

# DyGMamba: Efficiently Modeling Long-Term Temporal Dependency on Continuous-Time Dynamic Graphs with State Space Models

Anonymous Author(s)

## Abstract

Learning useful representations for continuous-time dynamic graphs (CTDGs) is challenging, due to the concurrent need to span long node interaction histories and grasp nuanced temporal details. In particular, two problems emerge: (1) Encoding longer histories requires more computational resources, making it crucial for CTDG models to maintain low computational complexity to ensure efficiency; (2) Meanwhile, more powerful models are needed to identify and select the most critical temporal information within the extended context provided by longer histories. To address these problems, we propose a CTDG representation learning model named DyGMamba, originating from the popular Mamba state space model (SSM). DyGMamba first leverages a node-level SSM to encode the sequence of historical node interactions. Another time-level SSM is then employed to exploit the temporal patterns hidden in the historical graph, where its output is used to dynamically select the critical information from the interaction history. We validate DyGMamba experimentally on the dynamic link prediction task. The results show that our model achieves state-of-the-art in most cases. DyGMamba also maintains high efficiency in terms of computational resources, making it possible to capture long temporal dependencies with a limited computation budget.

## CCS Concepts

• **Do Not Use This Code** → **Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper; • **Computing methodologies** → **Machine learning**; • **Information systems**;

## Keywords

temporal graph learning, data mining

## ACM Reference Format:

Anonymous Author(s). 2025. DyGMamba: Efficiently Modeling Long-Term Temporal Dependency on Continuous-Time Dynamic Graphs with State Space Models. In *Proceedings of Temporal Graph Learning Workshop, SIGKDD International Conference on Knowledge Discovery and Data Mining 2025 (TGL Workshop, KDD 2025)*. ACM, New York, NY, USA, 20 pages. <https://doi.org/XXXXXXX.XXXXXXX>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

*TGL Workshop, KDD 2025, Toronto, Canada*

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.  
ACM ISBN 978-1-4503-XXXX-X/2018/06  
<https://doi.org/XXXXXXX.XXXXXXX>

## 1 Introduction

Dynamic graphs store node interactions in the form of links labeled with timestamps [11]. In recent years, learning dynamic graphs has gained increasing interest since it can be used to facilitate various real-world applications. Dynamic graphs can be classified into two types, i.e., discrete-time dynamic graph (DTDG) and continuous-time dynamic graph (CTDG). A DTDG is represented as a sequence of graph snapshots that are observed at regular time intervals, where all the edges in a snapshot are taken as existing simultaneously, while a CTDG consists of a stream of events where each of them is observed individually with its own timestamp. Previous works [11, 23] have indicated that CTDGs have an advantage over DTDGs in preserving temporal details, and therefore, more attention is paid to developing novel CTDG modeling approaches for dynamic graph representation learning.

Recent effort in CTDG modeling has resulted in a wide range of models. However, most of them are unable to model long-term temporal dependencies of nodes, despite the existence of abundant historical information. To solve this problem, Yu et al. [35] propose a CTDG model DyGFormer that can handle long-term node interaction histories based on Transformer [28]. Despite its ability in modeling longer histories, employing a Transformer naturally introduces excessive usage of computational resources due to its quadratic complexity. Another recent work CTAN [5] tries to capture long-term temporal dependencies by propagating graph information in a non-dissipative way over time with a graph convolution-based model. Despite the model's high efficiency, Gravina et al. [5] show that CTAN cannot capture very long histories and is surpassed by DyGFormer on the CTDGs where learning from very far away temporal information is critical. Based on these observations, we summarize the first challenge in CTDG modeling: **How to develop a model that is scalable in modeling very long-term historical interactions?** Another point worth noting is that as longer histories introduce more temporal information, more powerful models are needed to identify and select the most critical parts. This reveals another challenge: **How to effectively select critical temporal information with long node interaction histories?**

To address the first challenge, we propose to leverage a popular state space model (SSM), i.e., Mamba SSM [6] to encode the long sequence of historical node interactions. Since Mamba is proven effective and efficient in long sequence modeling [6], it maintains low computational complexity and is scalable in modeling long-term temporal dependencies. For the second challenge, we address it by learning temporal patterns of node interactions and dynamically selecting the critical temporal information based on them. The motivation can be explained by the following example. Consider a CTDG with nodes as people or songs and edges representing a person playing a song at a specific time. If a person  $u$  frequently

plays a hit song  $v$  initially but decreases the frequency later on, the time intervals between plays increase. Ignoring this pattern can lead models to incorrectly predict that  $u$  will still play  $v$  at future timestamps due to their high appearances in each other's historical interactions. If a CTDG model recognizes this pattern, it can prioritize other temporal information, such as  $u$  increasingly listening to a new song  $v'$  before  $t$ , instead of focusing on  $u, v$  interactions. Since each pattern corresponds to a specific edge, e.g.,  $(u, v, t)$ , we name these patterns as edge-specific temporal patterns.

To this end, we propose a new CTDG model named DyGMamba. DyGMamba first leverages a node-level Mamba SSM to encode historical node interactions. Another time-level Mamba SSM is then employed to exploit the edge-specific temporal patterns, where its output is used to dynamically select the critical information from the interaction history. To summarize: (1) We present DyGMamba, the first model using SSMs for CTDG representation learning; (2) DyGMamba demonstrates high efficiency and strong effectiveness in modeling long-term temporal dependencies in CTDGs; (3) Experimental results show that DyGMamba achieves new state-of-the-art on dynamic link prediction over most common CTDG datasets.

## 2 Related Work and Preliminaries

### 2.1 Related Work

*Dynamic Graph Representation Learning.* Dynamic graph representation learning methods can be categorized into two groups, i.e., DTDG and CTDG methods. DTDG methods [4, 13, 17, 22, 34] can only model DTDGs where each of them is represented as a sequence of graph snapshots. Modeling a dynamic graph as graph snapshots requires time discretization and will inevitably cause information loss [11]. To overcome this problem, recent works focus more on developing CTDG methods that treat a dynamic graph as a stream of events, where each event has its own unique timestamp. Some works [1, 27] model CTDGs by using temporal point process. Another line of works [5, 16, 31, 33] designs advanced temporal graph neural networks for CTDGs. Besides, some other methods are developed based on memory networks [14, 21], temporal random walk [10, 32] and temporal sequence modeling [2, 25, 35]. Since some real-world CTDGs heavily rely on long-term temporal information for effective learning, a number of works start to develop CTDG models that can do long range propagation of information over time [5, 35].

*State Space Models.* Transformer [28] is a de facto backbone architecture in modern deep learning. However, its self-attention mechanism results in large space and time complexity, making it unsuitable for extremely long sequence modeling [3]. To address this, many works focus on building structured state space models that scale linearly or near-linearly with input sequence length [6–9, 15, 19, 24]. Most structured SSMs exhibit linear time invariance (LTI), meaning their parameters are not input-dependent and fixed for all time-steps. Gu and Dao [6] demonstrate that LTI prevents SSMs from effectively selecting relevant information from the input context, which is problematic for tasks requiring context-aware reasoning. To solve this issue, Gu and Dao [6] proposes S6, also known as Mamba, which uses a selection mechanism to dynamically choose important information from input sequence elements.

Selection mechanism involves learning functions that map input data to SSM's parameters, making Mamba both efficient and effective in modeling language, DNA sequences, and audio.

### 2.2 Preliminaries

*CTDG and Task Formulation.* We define CTDG and dynamic link prediction as follows.

**DEFINITION 1 (CONTINUOUS-TIME DYNAMIC GRAPH).** Let  $\mathcal{N}$  and  $\mathcal{T}$  denote a set of nodes and timestamps, respectively. A CTDG is a sequence of  $|\mathcal{G}|$  chronological interactions  $\mathcal{G} = \{(u_i, v_i, t_i)\}_{i=1}^{|\mathcal{G}|}$  with  $0 \leq t_1 \leq t_2 \leq \dots \leq t_{|\mathcal{G}|}$ , where  $u_i, v_i \in \mathcal{N}$  are the source and destination node of the  $i$ -th interaction happening at  $t_i \in \mathcal{T}$ , respectively. Each node  $u \in \mathcal{N}$  can be equipped with a node feature  $\mathbf{x}_u \in \mathbb{R}^{d_N}$ , and each interaction  $(u, v, t)$  can be associated with a link (edge) feature  $\mathbf{e}_{u,v}^t \in \mathbb{R}^{d_E}$ . If  $\mathcal{G}$  is not attributed, we set node and link features to zero vectors.

**DEFINITION 2 (DYNAMIC LINK PREDICTION).** Given a CTDG  $\mathcal{G}$ , a source node  $u \in \mathcal{N}$ , a destination node  $v \in \mathcal{N}$ , a timestamp  $t \in \mathcal{T}$ , and all the interactions before  $t$ , i.e.,  $\{(u_i, v_i, t_i) | t_i < t, (u_i, v_i, t_i) \in \mathcal{G}\}$ , dynamic link prediction aims to predict whether the interaction  $(u, v, t)$  exists.

*S4 and Mamba SSM.* S4 and Mamba [6, 8] are inspired by a continuous system which can be described as  $\mathbf{z}(\tau)' = \mathbf{A}\mathbf{z}(\tau) + \mathbf{B}q(\tau)$  and  $\mathbf{r}(\tau) = \mathbf{C}\mathbf{z}(\tau)$ .  $q(\tau) \in \mathbb{R}$  and  $\mathbf{r}(\tau) \in \mathbb{R}$  are the 1-dimensional input and output over time  $\tau^1$ , respectively.  $\mathbf{A} \in \mathbb{R}^{d_1 \times d_1}$ ,  $\mathbf{B} \in \mathbb{R}^{d_1 \times 1}$ ,  $\mathbf{C} \in \mathbb{R}^{1 \times d_1}$  are three parameters deciding the system. Based on it, both S4 and Mamba include a time-scale parameter  $\Delta \in \mathbb{R}$  and discretize all the parameters to adapt to a discretized system

$$\begin{aligned} \mathbf{z}_\tau &= \bar{\mathbf{A}}\mathbf{z}_{\tau-1} + \bar{\mathbf{B}}p_\tau, \quad q_\tau = \mathbf{C}\mathbf{z}_\tau; \quad \bar{\mathbf{A}} = \exp(\Delta\mathbf{A}), \\ \bar{\mathbf{B}} &= (\Delta\mathbf{A})^{-1}(\exp(\Delta\mathbf{A}) - \mathbf{I})\Delta\mathbf{B}. \end{aligned} \quad (1)$$

Here,  $\tau$  is also discretized to denote the position of a sequence element. Given Eq. 1, sequence processing with S4 and Mamba can be written as computing an output sequence with convolution

$$\mathbf{q} = \mathbf{p} * \bar{\mathbf{K}}_{\text{SSM}}, \quad \text{where } \bar{\mathbf{K}}_{\text{SSM}} = [\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^{|\mathbf{p}|-1}\bar{\mathbf{B}}] \in \mathbb{R}^{|\mathbf{p}|-1}. \quad (2)$$

$\mathbf{p} \in \mathbb{R}^{|\mathbf{p}|}$  and  $\mathbf{q} \in \mathbb{R}^{|\mathbf{p}|}$  are input and output sequences, where  $|\mathbf{p}|$  is the sequence length of  $\mathbf{p}$ .  $*$  denotes the element-wise multiplication. When the dimension size of each element  $p_\tau$  in  $\mathbf{p}$  becomes higher (i.e.,  $p_\tau \in \mathbb{R}^{d_2}$  is a vector and  $d_2 > 1$ ), both S4 and Mamba are in a Single-Input Single-Output (SISO) fashion, processing each input dimension in parallel with the same set of parameters. We follow Gu and Dao [6] and denote the computation in Eq. 2 on the input sequences with vector elements as a function  $\text{SSM}_{\bar{\mathbf{A}}, \bar{\mathbf{B}}, \mathbf{C}}(\cdot)^2$ . Different from S4 which uses same parameters to process each element, Mamba changes its parameters into input-dependent by employing several trainable linear layers to map input into  $\bar{\mathbf{B}}, \mathbf{C}$  and  $\Delta$ . The system is evolving as it processes different elements in the input sequence, making Mamba time-variant and suitable for modeling temporal sequences.

<sup>1</sup>We use  $\tau$  rather than  $t$  to indicate time in a continuous system to distinguish from the time in CTDGs.

<sup>2</sup>Input and output of  $\text{SSM}_{\bar{\mathbf{A}}, \bar{\mathbf{B}}, \mathbf{C}}(\cdot)$  are matrices where each row is a vector corresponding to an element. See App. I for more details of SISO and the function.

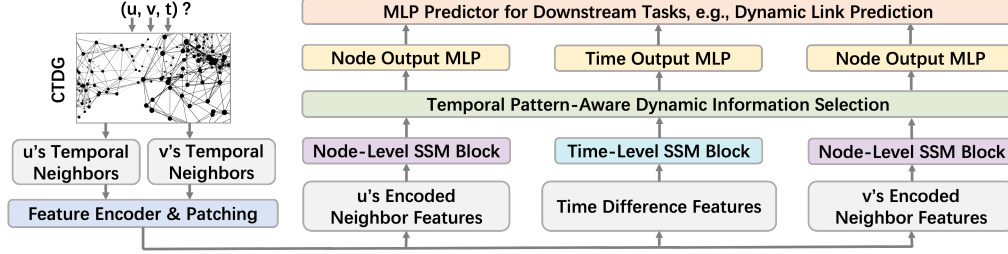


Figure 1: Model overview of DyGMamba.

### 3 DyGMamba

Fig. 1 illustrates the overview of DyGMamba. Given a potential interaction  $(u, v, t)$ , CTDG models are asked to predict whether it exists or not. DyGMamba extracts the historical one-hop interactions of node  $u$  and  $v$  before timestamp  $t$  from the CTDG  $\mathcal{G}$  and gets two interaction sequences  $S_u^t = \{(u, u', t') | t' < t, (u, u', t') \in \mathcal{G}\} \cup \{(u', u, t') | t' < t, (u', u, t') \in \mathcal{G}\}$  and  $S_v^t = \{(v, v', t') | t' < t, (v, v', t') \in \mathcal{G}\} \cup \{(v', v, t') | t' < t, (v', v, t') \in \mathcal{G}\}$  containing  $u$ 's and  $v$ 's one-hop temporal neighbors  $Nei_u^t = \{(u', t') | (u, u', t') \in \mathcal{G}, t' < t\}$  and  $Nei_v^t = \{(v', t') | (v, v', t') \in \mathcal{G}, t' < t\}$  (link features are omitted for clarity). Then it encodes the neighbors in  $Nei_u^t$  and  $Nei_v^t$  to get two sequences of encoded neighbor representations for  $u$  and  $v$ . To learn the edge-specific temporal pattern of  $(u, v, t)$ , we find the interactions between  $u$  and  $v$  before  $t$ , compute the time difference between each pair of neighboring interactions, and build a sequence of time differences  $S_{u,v}^t$ . Finally, DyGMamba dynamically selects critical information by assigning different weights to different encoded neighbors based on the learned temporal pattern, and uses the selected information to achieve link prediction.

#### 3.1 Learning One-Hop Temporal Neighbors

**Encode Neighbor Features.** Given one-hop temporal neighbors  $Nei_u^t$  of the source node  $u$ , we sort them in the chronological order and append  $(u, t)$  at the end to form a sequence of  $|Nei_u^t| + 1$  temporal nodes. We take their node features from the dataset and stack them into a feature matrix  $\tilde{X}_u^t \in \mathbb{R}^{(|Nei_u^t|+1) \times d_N}$ . Similarly, we build a link feature matrix  $\tilde{E}_u^t \in \mathbb{R}^{(|Nei_u^t|+1) \times d_E}$ . To incorporate temporal information, we encode the time difference between  $u$  and each one-hop temporal neighbor  $(u', t')$  using the time encoding function introduced in TGAT [33]:  $\sqrt{1/d_T}[\cos(\omega_1(t - t') + \phi_1), \dots, \cos(\omega_d(t - t') + \phi_d)]$ .  $d_T$  is the dimension of time representation.  $\omega_1 \dots \omega_{d_T}$  and  $\phi_1 \dots \phi_{d_T}$  are trainable parameters. The time feature of  $u$ 's temporal neighbors are denoted as  $\tilde{T}_u^t \in \mathbb{R}^{(|Nei_u^t|+1) \times d_T}$ . We follow the same way to get  $\tilde{X}_v^t \in \mathbb{R}^{(|Nei_v^t|+1) \times d_N}$ ,  $\tilde{E}_v^t \in \mathbb{R}^{(|Nei_v^t|+1) \times d_E}$  and  $\tilde{T}_v^t \in \mathbb{R}^{(|Nei_v^t|+1) \times d_T}$  for  $v$ 's temporal neighbors. Following Tian et al. [25], we also consider the historical node interaction frequencies in the interaction sequences  $S_u^t$  and  $S_v^t$  of source  $u$  and destination  $v$ . For example, assume the interacted nodes of  $u$  and  $v$  (arranged in chronological order) are  $\{a, v, a\}$  and  $\{b, b, u, a\}$ , the appearing frequencies of  $a, b$  in  $u/v$ 's historical interactions are 2/1, 0/2, respectively. And the frequency of the interaction involving  $u$  and  $v$  is 1. Thus, the node interaction frequency features

of  $u$  and  $v$  are written as  $\tilde{F}_u^t = [[2, 1], [1, 1], [2, 1], [0, 1]]^T$  and  $\tilde{F}_v^t = [[0, 2], [0, 2], [1, 1], [2, 1], [0, 1]]^T$ , respectively. Note that the last elements  $([0, 1]$  and  $[0, 1])$  in  $\tilde{F}_u^t$  and  $\tilde{F}_v^t$  correspond to the appended  $(u, t)$  and  $(v, t)$  not existing in the observed histories. We initialize them with  $[0, \text{number of historical interactions between } u, v]$ . An encoding multilayer perceptron (MLP)  $f(\cdot) : \mathbb{R} \rightarrow \mathbb{R}^{d_F}$  is employed to encode these features into representations:  $\tilde{F}_u^t = f(\tilde{F}_u^t[: , 0]) + f(\tilde{F}_u^t[: , 1]) \in \mathbb{R}^{(|Nei_u^t|+1) \times d_F}$ ,  $\tilde{F}_v^t = f(\tilde{F}_v^t[: , 0]) + f(\tilde{F}_v^t[: , 1]) \in \mathbb{R}^{(|Nei_v^t|+1) \times d_F}$ .

**Patching Neighbors.** We employ the patching technique proposed by [35] to save computational resources when dealing with a large number of temporal neighbors. We treat  $p$  temporally adjacent neighbors as a patch and flatten their features. For example, with patching,  $\tilde{X}_u^t \in \mathbb{R}^{(|Nei_u^t|+1) \times d_N}$  results in a new patched feature matrix  $X_u^t \in \mathbb{R}^{(|Nei_u^t|+1)/p \times (p \cdot d_N)}$  (we pad  $\tilde{X}_u^t$  with zero-valued features when  $|Nei_u^t| + 1$  cannot be divided by  $p$ ). Similarly, we get  $E_\theta^t \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times (p \cdot d_E)}$ ,  $T_\theta^t \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times (p \cdot d_T)}$  and  $F_\theta^t \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times (p \cdot d_F)}$  ( $\theta$  is either  $u$  or  $v$ ). Each row of a feature matrix corresponds to an element of the input sequence sent into an SSM later. Recall that SSMs process sequences in a recurrent way. Patching decreases the length of the sequence by roughly  $p$  times, making great contribution in saving computational resources.

**Node-Level SSM Block.** We first map the padded features of  $u$ 's and  $v$ 's one-hop temporal neighbors to the same dimension  $d$ , i.e.,  $X_\theta^t := f_N(X_\theta^t)$ ,  $E_\theta^t := f_E(E_\theta^t)$ ,  $T_\theta^t := f_T(T_\theta^t)$ ,  $F_\theta^t := f_F(F_\theta^t)$ .  $f_N(\cdot) : \mathbb{R}^{p \cdot d_N} \rightarrow \mathbb{R}^d$ ,  $f_E(\cdot) : \mathbb{R}^{p \cdot d_E} \rightarrow \mathbb{R}^d$ ,  $f_T(\cdot) : \mathbb{R}^{p \cdot d_T} \rightarrow \mathbb{R}^d$ ,  $f_F(\cdot) : \mathbb{R}^{p \cdot d_F} \rightarrow \mathbb{R}^d$  are four MLPs for different types of neighbor features. We take the concatenation of them as the encoded representations of the temporal neighbors, i.e.,  $H_\theta^t = X_\theta^t \| E_\theta^t \| T_\theta^t \| F_\theta^t \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times 4d}$ . We input  $H_u^t$  and  $H_v^t$  separately into a node-level SSM block to learn the temporal dependencies of temporal neighbors. The node-level SSM block consists of  $l_N$  layers, where each layer is defined as follows (Eq. 3-4). First, we input  $H_\theta^t$  into a

$$B_1 = H_\theta^t W_{B_1} \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times d_{SSM}}, \quad (3a)$$

$$C_1 = H_\theta^t W_{C_1} \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times d_{SSM}}, \quad (3b)$$

$$\Delta_1 = \sigma(\text{Broad}_{4d}(H_\theta^t W_{\Delta_1}) + \text{Par}_{\Delta_1}) \in \mathbb{R}^{(|Nei_\theta^t|+1)/p \times 4d}, \quad (3c)$$

$$\bar{A}_1 = \exp(\Delta_1 A_1), \quad \bar{B}_1 = (\Delta_1 A_1)^{-1} (\exp(\Delta_1 A_1) - I) \Delta_1 B_1; \quad (3d)$$

$$H_\theta^t := H_\theta^t + \text{SSM}_{\bar{A}_1, \bar{B}_1, C_1}(H_\theta^t). \quad (3e)$$



$\bar{\mathbf{A}}_1, \bar{\mathbf{B}}_1 \in \mathbb{R}^{\lceil(|Nei_\theta^t|+1)/p\rceil \times 4d \times d_{SSM}}$  are discretized parameters.  $\mathbf{W}_{B_1}, \mathbf{W}_{C_1} \in \mathbb{R}^{4d \times d_{SSM}}$  and  $\mathbf{W}_{\Delta_1} \in \mathbb{R}^{4d \times 1}$ .  $\text{Par}_{\Delta_1} \in \mathbb{R}^{\lceil(|Nei_\theta^t|+1)/p\rceil \times 4d}$  is a parameter defined by Gu and Dao [6].  $\text{Broad}_{4d}(\cdot)$  is a function that copies its vector input for  $4d$  times to form a matrix with  $4d$  identical columns (following the definition in [6]).  $\sigma(\cdot)$  is the Softplus function.  $\mathbf{I}$  is an identity matrix. Then we use an MLP  $f_{\text{node}}(\cdot) : \mathbb{R}^{4d} \rightarrow \mathbb{R}^{4d}$  on SSM's output

$$\mathbf{H}_\theta^t := \mathbf{H}_\theta^t + f_{\text{node}}(\text{LayerNorm}(\mathbf{H}_\theta^t)). \quad (4)$$

After  $l_N$  layers, we have  $\mathbf{H}_u^t$  and  $\mathbf{H}_v^t$  that contain the encoded information of all one-hop temporal neighbors for the entities  $u$  and  $v$  as well as the information of themselves. Since we sort temporal neighbors chronologically, our node-level SSM block can directly learn the temporal dynamics for graph forecasting.

### 3.2 Learning from Temporal Patterns

**Time-Level SSM Block.** To capture edge-specific temporal patterns, we use another time-level SSM block consisting of  $l_T$  layers. We first find out  $k$  temporally nearest historical interactions between  $u$  and  $v$  before  $t$  and sort them in the chronological order, i.e.,  $\{(u, v, t_0), \dots, (u, v, t_{k-1}) | t_0 < \dots < t_{k-1} < t\}$ . Then we construct a timestamp sequence  $\{t_0, t_1, \dots, t_{k-1}, t\}$  based on these interactions and the prediction timestamp  $t$ . We compute the time difference between each neighboring pair of them and further get a time difference sequence  $\{t_1 - t_0, t_2 - t_1, \dots, t - t_{k-1}\}$ , representing the change of time interval between two identical interactions. Each element in this sequence is input into the time encoding function stated above to get a edge-specific (specific to the edge  $(u, v, t)$ ) time feature. The features are stacked into a feature matrix  $\mathbf{H}_{u,v}^t \in \mathbb{R}^{k \times d_T}$  and mapped by an MLP  $f_{\text{map1}}(\cdot) : \mathbb{R}^{d_T} \rightarrow \mathbb{R}^{d_Y}$  ( $\gamma \in [0, 1]$  is a hyperparameter), i.e.,  $\mathbf{H}_{u,v}^t := f_{\text{map1}}(\mathbf{H}_{u,v}^t)$ . A time-level SSM layer takes  $\mathbf{H}_{u,v}^t$  as input and computes

$$\mathbf{B}_2 = \mathbf{H}_{u,v}^t \mathbf{W}_{B_2} \in \mathbb{R}^{k \times d_{SSM}}, \quad \mathbf{C}_2 = \mathbf{H}_{u,v}^t \mathbf{W}_{C_2} \in \mathbb{R}^{k \times d_{SSM}}, \quad (5a)$$

$$\Delta_2 = \text{Softplus}(\text{Broad}_{\gamma d}(\mathbf{H}_{u,v}^t \mathbf{W}_{\Delta_2}) + \text{Par}_{\Delta_2}) \in \mathbb{R}^{k \times \gamma d}, \quad (5b)$$

$$\bar{\mathbf{A}}_2 = \exp(\Delta_2 \mathbf{A}_2), \quad \bar{\mathbf{B}}_2 = (\Delta_2 \mathbf{A}_2)^{-1} (\exp(\Delta_2 \mathbf{A}_2) - \mathbf{I}) \Delta_2 \mathbf{B}_2; \quad (5c)$$

$$\mathbf{H}_{u,v}^t := \mathbf{H}_{u,v}^t + \text{SSM}_{\bar{\mathbf{A}}_2, \bar{\mathbf{B}}_2, \mathbf{C}_2}(\mathbf{H}_{u,v}^t). \quad (5d)$$

$\mathbf{W}_{B_2}, \mathbf{W}_{C_2} \in \mathbb{R}^{\gamma d \times d_{SSM}}$  and  $\mathbf{W}_{\Delta_2} \in \mathbb{R}^{\gamma d \times 1}$ .  $\bar{\mathbf{A}}_2, \bar{\mathbf{B}}_2 \in \mathbb{R}^{k \times \gamma d \times d_{SSM}}$  are discretized parameters.  $\text{Par}_{\Delta_2} \in \mathbb{R}^{k \times \gamma d}$  is a parameter defined as same as  $\text{Par}_{\Delta_1}$ . In practice, we set  $k$  to a number much smaller than  $|Nei_\theta^t|$ , e.g., 10. This ensures that time-level SSM will not incur huge computational burden and the model focuses more on the recent histories. Note that we cannot always find  $k$  recent historical interactions between each pair of nodes, leading to varying lengths of time difference sequences for different  $(u, v, t)$  in a batch of data. To enable batch processing, we set the time difference without a found historical interaction to a very large number  $10^{10}$ . For example, if  $k = 2$ , and for  $(u, v, t)$  we can only find  $(u, v, t_0)$ . The time difference sequence will be  $\{10^{10}, t - t_0\}$ .  $10^{10}$  is much larger than  $t - t_0$ , indicating that  $u$  and  $v$  have not had an interaction for an extremely long time, same as existing no historical interaction. We further explain why we use SSM to learn temporal patterns in App. J.

**Dynamic Information Selection with Temporal Patterns.** After the time-level SSM block, we compute a compressed representation to represent the edge-specific temporal pattern by averaging over  $k$  encoded time intervals:  $\mathbf{h}_{u,v}^t = \text{MeanPooling}(\mathbf{H}_{u,v}^t)$ . As a result, we have  $\mathbf{h}_{u,v}^t \in \mathbb{R}^{d_Y}$  to represent the temporal pattern specific to the edge  $(u, v, t)$ . To leverage learned temporal pattern, we use it to dynamically select the information from the encoded temporal neighbors  $\mathbf{H}_\theta^t$

$$\hat{\mathbf{h}}_{u,v}^t = f_{\text{map2}}(\mathbf{h}_{u,v}^t) \in \mathbb{R}^{4d}, \quad (6a)$$

$$\hat{\mathbf{h}}_\theta^t = \mathbf{w}_{\text{agg}}^\top \mathbf{H}_\theta^t \in \mathbb{R}^{4d}; \mathbf{w}_{\text{agg}} = f_{\text{map3}}(\mathbf{H}_\theta^t) \in \mathbb{R}^{\lceil(|Nei_\theta^t|+1)/p\rceil}, \quad (6b)$$

$$\alpha_u = f'(\hat{\mathbf{h}}_\theta^t) * \hat{\mathbf{h}}_{u,v}^t \in \mathbb{R}^{4d}, \quad \alpha_v = f'(\hat{\mathbf{h}}_\theta^t) * \hat{\mathbf{h}}_{u,v}^t \in \mathbb{R}^{4d}, \quad (6c)$$

$$\mathbf{h}_\theta^t = \beta_\theta^\top \mathbf{H}_\theta^t; \beta_\theta = \text{Softmax}(\mathbf{H}_\theta^t \alpha_\theta) \in \mathbb{R}^{\lceil(|Nei_\theta^t|+1)/p\rceil}. \quad (6d)$$

$f_{\text{map2}}(\cdot) : \mathbb{R}^{d_Y} \rightarrow \mathbb{R}^{4d}$  and  $f_{\text{map3}}(\cdot) : \mathbb{R}^{4d} \rightarrow \mathbb{R}^1$  are two mapping MLPs.  $f'(\cdot) : \mathbb{R}^{4d} \rightarrow \mathbb{R}^{4d}$  is another MLP introducing training parameters. Note that  $\alpha_u/\alpha_v$  is computed by considering both the edge-specific temporal pattern and the opposite node  $v/u$ . In the node-level SSM block, we separately model the one-hop temporal neighbors of each node  $\theta$ , making it hard to connect  $u$  and  $v$ . Computing  $\alpha_\theta$  as Eq. 6c helps to strengthen the connection between both nodes and meanwhile incorporates the learned temporal pattern.  $\beta_\theta$  is derived by transforming the queried results based on  $\alpha_\theta$  into weights. It is then used to compute a weighted-sum of all temporal neighbors for representing  $\theta$  at  $t$ , i.e.,  $\mathbf{h}_\theta^t$ . The neighbors assigned with greater weights from  $\beta_\theta$  are selected as more critical and will contribute more to  $\mathbf{h}_\theta^t$ . Finally, we output the representations of  $u, v$  and the edge-specific temporal pattern by employing two output MLPs  $f_{\text{out1}}(\cdot) : \mathbb{R}^{4d} \rightarrow \mathbb{R}^{d_N}$  and  $f_{\text{out2}}(\cdot) : \mathbb{R}^{d_Y} \rightarrow \mathbb{R}^{d_N}$ , i.e.,  $\mathbf{h}_\theta^t := f_{\text{out1}}(\mathbf{h}_\theta^t) \in \mathbb{R}^{d_N}$ ,  $\mathbf{h}_{u,v}^t := f_{\text{out2}}(\mathbf{h}_{u,v}^t) \in \mathbb{R}^{d_N}$ .

### 3.3 Leveraging Learned Representations for Link Prediction

We leverage  $\mathbf{h}_\theta^t$  and  $\mathbf{h}_{u,v}^t$  for dynamic link prediction. We employ a prediction MLP, i.e.,  $f_{\text{LP}}(\cdot) : \mathbb{R}^{3d_N} \rightarrow \mathbb{R}$ , as the predictor. The probability of existing a link  $(u, v, t)$  is computed as  $y'(u, v, t) = \text{Sigmoid}(f_{\text{LP}}(\mathbf{h}_u^t \|\mathbf{h}_v^t \|\mathbf{h}_{u,v}^t))$ . For model parameter learning, we use the following loss function

$$\mathcal{L} = -\frac{1}{2M} \sum_{2M} \left( y(u, v, t) \log(y'(u, v, t)) + (1 - y(u, v, t)) \log(1 - y'(u, v, t)) \right). \quad (7)$$

$y(u, v, t)$  is the ground truth label denoting the existence of  $(u, v, t)$  (1/0 means existing/non-existing).  $M$  is the total number of edges existing in the training data (positive edges). We follow previous work [35] and randomly sample one negative edge for each positive edge during training. Therefore, in total we have  $2M$  edges considered in our loss  $\mathcal{L}$ .

## 4 Experiments

In Sec. 4.2.1, we validate DyGMamba's ability in CTDG representation learning by comparing it with baseline methods on dynamic

**Table 1: AP of transductive dynamic link prediction. The best and the second best results are marked as bold and underlined, respectively. CTAN cannot be trained before 120 hours timeout on Social Evo. so is ranked bottom on this dataset.**

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	EdgeBank	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	LastFM	70.95 $\pm$ 2.94	71.85 $\pm$ 2.44	73.30 $\pm$ 0.18	75.31 $\pm$ 5.62	86.60 $\pm$ 0.11	79.29 $\pm$ 0.00	76.62 $\pm$ 1.83	75.56 $\pm$ 0.19	<u>92.95 <math>\pm</math> 0.14</u>	86.44 $\pm$ 0.80	<b>93.35 <math>\pm</math> 0.20</b>
	Enron	84.85 $\pm$ 3.13	79.80 $\pm$ 2.28	70.76 $\pm$ 1.05	86.98 $\pm$ 1.05	89.50 $\pm$ 0.10	83.53 $\pm$ 0.00	85.41 $\pm$ 0.71	82.13 $\pm$ 0.30	92.42 $\pm$ 0.11	<u>92.52 <math>\pm</math> 1.20</u>	<b>92.65 <math>\pm</math> 0.12</b>
	MOOC	81.04 $\pm$ 0.83	81.50 $\pm$ 0.77	85.71 $\pm$ 0.20	<u>89.15 <math>\pm</math> 1.69</u>	80.30 $\pm$ 0.43	57.97 $\pm$ 0.00	83.89 $\pm$ 0.86	82.80 $\pm$ 0.15	87.66 $\pm$ 0.48	84.71 $\pm$ 2.85	<b>89.21 <math>\pm</math> 0.08</b>
	Reddit	98.31 $\pm$ 0.06	98.18 $\pm$ 0.03	98.57 $\pm$ 0.01	98.65 $\pm$ 0.04	99.11 $\pm$ 0.01	94.86 $\pm$ 0.00	97.78 $\pm$ 0.02	97.31 $\pm$ 0.01	<u>99.22 <math>\pm</math> 0.01</u>	97.21 $\pm$ 0.84	<b>99.32 <math>\pm</math> 0.01</b>
	Wikipedia	96.51 $\pm$ 0.22	94.88 $\pm$ 0.29	96.88 $\pm$ 0.06	98.45 $\pm$ 0.10	98.77 $\pm$ 0.01	90.37 $\pm$ 0.00	97.75 $\pm$ 0.04	97.22 $\pm$ 0.02	<u>99.03 <math>\pm</math> 0.03</u>	96.61 $\pm$ 0.79	<b>99.15 <math>\pm</math> 0.02</b>
	UCI	89.28 $\pm$ 1.02	66.11 $\pm$ 2.75	79.40 $\pm$ 0.61	92.33 $\pm$ 0.64	95.13 $\pm$ 0.23	76.20 $\pm$ 0.00	86.63 $\pm$ 1.30	93.15 $\pm$ 0.41	<u>95.74 <math>\pm</math> 0.17</u>	76.64 $\pm$ 4.11	<b>95.91 <math>\pm</math> 0.15</b>
	Social Evo.	89.88 $\pm$ 0.40	88.39 $\pm$ 0.69	93.33 $\pm$ 0.06	93.45 $\pm$ 0.29	84.90 $\pm$ 0.11	74.95 $\pm$ 0.00	93.82 $\pm$ 0.19	93.36 $\pm$ 0.06	<u>94.63 <math>\pm</math> 0.07</u>	Timeout	<b>94.77 <math>\pm</math> 0.01</b>
Avg. Rank		8.29	9.29	7.00	4.29	6.00	9.43	5.57	6.43	<u>2.43</u>	6.29	<b>1.00</b>
Historical	LastFM	74.38 $\pm$ 6.27	71.85 $\pm$ 2.91	71.60 $\pm$ 0.36	75.03 $\pm$ 6.90	69.93 $\pm$ 0.33	73.03 $\pm$ 0.00	71.02 $\pm$ 2.07	72.28 $\pm$ 0.37	81.51 $\pm$ 0.14	<u>82.29 <math>\pm</math> 0.94</u>	<b>83.02 <math>\pm</math> 0.16</b>
	Enron	69.13 $\pm$ 1.66	72.58 $\pm$ 1.83	64.24 $\pm$ 1.24	74.31 $\pm$ 0.99	65.40 $\pm$ 0.36	76.53 $\pm$ 0.00	72.39 $\pm$ 0.61	77.35 $\pm$ 1.22	76.93 $\pm$ 0.76	<u>77.24 <math>\pm</math> 1.53</u>	<b>77.77 <math>\pm</math> 1.32</b>
	MOOC	78.62 $\pm$ 2.43	75.14 $\pm$ 2.86	82.83 $\pm$ 0.71	85.46 $\pm$ 2.32	74.46 $\pm$ 0.53	60.71 $\pm$ 0.00	78.09 $\pm$ 1.24	77.09 $\pm$ 0.83	<u>85.65 <math>\pm</math> 0.88</u>	67.73 $\pm$ 2.08	<b>85.89 <math>\pm</math> 0.94</b>
	Reddit	79.96 $\pm$ 0.30	79.40 $\pm$ 0.30	79.78 $\pm$ 0.25	81.05 $\pm$ 0.32	80.96 $\pm$ 0.28	73.59 $\pm$ 0.00	77.38 $\pm$ 0.20	78.39 $\pm$ 0.40	81.63 $\pm$ 1.08	<b>89.77 <math>\pm</math> 2.28</b>	<u>81.80 <math>\pm</math> 1.52</u>
	Wikipedia	81.16 $\pm$ 0.73	79.46 $\pm$ 0.95	87.31 $\pm$ 0.36	87.31 $\pm$ 0.25	66.77 $\pm$ 6.62	73.35 $\pm$ 0.00	86.12 $\pm$ 1.69	<u>90.74 <math>\pm</math> 0.06</u>	70.13 $\pm$ 11.02	<b>95.91 <math>\pm</math> 0.10</b>	81.77 $\pm$ 1.20
	UCI	74.77 $\pm$ 5.35	55.89 $\pm$ 2.83	66.78 $\pm$ 0.77	<u>81.32 <math>\pm</math> 1.26</u>	64.69 $\pm$ 1.78	65.50 $\pm$ 0.00	74.62 $\pm$ 2.70	<b>83.88 <math>\pm</math> 1.06</b>	80.44 $\pm$ 1.16	76.62 $\pm$ 0.33	81.03 $\pm$ 1.09
	Social Evo.	91.26 $\pm$ 2.47	92.86 $\pm$ 0.90	95.31 $\pm$ 0.30	93.84 $\pm$ 1.68	85.65 $\pm$ 0.11	80.57 $\pm$ 0.00	95.93 $\pm$ 0.63	95.30 $\pm$ 0.34	<u>97.05 <math>\pm</math> 0.16</u>	Timeout	<b>97.35 <math>\pm</math> 0.52</b>
Avg. Rank		6.57	8.14	6.57	4.14	9.29	8.71	7.00	4.71	<u>4.00</u>	4.71	<b>2.14</b>
Inductive	LastFM	62.63 $\pm$ 6.89	62.49 $\pm$ 3.04	71.16 $\pm$ 0.33	65.09 $\pm$ 7.05	67.38 $\pm$ 0.57	<u>75.49 <math>\pm</math> 0.00</u>	62.76 $\pm$ 0.81	67.87 $\pm$ 0.37	72.60 $\pm$ 0.06	<b>80.06 <math>\pm</math> 0.85</b>	73.63 $\pm$ 0.54
	Enron	69.51 $\pm$ 1.06	66.78 $\pm$ 2.21	63.16 $\pm$ 0.59	73.27 $\pm$ 0.58	75.08 $\pm$ 0.81	73.89 $\pm$ 0.00	70.98 $\pm$ 0.96	74.12 $\pm$ 0.65	<u>78.22 <math>\pm</math> 0.80</u>	72.02 $\pm$ 2.64	<b>80.86 <math>\pm</math> 1.24</b>
	MOOC	68.56 $\pm$ 1.49	61.48 $\pm$ 0.96	76.96 $\pm$ 0.89	77.59 $\pm$ 1.83	73.55 $\pm$ 0.36	49.43 $\pm$ 0.00	76.35 $\pm$ 1.41	74.24 $\pm$ 0.75	<u>80.99 <math>\pm</math> 0.88</u>	64.93 $\pm$ 3.31	<b>81.11 <math>\pm</math> 0.63</b>
	Reddit	86.93 $\pm$ 0.21	86.06 $\pm$ 0.36	89.93 $\pm$ 0.10	88.12 $\pm$ 0.13	<b>91.89 <math>\pm</math> 0.18</b>	85.48 $\pm$ 0.00	86.97 $\pm$ 0.26	85.37 $\pm$ 0.26	91.06 $\pm$ 0.60	90.99 $\pm$ 2.19	<u>91.15 <math>\pm</math> 0.54</u>
	Wikipedia	74.78 $\pm$ 0.56	70.55 $\pm$ 1.22	86.77 $\pm$ 0.29	85.80 $\pm$ 0.15	69.27 $\pm$ 7.07	80.63 $\pm$ 0.00	72.54 $\pm$ 4.69	<u>88.54 <math>\pm</math> 0.20</u>	62.00 $\pm$ 14.00	<b>94.15 <math>\pm</math> 0.08</b>	79.86 $\pm$ 2.18
	UCI	66.02 $\pm$ 1.28	54.64 $\pm$ 2.52	67.63 $\pm$ 0.51	70.34 $\pm$ 0.72	64.08 $\pm$ 1.06	57.43 $\pm$ 0.00	<u>73.49 <math>\pm</math> 2.21</u>	<b>79.57 <math>\pm</math> 0.61</b>	70.51 $\pm$ 1.83	66.25 $\pm$ 0.51	71.95 $\pm$ 2.51
	Social Evo.	91.08 $\pm$ 3.29	92.84 $\pm$ 0.98	95.20 $\pm$ 0.30	94.58 $\pm$ 1.52	88.50 $\pm$ 0.13	83.69 $\pm$ 0.00	96.14 $\pm$ 0.63	95.11 $\pm$ 0.32	<u>97.62 <math>\pm</math> 0.12</u>	Timeout	<b>97.68 <math>\pm</math> 0.42</b>
Avg. Rank		8.29	9.57	5.43	5.43	6.57	7.57	6.00	5.00	<u>4.00</u>	5.71	<b>2.43</b>

link prediction<sup>3</sup>. We show the effectiveness of model components by conducting ablation studies (Sec. 4.2.2) and analysis on synthetic datasets (Sec. 4.2.3). In Sec. 4.3.1, we show DyGMamba’s efficiency against various baselines. We also show that it achieves much stronger scalability in modeling long-term temporal information compared with the current state-of-the-art DyGFormer (Sec. 4.3.2 and Sec. 4.3.3).

## 4.1 Experimental Setting

**CTDG Datasets and Baselines.** We consider seven real-world CTDG datasets collected by [20], i.e., LastFM, Enron, MOOC, Reddit, Wikipedia, UCI and Social Evo.. Dataset statistics are presented in App. A.1. Among them, we take LastFM, Enron and MOOC as long-range temporal dependent datasets because according to Yu et al. [35], much longer histories are needed for optimal representation learning on them. We compare DyGMamba with ten recent CTDG baseline models, i.e., JODIE [12], DyRep [27], TGAT [33], TGN [21], CAWN [32], EdgeBank [20], TCL [30], GraphMixer [2], DyGFormer [35] and CTAN [5]. Among them, only DyGFormer and CTAN are designed for long-range temporal information propagation. Detailed descriptions of baseline methods are presented in App. B. We also implemented FreeDyG [25] by using its official code repository, however, on LastFM, we find that FreeDyG’s loss cannot converge and the reported results are not reproducible. So we do not report its performance in our paper.

<sup>3</sup>To supplement, we also validate on the dynamic node classification task. Since current mainstream datasets of this task requires no long-term temporal reasoning, we put the discussion in App. F. Additionally, we also benchmark DyGMamba on DTDGs in App. L. This serves as supplementary experiment and does not directly connect to our main focus.

**Implementation Details and Evaluation Settings.** We use the implementations and the best hyperparameters provided by Yu et al. [35] for all baseline models except CTAN. For CTAN, we use its official implementation, fixing the number of layers to 5. All models are trained with a batch size of 200 for fair efficiency analysis. For DyGMamba, we report the number of sampled one-hop temporal neighbors  $\rho$  and the patch size  $p$  here. On Wikipedia, Social Evo., and UCI,  $\rho$  &  $p = 32$  & 1. On Reddit,  $\rho$  &  $p = 64$  & 2. On MOOC,  $\rho$  &  $p = 128$  & 4. On Enron,  $\rho$  &  $p = 256$  & 8. On LastFM,  $\rho$  &  $p = 512$  & 16. Note that to fairly compare DyGMamba’s efficiency with DyGFormer, we keep the sequence length  $\rho/p$  input into the SSM as same as the length input into Transformer in Yu et al. [35], i.e.,  $\rho/p = 32$ . All experiments are implemented with PyTorch [18] on a server equipped with an AMD EPYC 7513 32-Core Processor and a single NVIDIA A40 with 45GB memory. We run each experiment for five times with five random seeds and report the mean results together with error bars. Further implementation details including complete hyperparameter configurations are presented in App. C. We employ two evaluation settings following previous works: the transductive and inductive settings. As suggested in [20], we do link prediction evaluation using three negative sampling strategies (NSSs): random, historical and inductive. Historical NSS is only considered under the transductive setting. See App. D for detailed explanations. We employ two metrics, i.e., average precision (AP) and area under the receiver operating characteristic curve (AUC-ROC)

## 4.2 Performance Analysis

**4.2.1 Comparative Study on Benchmark Datasets.** We report the AP of transductive and inductive link prediction in Table 1 and 2 (AUC-ROC reported in Table 12 and 13 in App. E). We find that: (1)

**Table 2: AP of inductive dynamic link prediction. EdgeBank cannot do inductive link prediction so is not reported.**

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	LastFM	83.13 ± 1.19	83.47 ± 1.06	78.40 ± 0.30	81.18 ± 3.27	89.33 ± 0.06	81.38 ± 1.53	82.07 ± 0.31	<u>94.17 ± 0.10</u>	60.40 ± 3.01	<b>94.42 ± 0.21</b>
	Enron	78.97 ± 1.59	73.97 ± 3.00	66.67 ± 1.07	78.76 ± 1.69	86.30 ± 0.56	82.61 ± 0.61	75.55 ± 0.81	<u>89.62 ± 0.27</u>	74.61 ± 1.64	<b>89.67 ± 0.27</b>
	MOOC	80.57 ± 0.52	80.50 ± 0.68	85.28 ± 0.30	<u>88.01 ± 1.48</u>	81.32 ± 0.42	82.28 ± 0.99	81.38 ± 0.17	87.05 ± 0.51	64.99 ± 2.24	<b>88.64 ± 0.08</b>
	Reddit	96.43 ± 0.16	95.89 ± 0.26	97.13 ± 0.04	97.41 ± 0.12	98.62 ± 0.01	95.01 ± 0.10	95.24 ± 0.08	<u>98.83 ± 0.02</u>	80.07 ± 2.53	<b>98.97 ± 0.01</b>
	Wikipedia	94.91 ± 0.32	92.21 ± 0.29	96.26 ± 0.12	97.81 ± 0.18	98.27 ± 0.02	97.48 ± 0.06	96.61 ± 0.04	<u>98.58 ± 0.01</u>	93.58 ± 0.65	<b>98.77 ± 0.03</b>
	UCI	79.73 ± 1.48	58.39 ± 2.38	79.10 ± 0.49	87.81 ± 1.32	92.61 ± 0.35	84.19 ± 1.37	91.17 ± 0.29	<u>94.45 ± 0.13</u>	49.78 ± 5.02	<b>94.76 ± 0.19</b>
	Social Evo.	91.72 ± 0.66	89.10 ± 1.90	91.47 ± 0.10	90.74 ± 1.40	79.83 ± 0.14	92.51 ± 0.11	91.89 ± 0.05	<u>93.05 ± 0.10</u>	Timeout	<b>93.13 ± 0.05</b>
	<b>Avg. Rank</b>	6.29	8.00	7.00	5.14	4.43	5.57	5.86	2.14	9.57	1.00
Inductive	LastFM	71.37 ± 3.45	69.75 ± 2.73	76.26 ± 0.34	68.47 ± 6.07	71.28 ± 0.43	68.79 ± 0.93	<u>76.27 ± 0.37</u>	75.07 ± 1.45	55.60 ± 3.91	<b>76.76 ± 0.43</b>
	Enron	66.99 ± 1.15	62.64 ± 2.33	59.95 ± 1.00	64.51 ± 1.66	60.61 ± 0.63	68.93 ± 1.34	<b>71.71 ± 1.33</b>	67.21 ± 0.72	68.66 ± 2.31	<u>68.77 ± 0.60</u>
	MOOC	64.67 ± 1.18	62.05 ± 2.11	77.43 ± 0.81	76.81 ± 2.83	74.36 ± 0.78	75.95 ± 1.46	73.87 ± 0.99	<u>80.66 ± 0.94</u>	57.49 ± 1.34	<b>80.75 ± 1.00</b>
	Reddit	62.54 ± 0.52	61.07 ± 0.86	63.96 ± 0.25	65.27 ± 0.57	64.10 ± 0.22	61.45 ± 0.25	64.82 ± 0.30	65.03 ± 1.20	<b>78.35 ± 5.03</b>	<u>65.30 ± 1.05</u>
	Wikipedia	68.22 ± 0.36	61.07 ± 0.82	84.19 ± 0.96	81.96 ± 0.62	62.34 ± 6.79	71.46 ± 4.95	<u>87.47 ± 0.25</u>	57.90 ± 11.05	<b>92.61 ± 0.90</b>	71.14 ± 2.44
	UCI	63.57 ± 2.15	52.63 ± 1.87	69.77 ± 0.43	69.94 ± 0.50	63.44 ± 1.52	<u>74.39 ± 1.81</u>	<b>81.40 ± 0.52</b>	70.25 ± 2.02	52.31 ± 2.67	72.17 ± 2.20
	Social Evo.	89.06 ± 1.23	87.30 ± 1.55	94.24 ± 0.36	90.67 ± 2.41	80.30 ± 0.21	95.94 ± 0.37	94.56 ± 0.24	<u>96.73 ± 0.11</u>	Timeout	<b>96.83 ± 0.56</b>
	<b>Avg. Rank</b>	6.86	8.57	5.29	5.43	7.43	4.86	<u>3.14</u>	4.43	6.57	<b>2.43</b>

**Table 3: Ablation studies under transductive setting. R/H/I means random/historical/inductive NSS. Metric is AP.**

Datasets	LastFM			Enron			MOOC			Reddit			Wikipedia			UCI			Social Evo.		
	R	H	I	R	H	I	R	H	I	R	H	I	R	H	I	R	H	I	R	H	I
Variant A	93.14	80.30	71.29	91.35	70.07	75.44	87.78	83.25	77.04	99.19	81.60	90.70	98.99	80.99	79.26	94.88	79.37	70.43	94.59	96.97	97.42
Variant B	93.07	82.53	72.97	92.46	76.88	78.87	86.95	83.78	75.81	97.97	73.47	84.16	94.17	81.37	79.24	91.69	71.13	60.45	92.90	96.61	97.14
DyGMamba	<b>93.35</b>	<b>83.02</b>	<b>73.63</b>	<b>92.65</b>	<b>77.77</b>	<b>80.86</b>	<b>89.21</b>	<b>85.89</b>	<b>81.11</b>	<b>99.32</b>	81.80	<b>91.15</b>	<b>99.15</b>	<b>81.77</b>	<b>79.86</b>	<b>95.91</b>	81.03	71.95	<b>94.77</b>	<b>97.35</b>	<b>97.68</b>

**Table 4: Ablation studies under inductive setting. R/I means random/inductive NSS. Metric is AP.**

Datasets	LastFM		Enron		MOOC		Reddit		Wikipedia		UCI		Social Evo.	
	R	I	R	I	R	I	R	I	R	I	R	I	R	I
Variant A	94.12	73.03	85.97	61.43	84.25	76.16	98.84	65.19	98.49	70.98	93.23	70.84	92.99	96.54
Variant B	94.25	75.26	89.13	67.87	86.21	75.08	97.32	58.22	92.41	70.76	90.42	60.43	91.11	96.32
DyGMamba	<b>94.42</b>	<b>76.76</b>	<b>89.67</b>	<b>68.77</b>	<b>88.64</b>	80.75	<b>98.97</b>	<b>65.30</b>	<b>98.77</b>	<b>71.14</b>	<b>94.76</b>	72.17	<b>93.13</b>	<b>96.83</b>

DyGMamba constantly ranks top 1 under the random NSS, showing a superior performance; (2) Under the historical and inductive NSS, DyGMamba can achieve the best average rank compared with all baselines. More importantly, it shows more superiority on the datasets where encoding longer-term temporal dependencies is necessary, e.g., on LastFM, Enron and MOOC. (3) Among the models that can do long range propagation of information over time (i.e., DyGFormer, CTAN and DyGMamba), DyGMamba achieves the best average rank under any NSS setting in both transductive and inductive link prediction. On the long-range temporal dependent datasets, DyGMamba outperforms DyGFormer and CTAN in most cases; (4) CTAN achieves much better results in transductive than in inductive link prediction. This is because CTAN requires multi-hop temporal neighbors to learn node representations, which is difficult for unseen nodes. By contrast, DyGMamba and DyGFormer require only one-hop temporal neighbors, thus performing much better in inductive link prediction.

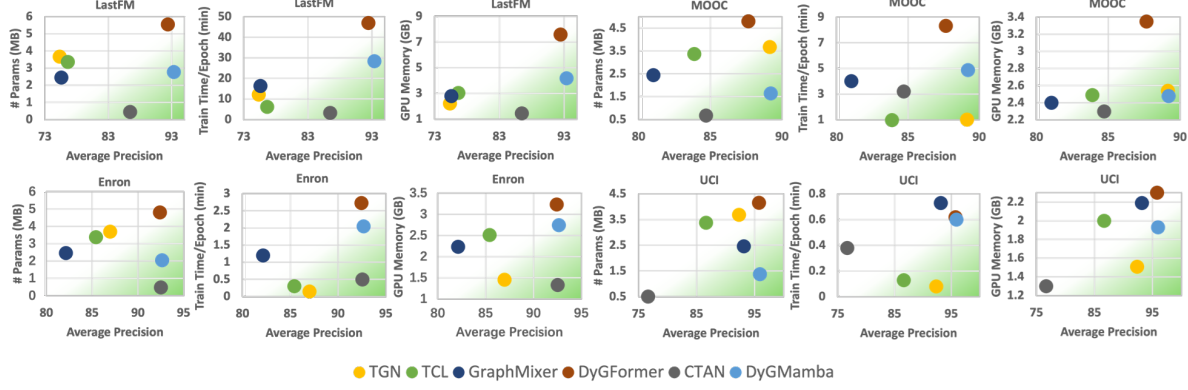
**4.2.2 Ablation Study.** We conduct four ablation studies to study the effectiveness of model components. In study A, we make a model variant (Variant A) by removing the time-level SSM block and restrain our model from learning temporal patterns (information selection is substituted by mean pooling over the output of Eq.

4). In study B, we make a model variant (Variant B) by removing the Mamba SSM layers (Eq. 3) in the node-level SSM block. From Table 3 and 4, we find that: (1) Variant A is constantly beaten by DyGMamba, showing the effectiveness of dynamic information selection based on edge-specific temporal patterns; (2) DyGMamba always outperforms Variant B, indicating the importance of encoding the one-hop temporal neighbors with SSM layers for capturing graph dynamics. See App. G for more ablation studies.

**4.2.3 A Closer Look into Temporal Pattern Modeling with Synthetic Datasets.** We observe from ablation studies that dynamic information selection based on temporal patterns contributes to better model performance on real-world datasets. To better quantify its benefits, we construct three synthetic datasets, i.e., S1, S2 and S3, that follow different patterns and compare our model with DyGFormer, CTAN as well as Variant A, C, D on them. Each synthetic dataset contains 7 nodes, where the interactions of each pair of two nodes follow a certain pattern along time. And for each node, we generate interactions with all the other nodes. Assume we have a pair of node  $u$  and  $v$  and they have interactions at  $\{t_i\}_{i=0}^N$ , in S1, the time intervals between neighboring interactions  $\{t_1 - t_0, \dots, t_N - t_{N-1}\}$  follow an increasing trend with a constant velocity of 0.05, i.e.,  $(t_{i+2} - t_{i+1}) - (t_{i+1} - t_i) = 0.05$ . In S2, we

**Table 5: Performance (Random NSS) on synthetic datasets.**

(a) AP on synthetic datasets.							(b) AUC-ROC on synthetic datasets.						
Datasets	DyGFormer	CTAN	Variant A	Variant C	Variant D	DyGMamba	Datasets	DyGFormer	CTAN	Variant A	Variant C	Variant D	DyGMamba
S1	55.19 $\pm$ 0.98	51.25 $\pm$ 2.11	53.72 $\pm$ 0.04	55.45 $\pm$ 0.32	54.52 $\pm$ 0.71	<b>81.58 <math>\pm</math> 1.31</b>	S1	56.27 $\pm$ 0.54	51.25 $\pm$ 2.38	53.16 $\pm$ 0.39	57.41 $\pm$ 0.04	55.83 $\pm$ 1.03	<b>86.61 <math>\pm</math> 1.30</b>
S2	57.80 $\pm$ 4.61	51.17 $\pm$ 0.93	60.16 $\pm$ 2.20	64.71 $\pm$ 2.33	61.51 $\pm$ 3.00	<b>85.36 <math>\pm</math> 2.55</b>	S2	59.06 $\pm$ 6.07	51.46 $\pm$ 0.91	62.50 $\pm$ 2.28	64.93 $\pm$ 2.59	62.06 $\pm$ 3.40	<b>89.94 <math>\pm</math> 2.70</b>
S3	79.20 $\pm$ 0.60	51.46 $\pm$ 0.19	77.61 $\pm$ 2.31	77.61 $\pm$ 2.31	79.41 $\pm$ 2.13	<b>86.59 <math>\pm</math> 0.09</b>	S3	82.89 $\pm$ 1.34	52.12 $\pm$ 0.44	81.78 $\pm$ 0.40	82.73 $\pm$ 1.97	84.20 $\pm$ 3.57	<b>91.72 <math>\pm</math> 0.11</b>



**Figure 2: Efficiency comparison on four datasets among DyGMamba and five baselines in terms of number (#) of parameters, training time per epoch and GPU memory. The performance metric here is AP of transductive link prediction under random NSS. The greener, the better performance/efficiency. In contrast to other methods, DyGMamba consistently shows strong overall capability across different datasets. More explanations in Sec. 4.3.1. Further comparison on Reddit and Social Evo. is presented in App. H.1**

set the time intervals to a decreasing trend with the same velocity, i.e.,  $(t_{i+1} - t_i) - (t_{i+2} - t_{i+1}) = 0.05$ . And in S3, we modify S1 by repeating several periods of increasing patterns taken from S1 to form a periodic dataset. In this way, we have three datasets demonstrating diverse temporal patterns: increasing/decreasing/periodic time intervals between neighboring interactions. Details of dataset construction and statistics are provided in App. A.2. From Table 5, we observe that DyGMamba greatly outperforms DyGFormer and CTAN. More importantly, Variant A, C and D show similar performance to DyGFormer, meaning that our time-level SSM block is able to capture temporal patterns and modeling such patterns for dynamic information selection is important in CTDG reasoning. For more implementation details on synthetic datasets, please refer to App. C.2.

### 4.3 Efficiency Analysis

We evaluate the models’ efficiency based on the following aspects: model size (number of trainable parameters), per-epoch training time, and GPU memory consumption during training. Since DyGMamba shares similarities with DyGFormer, i.e., both of them model large sequences of one-hop temporal neighbors and use patching to enhance scalability, we further analyze the impact of patch size on their scalability and performance. We also compare the complexity of DyGMamba and DyGFormer to highlight DyGMamba’s efficiency.

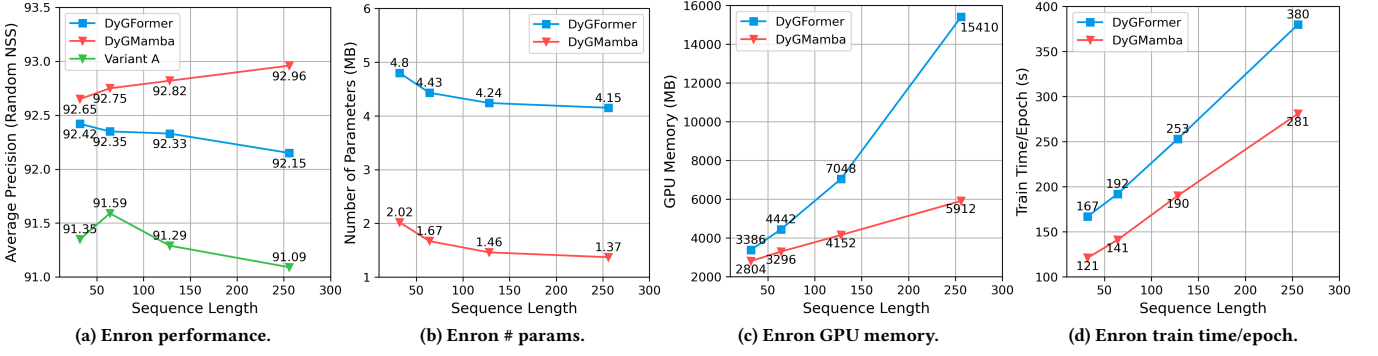
#### 4.3.1 Model Size, Per Epoch Training Time and GPU Memory Comparison Across Various Models.

Fig. 2 compares DyGMamba with five baselines in terms of number of parameters (model size), per

epoch training time, and GPU memory consumption during training<sup>4</sup>. We find that: (1) DyGMamba uses very few parameters while maintaining the best performance, showing a strong parameter efficiency. Only CTAN constantly uses fewer parameters than DyGMamba, however, its performance is significantly worse with the only exception of Enron; (2) DyGMamba is always more efficient than DyGFormer with the same length of input sequence ( $p/p=32$ ) while achieving the same performance; (3) Although DyGMamba generally consumes more GPU memory and takes more time to train per epoch compared with most baselines, the gap of consumption is modest. To model more temporal neighbors for long-range temporal dependent datasets, DyGMamba naturally requires more computational resources, thus enlarging the consumption gap. DyGFormer shows the same trend as DyGMamba since it also captures long-term temporal dependencies but at a higher cost; (4) CTAN requires very few computational resources. However, on long-range temporal dependent datasets, it is beaten by DyGFormer and DyGMamba by a large margin, e.g., on LastFM and Enron. Besides, CTAN is also hard to converge. Although it takes little time to train a single epoch, it needs more epochs to reach the best performance, leading to a long total training time. See App. H.2 for a total training time comparison among DyGFormer, CTAN and DyGMamba. To supplement, we also present in App. H.3 another experiment to compare DyGMamba with baselines in modeling an increasing number of temporal neighbors on Enron with limited total training time.

<sup>4</sup>The baselines not included here are either extremely inefficient (e.g., CAWN) or inferior in performance (e.g., DyRep). Complete statistics of all baseline models presented in App. H.1





**Figure 3: Impact of patch size on DyGFormer, DyGMamba and Variant A, given a fixed number of sampled temporal neighbors  $\rho$  on Enron. Patch size  $p$  varies from 8, 4, 2, 1. Sequence length  $\rho/p$  increases as patch size decreases. Performance is the transductive AP under random NSS.**

**4.3.2 Impact of Patch Size on Scalability and Performance.** Patching treats  $p$  temporal neighbors as one patch and thus decreases the sequence length by  $p$  times. This is very helpful in cutting the consumption of GPU memory and training/evaluation time. However, patching introduces excessive parameters because it is done through  $f_N$ ,  $f_E$ ,  $f_T$  and  $f_F$  whose sizes increase as the patch size grows. Fig. 3b shows the numbers of parameters of DyGFormer and DyGMamba with different patch sizes on a long-term temporal dependent dataset Enron. We find that patching greatly affects model sizes. To further study how patching affects DyGMamba, we decrease the patch size gradually from 8 to 1 and track DyGMamba’s performance (Fig. 3a) as well as efficiency (Fig. 3b to 3d) on Enron. Meanwhile, we also keep track on DyGFormer under the same patch size for comparison. We have several findings: (1) Whatever the patch size is, DyGMamba always consumes fewer parameters, less GPU memory and per epoch training time, showing its high efficiency; (2) While both models require increasing computational budgets as the patch size decreases, the speed of increase is much lower for DyGMamba, demonstrating its strong scalability in modeling longer sequences; (3) Different trends in performance change are observed between two models. While DyGFormer performs worse, DyGMamba can benefit from a smaller patch size, indicating its strong ability to capture nuanced temporal details even if the sequence becomes much longer. Note that the models use fewer parameters under smaller patch size. This also shows that DyGMamba can achieve much stronger parameter efficiency by reducing patch sizes. To further study the reason for finding (3), we plot the performance of Variant A under different patch sizes in Fig. 3a. We find that Variant A’s performance degrades when sequence length is more than 64. This means that dynamic information selection based on edge-specific temporal patterns is essential for DyGMamba to optimally process long sequences. To supplement, we provide additional analysis on MOOC in App. K.

**4.3.3 Complexity: DyGMamba vs. DyGFormer.** DyGMamba follows the current state-of-the-art DyGFormer by learning from one-hop temporal neighbors for temporal reasoning. We analyze the complexity of both models to show DyGMamba’s efficiency. Sequence length is the key factor affecting the consumption of computational resources in DyGFormer and DyGMamba. Following

the computation of previous work [36], the complexity of Transformer and Mamba in DyGFormer and DyGMamba can be written as  $O(T) = 4(\rho/p)(4d)^2 + 2(\rho/p)^2(4d) = 64(\rho/p)d^2 + 8d(\rho/p)^2$  and  $O(M) = 3(\rho/p)(4d)d_{SSM} + (\rho/p)(4d)d_{SSM} = 16d_{SSM}d(\rho/p)$ . This means that DyGMamba holds a computational complexity linear to  $\rho/p$ , while DyGFormer’s complexity is quadratic to  $\rho/p$ . As a result, as the sequence length grows (either  $\rho$  increases or  $p$  decreases), DyGFormer is less scalable compared with DyGMamba<sup>5</sup>. Some may argue that increasing the patch size  $p$  to a large enough value can offset the negative influence of the higher complexity of Transformer. However, as discussed in Sec. 4.3.2, increasing patch size will substantially increase model parameters, causing burden in parameter optimization, and meanwhile lose temporal details. This indicates that patching is not always beneficial and DyGMamba’s low complexity provides an alternative way to maintain great efficiency while considering more temporal information. With the same number of sampled temporal neighbors and equivalent computational resources, DyGMamba can leverage a smaller patch size, mitigating the negative effects of lost temporal details and simplifying parameter optimization, as implied in Sec. 4.3.2.

## 5 Conclusion

We propose DyGMamba, an efficient CTDG representation learning model that can capture long-term temporal dependencies. DyGMamba first leverages a node-level SSM to encode long sequences of historical node interactions. It then employs a time-level SSM to learn edge-specific temporal patterns. The learned patterns are used to select the critical part of the encoded temporal information. DyGMamba achieves superior performance on dynamic link prediction, and moreover, it shows high efficiency and strong scalability compared with previous CTDG methods, implying a great potential in modeling huge amounts of temporal information with a limited computational budget.

<sup>5</sup>As we set  $k$  (number of edge-specific historical interactions discussed in Sec. 3.2) to a number much smaller than the number of sampled one-hop temporal neighbors, i.e., sequence length  $\rho/p$ , we omit here the contribution of the time-level SSM in complexity analysis.



## References

- [1] Xiaofu Chang, Xuqin Liu, Jianfeng Wen, Shuang Li, Yanming Fang, Le Song, and Yuan Qi. 2020. Continuous-Time Dynamic Graph Learning via Neural Interaction Processes. In *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, Mathieu d'Aquin, Stefan Dietze, Claudia Hauff, Edward Curry, and Philippe Cudré-Mauroux (Eds.). ACM, 145–154. doi:10.1145/3340531.3411946
- [2] Weilin Cong, Si Zhang, Jian Kang, Baichuan Yuan, Hao Wu, Xin Zhou, Hanghang Tong, and Mehrdad Mahdavi. 2023. Do We Really Need Complicated Model Architectures For Temporal Networks?. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. <https://openreview.net/pdf?id=ayPPc0SyLv1>
- [3] Feysa Duman Keles, Pruthuvi Mahesakya Wijewardena, and Chinmay Hegde. 2023. On The Computational Complexity of Self-Attention. In *Proceedings of The 34th International Conference on Algorithmic Learning Theory (Proceedings of Machine Learning Research, Vol. 201)*, Shipra Agrawal and Francesco Orabona (Eds.). PMLR, 597–619. <https://proceedings.mlr.press/v201/duman-keles23a.html>
- [4] Palash Goyal, Sujit Rokka Chhetri, and Arquimedes Canedo. 2020. dyngraph2vec: Capturing network dynamics using dynamic graph representation learning. *Knowl. Based Syst.* 187 (2020). doi:10.1016/J.KNSYS.2019.06.024
- [5] Alessio Gravina, Giulio Lovisotto, Claudio Gallicchio, Davide Bacciu, and Claas Grohnfeldt. 2024. Long Range Propagation on Continuous-Time Dynamic Graphs. In *Forty-first International Conference on Machine Learning*. <https://openreview.net/forum?id=gVg8V9isul>
- [6] Albert Gu and Tri Dao. 2023. Mamba: Linear-Time Sequence Modeling with Selective State Spaces. *CoRR abs/2312.00752* (2023). doi:10.48550/ARXIV.2312.00752 arXiv:2312.00752
- [7] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. 2022. On the Parameterization and Initialization of Diagonal State Space Models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). [http://papers.nips.cc/paper\\_files/paper/2022/hash/e9a32fade47b906de908431991440f7c-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/e9a32fade47b906de908431991440f7c-Abstract-Conference.html)
- [8] Albert Gu, Karan Goel, and Christopher Ré. 2022. Efficiently Modeling Long Sequences with Structured State Spaces. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net. <https://openreview.net/forum?id=uYLFoz1v1AC>
- [9] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. 2021. Combining Recurrent, Convolutional, and Continuous-time Models with Linear State Space Layers. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 572–585. <https://proceedings.neurips.cc/paper/2021/hash/05546b0e38ab9175cd905eebc6eb76-Abstract.html>
- [10] Ming Jin, Yuan-Fang Li, and Shirui Pan. 2022. Neural Temporal Walks: Motif-Aware Representation Learning on Continuous-Time Dynamic Graphs. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). [http://papers.nips.cc/paper\\_files/paper/2022/hash/7dad855cef7494d5d956a8d28add871-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2022/hash/7dad855cef7494d5d956a8d28add871-Abstract-Conference.html)
- [11] Seyed Mehran Kazemi, Rishab Goel, Kshitij Jain, Ivan Kobyzev, Akshay Sethi, Peter Forsyth, and Pascal Poupart. 2020. Representation Learning for Dynamic Graphs: A Survey. *J. Mach. Learn. Res.* 21 (2020), 70:1–70:73. <http://jmlr.org/papers/v21/19-447.html>
- [12] Srijan Kumar, Xikun Zhang, and Jure Leskovec. 2019. Predicting Dynamic Embedding Trajectory in Temporal Interaction Networks. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*, Ankur Teredesai, Vipin Kumar, Ying Li, Rómer Rosales, Evimaria Terzi, and George Karypis (Eds.). ACM, 1269–1278. doi:10.1145/3292500.3330895
- [13] Jintang Li, Ruofan Wu, Xinzhou Jin, Boqun Ma, Liang Chen, and Zibin Zheng. 2024. State Space Models on Temporal Graphs: A First-Principles Study. *arXiv preprint arXiv:2406.00943* (2024).
- [14] Yuyu Liu, Jianzhu Ma, and Pan Li. 2022. Neural Predicting Higher-order Patterns in Temporal Networks. In *WWW '22: The ACM Web Conference 2022, Virtual Event, Lyon, France, April 25 - 29, 2022*, Frédérique Laforest, Raphaël Troncy, Elena Simperl, Deepak Agarwal, Aristides Gionis, Ivan Herman, and Lionel Médini (Eds.). ACM, 1340–1351. doi:10.1145/3485447.3512181
- [15] Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan May, and Luke Zettlemoyer. 2023. Mega: Moving Average Equipped Gated Attention. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. <https://openreview.net/pdf?id=qNLe3iq2El>
- [16] Yao Ma, Ziyi Guo, Zhaochun Ren, Jiliang Tang, and Dawei Yin. 2020. Streaming Graph Neural Networks. In *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, Jimmy X. Huang, Yi Chang, Xueqi Cheng, Jaap Kamps, Vanessa Murdock, Ji-Rong Wen, and Yiqun Liu (Eds.). ACM, 719–728. doi:10.1145/3397271.3401092
- [17] Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao B. Schardl, and Charles E. Leiserson. 2020. EvolveGCN: Evolving Graph Convolutional Networks for Dynamic Graphs. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 5363–5370. doi:10.1145/AAAILV34I04.5984
- [18] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (Eds.). 8024–8035. <https://proceedings.neurips.cc/paper/2019/hash/bdbca288fee7f92f2bfa9f7012727740-Abstract.html>
- [19] Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kranthi Kiran GV, Xuzheng He, Haowen Hou, Przemysław Kazienko, Jan Kocon, Jiaming Kong, Bartłomiej Koptyr, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, Guangyu Song, Xiangru Tang, Johan S. Wind, Stanisław Wozniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. 2023. RWKV: Reinventing RNNs for the Transformer Era. In *Findings of the Association for Computational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, Houda Bouamor, Juan Pino, and Kalika Bali (Eds.). Association for Computational Linguistics, 14048–14077. doi:10.18653/V1/2023.FINDINGS-EMNLP.936
- [20] Farimah Poursafaei, Shenyang Huang, Kellin Pelrine, and Reihaneh Rabbany. 2022. Towards Better Evaluation for Dynamic Link Prediction. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (Eds.). [http://papers.nips.cc/paper\\_files/paper/2022/hash/d49042a5d49818711c40d34172f9900-Abstract-Datasets\\_and\\_Benchmarks.html](http://papers.nips.cc/paper_files/paper/2022/hash/d49042a5d49818711c40d34172f9900-Abstract-Datasets_and_Benchmarks.html)
- [21] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael M. Bronstein. 2020. Temporal Graph Networks for Deep Learning on Dynamic Graphs. *CoRR abs/2006.10637* (2020). arXiv:2006.10637 <https://arxiv.org/abs/2006.10637>
- [22] Aravind Sankar, Yanhong Wu, Liang Gou, Wei Zhang, and Hao Yang. 2020. DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention Networks. In *WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020*, James Caverlee, Xia (Ben) Hu, Mounia Lalmas, and Wei Wang (Eds.). ACM, 519–527. doi:10.1145/3336191.3371845
- [23] Razieh Shirzadkhani, Shenyang Huang, Elahe Kooshafar, Reihaneh Rabbany, and Farimah Poursafaei. 2024. Temporal Graph Analysis with TGX. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining (Merida, Mexico) (WSDM '24)*. Association for Computing Machinery, New York, NY, USA, 1086–1089. doi:10.1145/3616855.3635694
- [24] Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. 2023. Simplified State Space Layers for Sequence Modeling. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net. <https://openreview.net/pdf?id=Ai8Hw3AXqks>
- [25] Yuxing Tian, Yiyan Qi, and Fan Guo. 2024. FreeDyG: Frequency Enhanced Continuous-Time Dynamic Graph Model for Link Prediction. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=82Mc5IlInM>
- [26] Ilya O. Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, Mario Lucic, and Alexey Dosovitskiy. 2021. MLP-Mixer: An all-MLP Architecture for Vision. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (Eds.). 24261–24272. <https://proceedings.neurips.cc/paper/2021/hash/cba0a4ee5ccd02fda0fe3f9a3e7b89fe-Abstract.html>
- [27] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. 2019. 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. <https://openreview.net/forum?id=HyePhR5KX>
- [28] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (Eds.). 5998–6008. <https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html>
- [29] Roger Waleffe, Wonmin Byeon, Duncan Riach, Brandon Norick, Vijay Korthikanti, Tri Dao, Albert Gu, Ali Hatamizadeh, Sudhakar Singh, Deepak Narayanan, Garvit Kulshreshtha, Vartika Singh, Jared Casper, Jan Kautz, Mohammad Shoeybi, and Bryan Catanzaro. 2024. An Empirical Study of Mamba-based Language Models. *CoRR* abs/2406.07887 (2024). doi:10.48550/ARXIV.2406.07887 arXiv:2406.07887
- [30] Lu Wang, Xiaofu Chang, Shuang Li, Yunfei Chu, Hui Li, Wei Zhang, Xiaofeng He, Le Song, Jingren Zhou, and Hongxia Yang. 2021. TCL: Transformer-based Dynamic Graph Modelling via Contrastive Learning. *CoRR* abs/2105.07944 (2021). arXiv:2105.07944 <https://arxiv.org/abs/2105.07944>
- [31] Xuhong Wang, Ding Lyu, Mengjian Li, Yang Xia, Qi Yang, Xinwen Wang, Xinggang Wang, Ping Cui, Yupu Gu, Bowen Sun, and Zhenyu Guo. 2021. APAN: Asynchronous Propagation Attention Network for Real-time Temporal Graph Embedding. In *SIGMOD '21: International Conference on Management of Data, Virtual Event, China, June 20-25, 2021*, Guoliang Li, Zhanhuai Li, Stratos Idreos, and Divesh Srivastava (Eds.). ACM, 2628–2638. doi:10.1145/3448016.3457564
- [32] Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. 2021. Inductive Representation Learning in Temporal Networks via Causal Anonymous Walks. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net. <https://openreview.net/forum?id=KYPz4YsCpJ>
- [33] Da Xu, Chuanwei Ruan, Evren Körpeoglu, Sushant Kumar, and Kannan Achan. 2020. Inductive representation learning on temporal graphs. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net. <https://openreview.net/forum?id=rJeW1yHYwH>
- [34] Jiaxuan You, Tianyu Du, and Jure Leskovec. 2022. ROLAND: Graph Learning Framework for Dynamic Graphs. In *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 - 18, 2022*, Aidong Zhang and Huzefa Rangwala (Eds.). ACM, 2358–2366. doi:10.1145/3534678.3539300
- [35] Le Yu, Leilei Sun, Bowen Du, and Weifeng Lv. 2023. Towards Better Dynamic Graph Learning: New Architecture and Unified Library. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (Eds.). [http://papers.nips.cc/paper\\_files/paper/2023/hash/d611019afba70d547bd595e8a4158f55-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/d611019afba70d547bd595e8a4158f55-Abstract-Conference.html)
- [36] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. 2024. Vision Mamba: Efficient Visual Representation Learning with Bidirectional State Space Model. In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*. OpenReview.net. <https://openreview.net/forum?id=YbHCqn4qF4>

## A CTDG Dataset Details

### A.1 Real-World Benchmark Datasets

We present the dataset statistics of all considered CTDG datasets in Table 6. All the datasets in our experiments are taken from Yu et al. [35]. We chronologically split each dataset with the ratio of 70%/15%/15% for training/validation/testing. Please refer to it for detailed dataset descriptions.

### A.2 Synthetic Datasets

For all of our three synthetic datasets, node and link features are not involved during dataset construction. The construction details are as follows:

- **S1**: For each interaction pair  $u$  and  $v$ , their first interaction is at timestamp 0 and the second interaction is generated randomly. Thus, the first time interval is also determined. Starting from the second interval, they follow an increasing

trend with a constant velocity of 0.05, i.e.,  $(t_{i+2} - t_{i+1}) - (t_{i+1} - t_i) = 0.05$ . The number of interactions for each node pair is randomly determined and the interaction numbers of all node pairs sum up to 100000.

- **S2**: For each interaction pair  $u$  and  $v$ , their first interaction is at timestamp 0 and the second interaction is generated randomly. However, it should be large enough so that the interval will not drop to zero or negative afterwards. Starting from the second interval, they follow a decreasing trend with a constant velocity of 0.05, i.e.,  $(t_{i+1} - t_i) - (t_{i+2} - t_{i+1}) = 0.05$ . The number of interactions for each node pair is randomly determined and the interaction numbers of all node pairs sum up to 100000.
- **S3**: S3 contains 8 periods. In each period, the interactions of each node pair  $u$  and  $v$  are generated following the same pattern in S1. The number of interactions for each node pair is randomly determined and the interaction numbers of all node pairs sum up to 12000 in the period.

We present the statistics of all synthetic datasets in Table 7. We chronologically split each dataset with the ratio of 70%/15%/15% for training/validation/testing. To better visualize the temporal patterns in each dataset, we pick one pair of interacting nodes and plot the time intervals between neighboring interactions in each dataset in Figure 4. Note that for the periodic dataset S3 (Figure 4c), each of the train, validation and test sets contains at least one start of a new period. This ensures that models have to capture periodic temporal patterns in order to achieve good performance during evaluation, rather than only learning the increasing time intervals as specified in S1.

Furthermore, we provide the information about the numbers of interactions regarding interacting node pairs in Table 8. We show that each node pair is equipped with a substantial number of interactions, meaning that temporal patterns in our synthetic datasets span across long time periods. This encourages models to consider long-term temporal dependencies for better graph reasoning.

## B Baseline Details

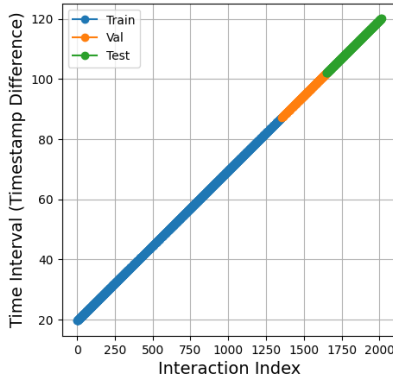
We provide the detailed descriptions of all baselines here. The baselines can be split into two groups: the methods designed/not designed for long-range temporal information propagation.

### B.1 Baselines Not Designed for Long-Range Temporal Information Propagation

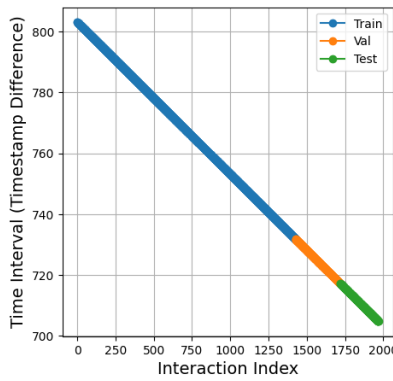
- **JODIE** [12]: JODIE employs a recurrent neural network (RNN) for each node and uses a projection operation to learn the future representation trajectory of each node.
- **DyRep** [27]: DyRep updates node representations as events appear. It designs a two-time scale deep temporal point process approach for source and destination nodes and couples the structural and temporal components with a temporal-attentive aggregation module.
- **TGAT** [33]: TGAT computes the node representations by aggregating each node’s temporal neighbors based on a self-attention module. A time encoding function is proposed to learn functional representations of time.

**Table 6: Dataset statistics. # N&E Feat means the numbers of node and edge features.**

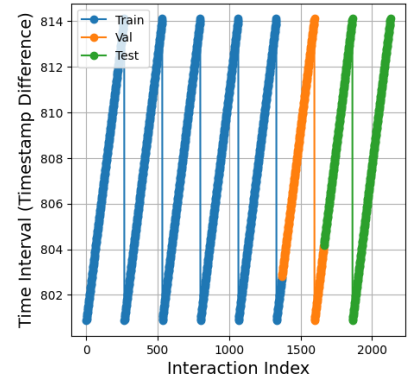
Datasets	# Nodes	# Edges	# N&E Feat	Bipartite	Duration	# Timestamps	Time Granularity
LastFM	1,980	1,293,103	0 & 0	True	1 month	1,283,614	Unix timestamps
Enron	184	125,235	0 & 0	False	3 years	22,632	Unix timestamps
MOOC	7,144	411,749	0 & 4	True	17 months	345,600	Unix timestamps
Reddit	10,984	672,447	0 & 172	True	1 month	669,065	Unix timestamps
Wikipedia	9,227	157,474	0 & 172	True	1 month	152,757	Unix timestamps
UCI	1,899	59,835	0 & 0	False	196 days	58,911	Unix timestamps
Social Evo.	74	2,099,519	0 & 2	False	8 months	565,932	Unix timestamps



(a) Synthetic dataset S1.



(b) Synthetic dataset S2.



(c) Synthetic dataset S3.

**Figure 4: Time intervals of a node pair  $u, v$  in synthetic datasets S1, S2 and S3.****Table 7: Synthetic dataset statistics.**

Datasets	# Nodes	# Edges	# Timestamps	Time Range
S1	7	100,000	96,869	0 - 163241.65
S2	7	100,000	98,004	0 - 1573561.52
S3	7	95,657	95,370	0 - 1771300.40

**Table 8: Interaction information of node pairs in synthetic datasets. Complete Dataset includes the numbers of interactions across the whole datasets, including training, validation and testing.**

	Datasets	Avg. # Interactions	Min # Interactions	Max # Interactions
Training Set	S1	1,428.57	1,205	1,663
	S2	1,428.57	1,424	1,433
	S3	1,367.91	1,364	1,372
Complete Dataset	S1	2,010.20	1,879	2,150
	S2	2,010.20	1,921	2,097
	S3	1,952.18	1,721	2,193
	S3 (each period)	244.02	215	274

- **TGN** [21]: TGN leverages an evolving memory for each node and updates the memory when a node-relevant interaction occurs by using a message function, a message

aggregator, and a memory updater. An embedding module is used to generate the temporal representations of nodes.

- **CAWN** [32]: CAWN is a random walk-based method. It does multiple causal anonymous walks for each node and extracts relative node identities from the walk results. RNNs are then introduced to encode the anonymous walks. The aggregated walk information forms the final node representation.
- **EdgeBank** [20]: EdgeBank is a non-parametric method purely based on memory. It stores the observed interactions in its memory and updates the memory through various strategies. An interaction, i.e., link, will be predicted as existing if it is stored in the memory, and non-existing otherwise. EdgeBank uses four memory update strategies: (1)  $\text{EdgeBank}_{\infty}$ , where all the observed edges are stored in the memory; (2)  $\text{EdgeBank}_{\text{tw-ts}}$ , where only the edges within the duration of the test set from the immediate past are kept in the memory; (3)  $\text{EdgeBank}_{\text{tw-re}}$ , where only the edges within the average time intervals of repeated edges from the immediate past are kept in the memory; (4)  $\text{EdgeBank}_{\text{th}}$ , where the edges with appearing counts higher than a threshold are stored in the memory. The results reported in our paper correspond to the best results achieved among the four memory update strategies.



- **TCL** [30]: TCL first extracts temporal dependency interaction sub-graphs for source and interaction nodes and then presents a graph transformer to aggregate node information from the sub-graphs. A cross-attention operation is implemented to enable information communication between two source and destination nodes.
- **GraphMixer** [2]: GraphMixer designs a link-encoder based on MLP-Mixer [26] to learn from the temporal interactions. A mean pooling-based node-encoder is used to aggregate the node features. Link prediction is done with a link classifier that leverages the representations output by link-encoder and node-encoder.

Note that TGN uses a memory network to store the whole graph history, making it able to preserve long-range temporal information. However, as discussed in Yu et al. [35], it faces a problem of vanishing/exploding gradients, preventing it from optimally capturing long-term temporal dependencies. EdgeBank can also preserve a very long graph history, but we can observe from the experimental results (Table 1, 12) that without learnable parameters, it is not strong enough on long-range temporal dependent datasets.

## B.2 Baselines Designed for Long-Range Temporal Information Propagation

- **DyGFormer** [35]: DyGFormer is a Transformer-based CTDG model. It takes the long-term one-hop temporal interactions of source and destination nodes and uses a Transformer to encode them. A patching technique is developed to cut the computational consumption and a node co-occurrence encoding scheme is used to exploit the correlations of nodes in each interaction. DyGFormer achieves long-range temporal information propagation by increasing the number of sampled one-hop historical interactions. The patching technique ensures that even with a huge number of sampled interactions, the length of the sequence input into Transformer will not be too long, making it possible to implement DyGFormer with a limited computational budget.
- **CTAN** [5]: CTAN is deep graph network for learning CTDGs based on non-dissipative ordinary differential equations. CTAN’s formulation allows for a scalable long-range temporal information propagation in CTDGs because its non-dissipative layer can retain the information from a specific event indefinitely, ensuring that the historical context of a node is preserved despite the occurrence of additional events involving this node.

## C Implementation Details

We train every CTDG model except for CTAN for a maximum number of 200 epochs. Maximum epochs for CTAN training is 1000. We evaluate each model on the validation set at the end of every training epoch and adopt an early stopping strategy with a patience of 20. We take the model that achieves the best validation result for testing. We use the implementations<sup>6</sup> provided by Yu et al. [35] for all baseline models except CTAN. For CTAN, we use its official implementation<sup>7</sup>. All models are trained with a batch size of 200

<sup>6</sup><https://github.com/yule-BUAA/DyGLib>

<sup>7</sup><https://github.com/gravins/non-dissipative-propagation-CTDGs>

for fair efficiency analysis. All experiments are implemented with PyTorch [18] on a server equipped with an AMD EPYC 7513 32-Core Processor and a single NVIDIA A40 with 45GB memory. We run each experiment for five times with five random seeds and report the mean results together with error bars.

## C.1 Hyperparameter Configurations on Real-World Datasets

For all the baselines except CTAN, please refer to Yu et al. [35] for the hyperparameter configurations on real-world datasets. For CTAN, we present its hyperparameter configurations in Table 9. We keep its hyperparameters unchanged for all real-world datasets. Note that we set the number of graph convolution layers (GCLs) in CTAN to its maximum, i.e., 5, in order to maximize its performance in capturing long-term temporal dependencies.

**Table 9: Hyperparameter configurations of CTAN on all real-world datasets.  $\gamma$  here denotes the discretization step size introduced in [5], different from the one in DyGMamba.**

Model	# GCL	$\epsilon$	$\gamma$	Embedding Dim
CTAN	5	0.5	0.5	128

We report the hyperparameter searching strategy of DyGMamba on real-world datasets and the best hyperparameters in Table 10. To achieve fair efficiency comparison with DyGFormer, we fix the length of the input sequence into the node-level SSM to 32, i.e.,  $\rho$  &  $p = 32$ . The results reported in Table 1, 2, 12, 13 are all achieved by DyGMamba with  $\rho$  &  $p = 32$ . In practice, we can decrease  $p$  to have a better performance given more computational resources (as discussed in Sec. 4.3). DyGMamba keeps the embedding size as same as DyGFormer on all real-world datasets, i.e.,  $d_N = d_E = 172$ ,  $d_T = 100$ ,  $d_F = 50$ . We also set  $\gamma = 0.5$  for all experiments of DyGMamba. The dimension of SSMs  $d_{SSM} = 16$  remains the default value of mamba SSM’s official repository<sup>8</sup>. For all experiments, we set the numbers of layers in both node-level and time-level SSMs as 2. We also set  $\gamma$  to 0.5 for all datasets. A smaller  $\gamma$  lowers the computational resource consumption in time-level SSM, potentially at the cost of performance. Raising  $\gamma$ ’s value does not necessarily lead to better performance but will lower efficiency. We search  $\gamma$ ’s value in  $\{0.1, 0.5, 0.7, 1\}$  and find that 0.5 brings a good balance between performance and efficiency.

## C.2 Hyperparameter Configurations on Synthetic Datasets

We use the same settings of CTAN on real-world datasets when we experiment it on synthetic datasets. For DyGFormer and DyGMamba, we fix the length of the input sequence into the Transformer and the node-level SSM to 32, i.e.,  $\rho/p = 32$ . For DyGFormer, we set the hyperparameters except  $\rho$  and  $p$  to the same default values as on real-world datasets, and search for the best  $\rho$  &  $p$  within  $\{512 \& 16, 256 \& 8, 128 \& 4, 64 \& 2, 32 \& 1\}$ . For DyGMamba, we search for the best  $\rho$  &  $p$  within the same search range and

<sup>8</sup><https://github.com/state-spaces/mamba>



**Table 10: DyGMamba hyperparameter searching strategy on real-world datasets. The best settings are marked as bold.**

Datasets	Dropout	$\rho$ & $p$	$k$
LastFM	{0.0, <b>0.1</b> , 0.2}	{1024 & 32, <b>512 &amp; 16</b> , 256 & 8}	{30, <b>10</b> , 5}
Enron	{ <b>0.0</b> , 0.1, 0.2}	{512 & 16, <b>256 &amp; 8</b> , 128 & 4}	{ <b>30</b> , 10, 5}
MOOC	{0.0, <b>0.1</b> , 0.2}	{512 & 16, 256 & 8, <b>128 &amp; 4</b> }	{30, <b>10</b> , 5}
Reddit	{0.0, 0.1, <b>0.2</b> }	{128 & 4, <b>64 &amp; 2</b> , 32 & 1}	{30, 10, <b>5</b> }
Wikipedia	{0.0, <b>0.1</b> , 0.2}	{64 & 2, <b>32 &amp; 1</b> }	{30, 10, <b>5</b> }
UCI	{0.0, <b>0.1</b> , 0.2}	{64 & 2, <b>32 &amp; 1</b> }	{30, 10, <b>5</b> }
Social Evo.	{0.0, <b>0.1</b> , 0.2}	{64 & 2, <b>32 &amp; 1</b> }	{30, 10, <b>5</b> }

further search for the best  $k$ . All the other hyperparameters are set as same as the setting on LastFM. We report the hyperparameter searching strategy as well as the best settings of DyGFormer and DyGMamba on synthetic datasets in Table 11. For all experiments with DyGMamba, we set the numbers of layers in both node-level and time-level SSMs as 2.

## D Negative Edge Sampling Strategies during Evaluation

We justify why we do not do historical NSS for inductive link prediction. As described in Poursafaei et al. [20], historical NSS focuses on sampling negative edges from the set of edges that have been observed during previous timestamps but are absent in the current step. In the setting of inductive link prediction, models are asked to predict the links between the nodes unseen in the training dataset. This means when doing historical NSS, models only need to care about the previously observed edges in the test set (or validation set during validation) for choosing negative edges. This makes historical NSS the same as inductive NSS in the inductive link prediction, where inductive NSS samples negative edges that have been observed only in the test set, but not training set. Empirical results shown in Appendix C.2 Table 13 and 14 of Yu et al. [35] also prove that there is no difference between historical and inductive NSS in inductive link prediction. So we omit the results of historical NSS in our paper.

## E AUC-ROC Results on Real-World Datasets

Table 12 and 13 presents the AUC-ROC results of all baselines and DyGMamba on real-world datasets. We have similar observations as the AP results shown in Table 1 and 2. DyGMamba still demonstrates superior performance and can achieve the best average rank under any NSS setting in both transductive and inductive link prediction.

## F Dynamic Node Classification

We first give the definition of the dynamic node classification task.

**DEFINITION 3 (DYNAMIC NODE CLASSIFICATION).** *Given a CTDG  $\mathcal{G}$ , a source node  $u \in \mathcal{N}$ , a destination node  $v \in \mathcal{N}$ , a timestamp  $t \in \mathcal{T}$ , and all the interactions before  $t$ , i.e.,  $\{(u_i, v_i, t_i) | t_i < t, (u_i, v_i, t_i) \in \mathcal{G}\}$ , dynamic node classification aims to predict the state (e.g., dynamic node label) of  $u$  or  $v$  at  $t$  in the condition that the interaction  $(u, v, t)$  exists.*

We follow Rossi et al. [21], Xu et al. [33], Yu et al. [35] to conduct dynamic node classification by estimating the state of a node in a given interaction at a specific timestamp. A classification MLP is employed to map the node representations as well as the learned temporal patterns to the labels. AUC-ROC is used as the evaluation metric and we follow the dataset splits introduced in Yu et al. [35] (70%15%/15% for training/validation/testing in chronological order) for node classification. Table 14 shows the node classification results on Wikipedia and Reddit (the only two CTDG datasets for dynamic node classification), we observe that DyGMamba can achieve the best average rank, showing its strong performance. Note that both Wikipedia and Reddit are not long-range temporal dependent datasets, therefore we do not include this part into the main body of the paper. Nonetheless, DyGMamba’s great results on these datasets further prove its strength in CTDG modeling, regardless of the type of the dataset (whether long-range temporal dependent or not).

## G Further Ablation Study

We do further ablation studies here. In study C, we switch how we compute  $\beta_\theta$  in Eq. 6d to  $\beta_\theta = \text{Softmax}(f_{\text{sel}}(\mathbf{H}_\theta^t))$  ( $f_{\text{sel}}(\cdot) : \mathbb{R}^{4d} \rightarrow \mathbb{R}^{4d}$ ) to create Variant C. In study D, we base on Variant C and develop Variant D that further enables information selection from opposite nodes, i.e.,  $\beta_u = \text{Softmax}(f_{\text{sel}}(\mathbf{H}_v^t)) / \beta_v = \text{Softmax}(f_{\text{sel}}(\mathbf{H}_u^t))$ . Both ablation C and D do information selection without learning temporal patterns. In study E and F, we devise Variant E and Variant F by removing the time features and the node interaction frequency features, respectively, during neighbor encoding. From Table 15 and 16 we find that: (1) Variant C and D perform better than Variant A in most cases, implying that selecting temporal information is generally contributive; (2) Variant C generally lags behind Variant D, meaning that information selection from opposite node is beneficial; (3) DyGMamba performs better than both Variant C and D in almost all cases, proving that information selection based on temporal patterns is more effective; (4) Variant E and F in general achieve worse results than DyGMamba, validating the contribution of both time and node interaction frequency features.

## H Efficiency Analysis Complete Results

We first provide the efficiency analysis results of all baselines in this section. We then provide a comparison of total training time among DyGFormer, CTAN and DyGMamba.

### H.1 Efficiency Statistics for all baselines

We provide the efficiency statistics for all baselines in Table 17. To supplement, in Figure 5, we further plot the comparison of model size, per epoch training time and GPU Memory across models on Reddit and Social Evo., analogous to the analysis in Sec. 4.3.1. CTAN cannot be trained on Social Evo. before timeout and therefore does not appear in the plot of Social Evo.. The complete statistics including the numbers in the plots can be found in Table 17.

**Table 11: DyGFormer and DyGMamba hyperparameter searching strategy on synthetic datasets. The best settings are marked as bold.**

Models	DyGFormer		DyGMamba	
Datasets	$\rho$ & $p$	$\rho$ & $p$	$k$	
S1	{512 & 16, 256 & 8, 128 & 4, <b>64 &amp; 2</b> , 32 & 1}	{512 & 16, <b>256 &amp; 8</b> , 128 & 4, 64 & 2, 32 & 1}	{30, <b>10</b> , 5}	
S2	{ <b>512 &amp; 16</b> , 256 & 8, 128 & 4, 64 & 2, 32 & 1}	{ <b>512 &amp; 16</b> , 256 & 8, 128 & 4, 64 & 2, 32 & 1}	{ <b>30</b> , 10, 5}	
S3	{512 & 16, 256 & 8, 128 & 4, 64 & 2, <b>32 &amp; 1</b> }	{ <b>512 &amp; 16</b> , 256 & 8, 128 & 4, 64 & 2, 32 & 1}	{30, <b>10</b> , 5}	

**Table 12: AUC-ROC of transductive dynamic link prediction.**

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	EdgeBank	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	LastFM	70.89 $\pm$ 1.97	71.40 $\pm$ 2.12	71.47 $\pm$ 0.14	76.64 $\pm$ 4.66	85.92 $\pm$ 0.16	83.77 $\pm$ 0.00	71.09 $\pm$ 1.48	73.51 $\pm$ 0.14	<u>93.03 <math>\pm</math> 0.11</u>	85.12 $\pm$ 0.77	<b>93.31 <math>\pm</math> 0.18</b>
	Enron	87.77 $\pm$ 2.43	83.09 $\pm$ 2.20	68.57 $\pm$ 1.46	88.72 $\pm$ 0.95	90.34 $\pm$ 0.23	87.05 $\pm$ 0.00	83.33 $\pm$ 0.93	84.16 $\pm$ 0.34	<u>93.20 <math>\pm</math> 0.12</u>	87.09 $\pm$ 1.51	<b>93.34 <math>\pm</math> 0.23</b>
	MOOC	84.50 $\pm$ 0.60	84.50 $\pm$ 0.87	87.01 $\pm$ 0.16	<b>91.91 <math>\pm</math> 0.82</b>	80.48 $\pm$ 0.41	60.86 $\pm$ 0.00	84.02 $\pm$ 0.59	84.04 $\pm$ 0.12	88.08 $\pm$ 0.50	85.40 $\pm$ 2.67	<u>89.58 <math>\pm</math> 0.12</u>
	Reddit	98.29 $\pm$ 0.05	98.13 $\pm$ 0.04	98.50 $\pm$ 0.01	98.61 $\pm$ 0.05	99.02 $\pm$ 0.00	95.37 $\pm$ 0.00	97.67 $\pm$ 0.01	97.17 $\pm$ 0.02	<u>99.15 <math>\pm</math> 0.01</u>	97.24 $\pm$ 0.75	<b>99.27 <math>\pm</math> 0.01</b>
	Wikipedia	96.36 $\pm$ 0.14	94.43 $\pm$ 0.32	96.60 $\pm$ 0.07	98.37 $\pm$ 0.10	98.54 $\pm$ 0.01	90.78 $\pm$ 0.00	97.27 $\pm$ 0.06	96.89 $\pm$ 0.04	<u>98.92 <math>\pm</math> 0.03</u>	97.00 $\pm$ 0.21	<b>99.08 <math>\pm</math> 0.02</b>
	UCI	90.35 $\pm$ 0.51	69.46 $\pm$ 2.66	78.76 $\pm$ 1.10	92.03 $\pm$ 0.69	93.81 $\pm$ 0.23	77.30 $\pm$ 0.00	85.49 $\pm$ 0.82	91.62 $\pm$ 0.52	<u>94.45 <math>\pm</math> 0.22</u>	76.25 $\pm$ 2.83	<b>94.77 <math>\pm</math> 0.18</b>
	Social Evo.	92.13 $\pm$ 0.20	90.37 $\pm$ 0.52	94.93 $\pm$ 0.06	95.31 $\pm$ 0.27	87.34 $\pm$ 0.10	81.60 $\pm$ 0.00	95.45 $\pm$ 0.21	95.21 $\pm$ 0.07	<u>96.25 <math>\pm</math> 0.04</u>	Timeout	<b>96.38 <math>\pm</math> 0.02</b>
Avg. Rank		7.14	8.86	7.14	3.86	4.86	9.14	7.29	7.14	<u>2.14</u>	7.29	<b>1.14</b>
Historical	LastFM	75.65 $\pm$ 4.43	70.63 $\pm$ 2.56	64.23 $\pm$ 0.45	78.00 $\pm$ 2.97	67.92 $\pm$ 0.32	78.09 $\pm$ 0.00	60.53 $\pm$ 2.54	64.06 $\pm$ 0.34	78.80 $\pm$ 0.02	<u>79.50 <math>\pm</math> 0.82</u>	<b>79.82 <math>\pm</math> 0.27</b>
	Enron	75.21 $\pm$ 1.27	76.36 $\pm$ 1.42	62.36 $\pm$ 1.07	76.75 $\pm$ 1.40	65.62 $\pm$ 0.49	<u>79.59 <math>\pm</math> 0.00</u>	71.72 $\pm$ 1.24	74.82 $\pm$ 2.04	77.35 $\pm$ 0.64	<b>81.95 <math>\pm</math> 1.64</b>	77.73 $\pm$ 0.61
	MOOC	82.38 $\pm$ 1.75	80.71 $\pm$ 2.08	81.53 $\pm$ 0.79	86.59 $\pm$ 2.03	71.74 $\pm$ 0.88	61.90 $\pm$ 0.00	73.22 $\pm$ 1.21	77.09 $\pm$ 0.83	<u>87.26 <math>\pm</math> 0.83</u>	73.87 $\pm$ 2.77	<b>87.91 <math>\pm</math> 0.93</b>
	Reddit	80.70 $\pm$ 0.20	79.96 $\pm$ 0.23	79.60 $\pm$ 0.09	81.04 $\pm$ 0.23	80.42 $\pm$ 0.20	78.58 $\pm$ 0.00	76.83 $\pm$ 0.12	77.83 $\pm$ 0.33	80.61 $\pm$ 0.48	<b>90.63 <math>\pm</math> 2.28</b>	81.71 $\pm$ 0.49
	Wikipedia	80.71 $\pm$ 0.64	77.49 $\pm$ 0.72	82.83 $\pm$ 0.27	83.28 $\pm$ 0.26	65.74 $\pm$ 3.46	77.27 $\pm$ 0.00	85.55 $\pm$ 0.47	<u>87.47 <math>\pm</math> 0.20</u>	72.78 $\pm$ 6.65	<b>95.43 <math>\pm</math> 0.07</b>	78.99 $\pm$ 1.24
	UCI	78.21 $\pm$ 3.18	58.65 $\pm$ 3.58	57.12 $\pm$ 0.98	<b>78.48 <math>\pm</math> 1.79</b>	57.67 $\pm$ 1.11	69.56 $\pm$ 0.00	65.42 $\pm$ 2.62	<u>77.46 <math>\pm</math> 1.63</u>	75.71 $\pm$ 0.57	75.05 $\pm$ 0.13	75.43 $\pm$ 1.99
	Social Evo.	91.83 $\pm$ 1.52	92.81 $\pm$ 0.60	93.63 $\pm$ 0.48	94.27 $\pm$ 1.33	87.61 $\pm$ 0.06	85.81 $\pm$ 0.00	95.03 $\pm$ 0.82	94.65 $\pm$ 0.28	<u>97.16 <math>\pm</math> 0.06</u>	Timeout	<b>97.27 <math>\pm</math> 0.30</b>
Avg. Rank		5.29	7.14	7.86	3.71	9.14	7.43	7.71	6.29	<u>4.29</u>	<u>4.29</u>	<b>2.86</b>
Inductive	LastFM	61.59 $\pm$ 5.72	60.62 $\pm$ 2.20	63.96 $\pm$ 0.41	65.48 $\pm$ 4.13	67.90 $\pm$ 0.44	77.37 $\pm$ 0.00	54.75 $\pm$ 1.31	59.98 $\pm$ 0.20	67.87 $\pm$ 0.53	<b>78.70 <math>\pm</math> 0.87</b>	68.74 $\pm$ 0.55
	Enron	70.75 $\pm$ 0.69	67.37 $\pm$ 2.21	59.78 $\pm$ 1.12	73.22 $\pm$ 0.42	75.29 $\pm$ 0.66	75.00 $\pm$ 0.00	69.74 $\pm$ 1.19	70.72 $\pm$ 1.08	74.67 $\pm$ 0.80	<u>75.40 <math>\pm</math> 1.92</u>	<b>75.47 <math>\pm</math> 1.41</b>
	MOOC	67.53 $\pm$ 1.76	62.60 $\pm$ 1.27	74.44 $\pm$ 0.81	76.89 $\pm$ 2.13	70.08 $\pm$ 0.33	48.18 $\pm$ 0.00	71.80 $\pm$ 1.09	72.25 $\pm$ 0.57	<u>80.78 <math>\pm</math> 0.89</u>	68.17 $\pm$ 3.73	<b>81.08 <math>\pm</math> 0.82</b>
	Reddit	83.40 $\pm$ 0.33	82.75 $\pm$ 0.36	87.46 $\pm$ 0.10	84.57 $\pm$ 0.19	88.19 $\pm$ 0.20	85.93 $\pm$ 0.00	84.41 $\pm$ 0.18	82.24 $\pm$ 0.24	86.25 $\pm$ 0.64	<b>91.42 <math>\pm</math> 2.18</b>	86.35 $\pm$ 0.52
	Wikipedia	70.41 $\pm$ 0.39	67.57 $\pm$ 0.94	81.54 $\pm$ 0.31	81.21 $\pm$ 0.30	68.48 $\pm$ 3.64	<u>81.73 <math>\pm</math> 0.00</u>	73.51 $\pm$ 1.88	84.20 $\pm$ 0.36	64.09 $\pm$ 9.75	<b>93.67 <math>\pm</math> 0.11</b>	75.64 $\pm$ 2.42
	UCI	64.14 $\pm$ 1.25	54.10 $\pm$ 2.74	59.60 $\pm$ 0.61	63.76 $\pm$ 0.99	57.85 $\pm$ 0.59	58.03 $\pm$ 0.00	65.46 $\pm$ 2.07	<b>74.25 <math>\pm</math> 0.71</b>	64.92 $\pm$ 0.83	66.51 $\pm$ 0.25	<u>66.83 <math>\pm</math> 2.83</u>
	Social Evo.	91.81 $\pm$ 1.69	92.77 $\pm$ 0.64	93.54 $\pm$ 0.48	94.86 $\pm$ 1.25	90.10 $\pm$ 0.11	87.88 $\pm$ 0.00	<u>95.13 <math>\pm</math> 0.83</u>	94.50 $\pm$ 0.26	95.01 $\pm$ 0.15	Timeout	<b>97.37 <math>\pm</math> 0.26</b>
Avg. Rank		7.86	9.57	6.14	5.43	6.29	6.43	6.71	6.00	5.14	<u>3.86</u>	<b>2.57</b>

**Table 13: AUC-ROC of inductive dynamic link prediction. EdgeBank cannot do inductive link prediction so is not reported.**

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	LastFM	82.49 $\pm$ 0.94	82.82 $\pm$ 1.17	76.76 $\pm$ 0.22	82.61 $\pm$ 2.62	87.92 $\pm$ 0.15	76.95 $\pm$ 1.34	80.34 $\pm$ 0.14	<u>94.10 <math>\pm</math> 0.09</u>	61.49 $\pm$ 2.78	<b>94.37 <math>\pm</math> 0.13</b>
	Enron	80.16 $\pm$ 1.50	75.82 $\pm$ 3.14	64.25 $\pm$ 1.29	79.40 $\pm$ 1.77	86.84 $\pm$ 0.89	81.03 $\pm$ 0.93	76.08 $\pm$ 0.92	<u>89.59 <math>\pm</math> 0.10</u>	75.23 $\pm$ 2.24	<b>89.76 <math>\pm</math> 0.21</b>
	MOOC	83.82 $\pm$ 0.30	83.42 $\pm$ 0.77	86.67 $\pm$ 0.24	<b>91.58 <math>\pm</math> 0.74</b>	81.76 $\pm$ 0.46	82.42 $\pm$ 0.71	82.76 $\pm$ 0.13	87.75 $\pm$ 0.42	66.38 $\pm$ 1.59	<u>89.34 <math>\pm</math> 0.12</u>
	Reddit	96.42 $\pm$ 0.13	95.87 $\pm$ 0.21	97.02 $\pm$ 0.04	97.30 $\pm$ 0.12	98.42 $\pm$ 0.01	94.63 $\pm$ 0.08	94.95 $\pm$ 0.08	<u>98.70 <math>\pm</math> 0.02</u>	82.35 $\pm$ 4.03	<b>98.88 <math>\pm</math> 0.01</b>
	Wikipedia	94.43 $\pm$ 0.28	91.31 $\pm$ 0.40	95.93 $\pm$ 0.19	97.71 $\pm$ 0.19	98.05 $\pm$ 0.03	97.03 $\pm$ 0.08	96.26 $\pm$ 0.04	<u>98.49 <math>\pm</math> 0.02</u>	92.59 $\pm$ 0.70	<b>98.72 <math>\pm</math> 0.03</b>
	UCI	78.78 $\pm$ 1.11	58.84 $\pm$ 2.54	77.41 $\pm$ 0.65	86.27 $\pm$ 1.49	90.27 $\pm$ 0.40	81.67 $\pm$ 1.01	89.26 $\pm$ 0.42	<u>92.43 <math>\pm</math> 0.20</u>	48.58 $\pm$ 6.02	<b>92.70 <math>\pm</math> 0.19</b>
	Social Evo.	93.62 $\pm$ 0.36	90.20 $\pm$ 2.05	93.52 $\pm$ 0.05	93.21 $\pm$ 0.90	84.73 $\pm$ 0.20	94.63 $\pm$ 0.06	94.09 $\pm$ 0.03	<u>95.30 <math>\pm</math> 0.05</u>	Timeout	<b>95.36 <math>\pm</math> 0.04</b>
Avg. Rank		6.00	7.43	7.00	4.57	4.71	6.14	6.14	<u>2.14</u>	9.71	<b>1.14</b>
Inductive	LastFM	69.85 $\pm$ 1.70	68.14 $\pm$ 1.61	69.89 $\pm$ 0.41	67.01 $\pm$ 5.77	67.72 $\pm$ 0.20	63.15 $\pm$ 1.17	69.93 $\pm$ 0.17	69.86 $\pm$ 0.80	57.85 $\pm$ 3.67	<b>70.59 <math>\pm</math> 0.57</b>
	Enron	65.95 $\pm$ 1.27	62.20 $\pm$ 2.15	56.52 $\pm$ 0.84	64.21 $\pm$ 0.94	62.07 $\pm$ 0.72	67.56 $\pm$ 1.34	67.39 $\pm$ 1.33	66.07 $\pm$ 0.65	<u>68.70 <math>\pm</math> 1.82</u>	<b>68.98 <math>\pm</math> 1.00</b>
	MOOC	65.37 $\pm$ 0.96	62.97 $\pm$ 2.05	74.94 $\pm$ 0.80	76.36 $\pm$ 2.91	71.18 $\pm$ 0.54	71.30 $\pm$ 1.21	72.15 $\pm$ 0.65	<u>80.42 <math>\pm</math> 0.72</u>	58.06 $\pm$ 0.89	<b>81.12 <math>\pm</math> 0.63</b>
	Reddit	61.84 $\pm$ 0.44	60.35 $\pm$ 0.53	64.92 $\pm$ 0.08	65.24 $\pm$ 0.08	<u>65.37 <math>\pm</math> 0.12</u>	61.85 $\pm$ 0.11	64.56 $\pm$ 0.26	64.80 $\pm$ 0.53	<b>81.70 <math>\pm</math> 4.71</b>	64.93 $\pm$ 0.89
	Wikipedia	61.66 $\pm$ 0.30	56.34 $\pm$ 0.67	78.40 $\pm$ 0.77	75.86 $\pm$ 0.50	59.00 $\pm$ 4.33	71.45 $\pm$ 2.23	<u>82.76 <math>\pm</math> 0.11</u>	58.21 $\pm$ 8.78	<b>91.12 <math>\pm</math> 0.13</b>	67.92 $\pm$ 2.23
	UCI	60.66 $\pm$ 1.82	51.50 $\pm$ 2.08	61.27 $\pm$ 0.78	62.07 $\pm$ 0.67	55.60 $\pm$ 1.22	65.87 $\pm$ 1.90	<b>75.72 <math>\pm</math> 0.70</b>	64.37 $\pm$ 0.98	51.68 $\pm$ 2.60	<u>66.95 <math>\pm</math> 2.22</u>
	Social Evo.	88.98 $\pm$ 0.81	86.43 $\pm$ 1.48	92.37 $\pm$ 0.50	91.66 $\pm$ 2.14	83.84 $\pm$ 0.21	<u>95.50 <math>\pm</math> 0.31</u>	93.88 $\pm$ 0.22	94.97 $\pm$ 0.36	Timeout	<b>96.65 <math>\pm</math> 0.29</b>
Avg. Rank		7.00	8.71	5.14	5.14	7.14	5.14	<u>3.57</u>	4.71	6.14	<b>2.29</b>

## H.2 Total Training Time Comparison among DyGFormer, CTAN and DyGMamba

We present the per epoch training time, number of epochs until the best performance and the total training time in Table 18. Total

**Table 14: AUC-ROC of dynamic node classification.**

Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Wikipedia	<b>88.10 <math>\pm</math> 1.57</b>	87.41 $\pm$ 1.94	83.42 $\pm$ 2.92	85.51 $\pm$ 3.28	84.59 $\pm$ 1.16	79.03 $\pm$ 1.18	85.60 $\pm$ 1.73	86.35 $\pm$ 2.19	87.38 $\pm$ 0.14	<u>87.44 <math>\pm</math> 0.82</u>
Reddit	59.53 $\pm$ 3.18	63.12 $\pm$ 0.51	<b>69.31 <math>\pm</math> 2.18</b>	63.21 $\pm$ 3.00	65.22 $\pm$ 0.79	<u>68.04 <math>\pm</math> 2.00</u>	64.42 $\pm$ 1.15	67.67 $\pm$ 1.39	67.29 $\pm$ 0.15	67.70 $\pm$ 1.32
<b>Avg. Rank</b>	5.50	6.00	5.00	7.50	7.00	6.00	6.50	<u>4.50</u>	<u>4.50</u>	<b>2.50</b>

**Table 15: Ablation studies under transductive setting. R/H/I means random/historical/inductive NSS. Metric is AP.**

Datasets	LastFM			Enron			MOOC			Reddit			Wikipedia			UCI			Social Evo.		
Models	R	H	I	R	H	I	R	H	I	R	H	I	R	H	I	R	H	I	R	H	I
Variant C	92.71	82.85	72.36	92.49	76.99	78.64	88.80	85.23	81.02	99.27	81.74	91.05	99.06	79.14	73.49	95.85	81.00	71.86	94.71	96.71	97.25
Variant D	92.74	82.87	72.68	92.52	77.07	78.05	88.71	85.76	81.09	99.27	82.10	91.07	99.08	81.75	79.79	95.87	<b>82.35</b>	<b>72.98</b>	94.74	97.17	97.60
Variant E	85.80	63.98	70.09	89.09	70.85	85.42	83.25	82.18	75.06	99.00	<b>82.21</b>	91.02	98.92	81.21	76.62	95.12	82.11	70.04	94.18	96.97	97.50
Variant F	87.47	67.11	64.10	88.32	69.16	83.98	82.42	76.18	74.12	99.08	76.09	88.75	98.76	72.41	70.03	94.74	63.94	63.56	94.17	96.15	96.29
DyGMamba	<b>93.35</b>	<b>83.02</b>	<b>73.63</b>	<b>92.65</b>	<b>77.77</b>	<b>80.86</b>	<b>89.21</b>	<b>85.89</b>	<b>81.11</b>	<b>99.32</b>	81.80	<b>91.15</b>	<b>99.15</b>	<b>81.77</b>	<b>79.86</b>	<b>95.91</b>	81.03	71.95	<b>94.77</b>	<b>97.35</b>	<b>97.68</b>

**Table 16: Ablation studies under inductive setting. R/I means random/inductive NSS. Metric is AP.**

Datasets	LastFM		Enron		MOOC		Reddit		Wikipedia		UCI		Social Evo.	
Models	R	I	R	I	R	I	R	I	R	I	R	I	R	I
Variant C	94.18	76.44	89.40	68.33	88.59	80.39	98.90	64.07	98.65	69.82	94.47	72.05	93.07	96.20
Variant D	94.21	76.64	89.44	67.91	88.29	<b>80.86</b>	98.91	65.10	98.69	71.10	94.51	73.50	93.10	96.75
Variant E	88.36	55.95	89.08	68.05	82.24	74.72	98.70	65.20	98.40	70.99	93.47	<b>74.64</b>	92.76	96.50
Variant F	87.71	66.45	84.02	66.85	81.13	73.85	98.54	61.74	98.37	67.59	92.81	64.59	92.42	95.38
DyGMamba	<b>94.42</b>	<b>76.76</b>	<b>89.67</b>	<b>68.77</b>	<b>88.64</b>	80.75	<b>98.97</b>	<b>65.30</b>	<b>98.77</b>	<b>71.14</b>	<b>94.76</b>	72.17	<b>93.13</b>	<b>96.83</b>

**Table 17: Efficiency statistics for all baselines. EdgeBank is non-parameterized and not a machine learning model so we omit it here. # Params means number of parameters (MB). Time and Mem denote per epoch training time (min) and GPU memory (GB), respectively. The numbers in this table are the average results of five runs with different random seeds.**

Datasets	LastFM			Enron			MOOC			UCI			Reddit			Social Evo.		
Models	# Params	Time	Mem	# Params	Time	Mem	# Params	Time	Mem	# Params	Time	Mem	# Params	Time	Mem	# Params	Time	Mem
JODIE	0.75	4.4	2.28	0.75	0.07	1.30	0.75	0.78	2.36	0.75	0.03	1.44	0.75	3.95	1.10	0.75	4.70	1.71
DyRep	2.64	6.6	2.29	2.64	0.10	1.34	2.64	0.88	2.38	2.64	0.05	1.51	2.64	5.75	1.21	2.64	7.55	1.76
TGAT	4.02	22.75	4.15	4.02	1.28	3.46	4.02	4.08	3.64	4.02	0.60	3.42	4.02	16.33	2.98	4.02	25.50	3.89
TGN	3.68	12.14	2.21	3.68	0.15	1.45	3.68	1.03	2.54	3.68	0.08	1.51	3.67	2.05	1.67	3.67	3.83	1.78
CAWN	15.35	99.00	14.92	15.35	2.62	4.03	15.35	13.45	8.02	15.35	1.95	9.40	15.35	20.16	5.89	15.35	85.66	8.14
TCL	3.37	6.23	3.04	3.37	0.30	2.51	3.37	1.00	2.49	3.37	0.13	2.00	3.37	2.25	1.82	3.37	5.05	2.48
GraphMixer	2.45	16.35	2.78	2.45	1.20	2.23	2.45	4.02	2.40	2.45	0.73	2.19	2.44	4.92	1.57	2.45	15.50	2.71
DyGFormer	5.56	47.00	7.57	4.80	2.73	3.23	4.80	8.32	3.35	4.15	0.62	2.30	4.24	7.00	2.42	4.14	20.00	2.77
CTAN	0.45	3.33	1.44	0.47	0.50	1.33	0.68	3.22	2.30	0.50	0.38	1.30	0.53	0.86	1.54	0.45	2.41	0.63
DyGMamba	2.78	28.45	4.17	2.03	2.05	2.74	1.65	4.88	2.48	1.37	0.60	1.93	3.32	6.30	2.07	3.22	17.80	2.59

training time computes the total amount of time a model requires to reach its maximum performance, without considering the patience during training. We observe that CTAN requires much more epochs to converge, e.g., on LastFM it uses almost 54 times of epochs than DyGMamba to reach its best performance.

### H.3 Modeling an Increasing Number of Temporal Neighbors with Limited Total Training Time

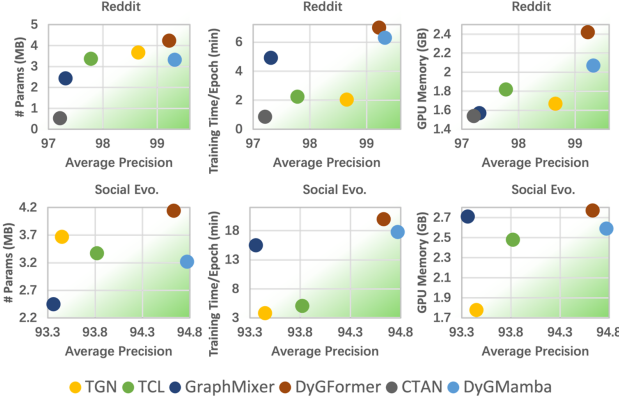
To further show DyGMamba’s superior efficiency against baseline methods, we do the following experiments. We train five best performing models (as shown in Table 1) on Enron with a gradually

increasing number of temporal neighbors<sup>9</sup> and report their performance. The number of sampled neighbors spans from 8, 16, 32, 64, 128 to 256 (Note that these numbers are different from the best hyperparameters reported in Yu et al. [35]). We fix the patch size  $p$  of DyGFormer and DyGMamba to 1 in order to maximize their input sequence lengths. We set a time limit of 120 minutes for the total training time. We let all the experiments finish the complete training process and note down the ones that exceed the time limit. In this way, we not only care about the per epoch training time, but also pay attention to how long it takes for models to converge. The experimental results are reported in Fig. 6. The points marked with crosses (×) mean that the training process cannot finish within

<sup>9</sup>For CTAN, by number of temporal neighbors we mean the sampler size in each graph convolutional layer, i.e., the size of the sampled temporal neighborhood for each node at a timestamp.

**Table 18: Comparison among DyGFormer, CTAN and DyGMamba on per epoch training time ( $T_{ep}$  (min)), number of epochs until the best performance (# Epoch) and the total training time ( $T_{tot}$  (min)).  $T_{tot} = T_{ep} \times \# \text{ Epoch}$ . The numbers in this table are the average results of five runs with different random seeds. CTAN cannot be trained on Social Evo. before timeout and thus without # Epoch and  $T_{tot}$ .**

Datasets	LastFM			Enron			MOOC			UCI			Reddit			Social Evo.		
Models	$T_{ep}$	# Epoch	$T_{tot}$	$T_{ep}$	# Epoch	$T_{tot}$	$T_{ep}$	# Epoch	$T_{tot}$	$T_{ep}$	# Epoch	$T_{tot}$	$T_{ep}$	# Epoch	$T_{tot}$	$T_{ep}$	# Epoch	$T_{tot}$
DyGFormer	47.00	49.60	2331.20	2.73	32.80	89.54	8.32	64.20	534.14	0.62	34.80	21.58	4.24	24.60	104.30	20.00	30.22	604.40
CTAN	3.33	635.00	2114.55	0.50	173.00	86.50	3.22	138.00	444.36	0.38	236.00	89.68	0.53	327.18	173.41	0.45	Timeout	Timeout
DyGMamba	28.45	11.80	335.71	2.05	33.00	67.65	4.88	38.00	185.44	0.60	28.00	16.80	3.32	26.80	88.98	17.80	24.40	434.32

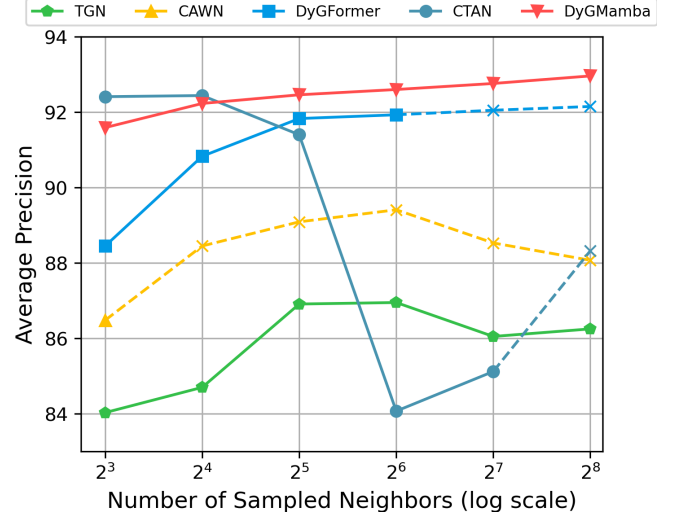


**Figure 5: Efficiency comparison on Reddit and Social Evo. among DyGMamba and five baselines in terms of number (#) of parameters, training time per epoch and GPU memory.**

the time limit (although we still plot their corresponding performance). We find that only TGN and DyGMamba can successfully converge within the time limit when the number of considered neighbors increases to 256. DyGMamba can constantly achieve performance gain from modeling more temporal neighbors while TGN cannot. CAWN is extremely time consuming so it cannot finish training within the time limit even when it is asked to model 16 temporal neighbors. As for the methods designed for long-range temporal information propagation, DyGFormer and CTAN consume much longer total training time than DyGMamba. They fail to converge within 120 minutes when the number of considered neighbors reaches 128 and 256, respectively. We also observe that CTAN’s performance fluctuates greatly with the increasing temporal neighbors, indicating that it is not stable to model a large number of temporal neighbors. This also implies that increasing the amount of historical information will gradually make CTAN harder to converge, which might cause trouble in modeling long range temporal dependent datasets.

## I Details of S4 and Mamba Operations

**Single-Input Single-Output.** Given a sequence of vector elements as input, SISO means that the SSM processes each input dimension in parallel with the same set of parameters. For example, a sequence of  $d_2$ -dimensional vectors will be split into  $d_2$  1-dimensional sequences with the same sequence length. Each of them will be computed in parallel as in Eq. 2 with a shared set of SSM parameters.



**Figure 6: Performance comparison among TGN, CAWN, DyGFormer, CTAN and DyGMamba on Enron, with an increasing number of encoded temporal neighbors. The metric is AP under random NSS on transductive link prediction. The time limit for training is 120 min. If a model fails to complete training within this limit, a cross  $\times$  is used to mark the data point. Dashed lines indicate that models start to exceed time limit as neighbor number increases.**

After computation, all these  $d_2$  sequences will be rearranged back into a sequence of  $d_2$ -dimensional vectors. SISO fails to mix the information across dimensions of each vector. To address this, S4 and Mamba employs a mixing linear layer  $f_{mix}(\cdot) : \mathbb{R}^{d_2} \rightarrow \mathbb{R}^{d_2}$  on each  $d_2$ -dimensional vector to mix the information across  $d_2$  dimensions. For more details, please refer to [8], [24] and [6].

**SSM Function.**  $SSM_{\bar{A}, \bar{B}, C}(\cdot)$  takes a matrix as input. The input matrix can be considered as a sequence of vector elements, where each row of the matrix corresponds to an element. The output of  $SSM_{\bar{A}, \bar{B}, C}(\cdot)$  is also a matrix, where each row of the output matrix is the output of its corresponding input vector element.  $SSM_{\bar{A}, \bar{B}, C}(\cdot)$  can be viewed as using S4 or Mamba to process a sequence of vectors in the SISO fashion, based on their parameters  $\bar{A}, \bar{B}, C$ .



## J Motivation of Using SSM for Temporal Pattern Modeling

The biggest motivation of using SSM for temporal pattern modeling is that it helps to maintain good efficiency. If in the future we want to deploy DyGMamba on larger datasets that require much longer historical histories for modeling, the value of  $k$  will also increase accordingly and the time difference sequence will not be short anymore. Besides, as we have chosen SSM to model historical one-hop temporal neighbors in the node-level SSM, it is natural to employ another SSM for temporal pattern modeling.

## K Impact of Patch Size on MOOC

Apart from the analysis on Enron, we further analyze the impact of patch size on another long-term temporal dependent dataset MOOC. Fig. 7b shows the numbers of parameters of DyGFormer and DyGMamba with different patch sizes on MOOC. We confirm that patching greatly affects model sizes. We decrease the patch size gradually from 4 to 1 (the optimal value of DyGMamba’s hyperparameter  $\rho$  &  $p$  is 128 & 4) and track DyGMamba’s performance (Fig. 7a) as well as efficiency (Fig. 7b to 7d) on MOOC. Meanwhile, we also keep track on DyGFormer under the same patch size for comparison. Same as our analysis on Enron, we plot the performance of Variant A under different patch sizes in Fig. 7a as well. **We can draw the same conclusions as in our analysis on Enron.**

## L Performance of DyGMamba on Discrete-Time Dynamic Graphs

To better benchmark DyGMamba, we test the performance of DyGMamba on 6 DTDG datasets collected in [35] (i.e., Flights, Can. Parl, US Legis., UN Trade, UN Vote and Contact) and compare it with the 10 recent baselines discussed in Sec. 4.1 (i.e., JODIE, DyRep, TGAT, TGN, CAWN, EdgeBank, TCL, GraphMixer, DyGFormer, CTAN)<sup>10</sup>. The value of the hyperparameter  $\rho$  &  $p$  of DyGMamba is set as same as the one for DyGFormer. The optimal  $\rho$  &  $p$  for Flights, Can. Parl, US Legis., UN Trade, UN Vote and Contact is {256 & 8, 2048 & 64, 256 & 8, 256 & 8, 256 & 8, 128 & 4, 32 & 1}. We report the best hyperparameter  $k$  of DyGMamba for each dataset in Table 19. We use the implementations and the best hyperparameters provided by Yu et al. [35] for all baseline models except CTAN. For CTAN, we use its official implementation, fixing the number of layers to 5. All models are trained with a batch size of 200 for fair efficiency analysis. We report the AP of both models on DTDG datasets in the transductive and inductive settings in Table 20 and Table 21, respectively. In addition, Table 22 and Table 23 present the AUC-ROC of both models in the transductive and inductive settings, respectively.

We find that DyGMamba achieves strong overall performance (Avg. Rank) on DTDGs in both transductive and inductive link prediction tasks. However, compared to its performance on CTDGs, DyGMamba’s performance on DTDGs is not always highly competitive. This is expected, as DyGMamba is not originally designed for DTDGs and lacks the ability to optimally handle concurrent edges.

<sup>10</sup> Although we have discussed in Sec. ?? that DyGMamba is not suitable to reason over DTDGs, we still benchmark our model on them to show its effectiveness.

One interesting finding is that on the long-range temporal dependent dataset Can. Parl which requires sampling 2048 neighbors for modeling, DyGMamba can benefit from such long neighbor sequence and achieve the best performance among all models, indicating the importance of modeling long-term temporal information as well as the strong capability of our model in capturing it. More importantly, we find that DyGMamba can benefit from greater value of  $k$  as the number of the sampled neighbors  $\rho$  increases. For example, on Can. Parl, as shown in Table 19, the optimal value of  $k$  is 100. This makes our selection of using Mamba for temporal pattern modeling more reasonable since Mamba can better demonstrate its advantage in efficiency when there is a growing time difference sequence corresponding to the temporal pattern.

**Table 19: DyGMamba hyperparameter searching strategy on DTDG datasets. The best settings are marked as bold.**

Datasets	$k$
Flights	{ <b>30</b> , 10, 5}
Can. Parl.	{200, <b>100</b> , 30}
US Legis.	{ <b>30</b> , 10, 5}
UN Trade	{ <b>30</b> , 10, 5}
UN Vote	{30, <b>10</b> , 5}
Contact	{30, 10, <b>5}</b> }

## M Detailed Discussion of Motivation

The goal of this work is to introduce a model capable of effective and efficient reasoning over CTDGs using extensive historical information. Our motivation is supported by the following considerations:

*Abundant Historical Information in Real-World Datasets.* We quantify the availability of historical neighbors across various datasets for transductive and inductive link prediction tasks under random, historical, and inductive NSS settings, as shown in Table 24. Observations indicate that real-world CTDG datasets, such as LastFM and Enron, typically possess extensive one-hop historical information. Although it is feasible to perform predictions based solely on recent interactions, we posit that effectively leveraging longer historical contexts can further enhance model performance.

*Limitations of Existing Models in Handling Extensive Historical Information.* Two recent state-of-the-art methods, DyGFormer and CTAN, explicitly designed to capture long-range historical dependencies, have been considered in our study. Our experimental results reveal distinct limitations: DyGFormer suffers from high computational complexity, while CTAN exhibits comparatively weaker predictive performance. Consequently, neither model achieves an optimal balance between effectiveness and computational efficiency. Our proposed model, DyGMamba, overcomes these limitations by demonstrating superior predictive capability and efficiency. Extensive experiments position DyGMamba as the top ranked method across multiple benchmark datasets.

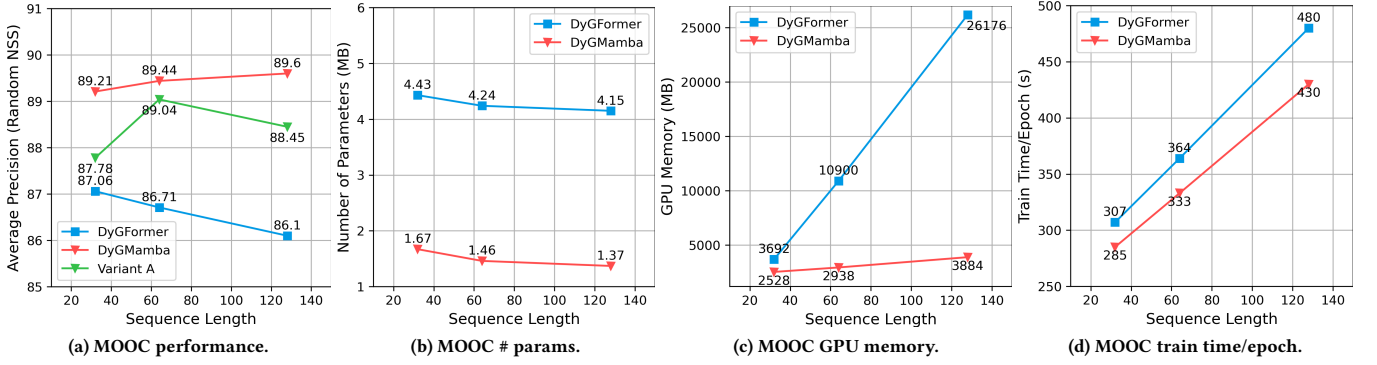


Figure 7: Impact of patch size on DyGFormer, DyGMamba and Variant A, given a fixed number of sampled temporal neighbors  $\rho$  on MOOC. Patch size  $p$  varies from 4, 2, 1. Sequence length  $\rho/p$  increases as patch size decreases. Performance is the transductive AP under random NSS.

Table 20: AP of transductive dynamic link prediction on DTDGs. The best and the second best results are marked as bold and underlined, respectively.

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	EdgeBank	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	Can. Parl.	69.58 ± 0.28	68.94 ± 3.55	71.57 ± 0.68	69.37 ± 1.86	69.07 ± 3.01	64.55 ± 0.00	71.49 ± 0.46	77.49 ± 0.62	97.63 ± 0.21	87.01 ± 1.10	<b>99.57 ± 0.08</b>
	Contact	94.75 ± 0.82	95.97 ± 0.16	96.55 ± 0.11	97.21 ± 0.44	90.65 ± 0.16	92.58 ± 0.00	94.67 ± 0.11	91.88 ± 0.07	98.30 ± 0.00	<b>98.53 ± 0.09</b>	<u>98.43 ± 0.12</u>
	Flights	96.03 ± 1.17	95.49 ± 0.49	93.92 ± 0.04	98.01 ± 0.39	98.50 ± 0.01	89.35 ± 0.00	91.30 ± 0.05	91.02 ± 0.04	<u>98.92 ± 0.01</u>	93.64 ± 0.54	<b>98.95 ± 0.05</b>
	UN Trade	64.43 ± 0.25	62.58 ± 0.93	61.53 ± 0.18	64.03 ± 0.69	65.41 ± 0.04	60.41 ± 0.00	62.19 ± 0.07	62.87 ± 0.12	66.69 ± 0.33	<b>69.01 ± 0.38</b>	<u>67.50 ± 0.24</u>
	UN Vote	63.35 ± 0.99	63.90 ± 0.08	52.89 ± 0.28	<u>65.18 ± 0.61</u>	52.77 ± 0.07	58.49 ± 0.00	51.89 ± 0.18	52.18 ± 0.13	55.90 ± 0.24	<b>70.19 ± 1.26</b>	56.39 ± 0.18
	US Legis.	74.29 ± 0.68	73.81 ± 2.29	67.85 ± 2.80	<u>76.46 ± 0.37</u>	70.60 ± 0.25	58.39 ± 0.00	70.99 ± 0.25	72.44 ± 0.29	71.06 ± 1.06	<b>79.44 ± 2.28</b>	71.75 ± 0.26
	Avg. Rank	5.17	6.17	7.50	4.33	7.50	9.67	8.50	7.67	4.00	2.50	3.00
Historical	Can. Parl.	51.66 ± 0.51	62.46 ± 2.73	68.87 ± 0.96	66.49 ± 0.47	65.58 ± 3.34	63.76 ± 0.00	68.92 ± 2.04	74.58 ± 1.42	<u>95.72 ± 1.84</u>	89.28 ± 1.24	<b>99.77 ± 0.12</b>
	Contact	94.33 ± 1.77	96.27 ± 0.17	96.75 ± 0.72	96.39 ± 0.47	84.61 ± 0.52	88.83 ± 0.00	95.26 ± 0.19	93.29 ± 0.47	97.59 ± 0.18	<b>97.86 ± 0.15</b>	<u>97.61 ± 0.04</u>
	Flights	66.56 ± 1.39	67.53 ± 1.79	<u>72.60 ± 0.33</u>	67.77 ± 1.08	64.60 ± 0.64	70.53 ± 0.00	71.13 ± 0.19	71.40 ± 0.32	65.63 ± 0.21	<b>75.75 ± 8.04</b>	67.80 ± 2.17
	UN Trade	61.38 ± 1.27	58.56 ± 1.04	54.64 ± 0.22	56.64 ± 3.43	57.09 ± 1.80	<b>81.08 ± 0.00</b>	56.80 ± 1.20	57.91 ± 0.81	62.86 ± 1.66	<u>67.38 ± 1.08</u>	65.10 ± 0.02
	UN Vote	71.11 ± 1.00	70.33 ± 1.21	53.67 ± 2.31	69.16 ± 1.82	51.29 ± 0.41	<b>84.76 ± 0.00</b>	53.78 ± 2.16	54.17 ± 1.08	60.59 ± 0.94	<u>74.75 ± 1.81</u>	61.07 ± 0.39
	US Legis.	47.20 ± 1.72	81.60 ± 7.14	76.28 ± 4.16	67.45 ± 6.79	73.25 ± 14.99	63.31 ± 0.00	83.27 ± 0.88	<b>86.05 ± 1.77</b>	63.99 ± 11.58	<u>84.20 ± 0.95</u>	82.15 ± 1.02
	Avg. Rank	7.83	6.50	6.50	7.00	9.33	6.00	6.17	5.33	5.83	1.83	3.67
Inductive	Can. Parl.	48.13 ± 0.42	57.90 ± 2.61	69.54 ± 0.66	63.58 ± 1.59	66.94 ± 1.60	62.20 ± 0.00	68.23 ± 1.42	70.00 ± 0.57	<u>95.01 ± 0.80</u>	87.42 ± 0.71	<b>98.32 ± 0.34</b>
	Contact	92.34 ± 1.91	93.99 ± 0.24	95.14 ± 1.19	94.36 ± 0.88	89.61 ± 0.36	85.19 ± 0.00	92.08 ± 0.32	90.64 ± 0.56	94.93 ± 0.54	<b>97.31 ± 0.19</b>	<u>95.43 ± 0.17</u>
	Flights	69.39 ± 2.19	71.32 ± 2.44	75.68 ± 0.37	72.27 ± 1.33	89.01 ± 0.35	<b>81.07 ± 0.00</b>	75.10 ± 0.08	74.68 ± 0.35	70.31 ± 1.00	<u>76.20 ± 7.96</u>	73.79 ± 5.69
	UN Trade	60.55 ± 1.20	60.15 ± 0.52	58.86 ± 0.38	59.52 ± 3.79	63.37 ± 2.14	<b>73.00 ± 0.00</b>	61.86 ± 1.42	62.03 ± 0.63	53.53 ± 0.27	<u>68.98 ± 0.73</u>	58.89 ± 0.98
	UN Vote	67.81 ± 0.75	<u>68.69 ± 0.24</u>	52.66 ± 1.18	67.58 ± 1.51	51.90 ± 1.40	66.31 ± 0.00	49.17 ± 4.04	51.28 ± 1.25	53.06 ± 0.23	<b>72.85 ± 1.36</b>	52.24 ± 0.95
	US Legis.	46.58 ± 1.19	79.16 ± 7.11	76.09 ± 4.00	61.85 ± 7.02	71.78 ± 14.33	64.83 ± 0.00	74.33 ± 2.79	<b>83.71 ± 0.69</b>	62.01 ± 10.94	79.13 ± 2.03	<u>81.67 ± 2.16</u>
	Avg. Rank	8.00	6.00	5.50	7.00	7.83	5.83	6.67	5.50	6.83	2.17	4.67

Table 21: AP of inductive dynamic link prediction on DTDGs. EdgeBank cannot do inductive link prediction so is not reported.

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	Can. Parl.	53.69 ± 1.27	54.28 ± 1.84	55.60 ± 0.57	53.47 ± 1.04	55.33 ± 0.98	55.06 ± 0.67	56.69 ± 0.04	<u>87.14 ± 1.40</u>	62.97 ± 0.34	<b>93.46 ± 2.62</b>
	Contact	95.20 ± 0.34	92.32 ± 0.20	96.18 ± 0.21	92.74 ± 2.86	89.98 ± 0.36	93.93 ± 0.14	90.51 ± 0.03	<u>98.04 ± 0.02</u>	76.19 ± 5.71	<b>98.16 ± 0.03</b>
	Flights	94.45 ± 0.89	92.85 ± 0.55	88.70 ± 0.02	95.03 ± 1.04	97.07 ± 0.02	83.54 ± 0.10	83.03 ± 0.10	<u>97.79 ± 0.04</u>	75.12 ± 5.52	<b>97.85 ± 0.22</b>
	UN Trade	58.94 ± 0.28	56.46 ± 0.60	61.20 ± 0.21	56.97 ± 1.53	<u>65.30 ± 0.14</u>	62.24 ± 0.27	62.55 ± 0.11	<u>64.62 ± 0.58</u>	53.83 ± 0.20	<b>70.55 ± 0.04</b>
	UN Vote	55.65 ± 1.32	55.53 ± 1.43	52.38 ± 0.92	<b>57.86 ± 1.72</b>	49.66 ± 0.77	53.38 ± 0.60	51.94 ± 1.32	56.40 ± 0.16	55.02 ± 3.04	56.61 ± 0.13
	US Legis.	54.36 ± 1.90	<u>57.54 ± 1.43</u>	50.21 ± 1.77	<b>59.39 ± 2.30</b>	52.79 ± 0.41	57.12 ± 0.22	54.68 ± 2.07	54.39 ± 1.63	52.55 ± 3.42	55.95 ± 1.16
	Avg. Rank	6.00	6.17	6.50	6.00	6.33	5.83	6.50	3.00	8.00	1.67
Inductive	Can. Parl.	51.91 ± 1.35	52.52 ± 1.85	57.42 ± 0.57	53.55 ± 1.36	57.56 ± 0.81	56.59 ± 1.22	56.91 ± 0.55	<u>84.72 ± 2.58</u>	63.27 ± 1.07	<b>92.68 ± 0.97</b>
	Contact	91.22 ± 0.18	89.18 ± 0.48	<b>94.93 ± 0.97</b>	88.74 ± 2.60	74.88 ± 0.81	91.41 ± 0.40	89.67 ± 0.66	93.78 ± 0.88	70.10 ± 5.03	<u>94.05 ± 0.32</u>
	Flights	60.79 ± 1.25	62.54 ± 1.73	65.23 ± 0.42	59.58 ± 1.56	56.76 ± 0.11	65.12 ± 0.09	<u>65.26 ± 0.44</u>	56.45 ± 0.17	<b>66.53 ± 2.41</b>	57.76 ± 2.06
	UN Trade	55.33 ± 0.88	54.52 ± 0.52	54.21 ± 0.24	51.84 ± 1.57	<b>56.29 ± 1.86</b>	55.82 ± 1.21	55.73 ± 0.20	51.92 ± 0.13	51.91 ± 0.57	52.81 ± 0.18
	UN Vote	59.24 ± 1.60	<u>61.97 ± 1.82</u>	52.52 ± 2.93	<b>66.41 ± 1.80</b>	47.35 ± 0.75	56.34 ± 2.25	48.15 ± 0.58	52.79 ± 0.84	55.49 ± 0.78	53.70 ± 0.24
	US Legis.	55.59 ± 2.38	<u>59.60 ± 1.64</u>	50.85 ± 2.59	<b>60.19 ± 2.60</b>	55.15 ± 1.55	58.56 ± 0.99	55.15 ± 0.27	54.76 ± 2.35	52.65 ± 1.99	57.85 ± 2.30
	Avg. Rank	5.50	5.00	5.50	5.83	6.50	4.00	5.33	6.33	6.17	4.67

Table 22: AUC-ROC of transductive dynamic link prediction on DTDGs.

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	EdgeBank	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	Can. Parl.	78.34 ± 0.24	76.69 ± 4.51	76.36 ± 0.95	75.66 ± 1.62	74.24 ± 3.64	64.14 ± 0.00	76.64 ± 0.39	83.22 ± 0.84	<u>98.00 ± 0.22</u>	86.77 ± 1.17	<b>99.69 ± 0.06</b>
	Contact	96.36 ± 0.48	96.46 ± 0.14	97.18 ± 0.09	97.76 ± 0.35	90.45 ± 0.24	94.34 ± 0.00	95.53 ± 0.06	93.92 ± 0.01	98.52 ± 0.00	<b>98.88 ± 0.06</b>	<u>98.68 ± 0.02</u>
	Flights	96.48 ± 1.17	96.10 ± 0.50	94.05 ± 0.08	98.26 ± 0.36	98.45 ± 0.01	90.23 ± 0.00	91.26 ± 0.03	91.14 ± 0.03	<u>98.93 ± 0.01</u>	94.57 ± 0.57	<b>98.98 ± 0.05</b>
	UN Trade	69.10 ± 0.45	66.92 ± 0.53	63.99 ± 0.18	67.80 ± 0.81	68.57 ± 0.14	66.75 ± 0.00	64.72 ± 0.01	65.97 ± 0.10	70.53 ± 0.34	<b>72.00 ± 0.24</b>	<u>71.41 ± 0.21</u>
	UN Vote	68.26 ± 1.19	68.72 ± 0.08	53.68 ± 0.37	<u>68.92 ± 0.72</u>	53.15 ± 0.03	62.97 ± 0.00	51.73 ± 0.25	52.60 ± 0.09	57.65 ± 0.48	<b>73.91 ± 1.08</b>	58.48 ± 0.12
	US Legis.	82.13 ± 0.30	80.98 ± 1.92	75.00 ± 2.19	<u>83.65 ± 0.27</u>	77.10 ± 0.14	62.57 ± 0.00	77.11 ± 0.59	79.19 ± 0.23	77.65 ± 0.95	<b>86.01 ± 1.68</b>	79.03 ± 0.26
Avg. Rank		4.83	5.33	8.00	5.67	8.83	9.83	9.00	7.83	3.00	2.67	2.00
Historical	Can. Parl.	62.08 ± 0.96	70.62 ± 3.57	72.11 ± 0.93	71.43 ± 0.99	70.44 ± 3.69	62.94 ± 0.00	74.07 ± 1.60	78.55 ± 0.72	96.46 ± 1.93	89.06 ± 1.14	<b>99.82 ± 0.10</b>
	Contact	96.23 ± 0.67	95.89 ± 0.14	96.01 ± 0.61	96.34 ± 0.24	83.63 ± 0.49	92.18 ± 0.00	94.70 ± 0.18	92.98 ± 0.39	97.20 ± 0.13	<b>98.37 ± 0.10</b>	<u>97.27 ± 0.06</u>
	Flights	68.84 ± 0.70	69.23 ± 1.44	72.33 ± 0.23	69.12 ± 0.81	65.93 ± 0.46	<u>74.64 ± 0.00</u>	70.81 ± 0.07	70.78 ± 0.32	67.50 ± 0.66	<b>78.87 ± 7.94</b>	68.98 ± 0.81
	UN Trade	68.89 ± 1.20	63.54 ± 1.02	59.33 ± 0.58	61.61 ± 2.98	64.74 ± 2.21	<b>86.44 ± 0.00</b>	62.11 ± 1.03	64.50 ± 0.17	72.17 ± 1.68	69.75 ± 0.90	71.41 ± 0.21
	UN Vote	77.54 ± 1.09	76.32 ± 1.04	55.61 ± 2.41	73.31 ± 1.65	50.95 ± 0.72	<b>89.53 ± 0.00</b>	54.70 ± 1.78	56.29 ± 1.46	64.79 ± 1.13	<u>78.62 ± 1.36</u>	65.17 ± 1.24
	US Legis.	54.26 ± 3.67	88.47 ± 4.92	83.49 ± 4.72	79.29 ± 4.52	80.85 ± 9.96	67.50 ± 0.00	86.20 ± 0.84	<b>90.38 ± 0.73</b>	75.60 ± 9.12	<u>89.88 ± 0.54</u>	88.36 ± 1.78
Avg. Rank		7.33	6.00	6.83	6.83	9.17	5.67	6.83	5.67	5.50	2.17	4.00
Inductive	Can. Parl.	52.09 ± 0.28	63.40 ± 3.74	72.42 ± 0.82	67.47 ± 1.82	71.66 ± 1.91	61.25 ± 0.00	72.29 ± 1.43	71.25 ± 0.46	<u>95.78 ± 0.85</u>	85.99 ± 0.95	<b>99.56 ± 0.21</b>
	Contact	94.24 ± 0.24	94.11 ± 0.17	94.76 ± 0.98	94.70 ± 0.39	88.26 ± 0.29	85.86 ± 0.00	92.28 ± 0.24	90.71 ± 0.43	95.12 ± 0.28	<b>97.93 ± 0.12</b>	95.68 ± 0.20
	Flights	70.09 ± 1.39	71.38 ± 1.84	73.56 ± 0.22	72.41 ± 0.71	69.61 ± 0.26	<b>81.10 ± 0.00</b>	72.77 ± 0.04	72.03 ± 0.35	69.33 ± 0.56	<u>79.82 ± 7.80</u>	71.16 ± 3.24
	UN Trade	66.91 ± 1.17	65.40 ± 0.56	64.50 ± 0.93	64.50 ± 3.18	<u>72.63 ± 2.17</u>	<b>74.25 ± 0.00</b>	68.28 ± 0.85	68.58 ± 0.13	59.98 ± 0.40	69.74 ± 1.54	67.60 ± 0.64
	UN Vote	73.82 ± 1.05	<u>74.67 ± 0.22</u>	52.88 ± 1.86	72.60 ± 1.62	52.05 ± 1.65	72.87 ± 0.00	50.25 ± 7.46	51.37 ± 1.05	55.36 ± 0.30	<b>76.64 ± 1.02</b>	54.09 ± 0.06
	US Legis.	53.13 ± 2.57	<u>86.74 ± 5.01</u>	83.16 ± 3.65	73.50 ± 5.46	79.80 ± 10.61	68.72 ± 0.00	80.53 ± 3.50	<b>88.82 ± 0.32</b>	73.31 ± 9.22	86.09 ± 1.11	86.06 ± 2.27
Avg. Rank		7.83	5.83	5.50	6.67	7.33	6.17	6.50	6.17	7.00	2.17	4.67

Table 23: AUC-ROC of inductive dynamic link prediction on DTDGs. EdgeBank cannot do inductive link prediction so is not reported.

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	Can. Parl.	53.62 ± 1.28	55.89 ± 1.77	56.67 ± 0.60	55.38 ± 2.11	57.91 ± 0.84	56.92 ± 0.97	59.29 ± 0.04	<u>89.06 ± 0.39</u>	59.13 ± 0.36	<b>94.02 ± 3.02</b>
	Contact	95.96 ± 0.20	92.00 ± 0.16	96.78 ± 0.17	94.04 ± 2.09	89.60 ± 0.32	94.82 ± 0.09	92.81 ± 0.02	<u>98.30 ± 0.00</u>	80.41 ± 4.30	<b>98.44 ± 0.05</b>
	Flights	94.82 ± 0.85	93.64 ± 0.44	88.61 ± 0.12	95.81 ± 0.87	96.87 ± 0.02	82.56 ± 0.13	82.28 ± 0.08	97.80 ± 0.05	79.41 ± 3.79	<b>97.98 ± 0.25</b>
	UN Trade	61.34 ± 0.38	58.41 ± 0.54	62.68 ± 0.21	58.71 ± 1.69	67.21 ± 0.15	63.83 ± 0.16	63.90 ± 0.02	<u>67.28 ± 0.54</u>	54.62 ± 0.93	<b>68.26 ± 0.26</b>
	UN Vote	56.24 ± 1.56	56.72 ± 1.90	52.44 ± 1.09	<b>59.42 ± 2.12</b>	47.92 ± 1.19	51.94 ± 0.69	51.65 ± 1.10	<u>57.47 ± 0.27</u>	53.40 ± 3.51	56.91 ± 0.12
	US Legis.	57.23 ± 1.92	<u>60.86 ± 1.59</u>	46.86 ± 3.44	<b>62.36 ± 2.04</b>	50.81 ± 1.02	58.02 ± 1.37	54.98 ± 2.39	52.96 ± 1.71	53.89 ± 3.55	57.17 ± 0.20
Avg. Rank		5.83	6.17	6.67	4.83	6.50	5.83	6.33	3.00	7.83	2.00
Inductive	Can. Parl.	50.56 ± 2.02	52.07 ± 1.53	58.91 ± 0.48	54.61 ± 2.40	60.17 ± 0.46	58.37 ± 1.15	58.59 ± 0.77	<u>86.13 ± 2.37</u>	59.68 ± 1.05	<b>92.37 ± 0.18</b>
	Contact	91.20 ± 0.17	88.90 ± 0.26	<b>94.44 ± 0.86</b>	89.78 ± 2.06	75.38 ± 0.50	91.78 ± 0.28	89.82 ± 0.51	94.31 ± 0.43	74.41 ± 3.94	94.35 ± 0.29
	Flights	60.58 ± 0.50	61.64 ± 1.60	63.61 ± 0.32	59.28 ± 0.95	56.51 ± 0.03	<u>63.64 ± 0.11</u>	63.07 ± 0.35	55.67 ± 0.48	<b>71.35 ± 1.86</b>	56.58 ± 2.12
	UN Trade	58.51 ± 1.05	56.96 ± 0.44	58.71 ± 0.43	54.45 ± 1.33	<b>62.65 ± 2.62</b>	<u>61.14 ± 1.09</u>	61.00 ± 0.35	55.80 ± 0.11	52.05 ± 0.59	57.58 ± 0.20
	UN Vote	62.88 ± 1.52	<u>65.80 ± 2.47</u>	51.79 ± 3.30	<b>71.29 ± 2.01</b>	46.61 ± 1.35	57.35 ± 2.24	45.75 ± 1.04	54.62 ± 0.50	54.31 ± 1.14	54.83 ± 2.17
	US Legis.	59.83 ± 2.06	<b>65.54 ± 1.42</b>	47.15 ± 4.26	<u>63.95 ± 2.30</u>	53.14 ± 2.15	59.73 ± 1.06	56.22 ± 2.36	53.57 ± 2.73	53.64 ± 2.12	57.91 ± 3.41
Avg. Rank		5.33	5.33	5.17	5.67	6.67	3.83	5.83	6.17	6.50	4.50

*Proven Performance Enhancement with Increased Historical Information.* In Appendix H.3, we have included an experiment showing different models’ performance with a varying number of temporal neighbors. From Figure 6, we observe that, for datasets with long-range temporal dependencies (such as Enron), an increase in sampled neighbors positively impacts model performance, particularly for DyGMamba and DyGFormer. We also find that for other baseline models, increasing the number of sampled neighbors does not consistently lead to performance improvements. This suggests that the key issue is not whether increasing sampled neighbors is always beneficial for dynamic graph modeling, but rather whether a model is robust enough to effectively process larger amounts of historical information. To further support this, we conducted additional experiments on Can. Parl<sup>11</sup> using DyGMamba with an

<sup>11</sup>Although Can. Parl is a DTDG, it requires huge number of sampled historical neighbors for optimal modeling. Therefore, it is considered here to prove the performance enhancement with increased historical information.

increasing number of sampled temporal neighbors (64, 128, 256, 512, 1024 and 2048). The results, presented in Table 25, confirm that DyGMamba consistently benefits from an increasing neighbor count until reaching the large number of 2048, further demonstrating its ability to leverage long-range temporal dependencies.

*Temporal Pattern Modeling is Game Changer.* As mentioned in introduction, we wish to capture edge-specific temporal patterns for better CTDG modeling. The motivation can be well supported with the following experiment. We provide a comparison between DyGMamba, DyGFormer, and a modified version of DyGFormer where Transformer is replaced with Mamba (referred to as Variant G). Our results in Table 26 show a performance drop for Variant G across all datasets compared to DyGFormer. Additionally, Variant G performs significantly worse than DyGMamba, with particularly notable declines on the synthetic datasets S1, S2, and S3. We attribute this to two reasons:

- Intrinsic Limitations of Mamba: Previous study [29] has shown that Mamba is generally less effective than Transformer in tasks requiring strong copying or in-context learning abilities. In CTDG modeling, strong copying ability is crucial, as many predictions are based on recalling repeated edges. Consequently, replacing Transformer with Mamba in DyGFormer inevitably leads to a performance drop, despite the gains in efficiency.
- The Importance of Temporal Pattern Modeling: Our dynamic information selection module plays a crucial role in enabling a Mamba-based model to outperform a Transformer-based model in CTDG modeling, especially when clear temporal patterns exist. This suggests that effectively capturing temporal patterns can significantly enhance the competitiveness of Mamba-based models in this context, further validating the novelty and the contribution of our work.

To summarize, the design our DyGMamba achieves a good balance between efficiency and effectiveness, thanks to temporal pattern modeling.

**Table 24: Number of available historical neighbors for all datasets on both transductive and inductive link prediction under random/historical/inductive NSS settings. Trans LP and Ind LP denote transductive and inductive link prediction, respectively.**

NSS	Random NSS				Historical NSS		Inductive NSS			
	Trans LP		Ind LP		Trans LP		Trans LP		Ind LP	
	Avg.	Max	Avg.	Max	Avg.	Max	Avg.	Max	Avg.	Max
Datasets										
LastFM	2,253.03	51,767	2,333.75	51,767	2,393.71	51,767	2,237.29	51,767	2,309.49	51,767
Enron	1,681.80	21,512	1,693.40	21,511	1,734.19	21,512	1,675.21	21,512	1,670.65	21,511
MOOC	2,304.04	19,473	2,312.62	19,473	3,265.62	19,473	2,293.64	19,473	2,275.23	19,473
Reddit	4,129.77	58,726	4,196.20	58,726	4,604.99	58,726	4,128.88	58,726	4,186.38	58,726
Wikipedia	144.39	1,937	149.67	1,937	164.49	1,937	144.17	1,937	149.22	1,937
UCI	172.82	1,546	193.87	1,546	261.31	1,546	173.08	1,546	192.72	1,546
Social Evo.	63,962.35	124,565	68,701.32	124,565	60,403.62	124,565	60,395.83	124,565	65,191.24	124,565

**Table 25: DyGMamba’s performance with increasing sampled neighbors on transductive link prediction under random NSS.**

Can. Parl	64	128	256	512	1024	2048
DyGMamba	72.58	74.46	76.85	82.67	94.88	99.57

**Table 26: Comparison among Variant G, DyGFormer and DyGMamba on transductive link prediction under random NSS.**

Models	LastFM	Enron	MOOC	Reddit	Wikipedia	UCI	Social Evo.	S1	S2	S3
Variant G	92.69	92.24	87.13	98.89	98.70	95.32	94.29	54.40	55.13	72.53
DyGFormer	92.95	92.42	87.66	99.22	99.03	95.74	94.63	55.19	57.80	79.20
DyGMmaba	93.35	92.65	89.21	99.32	99.15	95.91	94.77	81.85	85.36	86.59