DYGMAMBA: EFFICIENTLY MODELING LONG-TERM TEMPORAL DEPENDENCY ON CONTINUOUS-TIME DY NAMIC GRAPHS WITH STATE SPACE MODELS

Anonymous authors

Paper under double-blind review

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.

028 029

031 032

006

008 009 010

011 012 013

014

015

016

017

018

019

021

025

026

027

1 INTRODUCTION

033 Dynamic graphs store node interactions in the form of links labeled with timestamps (Kazemi et al., 2020). 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., 035 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 037 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 work (Kazemi et al., 2020) has indicated that CTDGs have an advantage over DTDGs in preserving temporal 040 details, and therefore, more attention is paid to developing novel CTDG modeling approaches for 041 dynamic graph representation learning. 042

Recent effort in CTDG modeling has resulted in a wide range of models. However, most of them are 043 unable to model long-term temporal dependencies of nodes. To solve this problem, Yu et al. (2023) 044 propose a CTDG model DyGFormer that can handle long-term node interaction histories based on Transformer (Vaswani et al., 2017). Despite its ability in modeling longer histories, employing a 046 Transformer naturally introduces excessive usage of computational resources due to its quadratic 047 complexity. Another recent work CTAN (Gravina et al., 2024) tries to capture long-term temporal 048 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. (2024) show that CTAN cannot capture very long histories and is surpassed by DyGFormer on the CTDGs where learning 051 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 052 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., 057 Mamba SSM (Gu & Dao, 2023) to encode the long sequence of historical node interactions. Since Mamba is proven effective and efficient in long sequence modeling (Gu & Dao, 2023), it maintains low computational complexity and is scalable in modeling long-term temporal dependencies. For the 060 second challenge, we address it by learning temporal patterns of node interactions and dynamically 061 selecting the critical temporal information based on them. The motivation can be explained by the 062 following example. Consider a CTDG with nodes as people or songs and edges representing a 063 person playing a song at a specific time. If a person u frequently plays a hit song v initially but 064 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 065 appearances in each other's historical interactions. If a CTDG model recognizes this pattern, it can 066 prioritize other temporal information, such as u increasingly listening to a new song v' before t, 067 instead of focusing on u, v interactions. Since each pattern corresponds to a specific edge, e.g., 068 (u, v, t), we name these patterns as edge-specific temporal patterns. 069

To this end, we propose a new CTDG model named DyGMamba. DyGMamba first leverages a nodelevel 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 efectiveness 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 commonly-used CTDG datasets.

077 078 079

2 RELATED WORK AND PRELIMINARIES

081 2.1 RELATED WORK

Dynamic Graph Representation Learning. Dynamic graph representation learning methods can 083 be categorized into two groups, i.e., DTDG and CTDG methods. DTDG methods (Pareja et al., 084 2020; Goyal et al., 2020; Sankar et al., 2020; You et al., 2022; Li et al., 2024a) can only model 085 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 087 (Kazemi et al., 2020). To overcome this problem, recent works focus more on developing CTDG 880 methods that treat a dynamic graph as a stream of events, where each event has its own unique 089 timestamp. Some works (Trivedi et al., 2019; Chang et al., 2020) model CTDGs by using temporal point process. Another line of works (Xu et al., 2020; Ma et al., 2020; Wang et al., 2021b; Gravina et al., 2024) designs advanced temporal graph neural networks for CTDGs. Besides, some other 091 methods are developed based on memory networks (Rossi et al., 2020; Liu et al., 2022), temporal 092 random walk (Wang et al., 2021c; Jin et al., 2022) and temporal sequence modeling (Cong et al., 093 2023; Yu et al., 2023; Tian et al., 2024). Since some real-world CTDGs heavily rely on long-term 094 temporal information for effective learning, a number of works start to develop CTDG models that 095 can do long range propagation of information over time (Yu et al., 2023; Gravina et al., 2024). 096

007

State Space Models. Transformer (Vaswani et al., 2017) is a de facto backbone architecture in 098 modern deep learning. However, its self-attention mechanism results in large space and time complexity, making it unsuitable for extremely long sequence modeling (Duman Keles et al., 2023). 100 To address this, many works focus on building structured state space models that scale linearly or 101 near-linearly with input sequence length (Gu et al., 2021; 2022b;a; Smith et al., 2023; Peng et al., 102 2023; Ma et al., 2023; Gu & Dao, 2023). Most structured SSMs exhibit linear time invariance (LTI), 103 meaning their parameters are not input-dependent and fixed for all time-steps. Gu & Dao (2023) 104 demonstrate that LTI prevents SSMs from effectively selecting relevant information from the input 105 context, which is problematic for tasks requiring context-aware reasoning. To solve this issue, Gu & Dao (2023) proposes S6, also known as Mamba, which uses a selection mechanism to dynam-106 ically choose important information from input sequence elements. Selection mechanism involves 107 learning functions that map input data to SSM's parameters, making Mamba both efficient and effective in modeling language, DNA sequences, and audio. Recently, there have been several works
 employing Mamba SSM for representation learning on static graphs (Wang et al., 2024; Behrouz & Hashemi, 2024) and spatial-temporal graphs (Li et al., 2024b). Unlike our work, they do not focus
 on capturing the dynamic evolution of graph structure and are thus not suitable for CTDG modeling.

113 114 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 $0 \le t_1 \le t_2 \le ... \le t_{|\mathcal{G}|}$, 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 (Gu et al., 2022b; Gu & Dao, 2023) are inspired by a continuous system which can be described as $\mathbf{z}(\tau)' = \mathbf{A}\mathbf{z}(\tau) + \mathbf{B}q(\tau)$ and $r(\tau) = \mathbf{C}\mathbf{z}(\tau)$. $q(\tau) \in \mathbb{R}$ and $r(\tau) \in \mathbb{R}$ are the 1-dimensional input and output over time τ^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

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

Here, τ 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

136 137 138

126

127

128

129 130

131

132 133

134

135

 $\mathbf{q} = \mathbf{p} * \bar{\mathbf{K}}_{\text{SSM}}, \text{ where } \bar{\mathbf{K}}_{\text{SSM}} = [\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, ..., \mathbf{C}\bar{\mathbf{A}}^{|\mathbf{p}|-1}\bar{\mathbf{B}}] \in \mathbb{R}^{|\mathbf{p}|-1}.$ (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} . 139 denotes the element-wise multiplication. When the dimension size of each element p_{τ} in p becomes 140 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-141 Output (SISO) fashion, processing each input dimension in parallel with the same set of parameters. 142 We follow Gu & Dao (2023) and denote the computation in Eq. 2 on the input sequences with vector 143 elements as a function $SSM_{\overline{A},\overline{B},C}(\cdot)^2$. Different from S4 which uses same parameters to process 144 each element, Mamba changes its parameters into input-dependent by employing several trainable 145 linear layers to map input into **B**, **C** and Δ . The system is evolving as it processes different elements 146 in the input sequence, making Mamba time-variant and suitable for modeling temporal sequences. 147

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 $\mathcal{S}_u^t = \{(u, u', t') | t' < t, (u, u', t') \in \mathcal{G}\} \cup \{(u', u, t') | t' < t, (u', u, t') \in \mathcal{G}\}$ and $\mathcal{S}_v^t = \{(v, v', t') | t' < t, (v, v', t') \in \mathcal{G}\} \cup \{(v', v, t') | t' < t, (v', v, t') \in \mathcal{G}\}$ containing u's and v's one-hop temporal neighbors $Nei_u^t = \{(u', t') | (u, u', t') \text{ or } (u', u, t') \in \mathcal{G}, t' < t\}$ and $Nei_v^t = \{(v', t') | (v, v', t') \text{ or } (v', v, t') \in \mathcal{G}\}$ (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

161

148

¹⁵⁹ 160

¹We use τ rather than t to indicate time in a continuous system to distinguish from the time in CTDGs.

²Input and output of SSM_{Å,B,C}(·) are matrices where each row is a vector corresponding to an element. See App. H for more details of SISO and the function.</sub>



Figure 1: Model overview of DyGMamba.

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.

178 3.1 LEARNING ONE-HOP TEMPORAL NEIGHBORS

170 171 172

173

174

175

176 177

179 **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 $Nei_u^t + 1$ 181 temporal nodes. We take their node features from the dataset and stack them into a feature ma-182 trix $\mathbf{\tilde{X}}_{u}^{t} \in \mathbb{R}^{(|Nei_{u}^{t}|+1) \times d_{N}}$. Similarly, we build a link feature matrix $\mathbf{\tilde{E}}_{u}^{t} \in \mathbb{R}^{(|Nei_{u}^{t}|+1) \times d_{E}}$. To 183 incorporate temporal information, we encode the time difference between u and each one-hop tem-184 poral neighbor (u', t') using the time encoding function introduced in TGAT (Xu et al., 2020): 185 $\sqrt{1/d_T}[\cos(\omega_1(t-t')+\phi_1),\ldots,\cos(\omega_d(t-t')+\phi_{d_T})]$. 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 186 are denoted as $\tilde{\mathbf{T}}_{u}^{t} \in \mathbb{R}^{(|Nel_{u}^{t}|+1) \times d_{T}}$. We follow the same way to get $\tilde{\mathbf{X}}_{v}^{t} \in \mathbb{R}^{(|Nel_{v}^{t}|+1) \times d_{N}}, \tilde{\mathbf{E}}_{v}^{t} \in \mathbb{R}^{(Nel_{v}^{t}|+1) \times d_{N}}$ 187 188 $\mathbb{R}^{(|Nei_v^t|+1) \times d_E}$ and $\tilde{\mathbf{T}}_v^t \in \mathbb{R}^{(|Nei_v^t|+1) \times d_T}$ for v's temporal neighbors. Following Tian et al. (2024), 189 we also consider the historical node interaction frequencies in the interaction sequences S_u^t and S_v^t of 190 source u and destination v. For example, assume the interacted nodes of u and v (arranged in chrono-191 logical order) are $\{a, v, a\}$ and $\{b, b, u, a\}$, the appearing frequencies of a, b in u/v's historical inter-192 actions 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 $\tilde{F}_u^t = [[2,1],[1,1],[2,1],[0,1]]^\top$ 193 and $\tilde{F}_v^t = [[0,2], [0,2], [1,1], [2,1], [0,1]]^\top$, respectively. Note that the last elements ([0,1] and 194 [0,1]) in \vec{F}_{u}^{t} and \vec{F}_{v}^{t} correspond to the appended (u,t) and (v,t) not existing in the observed histo-196 ries. We initialize them with [0, number of historical interactions between u, v]. An encoding multilayer perceptron (MLP) $f(\cdot) : \mathbb{R} \to \mathbb{R}^{d_F}$ is employed to encode these features into representations: $\tilde{\mathbf{F}}_u^t = f(\tilde{F}_u^t[:,0]) + f(\tilde{F}_u^t[:,1]) \in \mathbb{R}^{(|Nei_u^t|+1) \times d_F}, \tilde{\mathbf{F}}_v^t = f(\tilde{F}_v^t[:,0]) + f(\tilde{F}_v^t[:,1]) \in \mathbb{R}^{(|Nei_v^t|+1) \times d_F}.$ 197

200 **Patching Neighbors.** We employ the patching technique proposed by (Yu et al., 2023) to save 201 computational resources when dealing with a large number of temporal neighbors. We treat ptemporally adjacent neighbors as a patch and flatten their features. For example, with patching, 202 $\tilde{\mathbf{X}}_{u}^{t} \in \mathbb{R}^{(|Net_{u}^{t}|+1) \times d_{N}}$ results in a new patched feature matrix $\mathbf{X}_{u}^{t} \in \mathbb{R}^{\lceil (|Net_{u}^{t}|+1)/p \rceil \times (p \cdot d_{N})}$ (we 203 pad $\tilde{\mathbf{X}}_{u}^{t}$ with zero-valued features when $|Nei_{u}^{t}| + 1$ cannot be divided by p). Similarly, we get $\mathbf{E}_{\theta}^{t} \in \mathbb{R}^{\lceil (|Nei_{\theta}^{t}|+1)/p \rceil \times (p \cdot d_{E})}$, $\mathbf{T}_{\theta}^{t} \in \mathbb{R}^{\lceil (|Nei_{\theta}^{t}|+1)/p \rceil \times (p \cdot d_{F})}$ and $\mathbf{F}_{\theta}^{t} \in \mathbb{R}^{\lceil (|Nei_{\theta}^{t}|+1)/p \rceil \times (p \cdot d_{F})}$ (θ is ei-204 205 206 ther u or v). Each row of a feature matrix corresponds to an element of the input sequence sent into 207 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. 208

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., $\mathbf{X}_{\theta}^{t} := f_{N}(\mathbf{X}_{\theta}^{t}), \mathbf{E}_{\theta}^{t} := f_{E}(\mathbf{E}_{\theta}^{t}), \mathbf{T}_{\theta}^{t} := f_{T}(\mathbf{T}_{\theta}^{t}), \mathbf{F}_{\theta}^{t} := f_{F}(\mathbf{F}_{\theta}^{t})$. $f_{N}(\cdot) : \mathbb{R}^{p \cdot d_{N}} \to \mathbb{R}^{d}, f_{E}(\cdot) : \mathbb{R}^{p \cdot d_{E}} \to \mathbb{R}^{d}, f_{T}(\cdot) : \mathbb{R}^{p \cdot d_{T}} \to \mathbb{R}^{d}, f_{F}(\cdot) : \mathbb{R}^{p \cdot d_{F}} \to \mathbb{R}^{d}$ 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., $\mathbf{H}_{\theta}^{t} = \mathbf{X}_{\theta}^{t} \|\mathbf{E}_{\theta}^{t}\|\mathbf{T}_{\theta}^{t}\|\mathbf{F}_{\theta}^{t} \in \mathbb{R}^{\lceil (|Ne_{\theta}^{t}|+1)/p \rceil \times 4d}$. We input \mathbf{H}_{u}^{t} and \mathbf{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 \mathbf{H}_{θ}^{t} into a Mamba SSM

 $\mathbf{B}_{1} = \mathbf{H}_{\theta}^{t} \mathbf{W}_{\mathbf{B}_{1}} \in \mathbb{R}^{\left[(|Nei_{\theta}^{t}|+1)/p \right] \times d_{\text{SSM}}}, \quad \mathbf{C}_{1} = \mathbf{H}_{\theta}^{t} \mathbf{W}_{\mathbf{C}_{1}} \in \mathbb{R}^{\left[(|Nei_{\theta}^{t}|+1)/p \right] \times d_{\text{SSM}}}; \tag{3a}$ $\Delta_{\star} = \text{Softplus}(\text{Broadcast}_{\star}, (\mathbf{H}^{t} \mathbf{W}_{\star})) + \text{Param}_{\star}) \in \mathbb{R}^{\left[(|Nei_{\theta}^{t}|+1)/p \right] \times 4d}. \tag{3b}$

219 220

$$\Delta_{1} = \text{Softplus}(\text{Broadcast}_{4d}(\mathbf{H}_{\theta}^{c}\mathbf{W}_{\Delta_{1}}) + \text{Param}_{\Delta_{1}}) \in \mathbb{R}^{|\langle|\text{Pere}_{\theta}|+1\rangle/p|/\sqrt{4u}};$$
(3b)
$$\bar{\mathbf{A}}_{1} = \exp(\Delta_{1}\mathbf{A}_{1}), \ \bar{\mathbf{B}}_{1} = (\Delta_{1}\mathbf{A}_{1})^{-1}\exp((\Delta_{1}\mathbf{A}_{1}) - \mathbf{I})\Delta_{1}\mathbf{B}_{1};$$
(3c)

230 231

232

233

234

235 236

237

218

$$\mathbf{H}_{\theta}^{t} := \mathbf{H}_{\theta}^{t} + \mathrm{SSM}_{\bar{\mathbf{A}}_{1},\bar{\mathbf{B}}_{1},\mathbf{C}_{1}}(\mathbf{H}_{\theta}^{t}).$$
(3d)

224 $\mathbf{W}_{\mathbf{B}_{1}}, \mathbf{W}_{\mathbf{C}_{1}} \in \mathbb{R}^{4d \times d_{SSM}}$ and $\mathbf{W}_{\Delta_{1}} \in \mathbb{R}^{4d \times 1}$. $\bar{\mathbf{A}}_{1}, \bar{\mathbf{B}}_{1} \in \mathbb{R}^{\lceil (|Nei_{\theta}^{t}|+1)/p \rceil \times 4d \times d_{SSM}}$ are discretized pa-225 rameters. Param_{Δ_{1}} $\in \mathbb{R}^{\lceil (|Nei_{\theta}^{t}|+1)/p \rceil \times 4d}$ is a parameter defined by Gu & Dao (2023). Broadcast_{4d}(·) 226 is a function that copies its vector input for 4d times to form a matrix with 4d identical columns 227 (following the definition in Gu & Dao (2023)). I is an identity matrix. Then we use an MLP 228 $f_{node}(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d}$ on SSM's output

$$\mathbf{H}_{\theta}^{t} := \mathbf{H}_{\theta}^{t} + f_{\text{node}} \left(\text{LayerNorm}(\mathbf{H}_{\theta}^{t}) \right).$$
(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 EDGE-SPECIFIC TEMPORAL PATTERNS

238 Time-Level SSM Block. To capture edge-specific temporal patterns, we use another time-level 239 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., $\{(u, v, t_0), ..., (u, v, t_{k-1})|t_0 < 0\}$ 240 $\dots < t_{k-1} < t_k$. Then we construct a timestamp sequence $\{t_0, t_1, \dots, t_{k-1}, t_k\}$ based on these inter-241 actions and the prediction timestamp t. We compute the time difference between each neighboring 242 pair of them and further get a time difference sequence $\{t_1 - t_0, t_2 - t_1, ..., t - t_{k-1}\}$, representing 243 the change of time interval between two identical interactions. Each element in this sequence is in-244 put into the time encoding function stated above to get a edge-specific (specific to the edge (u, v, t)) 245 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 246 $f_{\text{map1}}(\cdot) : \mathbb{R}^{d_T} \to \mathbb{R}^{\gamma d} \ (\gamma \in [0, 1] \text{ is a hyperparameter}), \text{ i.e., } \mathbf{H}_{u,v}^t := f_{\text{map1}}(\mathbf{H}_{u,v}^t).$ A time-level 247 SSM layer takes $\mathbf{H}_{u,v}^{t}$ as input and computes 248

$$\mathbf{B}_{2} = \mathbf{H}_{u,v}^{t} \mathbf{W}_{\mathbf{B}_{2}} \in \mathbb{R}^{k \times d_{\text{SSM}}}, \ \mathbf{C}_{2} = \mathbf{H}_{u,v}^{t} \mathbf{W}_{\mathbf{C}_{2}} \in \mathbb{R}^{k \times d_{\text{SSM}}};$$
(5a)

249

253

254

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

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

(5d)

$$\mathbf{H}_{u,v}^{t} := \mathbf{H}_{u,v}^{t} + \mathrm{SSM}_{\bar{\mathbf{A}}_{2},\bar{\mathbf{B}}_{2},\mathbf{C}_{2}}(\mathbf{H}_{u,v}^{t}).$$

 $\mathbf{W}_{\mathbf{B}_2}, \mathbf{W}_{\mathbf{C}_2} \in \mathbb{R}^{\gamma d \times d_{\mathrm{SSM}}}$ and $\mathbf{W}_{\Delta_2} \in \mathbb{R}^{\gamma d \times 1}$. $\bar{\mathbf{A}}_2, \bar{\mathbf{B}}_2 \in \mathbb{R}^{k \times \gamma d \times d_{\mathrm{SSM}}}$ are discretized parameters. Param $_{\Delta_2} \in \mathbb{R}^{k \times \gamma d}$ is a parameter defined as same as Param $_{\Delta_1}$. In practice, we set k to a number 255 256 257 much smaller than $|Net_{\theta}^{t}|$, e.g., 10. This ensures that time-level SSM will not incur huge computa-258 tional burden and the model focuses more on the recent histories. Note that we cannot always find k259 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 260 difference without a found historical interaction to a very large number 10^{10} . For example, if k = 2, 261 and for (u, v, t) we can only find (u, v, t_0) . The time difference sequence will be $\{10^{10}, t - t_0\}$. 262 10^{10} is much larger than $t - t_0$, indicating that u and v have not had an interaction for an extremely 263 long time, same as existing no historical interaction. We further explain why we use SSM to learn 264 temporal patterns in App. I. 265

266

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 $\mathbf{h}_{u,v}^t \in \mathbb{R}^{\gamma d}$ 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}_{θ}^{t}

$$\hat{\mathbf{h}}_{u,v}^t = f_{ ext{map2}}(\mathbf{h}_{u,v}^t) \in \mathbb{R}^{4d};$$

272 273

274

275

276

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

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

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

(6a)

277 $f_{\text{map2}}(\cdot) : \mathbb{R}^{\gamma d} \to \mathbb{R}^{4d} \text{ and } f_{\text{map3}}(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^1 \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ is } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ and } f_{\text{map3}}(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d} \text{ are two mapping MLPs. } f'(\cdot) : \mathbb{R}^{4d} \to \mathbb{$ 278 another MLP introducing training parameters. Note that α_u/α_v is computed by considering both 279 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 θ , making it hard to connect u and 281 v. Computing α_{θ} as Eq. 6c helps to strengthen the connection between both nodes and meanwhile 282 incorporates the learned temporal pattern. β_{θ} is derived by transforming the queried results based on 283 α_{θ} into weights. It is then used to compute a weighted-sum of all temporal neighbors for represent-284 ing θ at t, i.e., \mathbf{h}_{θ}^{t} . The neighbors assigned with greater weights from β_{θ} are selected as more critical 285 and will contribute more to \mathbf{h}_{θ}^{t} . Finally, we output the representations of u, v and the edge-specific temporal pattern by employing two output MLPs $f_{out1}(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{d_N}$ and $f_{out2}(\cdot) : \mathbb{R}^{\gamma d} \to \mathbb{R}^{d_N}$, i.e., $\mathbf{h}_{\theta}^t := f_{out1}(\mathbf{h}_{\theta}^t) \in \mathbb{R}^{d_N}$, $\mathbf{h}_{u,v}^t := f_{out2}(\mathbf{h}_{u,v}^t) \in \mathbb{R}^{d_N}$. 286 287 288

3.3 LEVERAGING LEARNED REPRESENTATIONS FOR LINK PREDICTION

291 We leverage \mathbf{h}_{θ}^{t} and $\mathbf{h}_{u,v}^{t}$ for dynamic link prediction. We employ a prediction MLP, i.e., $f_{LP}(\cdot)$: 292 $\mathbb{R}^{3d_{N}} \to \mathbb{R}$, as the predictor. The probability of existing a link (u, v, t) is computed as y'(u, v, t) =293 Sigmoid $(f_{LP}(\mathbf{h}_{u}^{t} \| \mathbf{h}_{v}^{t} \| \mathbf{h}_{u,v}^{t}))$. For model parameter learning, we use the following loss function

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

y(u, v, t) is the ground truth label denoting the existence of (u, v, t) (1/0 means existing/nonexisting). M is the total number of edges existing in the training data (positive edges). We follow previous work (Yu et al., 2023) 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} .

300 301 302

303

308

309

289

290

295 296

297

298

299

4 EXPERIMENTS

In Sec. 4.2, we validate DyGMamba's ability in CTDG representation learning by comparing it with baseline methods on dynamic link prediction³. We do further analysis to show the effectiveness of model components. In Sec. 4.3, we show DyGMamba's efficiency. We also show that it achieves much stronger scalability in modeling long-term temporal information compared with DyGFormer.

4.1 EXPERIMENTAL SETTING

310 CTDG Datasets and Baselines. We consider seven real-world CTDG datasets collected by (Pour-311 safaei et al., 2022), i.e., LastFM, Enron, MOOC, Reddit, Wikipedia, UCI and Social Evo.. Dataset 312 statistics are presented in App. A.1. Among them, we take LastFM, Enron and MOOC as long-range 313 temporal dependent datasets because according to Yu et al. (2023), much longer histories are needed 314 for optimal representation learning on them. We compare DyGMamba with ten recent CTDG base-315 line models, i.e., JODIE (Kumar et al., 2019), DyRep (Trivedi et al., 2019), TGAT (Xu et al., 2020), TGN (Rossi et al., 2020), CAWN (Wang et al., 2021c), EdgeBank (Poursafaei et al., 2022), TCL 316 (Wang et al., 2021a), GraphMixer (Cong et al., 2023), DyGFormer (Yu et al., 2023) and CTAN 317 (Gravina et al., 2024). Among them, only DyGFormer and CTAN are designed for long-range tem-318 poral information propagation. Detailed descriptions of baseline methods are presented in App. B. 319 We also implemented FreeDyG (Tian et al., 2024) by using its official code repository, however, on 320 LastFM, we find that FreeDyG's loss cannot converge and the reported results are not reproducible. 321 So we do not report its performance in our paper.

³²² 323

³To 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

Implementation Details and Evaluation Settings. We use the implementations and the best hy-perparameters provided by Yu et al. (2023) 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 ρ 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 ρ/p input into the SSM as same as the length input into Transformer in Yu et al. (2023), i.e., $\rho/p = 32$. All experiments are implemented with PyTorch (Paszke et al., 2019) 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 hy-perparamter configurations are presented in App. C. We employ two evaluation settings following previous works: the transductive and inductive settings. As suggested in (Poursafaei et al., 2022), 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 de-tailed explanations. We employ two metrics, i.e., average precision (AP) and area under the receiver operating characteristic curve (AUC-ROC)

Table 1: AP of transductive dynamic link prediction. The best and the second best results are marked
 as **bold** and <u>underlined</u>, 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 Enron MOOC Reddit Wikipedia UCI Social Evo.	$\begin{array}{c} 70.95 \pm 2.94 \\ 84.85 \pm 3.13 \\ 81.04 \pm 0.83 \\ 98.31 \pm 0.06 \\ 96.51 \pm 0.22 \\ 89.28 \pm 1.02 \\ 89.88 \pm 0.40 \end{array}$	$\begin{array}{c} 71.85 \pm 2.44 \\ 79.80 \pm 2.28 \\ 81.50 \pm 0.77 \\ 98.18 \pm 0.03 \\ 94.88 \pm 0.29 \\ 66.11 \pm 2.75 \\ 88.39 \pm 0.69 \end{array}$	$\begin{array}{c} 73.30 \pm 0.18 \\ 70.76 \pm 1.05 \\ 85.71 \pm 0.20 \\ 98.57 \pm 0.01 \\ 96.88 \pm 0.06 \\ 79.40 \pm 0.61 \\ 93.33 \pm 0.06 \end{array}$	$\begin{array}{c} 75.31\pm 5.62\\ 86.98\pm 1.05\\ \underline{89.15\pm 1.69}\\ 98.65\pm 0.04\\ 98.45\pm 0.10\\ 92.33\pm 0.64\\ 93.45\pm 0.29\end{array}$	$\begin{array}{c} 86.60\pm0.11\\ 89.50\pm0.10\\ 80.30\pm0.43\\ 99.11\pm0.01\\ 98.77\pm0.01\\ 95.13\pm0.23\\ 84.90\pm0.11 \end{array}$	$\begin{array}{c} 79.29 \pm 0.00 \\ 83.53 \pm 0.00 \\ 57.97 \pm 0.00 \\ 94.86 \pm 0.00 \\ 90.37 \pm 0.00 \\ 76.20 \pm 0.00 \\ 74.95 \pm 0.00 \end{array}$	$\begin{array}{c} 76.62\pm1.83\\ 85.41\pm0.71\\ 83.89\pm0.86\\ 97.78\pm0.02\\ 97.75\pm0.04\\ 86.63\pm1.30\\ 93.82\pm0.19 \end{array}$	$\begin{array}{c} 75.56\pm0.19\\ 82.13\pm0.30\\ 82.80\pm0.15\\ 97.31\pm0.01\\ 97.22\pm0.02\\ 93.15\pm0.41\\ 93.36\pm0.06\end{array}$	$\begin{array}{c} \underline{92.95\pm0.14}\\ \underline{92.42\pm0.11}\\ 87.66\pm0.48\\ \underline{99.22\pm0.01}\\ \underline{99.03\pm0.03}\\ \underline{95.74\pm0.17}\\ \underline{94.63\pm0.07}\end{array}$	$\begin{array}{c} 86.44 \pm 0.80 \\ \underline{92.52 \pm 1.20} \\ 84.71 \pm 2.85 \\ 97.21 \pm 0.84 \\ 96.61 \pm 0.79 \\ 76.64 \pm 4.11 \\ Timeout \end{array}$	$\begin{array}{c} 93.35\pm 0.20\\ 92.65\pm 0.12\\ 89.21\pm 0.08\\ 99.32\pm 0.01\\ 99.15\pm 0.02\\ 95.91\pm 0.15\\ 94.77\pm 0.01\\ \end{array}$
_		Avg. Rank	8.29	9.29	7.00	4.29	6.00	9.43	5.57	6.43	<u>2.43</u>	6.29	1.00
-	Historical	LastFM Enron MOOC Reddit Wikipedia UCI Social Evo.	$74.38 \pm 6.27 \\ 69.13 \pm 1.66 \\ 78.62 \pm 2.43 \\ 79.96 \pm 0.30 \\ 81.16 \pm 0.73 \\ 74.77 \pm 5.35 \\ 91.26 \pm 2.47 \\ \end{cases}$	$\begin{array}{c} 71.85 \pm 2.91 \\ 72.58 \pm 1.83 \\ 75.14 \pm 2.86 \\ 79.40 \pm 0.30 \\ 79.46 \pm 0.95 \\ 55.89 \pm 2.83 \\ 92.86 \pm 0.90 \end{array}$	$\begin{array}{c} 71.60\pm0.36\\ 64.24\pm1.24\\ 82.83\pm0.71\\ 79.78\pm0.25\\ 87.31\pm0.36\\ 66.78\pm0.77\\ 95.31\pm0.30\\ \end{array}$	$\begin{array}{c} 75.03\pm 6.90\\ 74.31\pm 0.99\\ 85.46\pm 2.32\\ 81.05\pm 0.32\\ 87.31\pm 0.25\\ \underline{81.32\pm 1.26}\\ 93.84\pm 1.68 \end{array}$	$\begin{array}{c} 69.93 \pm 0.33 \\ 65.40 \pm 0.36 \\ 74.46 \pm 0.53 \\ 80.96 \pm 0.28 \\ 66.77 \pm 6.62 \\ 64.69 \pm 1.78 \\ 85.65 \pm 0.11 \end{array}$	$\begin{array}{c} 73.03 \pm 0.00 \\ 76.53 \pm 0.00 \\ 60.71 \pm 0.00 \\ 73.59 \pm 0.00 \\ 73.35 \pm 0.00 \\ 65.50 \pm 0.00 \\ 80.57 \pm 0.00 \end{array}$	$\begin{array}{c} 71.02\pm2.07\\ 72.39\pm0.61\\ 78.51\pm1.24\\ 77.38\pm0.20\\ 86.12\pm1.69\\ 74.62\pm2.70\\ 95.93\pm0.63\\ \end{array}$	$\begin{array}{c} 72.28 \pm 0.37 \\ 77.35 \pm 1.22 \\ 77.09 \pm 0.83 \\ 78.39 \pm 0.40 \\ \underline{90.74 \pm 0.06} \\ \textbf{83.88 \pm 1.06} \\ 95.30 \pm 0.34 \\ \end{array}$	$\begin{array}{c} 81.51\pm0.14\\ 76.93\pm0.76\\ 85.65\pm0.89\\ 81.63\pm1.08\\ 70.13\pm11.02\\ 80.44\pm1.16\\ 97.05\pm0.16\\ \end{array}$	$\frac{82.29 \pm 0.94}{77.24 \pm 1.53}$ $\overline{67.73 \pm 2.08}$ 89.77 ± 2.28 95.91 ± 0.10 76.62 ± 0.33 Timeout 4.71	$\begin{array}{c} \textbf{83.02} \pm \textbf{0.16} \\ \textbf{77.77} \pm \textbf{1.32} \\ \textbf{85.89} \pm \textbf{0.94} \\ \underline{\textbf{81.80} \pm \textbf{1.52}} \\ \textbf{81.77} \pm \textbf{1.20} \\ \textbf{81.03} \pm \textbf{1.09} \\ \textbf{97.35} \pm \textbf{0.52} \end{array}$
=	Inductive	Avg. Rank LastFM Enron MOOC Reddit Wikipedia UCI Social Evo.	$\begin{array}{c} 6.57\\ 62.63\pm 6.89\\ 69.51\pm 1.06\\ 66.56\pm 1.49\\ 86.93\pm 0.21\\ 74.78\pm 0.56\\ 66.02\pm 1.28\\ 91.08\pm 3.29\end{array}$	$\begin{array}{c} 8.14\\ 62.49\pm 3.04\\ 66.78\pm 2.21\\ 61.48\pm 0.96\\ 86.06\pm 0.36\\ 70.55\pm 1.22\\ 54.64\pm 2.52\\ 92.84\pm 0.98\end{array}$	$\begin{array}{c} \textbf{6.57} \\ \textbf{71.16} \pm \textbf{0.33} \\ \textbf{63.16} \pm \textbf{0.59} \\ \textbf{76.96} \pm \textbf{0.89} \\ \textbf{89.93} \pm \textbf{0.10} \\ \textbf{86.77} \pm \textbf{0.29} \\ \textbf{67.63} \pm \textbf{0.51} \\ \textbf{95.20} \pm \textbf{0.30} \end{array}$	$\begin{array}{r} 4.14\\ 65.09\pm7.05\\ 73.27\pm0.58\\ 77.59\pm1.83\\ 88.12\pm0.13\\ 85.80\pm0.15\\ 70.34\pm0.72\\ 94.58\pm1.52\end{array}$	$\begin{array}{c} 9.29\\ 67.38\pm 0.57\\ 75.08\pm 0.81\\ 73.55\pm 0.36\\ \textbf{91.89}\pm \textbf{0.18}\\ 69.27\pm 7.07\\ 64.08\pm 1.06\\ 88.50\pm 0.13 \end{array}$	$\frac{75.49 \pm 0.00}{73.89 \pm 0.00}$ $\frac{49.43 \pm 0.00}{85.48 \pm 0.00}$ 80.63 ± 0.00 57.43 ± 0.00 83.69 ± 0.00	$\begin{array}{c} 7.00\\ 62.76\pm0.81\\ 70.98\pm0.96\\ 76.35\pm1.41\\ 86.97\pm0.26\\ 72.54\pm4.69\\ \underline{73.49\pm2.21}\\ 96.14\pm0.63\\ \end{array}$	$\begin{array}{r} 4.71\\ \hline 67.87\pm 0.37\\ 74.12\pm 0.65\\ 74.24\pm 0.75\\ 85.37\pm 0.26\\ \hline 88.54\pm 0.20\\ \hline \textbf{79.57}\pm \textbf{0.61}\\ 95.11\pm 0.32\\ \end{array}$	$\begin{array}{r} \underline{4.00} \\ \hline 72.60 \pm 0.06 \\ \underline{78.22 \pm 0.80} \\ \underline{80.99 \pm 0.88} \\ 91.06 \pm 0.60 \\ 62.00 \pm 14.00 \\ 70.51 \pm 1.83 \\ \underline{97.62 \pm 0.12} \end{array}$	4.71 80.06 ± 0.85 72.02 ± 2.64 64.93 ± 3.31 90.99 ± 2.19 94.15 ± 0.08 66.25 ± 0.51 Timeout	$\begin{array}{r} \textbf{2.14} \\ \textbf{73.63} \pm \textbf{0.54} \\ \textbf{80.86} \pm \textbf{1.24} \\ \textbf{81.11} \pm \textbf{0.63} \\ \textbf{91.15} \pm \textbf{0.54} \\ \textbf{79.86} \pm \textbf{2.18} \\ \textbf{71.95} \pm \textbf{2.51} \\ \textbf{97.68} \pm \textbf{0.42} \end{array}$
		Avg. Rank	8.29	9.57	5.43	5.43	6.57	7.57	6.00	5.00	4.00	5.71	2.43

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
-	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	$\underline{94.17 \pm 0.10}$	60.40 ± 3.01	$\textbf{94.42} \pm \textbf{0.21}$
_	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	89.62 ± 0.27	74.61 ± 1.64	$\textbf{89.67} \pm \textbf{0.27}$
8	MOOC	80.57 ± 0.52	80.50 ± 0.68	85.28 ± 0.30	88.01 ± 1.48	81.32 ± 0.42	82.28 ± 0.99	81.38 ± 0.17	87.05 ± 0.51	64.99 ± 2.24	$\textbf{88.64} \pm \textbf{0.08}$
P	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	98.83 ± 0.02	80.07 ± 2.53	$\textbf{98.97} \pm \textbf{0.01}$
Ra	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	98.58 ± 0.01	93.58 ± 0.65	$\textbf{98.77} \pm \textbf{0.03}$
_	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	94.45 ± 0.13	49.78 ± 5.02	$\textbf{94.76} \pm \textbf{0.19}$
	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	93.05 ± 0.10	Timeout	$\textbf{93.13} \pm \textbf{0.05}$
	Avg. Rank	6.29	8.00	7.00	5.14	4.43	5.57	5.86	2.14	9.57	1.00
-	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	76.27 ± 0.37	75.07 ± 1.45	55.60 ± 3.91	$\textbf{76.76} \pm \textbf{0.43}$
	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	71.71 ± 1.33	67.21 ± 0.72	68.66 ± 2.31	68.77 ± 0.60
ŝ	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	80.66 ± 0.94	57.49 ± 1.34	$\textbf{80.75} \pm \textbf{1.00}$
2	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	$\textbf{78.35} \pm \textbf{5.03}$	65.30 ± 1.05
P	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	87.47 ± 0.25	57.90 ± 11.05	$\textbf{92.61} \pm \textbf{0.90}$	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	74.39 ± 1.81	$\overline{\textbf{81.40}\pm\textbf{0.52}}$	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	$\overline{95.94\pm0.37}$	94.56 ± 0.24	$\underline{96.73 \pm 0.11}$	Timeout	$\textbf{96.83} \pm \textbf{0.56}$
-	Avg. Rank	6.86	8.57	5.29	5.43	7.43	4.86	3.14	4.43	6.57	2.43

4.2 PERFORMANCE ANALYSIS

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) 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.

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

3	8	8
3	8	9

385

386

387

394

Models								MOOC			Reddit		v	Vikipedi	a		UCI		s	ocial Ev	0.
mouchs	R	Н	Ι	R	Н	Ι	R	Н	Ι	R	Н	Ι	R	Н	Ι	R	Н	Ι	R	Н	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 C	93.07 92.71	82.85 82.85	72.36	92.40 92.49	76.99	78.64	88.80 88.80	85.23	81.02	99.27	81.74	91.05	99.06	79.14	73.49	91.09 95.85	81.00	71.86	92.90 94.71	96.01 96.71	97.14 97.25
Variant D	<mark>92.74</mark>	<mark>82.87</mark>	<mark>72.68</mark>	<mark>92.52</mark>	<mark>77.07</mark>	<mark>78.05</mark>	<mark>88.71</mark>	<mark>85.76</mark>	<mark>81.09</mark>	<mark>99.27</mark>	82.10	<mark>91.07</mark>	<mark>99.08</mark>	<mark>81.75</mark>	<mark>79.79</mark>	<mark>95.87</mark>	<mark>82.35</mark>	<mark>72.98</mark>	<mark>94.74</mark>	<mark>97.17</mark>	<mark>97.60</mark>
Variant D	92.71 92.74	82.85 82.87	72.36 72.68	92.49 92.52	76.99 77.07	78.64 78.05	88.80 88.71	85.23 85.76	81.02 81.09	99.27 99.27	81.74 82.10	91.05 91.07	<mark>99.06</mark> 99.08	79.14 81.75	73.49 79.79	95.85 95.87	81 82	.00	.00 71.86 .35 72.98	.00 71.86 94.71 .35 72.98 94.74	.00 71.86 94.71 96.71 .35 72.98 94.74 97.17

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

Datasets	Las	tFM	En	ron	MO	OC	Re	ddit	Wiki	pedia	U	CI	Socia	l Evo.
Models	R	Ι	R	Ι	R	Ι	R	Ι	R	Ι	R	Ι	R	Ι
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
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	<mark>94.21</mark>	<mark>76.64</mark>	<mark>89.44</mark>	<mark>67.91</mark>	88.29	<mark>80.86</mark>	<mark>98.91</mark>	<mark>65.10</mark>	<mark>98.69</mark>	71.10	94.51	73.50	<mark>93.10</mark>	96.75
VyGMamba	94.42	76.76	89.67	68.77	88.64	80.75	98.97	65.30	98.77	71.14	94.76	72.17	93.13	96.83

400 **Ablation Study.** We conduct four ablation studies to study the effectiveness of model components. 401 In study A, we make a model variant (Variant A) by removing the time-level SSM block and restrain 402 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 403 Mamba SSM layers (Eq. 3) in the node-level SSM block. In study C, we switch the computation 404 of the selection weights β_{θ} in Eq. 6d to $\beta_{\theta} = \text{Softmax}(f_{\text{sel}}(\mathbf{H}_{\theta}^{t})) (f_{\text{sel}}(\cdot) : \mathbb{R}^{4d} \to \mathbb{R}^{4d})$ to create 405 Variant C. In study D, we base on Variant C and develop Variant D that further enables information 406 407 selection from opposite nodes, i.e., $\beta_u = \text{Softmax}(f_{\text{sel}}(\mathbf{H}_u^t)) / \beta_v = \text{Softmax}(f_{\text{sel}}(\mathbf{H}_u^t))$. Both ablation C and D do information selection without learning temporal patterns. From Table 3 and 4, we 408 find that: (1) Variant A is constantly beaten by DyGMamba, showing the effectiveness of dynamic 409 information selection based on edge-specific temporal patterns; (2) DyGMamba always outperforms 410 Variant B, indicating the importance of encoding the one-hop temporal neighbors with SSM layers 411 for capturing graph dynamics; (3) Variant C and D perform better than Variant A in most cases, 412 implying that selecting temporal information is generally contributive; (4) Variant C generally lags 413 behind Variant D, meaning that information selection from opposite node is beneficial; (5) DyG-414 Mamba performs better than both Variant C and D in almost all cases, proving that information 415 selection based on temporal patterns is more effective. 416

417 A Closer Look into Temporal Pattern Modeling. We observe from ablation studies that dy-418 namic information selection based on temporal patterns contributes to better model performance on 419 real-world datasets. To better quantify its benefits, we construct three synthetic datasets, i.e., S1, S2 420 and S3, that follow different patterns and compare our model with DyGFormer, CTAN as well as 421 Variant A on them. Each synthetic dataset contains 7 nodes, where the interactions of each pair of 422 two nodes follow a certain pattern along time. And for each node, we generate interactions with all 423 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, 424 the time intervals between neighboring interactions $\{t_1 - t_0, ..., t_N - t_{N-1}\}$ follow an increasing 425 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 set 426 the time intervals to a decreasing trend with the same velocity, i.e., $(t_{i+1} - t_i) - (t_{i+2} - t_{i+1}) =$ 427 0.05. And in S3, we modify S1 by repeating several periods of increasing patterns taken from 428 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 429 of dataset construction and statistics are provided in App. A.2. From Table 5, we observe that 430 DyGMamba greatly outperforms DyGFormer and CTAN. More importantly, Variant A, C and D 431 show similar performance to DyGFormer, meaning that our time-level SSM block is able to capture



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

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.

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

Model Size, Per Epoch Training Time and GPU Memory. Fig. 2 compares DyGMamba with five baselines in terms of number of parameters (model size), per epoch training time and GPU memory consumption during training⁴. 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 much worse; (2) DyGMamba is much more efficient than DyGFormer with the same length of input sequence $(\rho/p = 32)$; (3) Al-though DyGMamba generally consumes more GPU memory and per epoch training time compared with most baselines, the gap of consumption is not very large in most cases. To model more tem-poral 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 by considering more temporal neighbors; (4) CTAN requires very few computational resources. However, on long-range tempo-ral dependent datasets, it is beaten by DyGFormer and DyGMamba with a great margin. Besides, CTAN is also hard to converge. Although it requires very little per epoch training time, it requires much more epochs to reach the best performance, leading to a long total training time. See App. G.2 and G.3 for details.

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 exces-sive 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 gradu-ally 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 pa-

⁴The 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. G



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

500 rameters, 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 501 is much lower for DyGMamba, demonstrating its strong scalability in modeling longer sequences; 502 (3) Different trends in performance change are observed between two models. While DyGFormer 503 performs worse, DyGMamba can benefit from a smaller patch size, indicating its strong ability to 504 capture nuanced temporal details even if the sequence becomes much longer. Note that the models 505 use fewer parameters under smaller patch size. This also shows that DyGMamba can achieve much 506 stronger parameter efficiency by reducing patch sizes. To further study the reason for finding (3), 507 we plot the performance of Variant A under different patch sizes in Fig. 3a. We find that Variant A's 508 performance degrades when sequence length is more than 64. This means that dynamic information 509 selection based on edge-specific temporal patterns is essential for DyGMamba to optimally process 510 long sequences. We provide more explanations and additional analysis on MOOC in App. J.

511

499

512 Scalability in Modeling Increasing Number of Temporal Neighbors. Sequence length is decided by the sampled temporal neighbors ρ and the patch size p. If we want to model a huge

514 number of temporal neighbors, e.g., $\rho = 4096$, keeping 515 the sequence length unchanged as 32 means we need to set p to 128. This will greatly increase model parameters 516 and cause burden in parameter optimization. Besides, as 517 shown in Fig. 3a, bigger p does not necessarily lead to 518 better performance. Therefore, it is also important to see 519 if a model is scalable to ρ , with a fixed p. In Fig. 4, 520 we show the consumed GPU memory of DyGMamba and 521 DyGFormer with ρ varying from 64, 128, 256, 512, 1024, 522 2048, 4096 and 8192 on Enron, under a fixed patch size 523 p = 8. We find that DyGMamba shows superior scalabil-524 ity over DyGFormer. While DyGFormer can only process 2048 nodes, DyGMamba can deal with more than 8192 525 526 (at least 4 times) on a single 45GB NVIDIA A40 GPU. This implies our model's potential to process a huge num-527



Figure 4: GPU memory consumption on Enron with increasing neighbors. For clarity, memory is shown only when neighbors are more than 512.

ber of temporal neighbors with a limited GPU memory budget. We provide the complexity analysis
 of DyGMamba and DyGFormer in App. G.4 to further explain DyGMamba's scalability.

530 531

532

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 se-quences 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. To supplement, we discuss the limitation of our work in App. K.

540 REFERENCES

Ali Behrouz and Farnoosh Hashemi. Graph mamba: Towards learning on graphs with state space models. *CoRR*, abs/2402.08678, 2024. doi: 10.48550/ARXIV.2402.08678. URL https://doi.org/10.48550/arXiv.2402.08678.

Xiaofu Chang, Xuqin Liu, Jianfeng Wen, Shuang Li, Yanming Fang, Le Song, and Yuan Qi.
Continuous-time dynamic graph learning via neural interaction processes. In Mathieu d'Aquin,
Stefan Dietze, Claudia Hauff, Edward Curry, and Philippe Cudré-Mauroux (eds.), *CIKM '20: The 29th ACM International Conference on Information and Knowledge Management, Virtual Event, Ireland, October 19-23, 2020*, pp. 145–154. ACM, 2020. doi: 10.1145/3340531.3411946. URL
https://doi.org/10.1145/3340531.3411946.

- Weilin Cong, Si Zhang, Jian Kang, Baichuan Yuan, Hao Wu, Xin Zhou, Hanghang Tong, and Mehrdad Mahdavi. 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, 2023. URL https://openreview.net/pdf? id=ayPPc0SyLv1.
- Feyza Duman Keles, Pruthuvi Mahesakya Wijewardena, and Chinmay Hegde. On the computational complexity of self-attention. In Shipra Agrawal and Francesco Orabona (eds.), *Proceedings of The 34th International Conference on Algorithmic Learning Theory*, volume 201 of *Proceedings of Machine Learning Research*, pp. 597–619. PMLR, 20 Feb–23 Feb 2023. URL https://proceedings.mlr.press/v201/duman-keles23a.html.
- Palash Goyal, Sujit Rokka Chhetri, and Arquimedes Canedo. dyngraph2vec: Capturing network dynamics using dynamic graph representation learning. *Knowl. Based Syst.*, 187, 2020. doi: 10.1016/J.KNOSYS.2019.06.024. URL https://doi.org/10.1016/j.knosys. 2019.06.024.
 - Alessio Gravina, Giulio Lovisotto, Claudio Gallicchio, Davide Bacciu, and Claas Grohnfeldt. Long range propagation on continuous-time dynamic graphs. In *Forty-first International Conference on Machine Learning*, 2024. URL https://openreview.net/forum?id=gVg8V9isul.
 - Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *CoRR*, abs/2312.00752, 2023. doi: 10.48550/ARXIV.2312.00752. URL https://doi.org/10.48550/arXiv.2312.00752.
- Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 572–585, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/05546b0e38ab9175cd905eebcc6ebb76-Abstract.html.
- Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), 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, 2022a. URL http://papers.nips.cc/paper_files/paper/2022/hash/ e9a32fade47b906de908431991440f7c-Abstract-Conference.html.
- Albert Gu, Karan Goel, and Christopher Ré. 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, 2022b. URL https://openreview.net/forum?id=uYLFoz1vlAC.
- 592

566

567

568

569 570

571

572

573

Ming Jin, Yuan-Fang Li, and Shirui Pan. Neural temporal walks: Motif-aware representation learning on continuous-time dynamic graphs. In Sanmi Koyejo, S. Mohamed, 594 A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural In-595 formation Processing Systems 35: Annual Conference on Neural Information Process-596 ing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 597 2022, 2022. URL http://papers.nips.cc/paper_files/paper/2022/hash/ 598 7dadc855cef7494d5d956a8d28add871-Abstract-Conference.html. Seyed Mehran Kazemi, Rishab Goel, Kshitij Jain, Ivan Kobyzev, Akshay Sethi, Peter Forsyth, and 600 Pascal Poupart. Representation learning for dynamic graphs: A survey. J. Mach. Learn. Res., 21: 601 70:1-70:73, 2020. URL http://jmlr.org/papers/v21/19-447.html. 602 603 Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in tem-604 poral interaction networks. In Ankur Teredesai, Vipin Kumar, Ying Li, Rómer Rosales, Evi-605 maria Terzi, and George Karypis (eds.), Proceedings of the 25th ACM SIGKDD International 606 Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, Au-607 gust 4-8, 2019, pp. 1269-1278. ACM, 2019. doi: 10.1145/3292500.3330895. URL https: //doi.org/10.1145/3292500.3330895. 608 609 Jintang Li, Ruofan Wu, Xinzhou Jin, Boqun Ma, Liang Chen, and Zibin Zheng. State space models 610 on temporal graphs: A first-principles study. arXiv preprint arXiv:2406.00943, 2024a. 611 612 Lincan Li, Hanchen Wang, Wenjie Zhang, and Adelle Coster. Stg-mamba: Spatial-temporal graph 613 learning via selective state space model. CoRR, abs/2403.12418, 2024b. doi: 10.48550/ARXIV. 614 2403.12418. URL https://doi.org/10.48550/arXiv.2403.12418. 615 616 Yunyu Liu, Jianzhu Ma, and Pan Li. Neural predicting higher-order patterns in temporal networks. 617 In Frédérique Laforest, Raphaël Troncy, Elena Simperl, Deepak Agarwal, Aristides Gionis, Ivan Herman, and Lionel Médini (eds.), WWW '22: The ACM Web Conference 2022, Virtual Event, 618 Lyon, France, April 25 - 29, 2022, pp. 1340–1351. ACM, 2022. doi: 10.1145/3485447.3512181. 619 URL https://doi.org/10.1145/3485447.3512181. 620 621 Xuezhe Ma, Chunting Zhou, Xiang Kong, Junxian He, Liangke Gui, Graham Neubig, Jonathan 622 May, and Luke Zettlemoyer. Mega: Moving average equipped gated attention. In The Eleventh 623 International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 624 2023. OpenReview.net, 2023. URL https://openreview.net/pdf?id=qNLe3iq2El. 625 Yao Ma, Ziyi Guo, Zhaochun Ren, Jiliang Tang, and Dawei Yin. Streaming graph neural networks. 626 In Jimmy X. Huang, Yi Chang, Xueqi Cheng, Jaap Kamps, Vanessa Murdock, Ji-Rong Wen, 627 and Yiqun Liu (eds.), Proceedings of the 43rd International ACM SIGIR conference on research 628 and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020, 629 pp. 719-728. ACM, 2020. doi: 10.1145/3397271.3401092. URL https://doi.org/10. 630 1145/3397271.3401092. 631 632 Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, 633 Tim Kaler, Tao B. Schardl, and Charles E. Leiserson. Evolvegcn: Evolving graph convolutional 634 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 635 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, 636 New York, NY, USA, February 7-12, 2020, pp. 5363–5370. AAAI Press, 2020. doi: 10.1609/ 637 AAAI.V34I04.5984. URL https://doi.org/10.1609/aaai.v34i04.5984. 638 639 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, 640 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas 641 Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, 642 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, 643 high-performance deep learning library. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d'Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), Advances in 644

Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada, pp. 646 8024-8035, 2019. URL https://proceedings.neurips.cc/paper/2019/hash/ 647 bdbca288fee7f92f2bfa9f7012727740-Abstract.html.

- 648 Bo Peng, Eric Alcaide, Quentin Anthony, Alon Albalak, Samuel Arcadinho, Stella Biderman, 649 Huanqi Cao, Xin Cheng, Michael Chung, Leon Derczynski, Xingjian Du, Matteo Grella, Kran-650 thi Kiran GV, Xuzheng He, Haowen Hou, Przemyslaw Kazienko, Jan Kocon, Jiaming Kong, Bart-651 lomiej Koptyra, Hayden Lau, Jiaju Lin, Krishna Sri Ipsit Mantri, Ferdinand Mom, Atsushi Saito, 652 Guangyu Song, Xiangru Tang, Johan S. Wind, Stanislaw Wozniak, Zhenyuan Zhang, Qinghua Zhou, Jian Zhu, and Rui-Jie Zhu. RWKV: reinventing rnns for the transformer era. In Houda 653 Bouamor, Juan Pino, and Kalika Bali (eds.), Findings of the Association for Computational 654 Linguistics: EMNLP 2023, Singapore, December 6-10, 2023, pp. 14048–14077. Association 655 for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-EMNLP.936. URL 656 https://doi.org/10.18653/v1/2023.findings-emnlp.936. 657
- 658 Farimah Poursafaei, Shenyang Huang, Kellin Pelrine, and Reihaneh Rabbany. To-659 wards better evaluation for dynamic link prediction. In Sanmi Koyejo, S. Mohamed, 660 A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Infor-661 mation Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 662 2022. 2022. URL http://papers.nips.cc/paper_files/paper/2022/ 663 hash/d49042a5d49818711c401d34172f9900-Abstract-Datasets_and_ 664 Benchmarks.html. 665
- Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and
 Michael M. Bronstein. Temporal graph networks for deep learning on dynamic graphs. *CoRR*,
 abs/2006.10637, 2020. URL https://arxiv.org/abs/2006.10637.
- Aravind Sankar, Yanhong Wu, Liang Gou, Wei Zhang, and Hao Yang. Dysat: Deep neural representation learning on dynamic graphs via self-attention networks. In James Caverlee, Xia (Ben) Hu, Mounia Lalmas, and Wei Wang (eds.), *WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, Houston, TX, USA, February 3-7, 2020*, pp. 519–527. ACM, 2020. doi: 10.1145/3336191.3371845. URL https://doi.org/10.1145/3336191.3371845.
- Jimmy T. H. Smith, Andrew Warrington, and Scott W. Linderman. 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, 2023. URL https://openreview.net/pdf?id=Ai8Hw3AXqks.
- Yuxing Tian, Yiyan Qi, and Fan Guo. Freedyg: Frequency enhanced continuous-time dynamic graph
 model for link prediction. In *The Twelfth International Conference on Learning Representations*,
 2024. URL https://openreview.net/forum?id=82Mc5illnM.
- 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. Mlp-mixer: An all-mlp architecture for vision. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 24261–24272, 2021. URL https://proceedings.neurips.cc/paper/2021/ hash/cba0a4ee5ccd02fda0fe3f9a3e7b89fe-Abstract.html.
- Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. Dyrep: Learning representations over dynamic graphs. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. URL https://openreview.net/forum?id=HyePrhR5KX.

 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 5998–6008, 2017. URL https://proceedings.neurips.cc/paper/2017/hash/ 3f5ee243547dee91fbd053c1c4a845aa-Abstract.html.

702	Chloe Wang, Oleksii Tsepa, Jun Ma, and Bo Wang. Graph-mamba: Towards long-range graph
703	sequence modeling with selective state spaces. CoRR, abs/2402.00789, 2024. doi: 10.48550/
704	ARXIV.2402.00789. URL https://doi.org/10.48550/arXiv.2402.00789.

- Lu Wang, Xiaofu Chang, Shuang Li, Yunfei Chu, Hui Li, Wei Zhang, Xiaofeng He, Le Song, Jingren Zhou, and Hongxia Yang. TCL: transformer-based dynamic graph modelling via contrastive learning. *CoRR*, abs/2105.07944, 2021a. URL https://arxiv.org/abs/2105.07944.
- Xuhong Wang, Ding Lyu, Mengjian Li, Yang Xia, Qi Yang, Xinwen Wang, Xinguang Wang, Ping
 Cui, Yupu Yang, Bowen Sun, and Zhenyu Guo. APAN: asynchronous propagation attention
 network for real-time temporal graph embedding. In Guoliang Li, Zhanhuai Li, Stratos Idreos,
 and Divesh Srivastava (eds.), *SIGMOD '21: International Conference on Management of Data, Virtual Event, China, June 20-25, 2021*, pp. 2628–2638. ACM, 2021b. doi: 10.1145/3448016.
 3457564. URL https://doi.org/10.1145/3448016.3457564.
- Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. 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, 2021c. URL https://openreview.net/forum?id=KYPz4YsCPj.
- Da Xu, Chuanwei Ruan, Evren Körpeoglu, Sushant Kumar, and Kannan Achan. Inductive representation learning on temporal graphs. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net, 2020. URL https://openreview.net/forum?id=rJeWlyHYwH.
- Jiaxuan You, Tianyu Du, and Jure Leskovec. ROLAND: graph learning framework for dynamic graphs. In Aidong Zhang and Huzefa Rangwala (eds.), *KDD '22: The 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Washington, DC, USA, August 14 18, 2022*, pp. 2358–2366. ACM, 2022. doi: 10.1145/3534678.3539300. URL https://doi.org/10.1145/3534678.3539300.
- Le Yu, Leilei Sun, Bowen Du, and Weifeng Lv. Towards better dynamic graph learning: New architecture and unified library. In Alice Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), 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, 2023. URL http://papers.nips.cc/paper_files/paper/2023/hash/ d611019afba70d547bd595e8a4158f55-Abstract-Conference.html.
- Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. 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, 2024. URL https://openreview.net/forum?id=YbHCqn4qF4.
- 740 741 742

744 745

746

747

748

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. (2023). We chronologically split each dataset with the ratio of 70%/15%/15% for training/validation/testing. Please refer to it for detailed dataset descriptions.

749 A.2 SYNTHETIC DATASETS

For all of our three synthetic datasets, node and link features are not involved during dataset con struction. 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,

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.	/4	2,099,519	0&2	False	8 months	565,932	Unix timestamps
i.e., rand	domly det	t_{i+1}) – (t_i termined ar	$t_{i+1} - t_i = 0.$	05. The n tion numb	umber of in pers of all no	nteractions for ode pairs sum	each node pair a up to 100000.
• S2:	For each	interaction	pair u and v ,	, their first	interaction	is at timestam	p 0 and the secon
will	not dron	to zero or	ганцонну. п negative afte	rwards St	arting from	the second in	terval they follow
a de	ecreasing	trend with	a constant ve	locity of 0	$.05$, i.e., (t_i)	$(\pm 1 - t_i) - (t_i)$	$(12 - t_{i+1}) = 0.02$
The	number	of interacti	ons for each	node pair	is randoml	y determined a	and the interaction
nun	nhers of a	ll node nai	rs sum un to	100000		-	

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

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 5. Note that for the periodic dataset S3 (Figure 5c), 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.

Datasets	# Nodes	# Edges	# Timestamps	Time Range
S1	7	100,000	96,869	0 - 163241.65
S2	7	100,000	98,004	0 - 1573561.52
S 3	7	95,657	95,370	0 - 1771300.40

Table 7: Synthetic dataset statistics.

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
	S1	1,428.57	1,205	1,663
Training Set	S2	1,428.57	1,424	1,433
-	S 3	1,367.91	1,364	1,372
	S1	2,010.20	1,879	2,150
Complete Detect	S2	2,010.20	1,921	2,097
Complete Dataset	S3	1,952.18	1,721	2,193
	S3 (each period)	244.02	215	274



Figure 5: Time intervals of a node pair u, v in synthetic datasets S1, S2 and S3.

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** (Kumar et al., 2019): 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** (Trivedi et al., 2019): 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** (Xu et al., 2020): 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.
- **TGN** (Rossi et al., 2020): 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 (Wang et al., 2021c): 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 (Poursafaei et al., 2022): 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) EdgeBank_{∞}, where all the observed edges are stored in the memory; (2) EdgeBank_{tw-ts}, where only the edges within the duration of the test set from the immediate past are kept in the memory; (3) EdgeBank_{tw-re}, where only the edges within the average time intervals of repeated edges from the immediate past are kept in the memory; (4) EdgeBank_{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 (Wang et al., 2021a): 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** (Cong et al., 2023): GraphMixer designs a link-encoder based on MLP-Mixer (Tolstikhin et al., 2021) to learn from the temporal interactions. A mean poolingbased 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. (2023), 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** (Yu et al., 2023): 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** (Gravina et al., 2024): 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

894 We train every CTDG model except for CTAN for a maximum number of 200 epochs. Maximum 895 epochs for CTAN training is 1000. We evaluate each model on the validation set at the end of 896 every training epoch and adopt an early stopping strategy with a patience of 20. We take the model 897 that achieves the best validation result for testing. We use the implementations⁵ provided by Yu et al. (2023) for all baseline models except CTAN. For CTAN, we use its official implementation⁶. 899 All models are trained with a batch size of 200 for fair efficiency analysis. All experiments are implemented with PyTorch (Paszke et al., 2019) on a server equipped with an AMD EPYC 7513 900 32-Core Processor and a single NVIDIA A40 with 45GB memory. We run each experiment for five 901 times with five random seeds and report the mean results together with error bars. 902

903 904

864

865

866

867

868

875

876

877

878

879

881

882

883

885

886

888

889

890

891 892

893

C.1 HYPERPARAMETER CONFIGURATIONS ON REAL-WORLD DATASETS

For all the baselines except CTAN, please refer to Yu et al. (2023) for the hyperparameter configurations on real-world datasets. For CTAN, we present its hyperparameter configurations in Table
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.

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.,

⁶https://github.com/gravins/non-dissipative-propagation-CTDGs

^{916 &}lt;sup>5</sup>https://github.com/yule-BUAA/DyGLib

Table 9: Hyperparameter configurations of CTAN on all real-world datasets. γ here denotes the discretization step size introduced in Gravina et al. (2024), different from the one in DyGMamba.

Model	# GCL	ϵ	γ	Embedding Dim
CTAN	5	0.5	0.5	128

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

Datasets	Dropout	$ ho \ \& \ p$	k
LastFM	{0.0, 0.1 , 0.2}	{1024 & 32, 512 & 16 , 256 & 8}	{30, 10 , 5}
Enron	{ 0.0 , 0.1, 0.2}	{512 & 16, 256 & 8 , 128 & 4}	{ 30 , 10, 5}
MOOC	{0.0, 0.1 , 0.2}	{512 & 16, 256 & 8, 128 & 4 }	{30, 10 , 5}
Reddit	{0.0, 0.1, 0.2 }	{128 & 4, 64 & 2 , 32 & 1}	{30, 10, 5 }
Wikipedia	{0.0, 0.1 , 0.2}	{64 & 2, 32 & 1 }	{30, 10, 5 }
UCI	{0.0, 0.1 , 0.2}	{64 & 2, 32 & 1 }	{30, 10, 5 }
Social Evo.	{0.0, 0.1 , 0.2}	(64 & 2, 32 & 1)	{30, 10, 5 }

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, 64 & 2 , 32 & 1}	{512 & 16, 256 & 8 , 128 & 4, 64 & 2, 32 & 1}	{30, 10, 5}
S2	512 & 16 , 256 & 8, 128 & 4, 64 & 2, 32 & 1	512 & 16 , 256 & 8, 128 & 4, 64 & 2, 32 & 1	{ 30 , 10, 5}
S3	{512 & 16, 256 & 8, 128 & 4, 64 & 2, 32 & 1 }	512 & 16 , 256 & 8, 128 & 4, 64 & 2, 32 & 1	{30, 10 , 5}

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 ρ 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 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. (2022), 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

⁷https://github.com/state-spaces/mamba

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

NSS	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	EdgeBank	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Random	LastFM Enron MOOC Reddit Wikipedia UCI Social Evo.	$70.89 \pm 1.97 \\ 87.77 \pm 2.43 \\ 84.50 \pm 0.60 \\ 98.29 \pm 0.05 \\ 96.36 \pm 0.14 \\ 90.35 \pm 0.51 \\ 92.13 \pm 0.20 \\ \end{cases}$	$71.40 \pm 2.12 \\ 83.09 \pm 2.20 \\ 84.50 \pm 0.87 \\ 98.13 \pm 0.04 \\ 94.43 \pm 0.32 \\ 69.46 \pm 2.66 \\ 90.37 \pm 0.52 \\ \end{array}$	$71.47 \pm 0.14 \\ 68.57 \pm 1.46 \\ 87.01 \pm 0.16 \\ 98.50 \pm 0.01 \\ 96.60 \pm 0.07 \\ 78.76 \pm 1.10 \\ 94.93 \pm 0.06 \\ 84.93 \pm 0.06 \\ 84.94 \pm 0.06 \\ 84.94 \pm 0.06 \\ 84.9$	$76.64 \pm 4.66 \\ 88.72 \pm 0.95 \\ 91.91 \pm 0.82 \\ 98.61 \pm 0.05 \\ 98.37 \pm 0.10 \\ 92.03 \pm 0.69 \\ 95.31 \pm 0.27 \\ \end{tabular}$	$85.92 \pm 0.1690.34 \pm 0.2380.48 \pm 0.4199.02 \pm 0.0098.54 \pm 0.0193.81 \pm 0.2387.34 \pm 0.10$	$\begin{array}{c} 83.77 \pm 0.00 \\ 87.05 \pm 0.00 \\ 60.86 \pm 0.00 \\ 95.37 \pm 0.00 \\ 90.78 \pm 0.00 \\ 77.30 \pm 0.00 \\ 81.60 \pm 0.00 \end{array}$	$71.09 \pm 1.48 \\ 83.33 \pm 0.93 \\ 84.02 \pm 0.59 \\ 97.67 \pm 0.01 \\ 97.27 \pm 0.06 \\ 85.49 \pm 0.82 \\ 95.45 \pm 0.21 \\ 95.45 \pm 0.21 \\ 85.49 \pm 0.82 \\ 95.45 \pm 0.21 \\ 95.4$	$73.51 \pm 0.14 \\ 84.16 \pm 0.34 \\ 84.04 \pm 0.12 \\ 97.17 \pm 0.02 \\ 96.89 \pm 0.04 \\ 91.62 \pm 0.52 \\ 95.21 \pm 0.07 \\ 85.01 \pm 0.07 \\ 95.21 \pm 0.07 \\ 95.2$	$\begin{array}{c} 93.03 \pm 0.11\\ \hline 93.20 \pm 0.12\\ \hline 88.08 \pm 0.50\\ \hline 99.15 \pm 0.01\\ \hline 98.92 \pm 0.03\\ \hline 94.45 \pm 0.22\\ \hline 96.25 \pm 0.04\\ \end{array}$	$85.12 \pm 0.77 \\ 87.09 \pm 1.51 \\ 85.40 \pm 2.67 \\ 97.24 \pm 0.75 \\ 97.00 \pm 0.21 \\ 76.25 \pm 2.83 \\ Timeout$	$\begin{array}{c} \textbf{93.31} \pm \textbf{0.18} \\ \textbf{93.34} \pm \textbf{0.23} \\ \textbf{89.58} \pm \textbf{0.12} \\ \textbf{99.27} \pm \textbf{0.01} \\ \textbf{99.08} \pm \textbf{0.02} \\ \textbf{94.77} \pm \textbf{0.18} \\ \textbf{96.38} \pm \textbf{0.02} \end{array}$
	Avg. Rank	7.14	8.86	7.14	3.86	4.86	9.14	7.29	7.14	2.14	7.29	1.14
Historical	LastFM Enron MOOC Reddit Wikipedia UCI Social Evo.	$\begin{array}{c} 75.65\pm4.43\\ 75.21\pm1.27\\ 82.38\pm1.75\\ 80.70\pm0.20\\ 80.71\pm0.64\\ 78.21\pm3.18\\ 91.83\pm1.52 \end{array}$	$\begin{array}{c} 70.63 \pm 2.56 \\ 76.36 \pm 1.42 \\ 80.71 \pm 2.08 \\ 79.96 \pm 0.23 \\ 77.49 \pm 0.72 \\ 58.65 \pm 3.58 \\ 92.81 \pm 0.60 \end{array}$	$\begin{array}{c} 64.23 \pm 0.45 \\ 62.36 \pm 1.07 \\ 81.53 \pm 0.79 \\ 79.60 \pm 0.09 \\ 82.83 \pm 0.27 \\ 57.12 \pm 0.98 \\ 93.63 \pm 0.48 \end{array}$	$\begin{array}{c} 78.00\pm2.97\\ 76.75\pm1.40\\ 86.59\pm2.03\\ 81.04\pm0.23\\ 83.28\pm0.26\\ \textbf{78.48}\pm1.79\\ 94.27\pm1.33 \end{array}$	$\begin{array}{c} 67.92\pm0.32\\ 65.62\pm0.49\\ 71.74\pm0.88\\ 80.42\pm0.20\\ 65.74\pm3.46\\ 57.67\pm1.11\\ 87.61\pm0.06 \end{array}$	$\begin{array}{c} 78.09\pm 0.00\\ \hline 79.59\pm 0.00\\ \hline 61.90\pm 0.00\\ 78.58\pm 0.00\\ 77.27\pm 0.00\\ 69.56\pm 0.00\\ 85.81\pm 0.00\\ \end{array}$	$\begin{array}{c} 60.53 \pm 2.54 \\ 71.72 \pm 1.24 \\ 73.22 \pm 1.21 \\ 76.83 \pm 0.12 \\ 85.55 \pm 0.47 \\ 65.42 \pm 2.62 \\ 95.03 \pm 0.82 \end{array}$	$\begin{array}{c} 64.06\pm0.34\\ 74.82\pm2.04\\ 77.09\pm0.83\\ 77.83\pm0.33\\ \underline{87.47\pm0.20}\\ \overline{77.46\pm1.63}\\ \overline{94.65\pm0.28}\end{array}$	$\begin{array}{c} 78.80\pm0.02\\ 77.35\pm0.64\\ \underline{87.26\pm0.83}\\ 80.61\pm0.48\\ 72.78\pm6.65\\ 75.71\pm0.57\\ \underline{97.16\pm0.06} \end{array}$	$\begin{array}{r} \hline 79.50 \pm 0.82 \\ \hline \textbf{81.95} \pm \textbf{1.64} \\ 73.87 \pm 2.77 \\ \textbf{90.63} \pm \textbf{2.28} \\ \textbf{95.43} \pm \textbf{0.07} \\ 75.05 \pm 0.13 \\ \hline \text{Timeout} \end{array}$	$\begin{array}{c} \textbf{79.82} \pm \textbf{0.27} \\ \textbf{77.73} \pm \textbf{0.61} \\ \textbf{87.91} \pm \textbf{0.93} \\ \underline{81.71} \pm \textbf{0.49} \\ \textbf{78.99} \pm \textbf{1.24} \\ \textbf{75.43} \pm \textbf{1.99} \\ \textbf{97.27} \pm \textbf{0.30} \end{array}$
Inductive	Avg. Rank LastFM Enron MOOC Reddit Wikipedia UCI Social Evo.	$\begin{array}{c} 5.29\\ \hline 61.59\pm5.72\\ 70.75\pm0.69\\ 67.53\pm1.76\\ 83.40\pm0.33\\ 70.41\pm0.39\\ 64.14\pm1.25\\ 91.81\pm1.69\end{array}$	$\begin{array}{r} 7.14\\ \hline 60.62\pm2.20\\ 67.37\pm2.21\\ 62.60\pm1.27\\ 82.75\pm0.36\\ 67.57\pm0.94\\ 54.10\pm2.74\\ 92.77\pm0.64 \end{array}$	$\begin{array}{c} 7.86\\\hline\\63.96\pm0.41\\59.78\pm1.12\\74.44\pm0.81\\87.46\pm0.10\\81.54\pm0.31\\59.60\pm0.61\\93.54\pm0.48\end{array}$	$\begin{array}{r} 3.71 \\ \hline 65.48 \pm 4.13 \\ 73.22 \pm 0.42 \\ 76.89 \pm 2.13 \\ 84.57 \pm 0.19 \\ 81.21 \pm 0.30 \\ 63.76 \pm 0.99 \\ 94.86 \pm 1.25 \end{array}$	$\begin{array}{r} 9.14\\ \hline 67.90 \pm 0.44\\ 75.29 \pm 0.66\\ 70.08 \pm 0.33\\ \underline{88.19 \pm 0.20}\\ \overline{68.48 \pm 3.64}\\ 57.85 \pm 0.59\\ 90.10 \pm 0.11 \end{array}$	$\begin{array}{c} 7.43\\ \hline 77.37 \pm 0.00\\ 75.00 \pm 0.00\\ 48.18 \pm 0.00\\ 85.93 \pm 0.00\\ \underline{81.73 \pm 0.00}\\ 58.03 \pm 0.00\\ 87.88 \pm 0.00 \end{array}$	$\begin{array}{r} 7.71\\ \hline 54.75\pm1.31\\ 69.74\pm1.19\\ 71.80\pm1.09\\ 84.41\pm0.18\\ 73.51\pm1.88\\ 65.46\pm2.07\\ \underline{95.13\pm0.83}\end{array}$	$\begin{array}{c} 6.29\\ \hline 59.98 \pm 0.20\\ 70.72 \pm 1.08\\ 72.25 \pm 0.57\\ 82.24 \pm 0.24\\ 84.20 \pm 0.36\\ \textbf{74.25 \pm 0.71}\\ 94.50 \pm 0.26\\ \end{array}$	$\begin{array}{r} \underline{4.29} \\ \hline 67.87 \pm 0.53 \\ 74.67 \pm 0.80 \\ \underline{80.78 \pm 0.89} \\ 86.25 \pm 0.64 \\ 64.09 \pm 9.75 \\ 64.92 \pm 0.83 \\ 95.01 \pm 0.15 \end{array}$	4.29 78.70 ± 0.87 75.40 ± 1.92 68.17 ± 3.73 91.42 ± 2.18 93.67 ± 0.11 66.51 ± 0.25 Timeout	$\begin{array}{r} \textbf{2.86} \\ \hline \textbf{68.74 \pm 0.55} \\ \textbf{75.47 \pm 1.41} \\ \textbf{81.08 \pm 0.82} \\ \textbf{86.35 \pm 0.52} \\ \textbf{75.64 \pm 2.42} \\ \hline \textbf{66.83 \pm 2.83} \\ \textbf{97.37 \pm 0.26} \end{array}$
	Avg. Rank	7.86	9.57	6.14	5.43	6.29	6.43	6.71	6.00	5.14	3.86	2.57

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
	LastFM	82.49 ± 0.94	82.82 ± 1.17	76.76 ± 0.22	82.61 ± 2.62	87.92 ± 0.15	76.95 ± 1.34	80.34 ± 0.14	$\underline{94.10\pm0.09}$	61.49 ± 2.78	$\textbf{94.37} \pm \textbf{0.13}$
_	Enron	80.16 ± 1.50	75.82 ± 3.14	64.25 ± 1.29	79.40 ± 1.77	86.84 ± 0.89	81.03 ± 0.93	76.08 ± 0.92	89.59 ± 0.10	75.23 ± 2.24	$\textbf{89.76} \pm \textbf{0.21}$
10	MOOC	83.82 ± 0.30	83.42 ± 0.77	86.67 ± 0.24	$\textbf{91.58} \pm \textbf{0.74}$	81.76 ± 0.46	82.42 ± 0.71	82.76 ± 0.13	87.75 ± 0.42	66.38 ± 1.59	89.34 ± 0.12
р	Reddit	96.42 ± 0.13	95.87 ± 0.21	97.02 ± 0.04	97.30 ± 0.12	98.42 ± 0.01	94.63 ± 0.08	94.95 ± 0.08	98.70 ± 0.02	82.35 ± 4.03	$\textbf{98.88} \pm \textbf{0.01}$
Ra	Wikipedia	94.43 ± 0.28	91.31 ± 0.40	95.93 ± 0.19	97.71 ± 0.19	98.05 ± 0.03	97.03 ± 0.08	96.26 ± 0.04	98.49 ± 0.02	92.59 ± 0.70	$\textbf{98.72} \pm \textbf{0.03}$
_	UCI	78.78 ± 1.11	58.84 ± 2.54	77.41 ± 0.65	86.27 ± 1.49	90.27 ± 0.40	81.67 ± 1.01	89.26 ± 0.42	92.43 ± 0.20	48.58 ± 6.02	$\textbf{92.70} \pm \textbf{0.19}$
	Social Evo.	93.62 ± 0.36	90.20 ± 2.05	93.52 ± 0.05	93.21 ± 0.90	84.73 ± 0.20	94.63 ± 0.06	94.09 ± 0.03	95.30 ± 0.05	Timeout	$\textbf{95.36} \pm \textbf{0.04}$
	Avg. Rank	6.00	7.43	7.00	4.57	4.71	6.14	6.14	2.14	9.71	1.14
	LastFM	69.85 ± 1.70	68.14 ± 1.61	69.89 ± 0.41	67.01 ± 5.77	67.72 ± 0.20	63.15 ± 1.17	69.93 ± 0.17	69.86 ± 0.80	57.85 ± 3.67	$\textbf{70.59} \pm \textbf{0.57}$
	Enron	65.95 ± 1.27	62.20 ± 2.15	56.52 ± 0.84	64.21 ± 0.94	62.07 ± 0.72	67.56 ± 1.34	67.39 ± 1.33	66.07 ± 0.65	68.70 ± 1.82	$\textbf{68.98} \pm \textbf{1.00}$
Ę.	MOOC	65.37 ± 0.96	62.97 ± 2.05	74.94 ± 0.80	76.36 ± 2.91	71.18 ± 0.54	71.30 ± 1.21	72.15 ± 0.65	80.42 ± 0.72	58.06 ± 0.89	$\textbf{81.12} \pm \textbf{0.63}$
ĩ	Reddit	61.84 ± 0.44	60.35 ± 0.53	64.92 ± 0.08	65.24 ± 0.08	65.37 ± 0.12	61.85 ± 0.11	64.56 ± 0.26	64.80 ± 0.53	$\textbf{81.70} \pm \textbf{4.71}$	64.93 ± 0.89
pu	Wikipedia	61.66 ± 0.30	56.34 ± 0.67	78.40 ± 0.77	75.86 ± 0.50	59.00 ± 4.33	71.45 ± 2.23	82.76 ± 0.11	58.21 ± 8.78	$\textbf{91.12} \pm \textbf{0.13}$	67.92 ± 2.23
-	UCI	60.66 ± 1.82	51.50 ± 2.08	61.27 ± 0.78	62.07 ± 0.67	55.60 ± 1.22	65.87 ± 1.90	$\textbf{75.72} \pm \textbf{0.70}$	64.37 ± 0.98	51.68 ± 2.60	66.95 ± 2.22
	Social Evo.	88.98 ± 0.81	86.43 ± 1.48	92.37 ± 0.50	91.66 ± 2.14	83.84 ± 0.21	$\underline{95.50\pm0.31}$	93.88 ± 0.22	94.97 ± 0.36	Timeout	$\overline{\textbf{96.65}\pm\textbf{0.29}}$
	Avg. Rank	7.00	8.71	5.14	5.14	7.14	5.14	3.57	4.71	6.14	2.29

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. (2023) 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.

1010 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.

1017 F DYNAMIC NODE CLASSIFICATION

¹⁰¹⁹ 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.

- 1025 We follow Rossi et al. (2020); Xu et al. (2020); Yu et al. (2023) to conduct dynamic node classification by estimating the state of a node in a given interaction at a specific timestamp. A classification

1026 MLP is employed to map the node representations as well as the learned temporal patterns to the 1027 labels. AUC-ROC is used as the evaluation metric and we follow the dataset splits introduced in 1028 Yu et al. (2023) (70%15%/15% for training/validation/testing in chronological order) for node clas-1029 sification. Table 14 shows the node classification results on Wikipedia and Reddit (the only two 1030 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-1031 range temporal dependent datasets, therefore we do not include this part into the main body of the 1032 paper. Nonetheless, DyGMamba's great results on these datasets further prove its strength in CTDG 1033 modeling, regardless of the type of the dataset (whether long-range temporal dependent or not). 1034

Table	14:	AUC-ROC	of o	dvnamic	node	classification.
I GOIC		1100100	U 1 v	a j manne	noue	orabbiliteation.

Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	DyGFormer	CTAN	DyGMamba
Wikipedia Reddit	$\begin{array}{c} \textbf{88.10} \pm \textbf{1.57} \\ 59.53 \pm 3.18 \end{array}$	$\begin{array}{c} 87.41 \pm 1.94 \\ 63.12 \pm 0.51 \end{array}$	$\begin{array}{c} 83.42\pm2.92\\\textbf{69.31}\pm\textbf{2.18}\end{array}$	$\begin{array}{c} 85.51 \pm 3.28 \\ 63.21 \pm 3.00 \end{array}$	$\begin{array}{c} 84.59 \pm 1.16 \\ 65.22 \pm 0.79 \end{array}$	$\frac{79.03 \pm 1.18}{68.04 \pm 2.00}$	$\begin{array}{c} 85.60 \pm 1.73 \\ 64.42 \pm 1.15 \end{array}$	$\begin{array}{c} 86.35 \pm 2.19 \\ 67.67 \pm 1.39 \end{array}$	$\begin{array}{c} 87.38 \pm 0.14 \\ 67.29 \pm 0.15 \end{array}$	$\frac{87.44\pm0.82}{67.70\pm1.32}$
Avg. Rank	5.50	6.00	5.00	7.50	7.00	6.00	6.50	4.50	4.50	2.50

EFFICIENY ANALYSIS COMPLETE RESULTS G 1043

1044 We first provide the efficiency analysis results of all baselines in this section. We then provide a 1045 comparison of total training time among DyGFormer, CTAN and DyGMamba.

1047 G.1 EFFICIENCY STATISTICS FOR ALL BASELINES

1049 We provide the efficiency statistics for all baselines in Table 15.

1050 Table 15: Efficiency statistics for all baselines. EdgeBank is non-parameterized and not a machine 1051 learning model so we omit it here. # Params means number of parameters (MB). Time and Mem 1052 denote per epoch training time (min) and GPU memory (GB), respectively. The numbers in this 1053 table are the average results of five runs with different random seeds. 1054

Datasets	L	astFM		H	Enron		Ν	100C			UCI		
Models	# 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	
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	
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	
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	
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	
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	
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	
DyĜFormer	5.56	47.00	7.57	4.80	2.73	3.23	4.80	8.32	3.35	4.15	0.62	2.30	
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	
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	

Table 16: Comparison among DyGFormer, CTAN and DyGMamba on per epoch training time (Tep (min)), number of epochs until the best performance (# Epoch) and the total training time (T_{tot} 1067 (min)). $T_{tot} = T_{ep} \times \#$ Epoch. The numbers in this table are the average results of five runs with 1068 different random seeds. 1069

Datasets LastFM			Enron				MOOC			UCI		
Models	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	21.5
CTAN	3.33	635.00	2114.55	0.50	173.00	86.50	3.22	138.00	444.36	0.38	236.00	89.6
DyGMamba	28.45	11.80	335.71	2.05	33.00	67.65	4.88	38.00	185.44	0.60	28.00	16.8

1074 1075 1077

1078

1070 1071

1035 1036

1039 1040 1041

1046

1048

G.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 1079 training time in Table 16. Total training time computes the total amount of time a model requires to



Figure 6: Performance comparison among TGN, CAWN, DyGFormer, CTAN and DyGMamba on 1093 Enron, with an increasing number of encoded temporal neighbors. The metric is AP under random 1094 NSS on transductive link prediction. The time limit for training is 120 min. If a model fails to 1095 complete training within this limit, a cross \times is used to mark the data point. Dashed lines indicate 1096 that models start to exceed time limit as neighbor number increases. 1097

1099 reach its maximum performance, without considering the patience during training. We observe that 1100 CTAN requires much more epochs to converge, e.g., on LastFM it uses almost 54 times of epochs 1101 than DyGMamba to reach its best performance. 1102

1103 MODELING AN INCREASING NUMBER OF TEMPORAL NEIGHBORS WITH LIMITED G.3 1104 TOTAL TRAINING TIME

1106 To further show DyGMamba's superior efficiency against baseline methods, we do the following 1107 experiments. We train five best performing models (as shown in Table 1) on Enron with a gradually 1108 increasing number of temporal neighbors⁸ and report their performance. The number of sampled 1109 neighbors spans from 8, 16, 32, 64, 128 to 256 (Note that these numbers are different from the 1110 best hyperparameters reported in Yu et al. (2023)). 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 1111 minutes for the total training time. We let all the experiments finish the complete training process 1112 and note down the ones that exceed the time limit. In this way, we not only care about the per epoch 1113 training time, but also pay attention to how long it takes for models to converge. The experimental 1114 results are reported in Fig. 6. The points marked with crosses (\times) mean that the training process 1115 cannot finish within the time limit (although we still plot their corresponding performance). We find 1116 that only TGN and DyGMamba can successfully converge within the time limit when the number of 1117 considered neighbors increases to 256. DyGMamba can constantly achieve performance gain from 1118 modeling more temporal neighbors while TGN cannot. CAWN is extremely time consuming so it 1119 cannot finish training within the time limit even when it is asked to model 16 temporal neighbors. As 1120 for the methods designed for long-range temporal information propagation, DyGFormer and CTAN 1121 consume much longer total training time than DyGMamba. They fail to converge within 120 minutes 1122 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 1123 is not stable to model a large number of temporal neighbors. This also implies that increasing the 1124 amount of historical information will gradually make CTAN harder to converge, which might cause 1125 trouble in modeling long range temporal dependent datasets. 1126

1127 1128

1129

1080

1081 1082

1083

1090

1091

1098

1105

G.4 DYGMAMBA VS. DYGFORMER ON COMPLEXITY.

Sequence length is the key factor affecting the consumption of computational resources in DyG-1130 Former and DyGMamba. Following the computation of previous work Zhu et al. (2024), the com-1131

1132 1133

⁸For 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.

plexity of Transformer and Mamba in DyGFormer and DyGMamba can be written as

$$\Omega(\text{Transformer}) = 4(\rho/p)(4d)^2 + 2(\rho/p)^2(4d) = 64(\rho/p)d^2 + 8d(\rho/p)^2,$$
(8a)

$$\Omega(\text{Mamba}) = 3(\rho/p)(4d)d_{\text{SSM}} + (\rho/p)(4d)d_{\text{SSM}} = 16d_{\text{SSM}}d(\rho/p).$$
(8b)

This means that DyGMamba holds a computational complexity linear to ρ/p , while DyGFormer's complexity is quadratic to ρ/p . As a result, as the sequence length grows (either ρ increases or p decreases), DyGFormer is less scalable compared with DyGMamba.

1136 1137

1138

1139 1140

H DETAILS OF S4 AND MAMBA OPERATIONS

1145 Single-Input Single-Output. Given a sequence of vector elements as input, SISO means that the 1146 SSM processes each input dimension in parallel with the same set of parameters. For example, a 1147 sequence of d_2 -dimensional vectors will be split into d_2 1-dimensional sequences with the same 1148 sequence length. Each of them will be computed in parallel as in Eq. 2 with a shared set of SSM 1149 parameters. After computation, all these d_2 sequences will be rearranged back into a sequence of d_2 -1150 dimensional vectors. SISO fails to mix the information across dimensions of each vector. To address 1151 this, S4 and Mamba employs a mixing linear layer $f_{\min}(\cdot) : \mathbb{R}^{d_2} \to \mathbb{R}^{d_2}$ on each d_2 -dimensional 1152 vector to mix the information across d_2 dimensions. For more details, please refer to (Gu et al., 1153 2022b), (Smith et al., 2023) and (Gu & Dao, 2023). 1154

1155

1156 **SSM Function.** $\overline{\text{SSM}_{\bar{\mathbf{A}},\bar{\mathbf{B}},\mathbf{C}}(\cdot)}$ takes a matrix as input. The input matrix can be considered as a 1157 sequence of vector elements, where each row of the matrix corresponds to an element. The output of 1158 $\overline{\text{SSM}_{\bar{\mathbf{A}},\bar{\mathbf{B}},\mathbf{C}}(\cdot)}$ is also a matrix, where each row of the output matrix is the output of its corresponding 1159 input vector element. $\overline{\text{SSM}_{\bar{\mathbf{A}},\bar{\mathbf{B}},\mathbf{C}}(\cdot)}$ can be viewed as using S4 or Mamba to process a sequence of 1160 vectors in the SISO fashion, based on their parameters $\bar{\mathbf{A}}, \bar{\mathbf{B}}, \mathbf{C}$.

- 1161
- 1162 1163 1164

I 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 really 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.

1171 1172 1173

J MORE DETAILS REGARDING THE IMPACT OF PATCH SIZE

1174 More Explanations about Different Performance Trends. Patching decreases sequence length 1175 by applying linear transformation on a patch of node embeddings, giving a model much more pa-1176 rameters to tune (as indicated in Fig. 3b). Larger patch size mixes more sampled neighbors in each 1177 patch, introducing more training parameters while losing nuanced temporal details brought by the 1178 temporal order of the neighbors within each patch. For DyGFormer, the negative influence brought 1179 by mixing neighbors within patches is smaller than the positive influence brought by more train-1180 able parameters, so the performance constantly increases with a growing patch size in Fig. 3a. By 1181 contrast, for DyGMamba, the negative influence brought by mixing neighbors is much greater than 1182 the positive influence brought by more trainable parameters, so it shows better performance when 1183 the patch size is smaller. For Variant A, given an increasing patch size, it follows the trend of DyG-1184 Former when the patch size is below a threshold and shows degrading performance after that. This 1185 means that there is a trade-off between the lost temporal details and the additional parameters when 1186 we modify patch size. Also, by comparing Variant A and DyGMamba, we can tell that the differ-1187 ence in performance trend roots from the dynamic information selection module. Smaller patch size



¹²⁴¹ ⁹Although we have discussed in App. K that DyGMamba is not suitable to reason over DTDGs, we still benchmark our model on them to show its effectiveness.

inductive settings, respectively. All the results of DyGFormer are directly taken from the original paper Yu et al. (2023).

We find that DyGMamba outperforms DyGFormer on DTDGs in almost all cases. We also find 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, indicating the importance of modeing 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 ρ increases. For example, on Can. Parl, as shown in Table 17, 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 17: DyGMamba hyperparameter searching strategy on DTDG datasets. The best settings are marked as bold.

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

Table 18: AP of DyGFormer and DyGMamba on DTDGs under the transductive setting.

NSS	Datasets	DyGFormer	DyGMamba
Random	Flights Can. Parl. US Legis. UN Trade UN Vote Contact	$\begin{array}{c} 98.91 \pm 0.01 \\ 97.36 \pm 0.45 \\ 71.11 \pm 0.59 \\ 66.46 \pm 1.29 \\ 55.55 \pm 0.42 \\ 98.29 \pm 0.01 \end{array}$	$\begin{array}{c} 98.95 \pm 0.05 \\ 99.57 \pm 0.08 \\ 71.75 \pm 0.26 \\ 67.50 \pm 0.24 \\ 56.39 \pm 0.18 \\ 98.43 \pm 0.12 \end{array}$
Historical	Flights Can. Parl. US Legis. UN Trade UN Vote Contact	$\begin{array}{c} 66.59 \pm 0.49 \\ 97.00 \pm 0.31 \\ \textbf{85.30} \pm \textbf{3.88} \\ 64.41 \pm 1.40 \\ 60.84 \pm 1.58 \\ 97.57 \pm 0.06 \end{array}$	$\begin{array}{c} \textbf{67.80} \pm \textbf{2.17} \\ \textbf{99.77} \pm \textbf{0.12} \\ \textbf{82.15} \pm \textbf{1.02} \\ \textbf{65.10} \pm \textbf{0.02} \\ \textbf{61.07} \pm \textbf{1.39} \\ \textbf{97.61} \pm \textbf{0.04} \end{array}$
Inductive	Flights Can. Parl. US Legis. UN Trade UN Vote Contact	$\begin{array}{c} 70.92 \pm 1.78 \\ 95.44 \pm 0.57 \\ 81.25 \pm 3.62 \\ 55.79 \pm 1.02 \\ 51.91 \pm 0.84 \\ 94.75 \pm 0.28 \end{array}$	$\begin{array}{c} 73.79\pm 5.69\\ 98.32\pm 0.34\\ 81.67\pm 2.16\\ 58.89\pm 0.98\\ 52.24\pm 0.95\\ 95.43\pm 0.17\end{array}$

NSS	Datasets	DyGFormer	DyGMamba
Random	Flights Can. Parl. US Legis. UN Trade UN Vote Contact	$\begin{array}{c} 97.79 \pm 0.02 \\ 87.74 \pm 0.71 \\ 54.28 \pm 2.87 \\ 64.55 \pm 0.62 \\ 55.93 \pm 0.39 \\ 98.03 \pm 0.02 \end{array}$	$\begin{array}{c} 97.85 \pm 0.22 \\ 93.46 \pm 2.62 \\ 55.95 \pm 1.16 \\ 70.55 \pm 0.04 \\ 56.61 \pm 0.13 \\ 98.16 \pm 0.03 \end{array}$
Inductive	Flights Can. Parl. US Legis. UN Trade UN Vote Contact	$\begin{array}{c} 57.11 \pm 0.20 \\ 87.22 \pm 0.82 \\ 56.31 \pm 3.46 \\ 52.56 \pm 1.70 \\ 52.61 \pm 1.25 \\ 93.55 \pm 0.52 \end{array}$	$\begin{array}{c} 57.76 \pm 2.06 \\ 92.68 \pm 0.97 \\ 57.85 \pm 0.23 \\ 52.81 \pm 0.18 \\ 53.70 \pm 2.40 \\ 94.05 \pm 0.32 \end{array}$

Table 19: AP of DyGFormer and DyGMamba on DTDGs under the inductive setting.

Table 20: AUC-ROC of DyGFormer and DyGMamba on DTDGs under the transductive setting.

NSS	Datasets	DyGFormer	DyGMamba
	Flights	98.93 ± 0.01	$\textbf{98.98} \pm \textbf{0.05}$
Ξ	Can. Parl.	97.76 ± 0.41	$\textbf{99.69} \pm \textbf{0.06}$
qo	US Legis.	77.90 ± 0.58	$\textbf{79.03} \pm \textbf{0.26}$
an	UN Trade	70.20 ± 1.44	$\textbf{71.41} \pm \textbf{0.21}$
R	UN Vote	57.12 ± 0.62	$\textbf{58.48} \pm \textbf{0.12}$
	Contact	98.53 ± 0.01	$\textbf{98.68} \pm \textbf{0.02}$
	Flights	68.09 ± 0.43	68.98 ± 1.81
al	Can. Parl.	97.61 ± 0.40	$\textbf{99.82} \pm \textbf{0.10}$
Ľ	US Legis.	$\textbf{90.77} \pm \textbf{1.96}$	88.36 ± 1.78
sto	UN Trade	73.86 ± 1.13	$\textbf{74.10} \pm \textbf{2.02}$
Hi	UN Vote	64.27 ± 1.78	$\textbf{65.17} \pm \textbf{1.24}$
	Contact	97.17 ± 0.05	$\textbf{97.27} \pm \textbf{0.06}$
	Flights	69.53 ± 1.17	71.16 ± 3.24
ve	Can. Parl.	96.70 ± 0.59	$\textbf{99.56} \pm \textbf{0.21}$
cti	US Legis.	$\textbf{87.96} \pm \textbf{1.80}$	86.08 ± 2.27
qu	UN Trade	62.56 ± 1.51	$\textbf{67.60} \pm \textbf{0.64}$
In	UN Vote	53.37 ± 1.26	$\textbf{54.09} \pm \textbf{0.06}$
	Contact	95.01 ± 0.15	$\textbf{95.68} \pm \textbf{0.20}$

Table 21: AUC-ROC of DyGFormer and DyGMamba on DTDGs under the inductive setting.

NSS	Datasets	DyGFormer	DyGMamba
	Flights	97.80 ± 0.02	$\textbf{97.98} \pm \textbf{0.25}$
Ξ	Can. Parl.	89.33 ± 0.48	$\textbf{94.02} \pm \textbf{3.42}$
lop	US Legis.	53.21 ± 3.04	$\textbf{57.17} \pm \textbf{0.20}$
an	UN Trade	67.25 ± 1.05	$\textbf{68.26} \pm \textbf{0.26}$
R	UN Vote	56.73 ± 0.69	$\textbf{56.91} \pm \textbf{0.12}$
	Contact	98.30 ± 0.02	$\textbf{98.44} \pm \textbf{0.05}$
	Flights	56.05 ± 0.22	56.58 ± 2.12
ve	Can. Parl.	88.51 ± 0.73	$\textbf{92.37} \pm \textbf{0.18}$
cti	US Legis.	56.57 ± 3.22	$\textbf{57.91} \pm \textbf{3.41}$
np	UN Trade	57.28 ± 3.06	$\textbf{57.58} \pm \textbf{0.20}$
In	UN Vote	53.87 ± 2.01	$\textbf{54.83} \pm \textbf{2.17}$
	Contact	94.14 ± 0.26	$\textbf{94.35} \pm \textbf{0.29}$