 1: An illustrative example. Figure (a) shows an illegal cash-out event. It can be revealed by high-frequency transactions with multiple merchants. However, if we merge the transaction data into a static graph, we cannot distinguish it from the static graph generated from normal online shopping activities. Thus, learning from such a static graph will fail to detect the cash-out event.

et al. (2019) extended this approach by modeling two-time scale and adopting temporal-attention mechanism.

In this paper we present a powerful deep neural approach that models continuous time-evolving graphs at multiple time resolutions based on a temporal point process framework. We name the new approach *deep interaction processes* (DIPs). To model the dependency between latent dynamic representations of each node, we define a mixture of temporal cascades in which a node’s neural representation depends on not only this node’s previous representations but also the previous representations of related nodes that have interacted with this node. We generalize LSTM on this temporal cascade mixture and introduce novel time gates to model time intervals between interactions. Furthermore, We introduce a selection mechanism that gives important nodes large influence in both k -hop subgraphs of nodes in an interaction. To obtain representations from fine-to-coarse time-resolutions, we stack our neural representations in several layers and fuse them based on attention. Based on the temporal point process framework, our approach can naturally handle growth of graph nodes and interactions, making it inductive.

The rest of the paper is organized as follows. In Section 2 we give background on temporal point processes and in Section 3 we present the new DIP approach. In Section 4 we discuss related works. In Section 5 we report experimental results on multiple interaction prediction and classification tasks including an important real-world anti-fraud financial application, demonstrating superior performance of the new approach over the alternatives.

2 TEMPORAL POINT PROCESSES

We first describe temporal point processes (a class of nonparametric Bayesian models) that our approach is based on. Specifically, a temporal point process is a stochastic process that generates a sequence of discrete events localized at times $\{t_i\}_{i=1}^N$ in any given observed time window $[0, T]$. An important way to characterize temporal point processes is via the conditional intensity function $\lambda(t|H_t)$ -the stochastic model for the next event time t given all historical events before time t , denoted as $H_t = \{t_i|t_i < t\}$. Formally, within a small time window $[t, t + dt)$, $\lambda(t|H_t) dt$ is the probability for the occurrence for a new event given the H_t : $\lambda(t|H_t) dt = P\{\text{event in } [t, t + dt)|H_t\}$. From the survival analysis theory(Aalen et al., 2008), given the times of the past events $\{t_1, t_2, \dots, t_i\}$, the conditional density that an event occurs at t_{i+1} is given as follows: $p(t_{i+1}|H_{t_{i+1}}) = \lambda(t_{i+1}|H_{t_{i+1}}) \exp\left\{-\int_{t_i}^{t_{i+1}} \lambda(t|H_t) dt\right\}$, where the exponential part in the above equation means the conditional probability that no event happens during $[t_i, t_{i+1})$. The functional forms of the conditional intensity function $\lambda(t|H_t)$ can represent certain forms of dependencies of the historical events. For instance, for Poisson processes(Kingman, 2005) we set λ to be constant – making the assumption that the process is stationary and the temporal events in history are independent of each other. For classical Hawkes processes(Hawkes, 1971), the intensity function λ is often set to be a sum of multiple exponential functions, assuming that the mutual excitation among events is positive, additive over the past events, and exponentially decaying with time. Mei & Eisner (2017a) removed these limiting assumptions using LSTM to learn λ from data.

3 DEEP INTERACTION PROCESSES

In this section, we present the new neural nonparametric Bayesian approach over continuous-time evolving graphs. First, we present a temporal dependency graph that is a mixture of the temporal cascades, to model interdependence between graph nodes (as well as latent node representations). Then we present a novel deep model to learn dynamic node representations in the temporal dependency graph. This model naturally generalizes LSTM on the traditional chain-structured data. Given the dynamic node representations, we define deep interaction processes that model potential interactions between any two nodes over time. Finally we layout the maximum likelihood estimation method.

3.1 TEMPORAL DEPENDENCY GRAPH

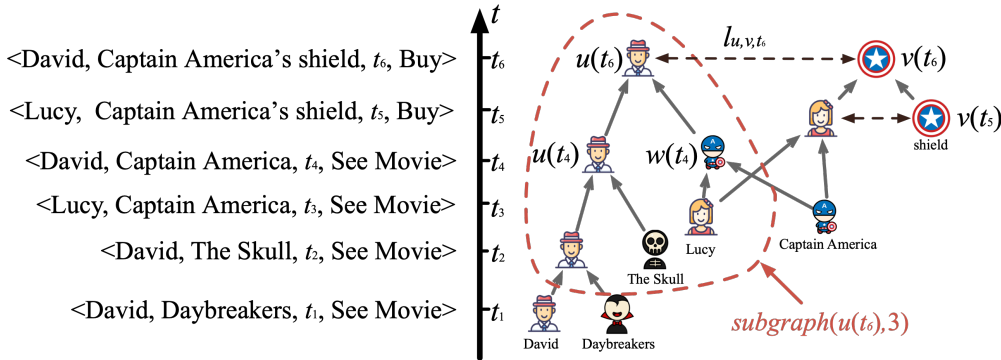


Figure 2: Dynamic interactions and the corresponding dependency graph

Consider a collections of people-movie records at different time points (e.g., Dave buy Captain America’s Shield toy at t_6 .) as shown in figure 2. The people and movies form a dynamic graph in which each person or movie is a node and interactions happen over time. After one interaction occurs, we update the neural representation of the two nodes linked to this interaction; e.g., right after time t_6 , we update the representations for David and the Captain America’s Shield toy. The new neural representation of David depends on both his current and previous interactions – as a result, depending on the representations of the two nodes associated with the previous interaction. This naturally forms a dependency cascade. Similarly we can obtain a dependency cascade for Lucy’s representations. Because of the common movies David and Lucy saw and toys they bought, their dependency cascades overlap and form a cascade mixture. Formally, we denote a dynamic interaction or link at time t by $l_{u,v,t}$ where u are v are two nodes associated with this interaction. We denote the node u at time t by $u(t)$ and the two nodes associated with u ’s precedent interaction at time t^- as $u^1(t)$ and $u^2(t)$. Note that one of $u^1(t)$ and $u^2(t)$ is simply $u(t^-)$. For example, $u^1(t_6)$ and $u^2(t_6)$ in figure 2 are $u(t_4)$ and $w(t_4)$, respectively. For later usage, We denote the subgraph rooted at $u(t)$ with $(k-1)$ depth as $subgraph(u(t),k)$ shown in Figure 2.

3.2 DIP NEURAL UNIT

Now we present the novel neural unit to update dynamic latent representations of nodes over the temporal dependency graph. First, let us denote node u ’s features or embedding (i.e., a static representation jointly learned from data) at time t by $\mathbf{x}_{u(t)}$ and denote features of interaction l by x_l . The interaction feature can be empty if the interaction contains only the temporal information. The concatenation of $\mathbf{x}_{u(t)}$ and x_l is denoted by $\hat{x}_{u(t)}$. Let $\Delta_{(u,t)} = t - t^-$ be the time interval between two consecutive interactions involving u at time t and t^- .

Our neural unit generalizes LSTM unit on the temporal dependency graph; we use an input gate, an output gate and two forget gates over $\hat{x}_{u(t)}$, dynamic representation of $\mathbf{h}_{u^i(t)}$ and cell states $\mathbf{c}_{u^i(t)}$ ($i = 1, 2$) as shown in figure 3. In addition to these gates, we introduce time gates to capture the impact of time interval $\Delta_{(u,t)}$.

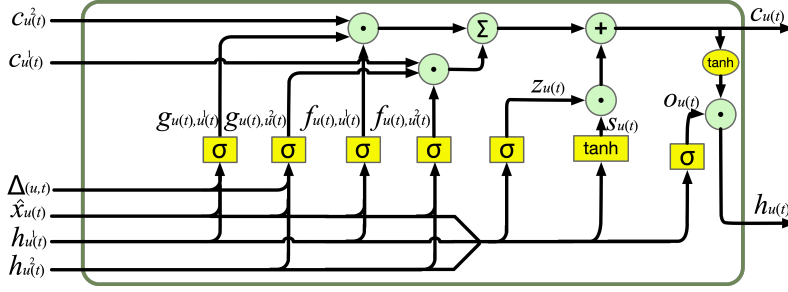


Figure 3: Time-evolving graph unit updates the representation $\mathbf{h}_{u(t)}$ and cell state $\mathbf{c}_{u(t)}$ for node u at time t based on both the features $\hat{\mathbf{x}}_{u(t)}$ and the representations and cell states of the nodes $u^1(t)$ and $u^2(t)$ associated with u 's precedent interaction.

Specifically, $\mathbf{h}_{u(t)}$ and $\mathbf{c}_{u(t)}$ are updated as follows:

$$\begin{aligned} \mathbf{z}_{u(t)} &= \sigma \left(\mathbf{W}_z \hat{\mathbf{x}}_{u(t)} + \sum_{i=1}^2 \mathbf{R}_{z_i} \mathbf{h}_{u^i(t)} + \mathbf{b}_z \right) & \mathbf{o}_{u(t)} &= \sigma \left(\mathbf{W}_o \hat{\mathbf{x}}_{u(t)} + \sum_{i=1}^2 \mathbf{R}_{o_i} \mathbf{h}_{u^i(t)} + \mathbf{b}_o \right) \\ \mathbf{s}_{u(t)} &= \tanh \left(\mathbf{W}_s \hat{\mathbf{x}}_{u(t)} + \sum_{i=1}^2 \mathbf{R}_{s_i} \mathbf{h}_{u^i(t)} + \mathbf{b}_s \right) & \mathbf{f}_{u(t), u^i(t)} &= \sigma \left(\mathbf{W}_{f_i} \hat{\mathbf{x}}_{u(t)} + \mathbf{R}_{f_i} \mathbf{h}_{u^i(t)} + \mathbf{b}_{f_i} \right) \\ \mathbf{g}_{u(t), u^i(t)} &= \sigma \left(\mathbf{W}_{g_i} \hat{\mathbf{x}}_{u(t)} + \mathbf{R}_{g_i} \mathbf{h}_{u^i(t)} + \mathbf{M}_{g_i} \Delta_u + \mathbf{b}_{g_i} \right) \\ \mathbf{c}_{u(t)} &= \mathbf{z}_{u(t)} \odot \mathbf{s}_{u(t)} + \sum_{i=1}^2 \mathbf{f}_{u(t), u^i(t)} \odot \mathbf{c}_{u^i(t)} \odot \mathbf{g}_{u(t), u^i(t)} \\ \mathbf{h}_{u(t)} &= \mathbf{o}_{u(t)} \odot \tanh(\mathbf{c}_{u(t)}) \end{aligned}$$

where σ , \tanh and \odot represent the sigmoid function, the hyperbolic tangent function, and the Hadamard product (pointwise multiplication), respectively, and parameters in the unit including the recurrent weights \mathbf{R}_{z_i} , \mathbf{R}_{o_i} , \mathbf{R}_{s_i} , \mathbf{R}_{f_i} and \mathbf{R}_{g_i} , the projection matrices \mathbf{W}_z , \mathbf{W}_o , \mathbf{W}_s , \mathbf{W}_{f_i} and \mathbf{W}_{g_i} , the bias vectors \mathbf{b}_z , \mathbf{b}_o , \mathbf{b}_s , \mathbf{b}_{f_i} and \mathbf{b}_{g_i} and the time weight matrix \mathbf{M}_{g_i} are learned from data.

For convenience, we use DIP-UNIT (\cdot) to summarize the above equations,

$$\mathbf{h}_{u(t)}, \mathbf{c}_{u(t)} = \text{DIP-UNIT}(\hat{\mathbf{x}}_{u(t)}, \mathbf{c}_{u^1(t)}, \mathbf{c}_{u^2(t)}, \mathbf{h}_{u^1(t)}, \mathbf{h}_{u^2(t)}, \Delta_u, \Theta)$$

where Θ represent all the parameters. Similar to LSTM training where a k -hop neighborhood is often used in practice, we limit the backtracking in $\text{subgraph}(u(t), k)$ to the computational cost.

3.3 DIP-UNIT AND FUSION

To model nonlinear dependency relationships at different temporal resolutions, we stack L layers of DIP-UNIT together. The output of the j -th layer is computed recursively as follows:

$$(\mathbf{h}_{u(t)}^j, \mathbf{c}_{u(t)}^j) = \text{DIP-UNIT}^j \left(\mathbf{h}_{u(t)}^{j-1}, \mathbf{c}_{u^1(t)}^j, \mathbf{c}_{u^2(t)}^j, \mathbf{h}_{u^1(t)}^j, \mathbf{h}_{u^2(t)}^j, \Delta_u, \Theta_j \right)$$

where $\mathbf{h}_{u(t)}^0 = \hat{\mathbf{x}}_{u(t)}$, $j = 1, \dots, L$. To train deeper dynamic neural networks easily, we employ the residual connection as the following form: $\text{skip}(h_{u(t)}^{j-1}, h_{u(t)}^j) = \mathbf{W}_{\text{skip}} h_{u(t)}^{j-1} + h_{u(t)}^j$ where \mathbf{W}_{skip} is a weight matrix. Motivated by *ELMo* (Peters et al., 2018), we fuse all internal dynamic representations from all the layers to achieve rich dynamic representations. The fusion is a weighted summation of all layers defined as follows: $\mathbf{h}_{u(t)} = \gamma^{\text{task}} \sum_{j=0}^L \alpha_j^{\text{task}} \mathbf{h}_{u(t)}^j$, where α_j^{task} are task-related softmax-normalized weights and γ^{task} is a scaling parameter.

3.4 SELECTION

Given an interaction $l_{u,v,t}$, it is reasonable to assume that not all the historical interactive nodes have the equal importance for formalizing this interaction. Thus we use an attention mechanism to *select*

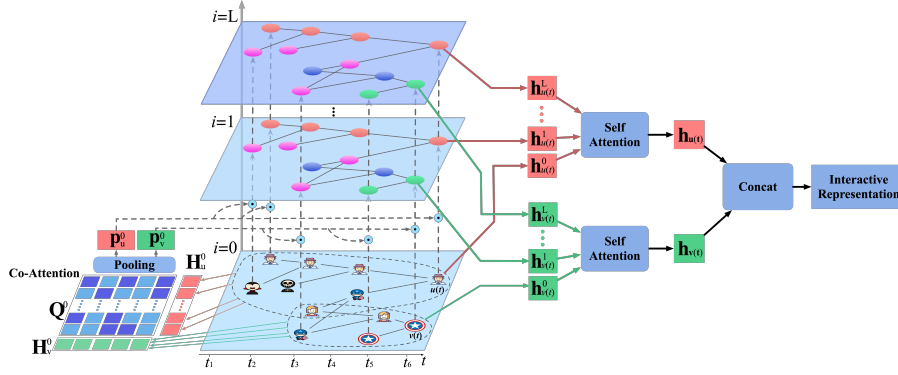


Figure 4: Neural stacking, fusion and selection.

relevant nodes to learn dynamic representations and cell states of the current node. Specifically, a co-attention mechanism is used to measure relevance of historical time-evolving patterns between $subgraph(u(t), k)$ and $subgraph(v(t), k)$,

$$\mathbf{Q}^j = \tanh(\mathbf{H}_u^{j\top} \mathbf{W}_Q \mathbf{H}_v^j)$$

where $\mathbf{H}_u^j = [\mathbf{h}_1^j, \dots, \mathbf{h}_a^j, \dots, \mathbf{h}_m^j]$, $a \in subgraph(u(t), k)$, $\mathbf{H}_v^j = [\mathbf{h}_1^j, \dots, \mathbf{h}_b^j, \dots, \mathbf{h}_n^j]$, $b \in subgraph(v(t), k)$, and $\mathbf{W}_Q \in \mathbb{R}^{d \times d}$ are the weight parameters. The \mathbf{Q}^j is a co-attention affinity matrix which captures the relevance information in $subgraph(u(t), k)$ and $subgraph(v(t), k)$. The co-dependent global embedding p_u^j, p_v^j are obtained by the following equations.

$$\mathbf{p}_u^j = \mathbf{H}_u^j \text{SoftMax} \left(\underset{ColWise}{\text{Max}} \mathbf{Q}^j \right) \quad \mathbf{p}_v^j = \mathbf{H}_v^j \text{SoftMax} \left(\underset{RowWise}{\text{Max}} (\mathbf{Q}^j)^\top \right)$$

where Max means max-pooling operation which is used to choose the most relevant information for the maximum influence (or affinity) on nodes in the corresponding subgraph. In addition, to adjust the importance of historical nodes, two adaptive gate functions are designed for previous nodes in $subgraph(u(t), k)$ and $subgraph(v(t), k)$ respectively,

$$g_u(p_u^j, h_a^j) = \sigma(w_p p_u^j + w_h h_a^j) \quad g_v(p_v^j, h_b^j) = \sigma(w_p p_v^j + w_h h_b^j)$$

where the weights w_p , and w_h are shared by all the stacked layers. Using these gates, we adjust dynamic node representations as follows:

$$(\mathbf{h}_{u(t)}^j, \mathbf{c}_{u(t)}^j) = \text{DIP-UNIT}^j \left(\mathbf{h}_{u(t)}^{j-1} \odot g_u(\mathbf{p}_{u(t)}^{j-1}, \mathbf{h}_{u(t)}^{j-1}), \mathbf{c}_{u^1(t)}^j, \mathbf{c}_{u^2(t)}^j, \mathbf{h}_{u^1(t)}^j, \mathbf{h}_{u^2(t)}^j, \Delta_u, \Theta_j \right)$$

Similarly, we can update $(\mathbf{h}_{v(t)}^j, \mathbf{c}_{v(t)}^j)$ based on the selection mechanism.

3.5 CONDITIONAL INTENSITY FUNCTION

We model the dynamic interactions as a multi-dimensional temporal point process. Specifically, we define the conditional intensity function of the temporal point process at the dimension indexed by (u, v) , given the dependant histories of non-chain structures $H_t^{u,v}$ where $H_t^{u,v} = subgraph(u^1(t), k) \cup subgraph(u^2(t), k) \cup subgraph(v^1(t), k) \cup subgraph(v^2(t), k)$, as follows:

$$\lambda^{u,v}(t | H_t^{u,v}) = \text{SoftPlus}(\mathbf{h}_t^{u,v} \mathbf{w}_\lambda + \mathbf{w}_t^\top \tau + b_\lambda)$$

where

$$\mathbf{h}_t^{u,v} = [\mathbf{h}_{u^1(t)}^\top, \mathbf{h}_{u^2(t)}^\top, \mathbf{h}_{v^1(t)}^\top, \mathbf{h}_{v^2(t)}^\top], \tau = [\Delta_{u,t}, \Delta_{v,t}]^\top,$$

the scalar b_λ can be viewed as a base intensity level for the occurrence of the next interaction, and the *SoftPlus* function is used to ensure the non-negativity of the intensity. A key step for obtaining $H_t^{u,v}$ is to get the k-hop *subgraphs* of u and v 's direct dependants Please see **Appendix.A** for more details about fast obtaining k-hop *subgraphs*

3.6 PARAMETER ESTIMATION

3.6.1 INTERACTION PREDICTION

Given a set of interactions as $I = \{(u_i, v_i, t_i)\}_{i=1}^{i=N}$, we can learn the model by minimizing the negative joint log-likelihood of I as follows: $\mathcal{L}_1 = -\sum_i \log P^{u_i, v_i}(t_i | H_{t_i}^{u_i, v_i})$ where $P^{u_i, v_i}(t_i | H_{t_i}^{u_i, v_i})$ represents the probability of formalizing an interaction between u_i and v_i at time t_i given the dependant history of non-chain structures $H_{t_i}^{u_i, v_i}$. Based on the intensity definition, we have $\mathcal{L}_1 = \sum_i -\lambda^{u_i, v_i}(t_i | H_{t_i}^{u_i, v_i}) + \int_{t_i^-}^{t_i} \Lambda(t) dt$, where t_i^- is the most recent time point when either u_i or v_i was involved in an interaction. $\Lambda(t) = \sum_{u, v} \lambda^{u, v}(t)$ which represents total survival probabilities for interactions that do not happen. Since the survival part does not have an analytic solution, we apply Monte Carlo to do numerical integrations. We follow the negative sampling approaches proposed by Dai et al. (2016); Trivedi et al. (2019) to accelerate the survival term calculation.

3.6.2 INTERACTION CLASSIFICATION

An interaction sequence with markers is denoted as $I' = \{(u_i, v_i, t_i, y_i)\}_{i=1}^{i=N}$, where y_i is a marker at time t_i and usually is a discrete variable. In practice, the markers have different meanings in distinct scenes. A marker can be treated as a magnitude in modeling earthquakes and aftershocks, while in financial transaction scenes, the marker can be used to label whether a transaction is a fraudulent trading or not. The joint conditional density of an interaction (u_i, v_i, t_i) with marker y_i is given as $P^{u_i, v_i}(t_i, y_i | \hat{H}_{t_i}^{u_i, v_i})$. By applying the Bayesian rule, the joint conditional density can be written as: $P^{u_i, v_i}(t_i, y_i | \hat{H}_{t_i}^{u_i, v_i}) = P^{u_i, v_i}(t_i | \hat{H}_{t_i}^{u_i, v_i}) P(y_i | t_i, \hat{H}_{t_i}^{u_i, v_i})$, where $P^{u_i, v_i}(t_i | \hat{H}_{t_i}^{u_i, v_i})$ has the same meaning as given in subsection 3.6.1, while $P(y_i | t_i, \hat{H}_{t_i}^{u_i, v_i})$ means the distribution of y_i given the interaction happened at t_i with interaction history $\hat{H}_{t_i}^{u_i, v_i}$. It should be noted that the history $\hat{H}_{t_i}^{u_i, v_i}$ contains the information of history markers and one need to design marker-specific intensity function (Mei & Eisner, 2017b). For simplicity, we assume $P^{u_i, v_i}(t_i, y_i | \hat{H}_{t_i}^{u_i, v_i})$ is independent on history markers. Meanwhile, $P(y_i | t_i, \hat{H}_{t_i}^{u_i, v_i})$ can be obtained by a multinomial function:

$$P(y_i = c | \mathbf{h}_{u_i, v_i, t_i}) = \frac{\exp(\mathbf{V}_c^y \mathbf{h}_{u_i, v_i, t_i})}{\sum_{c=1}^C \exp(\mathbf{V}_c^y \mathbf{h}_{u_i, v_i, t_i})}$$

where $\mathbf{h}_{u_i, v_i, t_i}$ is the concatenation of \mathbf{h}_{u_i} and \mathbf{h}_{v_i} which can be regarded as dynamic representation for an interaction between u_i and v_i at t_i . C is the number of markers, \mathbf{V}_c^y is the c -th row of matrix \mathbf{V}^y . Then the final objective function for interaction classification can be obtained as follows: $\mathcal{L}_2 = \mathcal{L}_1 + \mathcal{L}_{\text{cross-entropy}}$, where $\mathcal{L}_{\text{cross-entropy}}$ is a cross-entropy loss over marks:

$$\mathcal{L}_{\text{cross-entropy}} = -\sum_{i=1}^{i=N} \sum_{c=1}^C y_i \cdot \log P(y_i = c | \mathbf{h}_{u_i, v_i, t_i})$$

4 RELATED WORK

Inspired by the Skip-gram (Mikolov et al., 2013) for word embedding, a series of node embedding methods based on the random walks on graphs have been proposed (Perozzi et al., 2014; Tang et al., 2015; Grover & Leskovec, 2016; Wang et al., 2016; 2017). GCN and its variants (Bruna et al., 2013; Hamilton et al., 2017a; Kipf & Welling, 2017) are a recent class of algorithms which extend convolutions from spatial domains to graph-structured domains. Meanwhile they can efficiently generate node embeddings for previously unseen data. All models above are designed for static graphs. The intuitive and popular approaches for modeling dynamic graphs are based on a sequence for graph snapshots (Goyal et al., 2018; Zhou et al., 2018; Seo et al., 2018; Yu et al., 2018), while it can be difficult to specify the appropriate aggregation granularity. Nguyen et al. (2018) adds a temporal constraint on random walk sampling, but it can't model the rich temporal information explicitly. Temporal point processes (TPPs) are an another alternative to model dynamics (Daley & Vere-Jones, 2007). Several dynamic graph modeling methods based on the TPPs (Dai et al., 2016; Trivedi et al., 2019) have been proposed. Our method DIP differs from these TPP-based methods by

the extension of the LSTM model over temporal dependency graphs, the multiple time resolution modeling via stacking and fusing, and the selection mechanism. Detailed related work are included in **Appendix.C**.

5 EXPERIMENTS

We evaluate the proposed DIP model on the task of **Interaction Prediction** and **Interaction classification** on several real-world datasets.

5.1 BASELINES

GraphSage(Hamilton et al., 2017a) is an inductive graph neural network framework consisting of three different aggregators which are **GCN**, **Mean** and **LSTM** aggregators respectively. We report the best result among these three aggregators noted as Graphsage*. What’s more, for comparing with **GAT**(Veličković et al., 2017) we also implement a graph attention aggregator based on GraphSage. **CTDNE**(Nguyen et al., 2018) is a newly-proposed temporal network embedding method which is an extension of DeepWalk(Perozzi et al., 2014) by incorporating temporal order constraint when sampling sequences of walks from time-continuous graphs. **DynGEM**(Goyal et al., 2018) takes a sequence of static graph snapshots as inputs to learn node embeddings by a deep auto-encoder network. **DeepCoevolve** (Dai et al., 2016) models dynamic interaction sequences with two co-evolution recurrent neural networks. Hidden embeddings are learned for interactive nodes after each interaction. **DyREP** (Trivedi et al., 2019) uses a two-time scale deep temporal point process model to capture dynamics of graphs.

5.2 EXPERIMENTAL SETTING

We conduct all the experiments with a hyper-parameter grid search strategy. For all methods, we search the hidden vector dimension from {32, 64, 128, 256} and the learning rate from {0.01, 0.001, 0.0005, 0.0001, 0.00001}. For our DIP model, we go through {1, 2, 3, 4} for K and L . For Graphsage, the maximum number of 1-hop and 2-hop neighbor nodes are set to 25 and 20 respectively. The batch sizes for all candidates are {100, 300, 500}. All the models are trained for at most 50 epochs with an early-stop if the performance does not improve for 5 epochs. For Graphsage, DynGEM and DeepCoevolve, we use the open source codes provided by the authors. We implement the CTNDE and GAT based on the Graphsage framework, and implement DyREP based on the pytorch implementation of DeepCoevolve. After the best configuration has been found, we repeat the full experiments 5 times and report the mean results and standard deviation.

5.3 INTERACTION PREDICTION

5.3.1 DATASETS

CollegeMsg(Leskovec & Krevl, 2014) consists of sending message interactions on an online social network at the University of California, Irvine during 193 days. **Ubuntu**(Leskovec & Krevl, 2014) is a temporal interaction dataset extracted from the stack exchange website. An interaction between two users means one answered another’s questions or replied to his/her posts. **Amazon**(McAuley et al., 2015) is composed of commodity rating data from amazon users. We use the *Clothing* subset of this dataset. **MathOverflow**(Leskovec & Krevl, 2014) is comprised of interactions of commenting an existing answer on the *Math Overflow* website. **Table 1** shows the detailed dataset statistics. In this table, *Sparsity* indicates $|Edge|/(|u(t)| * |v(t)|)$. And the *Duplication* is the quotient of # of temporal edges divided by # of statics edges. Duplication equals 1.0 means each unique pair of participants interact only once. We also give the average time interval of two consequent interactions of a certain participant as $Avg \Delta t$. For each dataset, we first construct the corresponding dependency graph as described in section 3.1. Then we sort these interactions by occurrence time. The first 60% interactions are adopted as training set and the next 20% interactions are used for validation. The last 20% interactions are left as test set. The Cold-start participants which only exist in validation set or test set are removed.

Table 1: Dataset Statistics

	CollegeMsg	Ubuntu	Amazon-Clothing	Math Overflow
# Train	35902	204846	50209	58596
# Valid	7814	39913	9195	24045
# Test	5055	35271	7598	32705
Duplication	3.0569	1.886	1.000	2.819
Avg Δt (Hours)	27.79	575.89	3399.82	415.46
Sparsity	0.01181	0.00074	0.00038	0.00567

5.3.2 EVALUATION PROTOCOL AND RESULTS

For each interaction $l_{u,v,t}$ in test set, we fix the first participant u , then replace the second participant v by all possible participant candidates. The conditional density $p^{u,v}(t) = \lambda^{u,v}(t)S^{u,v}(t)$ for all candidates are first computed and then sorted by a descending order. The rank of the correct participant is finally stored and denoted as $rank_i$ for the i -th test interaction. We report the **Mean Rank** defined as $Mean Rank = \frac{\sum rank_i}{\#of Test Interaction}$, which can represent the overall performance. Figure 5 summarizes the **Mean Rank** performance of all the baselines and our DIP method. DIP outperforms all baselines by 65.84%, 41.64%, 10.69% and 43.99% over the four datasets. When the average Δt is small, which means interactions are sensitive to temporal information, our model achieves the best performance. This gives the evidence of the effectiveness of explicitly modeling the temporal information. What’s more, for the **CollegeMsg** dataset which is densest and has the most duplicated interactions, our model outperforms best again.

5.4 INTERACTION CLASSIFICATION

We conduct this task on an industrial dataset: **Huabei Trade Data**. For the performance measures, we employ KS(Kolmogorov-Smirnov) value which is a big concern of the loans provider in anti-fraud detection, as well as AUC(Area under the ROC Curve) score. This dataset consists of about 150,000 transaction records processed by Huabei during August 2018. Each transaction is initiated with three parties: the buyer, the seller and transaction details such as merchant category and transaction amount. Around 15% of the transaction are fraudulent and is labeled by a complicated *Ex-Post* method. In order to conducting an *Ex-ante* detection, we are required to find out fraudulent transaction at the time of initiating using those basic transaction features. For each interaction event, there are 11 context features that can be obtained when the trade request is created, including information about buyer types, seller types, purchased items’ categories and trading platform. We use the first 10 days data as training set, the following 10 days data as validation set, the rest as test set. Note that, in this scenario, there are always users who only appear in validation/testing dataset. Thus, traditional transductive methods are not applicable on this task. Alternatively, we employ the **XGBoost** (Chen & Guestrin, 2016) as an additional baseline which is a popular model in the cash-out detection task(Hu et al., 2019). Table 2 compares the results. Obviously, our model outperforms all the baseline methods again. The contributions of different modules of the proposed DIP models are given in figure 7.

Table 2: Interaction classification results

	Xgboost	GraphSage*	GAT	DIP
AUC	0.6818 \pm 0.0023	0.8603 \pm 0.0005	0.8597 \pm 0.0004	0.9017 \pm 0.0004
KS	0.2536 \pm 0.0015	0.5934 \pm 0.0012	0.6018 \pm 0.0005	0.6703 \pm 0.0060

5.5 ABLATION STUDY

As we described in Section 3, the DIP model consists of three important components: First, it uses a Time Gate in the DIP neural unit to explicitly model the temporal information. Second, the selection mechanism enables our model to select more important historical information for interactions. Third,

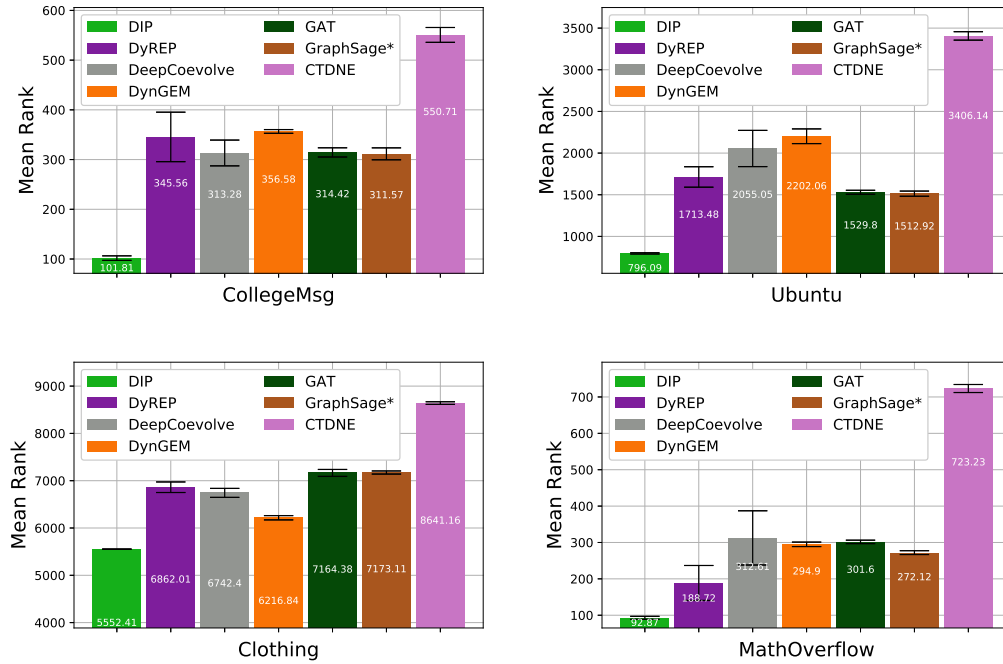


Figure 5: Mean rank results

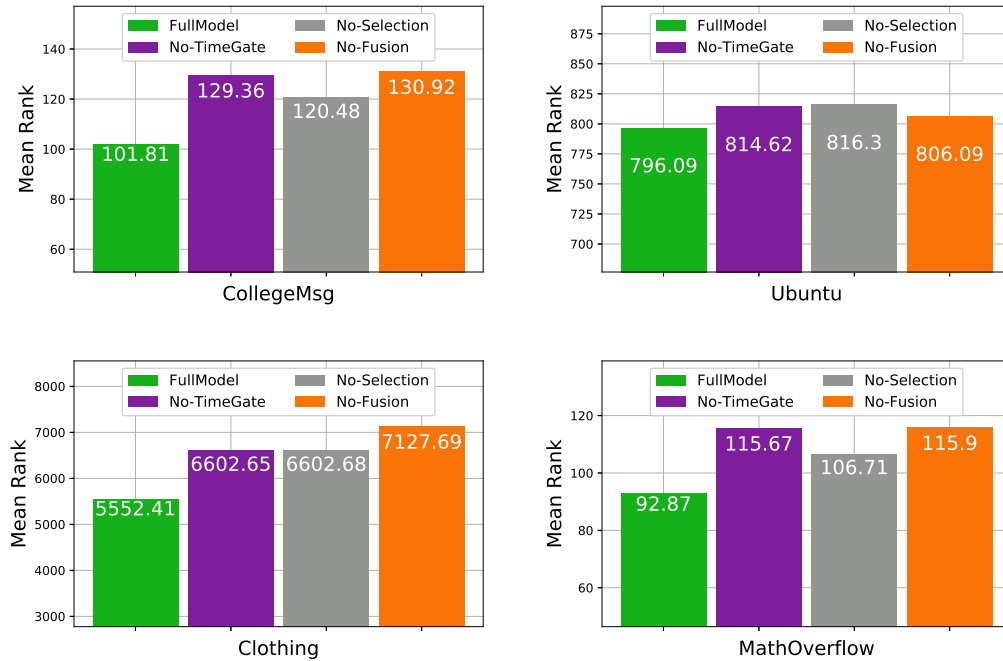


Figure 6: Ablation study results of the interaction prediction task

the Fusion of multi-layer DIP-UNIT’s hidden state vector helps to extract high level feature. We investigate the contribution of each component by disabling each of them one by one, and compare the corresponding result to the full model. figure 6 and figure 7 give the detailed ablation results. **FullModel** in the two figures means all the three components are enabled.

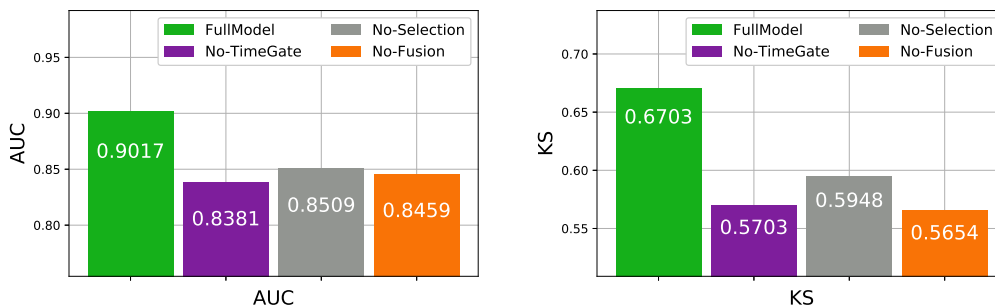


Figure 7: Ablation study results of the interaction classification task

- **No Time Gate:** In this configuration, the time gate in DIP-UNIT is disabled. This leads to a significant drop of the Mean Rank performance. It provides a strong evidence for the effectiveness of the time gate.
- **No Selection:** In this configuration the selection mechanism is disabled. Accordingly, all the historical node representations contribute equally, thus again leading to a performance drop.
- **No Fusion:** In this variant, we directly use the hidden state vector of the last layer. Again, the performance degrades significantly. This demonstrates that a fusion of different layers’ representations gives richer information than the last layer only.

6 CONCLUSIONS

In this paper, we have proposed a deep multidimensional point process approach, DIP, to learn dynamic graph representations. We generalize LSTM over temporal dependency graphs and model multiple time resolutions via stacking, selection and fusion. Experimental results show the effectiveness of the components of our neural unit and the superior performance on several datasets.

REFERENCES

- Odd Aalen, Ornulf Borgan, and Hakon Gjessing. *Survival and event history analysis: a process point of view*. Springer Science & Business Media, 2008.
- Amr Ahmed, Nino Shervashidze, Shравan Narayanamurthy, Vanja Josifovski, and Alexander J Smola. Distributed large-scale natural graph factorization. In *Proceedings of the 22nd international conference on World Wide Web*, pp. 37–48. ACM, 2013.
- Mikhail Belkin and Partha Niyogi. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *Advances in neural information processing systems*, pp. 585–591, 2002.
- Joan Bruna, Wojciech Zaremba, Arthur Szlam, and Yann LeCun. Spectral networks and locally connected networks on graphs. *arXiv preprint arXiv:1312.6203*, 2013.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794. ACM, 2016.
- Hanjun Dai, Yichen Wang, Rakshit Trivedi, and Le Song. Deep coevolutionary network: Embedding user and item features for recommendation. *arXiv preprint arXiv:1609.03675*, 2016.
- Daryl J Daley and David Vere-Jones. *An introduction to the theory of point processes: volume II: general theory and structure*. Springer Science & Business Media, 2007.
- Michael Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks on graphs with fast localized spectral filtering. *neural information processing systems*, pp. 3844–3852, 2016.
- Rianne Van Den Berg, Thomas N Kipf, and Max Welling. Graph convolutional matrix completion. *arXiv: Machine Learning*, 2017.

- Nan Du, Hanjun Dai, Rakshit Trivedi, Utkarsh Upadhyay, Manuel Gomez-Rodriguez, and Le Song. Recurrent marked temporal point processes: Embedding event history to vector. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1555–1564. ACM, 2016.
- Palash Goyal, Nitin Kamra, Xinran He, and Yan Liu. Dyngem: Deep embedding method for dynamic graphs. *arXiv preprint arXiv:1805.11273*, 2018.
- Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. *knowledge discovery and data mining*, pp. 855–864, 2016.
- Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in Neural Information Processing Systems*, pp. 1024–1034, 2017a.
- William L Hamilton, Jure Leskovec, and Dan Jurafsky. Diachronic word embeddings reveal statistical laws of semantic change. *arXiv preprint arXiv:1605.09096*, 2016.
- William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *IEEE Data(base) Engineering Bulletin*, 40:52–74, 2017b.
- Alan G Hawkes. Spectra of some self-exciting and mutually exciting point processes. *Biometrika*, 58 (1):83–90, 1971.
- Binbin Hu, Zhiqiang Zhang, Chuan Shi, Jun Zhou, Xiaolong Li, and Yuan Qi. Cash-out user detection based on attributed heterogeneous information network with a hierarchical attention mechanism. 2019.
- John Frank Charles Kingman. P oisson processes. *Encyclopedia of biostatistics*, 6, 2005.
- Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. *international conference on learning representations*, 2017.
- Vivek Kulkarni, Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. Statistically significant detection of linguistic change. In *Proceedings of the 24th International Conference on World Wide Web*, pp. 625–635. International World Wide Web Conferences Steering Committee, 2015.
- Jure Leskovec and Andrej Krevl. SNAP Datasets: Stanford large network dataset collection. <http://snap.stanford.edu/data>, June 2014.
- Jure Leskovec, Jon Kleinberg, and Christos Faloutsos. Graph evolution: Densification and shrinking diameters. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(1):2, 2007.
- Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 43–52. ACM, 2015.
- Hongyuan Mei and Jason M Eisner. The neural hawkes process: A neurally self-modulating multivariate point process. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems 30*, pp. 6754–6764. Curran Associates, Inc., 2017a.
- Hongyuan Mei and Jason M Eisner. The neural hawkes process: A neurally self-modulating multivariate point process. In *Advances in Neural Information Processing Systems*, pp. 6754–6764, 2017b.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pp. 3111–3119, 2013.
- Federico Monti, Michael M Bronstein, and Xavier Bresson. Geometric matrix completion with recurrent multi-graph neural networks. *neural information processing systems*, pp. 3697–3707, 2017.

- Giang Hoang Nguyen, John Boaz Lee, Ryan A. Rossi, Nesreen K. Ahmed, Eunye Koh, and Sungchul Kim. Continuous-time dynamic network embeddings. In *Companion of the The Web Conference 2018 on The Web Conference 2018, WWW 2018, Lyon , France, April 23-27, 2018*, pp. 969–976, 2018.
- Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk: Online learning of social representations. In *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 701–710. ACM, 2014.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proc. of NAACL*, 2018.
- Sam T Roweis and Lawrence K Saul. Nonlinear dimensionality reduction by locally linear embedding. *science*, 290(5500):2323–2326, 2000.
- Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- Youngjoo Seo, Michaël Defferrard, Pierre Vandergheynst, and Xavier Bresson. Structured sequence modeling with graph convolutional recurrent networks. In *International Conference on Neural Information Processing*, pp. 362–373. Springer, 2018.
- Sucheta Soundarajan, Acar Tamersoy, Elias B Khalil, Tina Eliassi-Rad, Duen Horng Chau, Brian Gallagher, and Kevin Roundy. Generating graph snapshots from streaming edge data. In *Proceedings of the 25th International Conference Companion on World Wide Web*, pp. 109–110. International World Wide Web Conferences Steering Committee, 2016.
- Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *Proceedings of the 24th International Conference on World Wide Web*, pp. 1067–1077. International World Wide Web Conferences Steering Committee, 2015.
- Joshua B Tenenbaum, Vin De Silva, and John C Langford. A global geometric framework for nonlinear dimensionality reduction. *science*, 290(5500):2319–2323, 2000.
- 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*, 2019. URL <https://openreview.net/forum?id=HyEPrhR5KX>.
- Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- Daixin Wang, Peng Cui, and Wenwu Zhu. Structural deep network embedding. pp. 1225–1234, 2016.
- Xiao Wang, Peng Cui, Jing Wang, Jian Pei, Wenwu Zhu, and Shiqiang Yang. Community preserving network embedding. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Shuai Xiao, Mehrdad Farajtabar, Xiaojing Ye, Junchi Yan, Le Song, and Hongyuan Zha. Wasserstein learning of deep generative point process models. In *Advances in Neural Information Processing Systems*, pp. 3247–3257, 2017.
- Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L Hamilton, and Jure Leskovec. Graph convolutional neural networks for web-scale recommender systems. *knowledge discovery and data mining*, pp. 974–983, 2018.
- Wenchao Yu, Wei Cheng, Charu C Aggarwal, Kai Zhang, Haifeng Chen, and Wei Wang. Netwalk: A flexible deep embedding approach for anomaly detection in dynamic networks. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2672–2681. ACM, 2018.
- Lekui Zhou, Yang Yang, Xiang Ren, Fei Wu, and Yueting Zhuang. Dynamic network embedding by modeling triadic closure process. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

Linhong Zhu, Dong Guo, Junming Yin, Greg Ver Steeg, and Aram Galstyan. Scalable temporal latent space inference for link prediction in dynamic social networks. *IEEE Transactions on Knowledge and Data Engineering*, 28(10):2765–2777, 2016.

A TEMPORAL DEPENDENCY GRAPH

The key step for efficient training or inference is to fast obtain $k_subgraph(u(t))$ or $k_subgraph(v(t))$ for an interaction (u, v, t) . Obviously, our definition of temporal dependency graph(TDG) in Subsection.3.1 provides a recursive and incremental way for its construction. As illustrated in figure 2, we can easily find that $k_subgraph(u(t))$ is the union of $k_subgraph(u(t^-))$, $k_subgraph(w(t^-))$ and the newly added two edges e_1, e_2 . Based on this feature, to fast obtain k -hop *subgraph*, we design two algorithms, TDG-COLORING and CONSTRUCTION OF K-HOP TDG orderly.

A.1 TDG-COLORING

The TDG-COLORING algorithm is shown in Algorithm 1, which is the pre-step of Algorithm 2. This algorithm takes $I = \{(u_i, v_i, t_i)\}_{i=1}^{i=N}$ as input and generate a new sorted sequence of interactions ordered by color numbers in an ascending way. The TDG-COLORING algorithm ensures that, interactions in the same color group are independent and interactions from groups with smaller color numbers are precedents of the larger ones.

Algorithm 1 TDG-COLORING

Require: I : A sequence of interaction with a chronological order.
Ensure: $ColorGpSeq$: A sequence of sorted interactions by color no.

- 1: Initialize $ColorGroupArray[x] \leftarrow -1$ \triangleright x represents nodes in I and assign an initial color no -1 for all nodes
- 2: Initialize $LastNodeTime[x] \leftarrow -1$ \triangleright Record the latest time when node x was involved in an interaction and initialize with -1
- 3: **for** $event$ in I **do**
- 4: $cur_u, cur_v, cur_time = event.u, event.v, event.t$
- 5: **if** $cur_time > LastNodeTime[cur_u]$ **then**
- 6: $ColorGroupArray[cur_u] ++$
- 7: $LastNodeTime[cur_u] = cur_time$
- 8: **end if**
- 9: **if** $cur_time > LastNodeTime[cur_v]$ **then**
- 10: $ColorGroupArray[cur_v] ++$
- 11: $LastNodeTime[cur_v] = cur_time$
- 12: **end if**
- 13: $event.gn = \max(ColorGroupArray[cur_u], ColorGroupArray[cur_v])$ \triangleright gn is short for group no
- 14: **end for**
- 15: $ColorGpSeq = Sort(I)$ \triangleright sort I by an ascending order with assigned group color no.
- 16: **return** $ColorGpSeq$

A.2 CONSTRUCTION OF K-HOP TDG

The details of CONSTRUCTION OF K-HOP TDG are given in Algorithm 2. In Algorithm 2, based on the definition of TDG and $ColorGpSeq$ output by Algorithm 1, we can incrementally construct the k -step subgraphs for any two nodes in a new interaction. Specifically, we first call PREVIOUS function to get each node’s adjacent interaction in which it was involved before current interaction. Then incrementally construct a graph($ugraph$) rooted at cur_u (lines 6 to 10). Here $H_{k_subgraph}[u_preInteraction]$ stores the *subgraph* of $ugraph$ which was obtained before $ugraph$ because $u_preInteraction$ ranks ahead of current $event$ in $ColorGpSeq$. Likely, we can get $vgraph$ incrementally. At last, we call the Breadth-First-Search algorithm with the traverse $depth = K$ to obtain the k -step subgraphs for cur_event and then store it into $H_{k_subgraph}[cur_event]$.

A.3 OBTAINING AND UPDATING K-HOP TDG

As for efficiency, fast obtaining k -step subgraph of TDG for an interaction is both important for offline training and online inference. Based on Algorithm 2, for an interaction l from the collected

Algorithm 2 CONSTRUCTION OF K-HOP TDG

Require: *ColorGpSeq*: A seq of sorted interaction by group color no.
Ensure: HashTable $H_{k_subgraph}$: Map an interaction to its corresponding $k_subgraph$.

```

1: Initialize an empty HashTable  $H_{k\_subgraph}$ 
2: for  $cur\_event$  in  $ColorGpSeq$  do
3:    $cur\_u, cur\_v = cur\_event.u, cur\_event.v$ 
4:    $u\_preInteraction = Previous[cur\_u]$   $\triangleright$  Previous Function returns last interaction in which
 $cur\_u$  was involved
5:    $v\_preInteraction = Previous[cur\_v]$ 
6:   if  $u\_preInteraction$  exists then
7:      $edge_1 = (u\_preInteraction.u, cur\_u)$   $\triangleright$  The edge is defined according to the
definition of dependency graph
8:      $edge_2 = (u\_preInteraction.v, cur\_u)$ 
9:      $ugraph = edge_1 \cup edge_2$ 
10:     $ugraph = H_{k\_subgraph}[u\_preInteraction] \cup ugraph$   $\triangleright$  Incremental Update
11:   else
12:      $ugraph = cur\_u$ 
13:   end if
14:   if  $v\_preInteraction$  exists then
15:      $edge_3 = (v\_preInteraction.u, cur\_v)$ 
16:      $edge_4 = (v\_preInteraction.v, cur\_v)$ 
17:      $vgraph = edge_3 \cup edge_4$ 
18:      $vgraph = H_{k\_subgraph}[v\_preInteraction] \cup vgraph$   $\triangleright$  Incremental Update
19:   else
20:      $vgraph = cur\_v$ 
21:   end if
22:    $G_{K-step} = BFS(ugraph \cup vgraph, depth = K)$   $\triangleright$  Call Breadth-First-Search with max
 $depth = K$ 
23:    $H_{k\_subgraph}[cur\_event] = G_{K-step}$ 
24: end for
25: return  $H_{k\_subgraph}$ 

```

training interaction data, we can obtain its corresponding k-hop subgraphs directly using $H_{k_subgraph}[l]$ with time complexity $O(1)$. For a new incoming interaction i with two nodes cur_u and cur_v , we can easily reuse Algorithm 2’s code from line 3 to 22 to find $ugraph$ and $vgraph$, respectively. Then we merge them to get k-step subgraph for doing inference online. At the same time, we can incrementally update TDG by $H_{k_subgraph}[l]$. The time complexity here for updating TDG mainly depends on union between $ugraph$ and $vgraph$ in line 22 of Algorithm 2 and is $O(n + e)$ where n and e are the total number of nodes and edges for $ugraph$ and $vgraph$.

B ADDITIONAL EXPERIMENT RESULTS

Figure 8 and Figure 9 provides HIT@5,1 results in addition to the Mean Rank results in Section 5.3.2 of the main paper. HIT@n is defined as $HIT@n = \frac{\sum \delta_i}{\#of\ Test\ Interaction}$, where $\delta_i = 1$ if $rank_i \leq n$ else 0, which measures the ability of top prediction precision.

C DETAILED RELATED WORK

Graph representation learning, which is also known as graph embedding or network embedding, is a task aiming to learn low-dimensional dense vectors for vertex and edges that preserve the original graph structural information and network properties. Before the era of deep learning, conventional graph representation learning methods usually adopts dimension reduction (Belkin & Niyogi, 2002; Tenenbaum et al., 2000; Roweis & Saul, 2000) or matrix factorization (Ahmed et al., 2013) techniques. While these methods usually suffer a heavy computational cost. After the great success of deep learning methods in the field of computer vision and natural language processing, especially Skip-

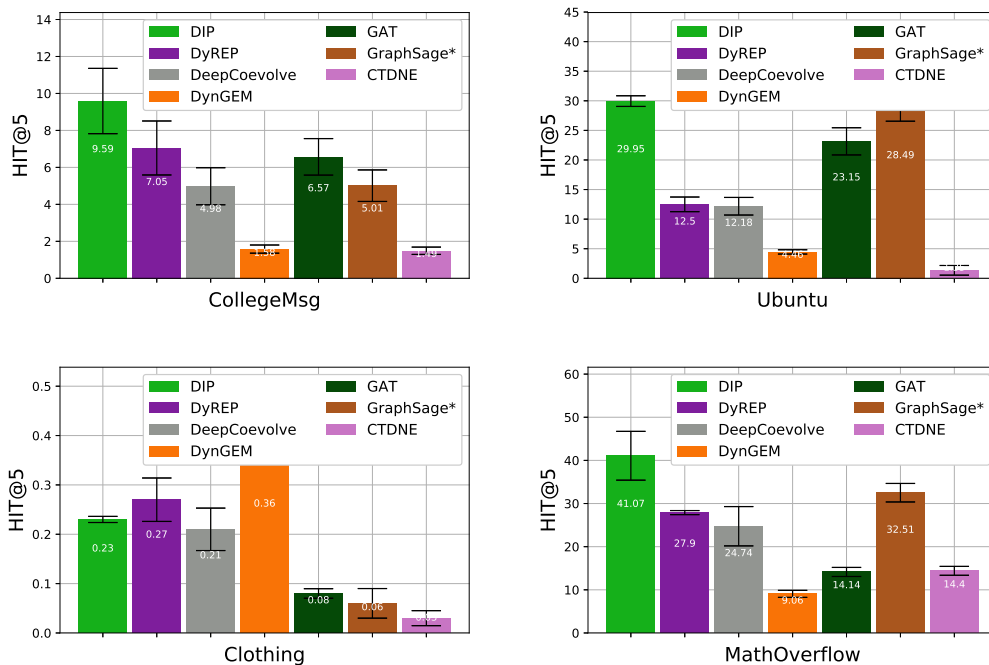


Figure 8: Hit@5 Results

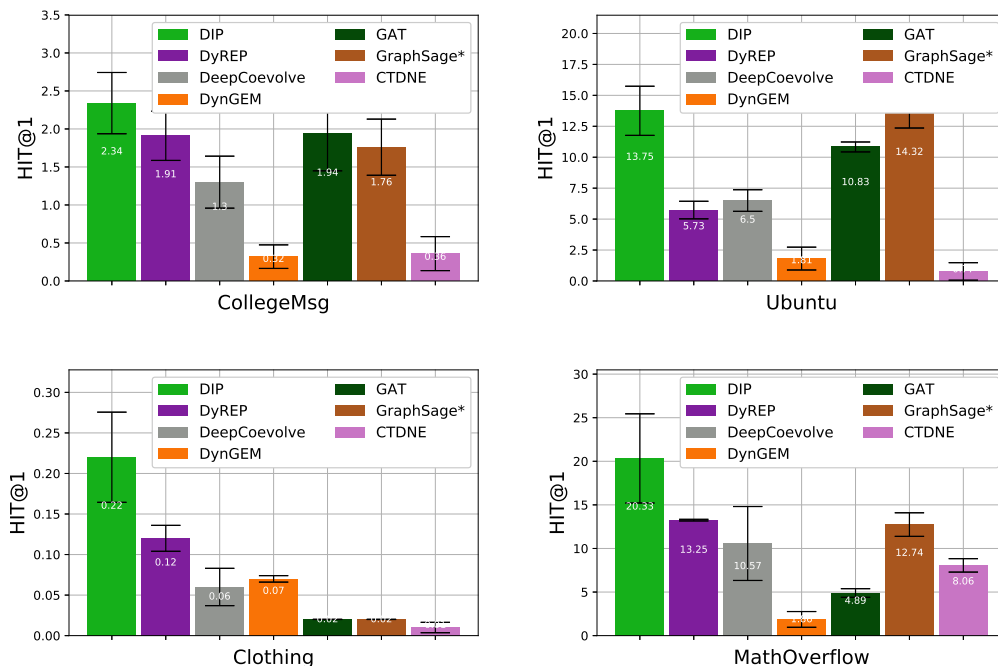


Figure 9: HIT@1 Results

gram(Mikolov et al., 2013) in word representation learning, this problem has been widely investigated by neural network methods. We review the development of this research direction from the following perspectives.

Static Graph Representation Learning: Perozzi et al. (2014) proposed the DeepWalk model utilizing random walks on graph to generate sequences of nodes, and feed them into the Skip-gram

model to learn nodes' low dimensional representation vector. Then, Tang et al. (2015) proposed the definition of first-order proximity and second-order proximity in the model named LINE, that jointly models this two different level structure information. Grover & Leskovec (2016) designed a biased random-walk sampling method to capture more flexible neighborhood information which is helpful to learn richer representations. And Wang et al. (2016; 2017) extended these models by considering high-order proximity and community structure. But these mentioned methods are somehow shallow. Another type of method is the Graph Convolution Network (Bruna et al., 2013; Kipf & Welling, 2017). They extends the convolution neural network to the graph spectral space, thus realized DEEP model on graph data. While these methods are still transductive, and can not jointly models the network structure information and the attributed on nodes or edges of graph. To overcome this problem, Hamilton et al. (2017a) proposed an inductive graph embedding framework, which divided the representation learning into two phrase: sampling and aggregation, such that it could incorporates node feature information and thus generate embedding from these features using the well-trained graph neural network for unseen nodes. Veličković et al. (2017) uses self-attention to learn weights for neighborhood, then aggregates them by the self-adapted weight. While, most of these approaches can only model static graph data.

Dynamic Graph Representation Learning: A popular approach for modeling dynamic graph data is considering the dynamics as a sequence of graph snapshot(Soundarajan et al., 2016). Zhu et al. (2016) uses an non-negative matrix factorization technique to embed social network in a temporal latent space. Zhou et al. (2018) models the specific triadics closure formation procedure over the snapshots. Goyal et al. (2018) adapts an graph autoencoder and learns a stable embedding over time. Seo et al. (2018) combines a CNN module and a RNN module to capture spatial characteristics and temporal characteristics. NetWalk(Yu et al., 2018) also learns vertex representations from sequences of snapshots, and it is designed specially for anomaly detection. In contrast, Nguyen et al. (2018) adds temporal order constraint on random walk sampling to capture evolving neighborhood. But, it can't explicitly models the rich temporal information.

Temporal Point Process: Temporal Point Process is a powerful statistical tool for modeling sequences of events with unequal interval. It has been wildly used in recent research with different intensity function: from parametric(Du et al., 2016) to recurrent neural networks(Dai et al., 2016; Mei & Eisner, 2017b), event reinforcement learning based(Xiao et al., 2017). DeepCoevolve(Dai et al., 2016) designs a recurrent neural network to capture the co-evolution dynamics of interaction data. But it just takes the interaction counterpart into the co-evolution module, which may loss temporal structural information. DyREP(Trivedi et al., 2019) divides the interactions between nodes into two different types as *communication* and *association* and models them separately. Although it uses the neighborhood of interaction counterpart to update node representation, the relevance between the different neighbors of counterpart and the node itself has not been modeled.