DyG-Mamba: Continuous State Space Modeling on Dynamic Graphs

Dongyuan Li¹, Shiyin Tan², Ying Zhang³, Ming Jin⁴, Shirui Pan⁴, Manabu Okumura², Renhe Jiang^{1*}

¹The University of Tokyo, ²Institute of Science Tokyo, ³RIKEN Center for Advanced Intelligence Project, ⁴Griffith University lidy@csis.u-tokyo.ac.jp, tanshiyin@lr.pi.titech.ac.jp, ying.zhang@riken.jp, mingjinedu@gmail.com, s.pan@griffith.edu.au, oku@pi.titech.ac.jp, jiangrh@csis.u-tokyo.ac.jp

Abstract

Dynamic graph modeling aims to uncover evolutionary patterns in real-world systems, enabling accurate social recommendation and early detection of cancer cells. Inspired by the success of recent state space models in efficiently capturing long-term dependencies, we propose DyG-Mamba by translating dynamic graph modeling into a long-term sequence modeling problem. Specifically, inspired by Ebbinghaus' forgetting curve, we treat the irregular timespans between events as control signals, allowing DyG-Mamba to dynamically adjust the forgetting of historical information. This mechanism ensures effective usage of irregular timespans, thereby improving both model effectiveness and inductive capability. In addition, inspired by Ebbinghaus' review cycle, we redefine core parameters to ensure that DyG-Mamba selectively reviews historical information and filters out noisy inputs, further enhancing the model's robustness. Through exhaustive experiments on 12 datasets covering dynamic link prediction and node classification tasks, we show that DyG-Mamba achieves state-of-the-art performance on most datasets, while demonstrating significantly improved computational and memory efficiency. Code is available at [https://github.com/Clearloveyuan/DyG-Mamba].

1 Introduction

Dynamic graph modeling represents entities as nodes and timestamped relationships as edges, aiming to explore the underlying evolution patterns of real-world systems [1]. It has attracted great attention in various fields, *e.g.*, social networks [2], traffic systems [3], and recommender systems [4].

Despite the great success of current methods, there are still two limitations. *Firstly, existing methods lack the ability to effectively and efficiently track long-term temporal dependencies in dynamic graphs.* Specifically, RNN-based methods, *e.g.*, JODIE [2] and TGN [5], model temporal evolution through recurrent updates of node embeddings. Although theoretically capable of capturing long-term dependencies, they suffer from vanishing/exploding gradients in practice, limiting their effectiveness on long sequences. On the other hand, Transformer-based models, *e.g.*, DyGFormer [6] and SimpleDyG [7], address gradient issues through the self-attention mechanism but require prohibitive quadratic $\mathcal{O}(N^2)$ computational complexity for sequences of length N. Recent efficiency improvements through patching [6] or temporal convolutions [8] inevitably sacrifice temporal resolution, forcing an effectiveness-efficiency trade-off. Other recent methods, using multi-layer perceptions (MLP), *e.g.*, GraphMixer [9], FreeDyG [10], or graph neural networks (GNN), *e.g.*, TGAT [11], primarily focus on short-term dependencies, and their performance often decreases as the sequence

^{*}Corresponding author.

length increases [12]. *Secondly, existing methods lack robustness against noise*. Real-world dynamic graphs frequently contain various types of noisy events [13]. RNN-based and GNN-based methods are naturally susceptible to noise interference, leading to unstable performance [14, 15]. Although Transformers partially mitigate historical noise via self-attention, they remain susceptible to noisy data and cannot fully eliminate its impact [16]. How to filter out noisy history information more effectively and efficiently remains a challenge [17].

To address these issues, we propose DyG-Mamba, a novel timespan-informed continuous state space model (SSM), for dynamic graph modeling. Firstly, compared to Transformer-based methods that rely on large number of trainable parameters, DyG-Mamba employs only one trainable step size parameter Δt to capture forgetting laws of historical information, along with a small set of parameters in the encoder and decoder layers. Under the same GPU memory constraints, DyG-Mamba can directly process the entire long-term sequence without pooling, thereby preserving temporal details and effectively modeling long-term dependencies. Furthermore, inspired by Ebbinghaus' forgetting curve [18], which posits that forgetting is primarily correlated with timespans rather than content, we aim to equip DyG-Mamba with the same forgetting mechanism. Specifically, we design a monotonically increasing and learnable timespan function to redefine Δt , enabling the dynamic system to automatically learn how to compress historical information across different timespans, i.e., the model forgets historical information in a "fast-then-slow" pattern as the timespan increases, thereby enhancing both its effectiveness and inductiveness. Additionally, compared to Transformer's quadratic time complexity, DyG-Mamba adopts the same hardware-aware parallel scan optimization as Mamba [19], enabling it to efficiently capture long-term dependencies with linear time complexity. Secondly, inspired by Ebbinghaus' review cycle that periodic review can counteract forgetting [20], to further enhance robustness, we redefine SSM's core parameters B and C to be input-dependent and add spectral norm constraints to ensure Lipschitz continuity. This strategy enables DyG-Mamba to selectively review historical information and thus remain robust against noise. Main contributions:

- To the best of our knowledge, we are the first to introduce SSMs for continuous-time dynamic graph modeling. By redefining the core SSM parameters, DyG-Mamba achieves high efficiency and effectiveness in capturing long-term temporal dependencies.
- Inspired by both the forgetting curve and the review cycle that counters it, we propose a timespaninformed continuous SSM that adopts timespans to control system forgetting while incorporating
 input-dependent parameterization. This design improves DyG-Mamba's capability to model longterm sequences with irregular timespans, enhancing its effectiveness, inductiveness and robustness.
- Extensive experiments on 12 benchmarks show that DyG-Mamba achieves state-of-the-art performance with superior effectiveness and robustness.

Table 1: Comparison of continuous-time dynamic graph baselines from six aspects. With a batch size of 200 and a sequence length of 512, a model is considered time and memory efficient if the running time and memory usage are less than GraphMixer, *i.e.*, running time 250 seconds and memory usage 30,000 MB. Adding 50% noisy temporal edges, the performance drop < 10% indicates robustness.

	JODIE	DyRep	TGN	CAWN	TGAT	EdgeBank	GraphMixer	TCL	DyGFormer	r DyG-Mamba
Long-Term Capability	×	×	×	×	×	×	×	×	~	~
Time Efficient	~	×	×	×	×	~	✓	~	×	✓
Memory Efficient	~	×	×	×	×	~	✓	×	×	✓
Irregular timespan Supportive	×	×	×	×	×	×	×	×	×	✓
Inductive Supportive	~	~	~	~	~	×	✓	~	~	✓
Noise Robust	×	×	~	×	×	~	×	×	✓	✓

2 Related Work

Dynamic Graph Modeling. Discrete-time methods segment the dynamic graph into snapshots at a predetermined time granularity, then employ a GNN (snapshot encoder) with a recurrent module (dynamic tracker) to learn node embedding [21–24]. However, fixing the time granularity in advance ignores the fine-grained temporal order within each snapshot. In contrast, continuous-time methods directly use timestamps to learn node embedding. Based on their neural architectures, they can be categorized into four types, including RNN-based methods, *e.g.*, JODIE [2], GNN-based methods,

e.g., DySAT [22], MLP-based methods, e.g., FreeDyG [10], and Transformer-based methods, e.g., SimpleDyG [7]. Additional techniques, such as ordinary differential equations [25, 26], random walks [27], and temporal point processes [28], have also been introduced to capture continuous temporal information. Table 1 provides a detailed comparison between DyG-Mamba and SOTAs, including JODIE [2], DyRep [29], TGN [5], CAWN [11], TGAT [11], EdgeBank [30], GraphMixer [9], TCL [23], and DyGFormer [6] from the following angles: if the method can effectively handle unseen nodes during training (i.e., inductive), capture long-term dependencies with both time and memory efficiency, exhibit robustness against noise, and effectively leverage irregular timespans.

State Space Models. SSMs have attracted great attention for long sequence modeling [31, 32]. Mamba [19] designs a data-dependent selection mechanism with parallel scan optimization, achieving SOTA performance on many fields [33]. Graph Mamba [34, 35] applies SSMs to static graphs for embedding learning. DG-Mamba [17] and GraphSSM [36] extend Mamba to discrete-time dynamic graphs by modeling snapshot sequences with fixed time intervals. And STG-Mamba [37] adopts Mamba layers on spatial-temporal graphs. However, these methods are not applicable to continuous-time dynamic graphs with irregular timestamps, and thus are not directly comparable to our setting. PIVEM [38] learns dynamic node embeddings by approximating temporal evolution through piecewise linear interpolation, based on a latent distance model with piecewise constant and node-specific velocities. It can be viewed as a special case of a first-order SSM, where the hidden state corresponds to node velocity and evolves linearly over time. In contrast, our method generalizes this idea by introducing learnable memory decay and input-adaptive updates, allowing it to better capture irregular temporal dynamics and model more complex patterns in dynamic graphs.

3 Preliminary

Dynamic Graph Modeling. Dynamic graphs can be modeled as a sequence of non-decreasing chronological interactions $\mathcal{G} = \{(u_1, v_1, t_1), \dots, (u_\tau, v_\tau, \tau)\}$ with $0 \le t_1 \le \tau$, where $u_i, v_i \in \mathcal{V}$ denote the source and destination nodes of the i-th link and \mathcal{V} denote all nodes. Each node is associated with a node feature $x \in \mathbb{R}^{d_N}$ and each interaction has a link feature $e^t \in \mathbb{R}^{d_E}$, where d_N and d_E denote dimensionality. Given the source node u, destination node v, timestamp t, and all their historical interactions before t, dynamic graph modeling aims to learn time-aware node embedding for them. We validate the learned node embedding via two common tasks: (i) dynamic link prediction, which predicts whether two nodes are connected in future; and (ii) dynamic node classification, which infers the class of nodes.

Continuous SSMs. They define a linear mapping from t-th input $u(t) \in \mathbb{R}^{1 \times d}$ to output $y(t) \in \mathbb{R}^d$ via a hidden state variable $h(t) \in \mathbb{R}^{m \times d}$, formulated by:

$$h'(t) = Ah(t) + Bu(t), \tag{1}$$

$$y(t) = Ch(t) + Du(t), (2)$$

where $A \in \mathbb{R}^{m \times m}$, $B \in \mathbb{R}^{m \times 1}$, $C \in \mathbb{R}^{1 \times m}$ are trainable parameters, and D = 0 since Du(t) can be viewed as a skip connection. Eq.(1,2) could be discretized for controllable optimization via the zero-order hold (ZOH), formulated by:

$$h_t = \overline{A}h_{t-1} + \overline{B}u_t, \tag{3}$$

$$y_t = \overline{C}h_t, \tag{4}$$

where $\overline{A} = \exp(\Delta t A)$, $\overline{B} = (\Delta t A)^{-1} (\overline{A} - I)(\Delta t B)$, $\overline{C} = C$, and Δt is predefined step size.

4 Methodology

The overview of DyG-Mamba is shown in Figure 1. First, in Section 4.1, we introduce dynamic graph encoding and encoding alignment. Then, in Section 4.2, we introduce two main limitations of current SSMs, and DyG-Mamba can alleviate these issues by redefining four core parameters of SSMs. Finally, in Section 4.3, we apply DyG-Mamba on downstream tasks and show its complexity.

4.1 Dynamic Graph Encoding

In Figure 1, we first extract the first-hop interaction sequence S_u^{τ} of node u before timestamp τ from dynamic graph, where $S_u^{\tau} = \{(u, k_1, t_1), \dots, (u, k_{|u|}, t_{|u|})\}$ with |u| denoting the sequence length.

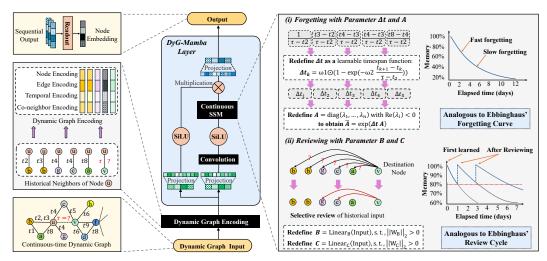


Figure 1: Overview of our proposed DyG-Mamba with four redefined core parameters Δ , A, B and C. Pseudocodes are in Appendix C.

Node and Edge Encoding. We directly adopt the node and edge features provided by datasets as node encoding $X_{u,V}^{\tau} \in \mathbb{R}^{|u| \times d_N}$ and edge encoding $X_{u,E}^{\tau} \in \mathbb{R}^{|u| \times d_E}$ for S_u^{τ} , respectively. If the graph is non-attributed, we simply set the node or edge encoding to zero vectors.

Absolute Temporal Encoding. We encode the absolute timespans between timestamp t_j and the final prediction timestamp τ by using an encoding function $\cos(\omega(\tau-t_j))$ to obtain the absolute temporal encoding $\boldsymbol{X}_{u,T}^{\tau} \in \mathbb{R}^{|u| \times d_T}$, where $\omega = \{\alpha^{-(i-1)/\beta}\}_{i=1}^{d_T}$ with α and β as trainable parameters. Following [9], we keep ω constant during training to facilitate easier model optimization.

Co-occurrence Frequency Encoding. Two nodes that frequently interact with the same neighbors tend to have similar embeddings. Thus, we capture this feature by adopting co-occurrence frequency encoding. Formally, let the neighbors of u and v be $S_u = \{a,b\}$ and $S_v = \{b,b,c,a\}$, the co-occurrence features of u could be denoted by $C_u^{\tau} = [[1,1],[1,2]]$, where [1,1] denotes the occurrence frequency of a in S_a and S_b , respectively. Then, we define a function $f(\cdot): \mathbb{R}^1 \to \mathbb{R}^{d_C}$ to encode the co-occurrence features by:

$$X_{u,C}^{\tau} = (f(C_u^{\tau}[:,0]) + f(C_u^{\tau}[:,1]))W_C + b_C,$$
 (5)

where $X_{u,C}^{\tau} \in \mathbb{R}^{|u| \times d_C}$ with d_C denotes dimensionality, W_C and b_C are trainable parameters. We implement $f(\cdot) : \mathbb{R}^1 \to \mathbb{R}^{d_C}$ by two-layer perception with ReLU activation.

Encoding Alignment. We align the above-mentioned encoding to the same dimension d:

$$Z_{u,*}^{\tau} = X_{u,*}^{\tau} W_* + b_*, \text{ where } * \in \{N, E, T, C\},$$
 (6)

where $W_* \in \mathbb{R}^{d_* \times d}$ and $b_* \in \mathbb{R}^d$ are trainable parameters. Finally, we concatenate aligned encoding for S_u^{τ} as $Z_u^{\tau} = Z_{u,N}^{\tau} \| Z_{u,E}^{\tau} \| Z_{u,T}^{\tau} \| Z_{u,C}^{\tau}$ and $Z_u^{\tau} \in \mathbb{R}^{|u| \times 4d}$.

4.2 DyG-Mamba: Dynamic Graph Mamba

4.2.1 Rethinking SSM on Dynamic Graph

To learn node embedding of u, SSM first encodes \mathbf{Z}_u^{τ} using a linear layer followed by a 1D convolution layer and SiLU activation function, which could be formulated by

$$M_u^{\tau} = \text{SiLU}\left(\text{Conv1D}\left(\text{Linear}\left(\mathbf{Z}_u^{\tau}\right)\right)\right) \in \mathbb{R}^{|u| \times 8d}.$$
 (7)

Then, SSM initializes four core trainable parameters: $A \in \mathbb{R}^{|u| \times 8d \times 8d}$ governs state transition, $B \in \mathbb{R}^{|u| \times 8d}$ and $C \in \mathbb{R}^{|u| \times 8d}$ governs input/output projections, and Δt controls system's update step size [19]. And the k-th output of the SSM is given by:

$$\boldsymbol{h}_k = \overline{\boldsymbol{A}}_k \boldsymbol{h}_{k-1} + \overline{\boldsymbol{B}}_k \boldsymbol{m}_k, \tag{8}$$

$$\widehat{\boldsymbol{m}}_{k}^{\tau} = \overline{\boldsymbol{C}}_{k} \boldsymbol{h}_{k}, \tag{9}$$

where m_k represents the k-th input, $\overline{A}_k = \exp(\Delta t_k A_k)$, $\overline{B}_k = (\Delta t_k A_k)^{-1} (\overline{A}_k - I)(\Delta t_k B_k)$, $\overline{C}_k = C_k$, and h_k denotes node's k-th hidden state representation.

Finally, SSM adopts skip connection to avoid gradient vanishing and generates the sequential output:

$$\widehat{Z}_{u,\text{out}}^{\tau} = (\widehat{M}_{u}^{\tau} \odot \text{SiLU}(\text{Linear}(Z_{u}^{\tau}))) W_{\text{out}} + b_{\text{out}}, \tag{10}$$

where \odot denotes element-wise product, $W_{\text{out}} \in \mathbb{R}^{8d \times 4d}$ and $b_{\text{out}} \in \mathbb{R}^{4d}$ are trainable parameters.

As shown in Eq.(8,9), SSM contains three core parameters, \overline{A}_k , \overline{B}_k and \overline{C}_k , which determine the effectiveness for long-term sequence modeling. Specifically, (i) \overline{A}_k controls the forgetting of historical information, determined by Δt_k and A_k . Existing SSMs, such as Hippo [39] and Mamba, typically initialize A_k randomly and either fix the step size Δt as a constant or adopt a data-dependent strategy to set $\Delta t = \text{SiLU}(\text{Linear}(Z_u^\tau))$. However, these SSMs do not account for the crucial role of irregular timespans in real-world sequential input, leading to suboptimal performance. Moreover, directly tying the input Z_u^τ to Δt further weakens SSM's effectiveness and inductiveness, as it will encounter a large variety of unseen input during testing. (ii) Existing SSMs struggle to effectively filter out noisy historical information. Although Mamba initializes B and C as data-dependent parameters, allowing \overline{B}_k and \overline{C}_k to selectively copy important past information, input noise can still affect their initialization, thereby weakening their robustness [17]. To address these issues, we propose DyG-Mamba, a timespan-informed continuous SSM designed for dynamic graph modeling.

4.2.2 Timespan-Informed Continuous SSM

To better utilize irregular timespans and enhance the effectiveness and inductiveness of SSMs, we first redefine Δt and A. Ebbinghaus's forgetting curve describes how memory retention decreases exponentially over time and can be formulated as $R = \exp(-t/S)$ [40, 41], where R denotes memory retention, t is the time interval, and S is a decay constant. Inspired by this formulation, we reinterpret R as a timespan-dependent decay coefficient, enabling the model to apply temporal decay to historical states proportionally to the elapsed timespan. Accordingly, we design DyG-Mamba with a similar exponential forgetting mechanism, where the core parameter \overline{A}_k , which governs the degree of forgetting, decays exponentially as the k-th timespan $(t_{k+1} - t_k)$ increases. Since \overline{A}_k is jointly determined by both Δt and A, this mechanism is realized by redefining these two variables.

Redefining Parameter Δt . To establish the forgetting curve relationship between timespans and \overline{A}_k , we first define the connection between timespans and the step size Δt . Since Δt could directly influence \overline{A}_k . Specifically, we define a monotonically increasing, learnable timespan function to redefine the step size parameter Δt as

$$\Delta t_k = \mathbf{w}_1 \odot \left(\mathbf{1} - \exp\left(-\mathbf{w}_2 \odot \frac{t_{k+1} - t_k}{\tau - t_1} \right) \right), \tag{11}$$

where Δt_k denotes the k-th step size, $(t_{k+1}-t_k)$ denotes the k-th timespan, $\boldsymbol{w}_1 \in \mathbb{R}^{8d}$ and $\boldsymbol{w}_2 \in \mathbb{R}^{8d}$ are trainable vectors designed to capture fine-grained timespan features, with each element constrained to be positive. 1 is an all-ones vector, \odot denotes element-wise multiplication, and τ and t_1 are the last and first appearing timestamps, respectively.

This redefinition of Δt establishes a direct relationship between timespan length and step size scaling. Its monotonically increasing property ensures that longer timespans induce stronger decay in $\overline{A}_k = \exp(\Delta t_k A_k)$. Consequently, careful initialization of A is required to maintain a balance between effective forgetting and numerical stability.

Redefining Parameter A. We consider two factors when redefining the initialization of A. (i) A should maintain the forgetting curve relationship, i.e., $\exp(\Delta t_k A_k)$ should diminish exponentially as Δt_k increases. (ii) A determines the stability of recurrent updating in long-term sequence modeling [42], i.e., it should prevent gradient vanishing or explosion over time. To satisfy these two requirements, we initialize A as a diagonal matrix whose eigenvalues all have negative real parts. Theorem 4.1 confirms that this initialization strategy satisfies both conditions.

Theorem 4.1. Let $A_k = \operatorname{diag}(\lambda_1, \ldots, \lambda_n)$, where the real parts of the eigenvalues satisfy $\operatorname{Re}(\lambda_i) < 0$. For any timespan Δt_k , we have $\overline{A}_k = \operatorname{diag}(e^{\lambda_1 \Delta t_{k,1}}, \ldots, e^{\lambda_n \Delta t_{k,n}})$ and $\overline{B}_k = \operatorname{diag}(\lambda_1^{-1}(e^{\lambda_1 \Delta t_{k,1}} - 1)B_{k,1}, \ldots, \lambda_n^{-1}(e^{\lambda_n \Delta t_{k,n}} - 1)B_{k,n})$, where $\Delta t_{k,i}$ and $B_{k,i}$ are the i-th elements of Δt_k and B_k . The i-th coordinate of h_k is denoted as $h_{k,i} = e^{\lambda_i \Delta t_{k,i}} h_{k-1,i} + \lambda_i^{-1}(e^{\lambda_i \Delta t_{k,i}} - 1)B_{k,i} m_{k,i}$.

(i) Theorem 4.1 guarantees the forgetting curve relationship. When $\Delta t_{k,i}$ is sufficiently small, then $e^{\lambda_i \, \Delta t_{k,i}} \approx 1$, i.e., $h_{k,i} \approx h_{k-1,i}$. This indicates that for small timespans, the model retains historical states and disregards the current input. On the other hand, as the timespan increases, $e^{\lambda_i \, \Delta t_{k,i}}$ gradually approaches 0 at a decreasing rate over time, causing the system to forget previous states and place greater emphasis on current input $m_{k,i}$. (ii) Since $\mathrm{Re}(\lambda_i) < 0$, the term $|e^{\lambda_i \, \Delta t_{k,i}}|$ is bounded by 1 and decreases as $\Delta t_{k,i}$ increases. This prevents the hidden state from diverging over extended sequences and mitigates gradient explosion, ensuring the stability of recurrent updates.

Traditional SSMs struggle to effectively filter out noise from the input sequence. To solve this issue, we redefine \boldsymbol{B} and \boldsymbol{C} as input-dependent parameters and introduce spectral norm constraints to enhance robustness. Specifically, Ebbinghaus' review cycle indicates that periodic review of previously learned information helps counteract memory decay [20]. Inspired by this, we design DyG-Mamba to continuously review important node while forgetting irrelevant or noisy inputs.

Redefining Parameter B **and** C**.** To align the SSM with the review cycle, we first define B and C as input-dependent parameters, e.g., $B = \text{Linear}_B(M_n^{\tau})$. Then we can filter out noise by Theorem 4.2.

Theorem 4.2. Let B and C be input-dependent parameters and m_k denote the k-th input in M_u^{τ} . Update process for SSMs, as shown in Eq.(8), could be further decomposed as

$$\widehat{\boldsymbol{m}}_{k}^{\tau} = \overline{\boldsymbol{C}}_{k} \prod_{i=0}^{k-2} \overline{\boldsymbol{A}}_{k-i} \overline{\boldsymbol{B}}_{1} \boldsymbol{m}_{1} + \dots + \overline{\boldsymbol{C}}_{k} \overline{\boldsymbol{B}}_{k} \boldsymbol{m}_{k},$$

$$= e^{(\sum_{i=0}^{k-2} \Delta t_{k-i} \boldsymbol{A}_{k-i})} \overline{\boldsymbol{C}}_{k} \overline{\boldsymbol{B}}_{1} \boldsymbol{m}_{1} + \dots + \overline{\boldsymbol{C}}_{k} \overline{\boldsymbol{B}}_{k} \boldsymbol{m}_{k}.$$
(12)

According to Theorem 4.2, \overline{C}_k can be interpreted as the query corresponding to the k-th input, while \overline{B}_k and m_k serve as the key and value, respectively. Thus, the product $\overline{C}_k \overline{B}_j$ represents the importance of the j-th historical input m_j to the k-th input. While Theorem 4.2 enables automatic filtering of irrelevant and noisy historical inputs, the construction of B and C using only linear layers makes them susceptible to input noise. To address this issue, we introduce spectral norm constraints on the initialization of B and C to achieve dual objectives, formulated by

$$B = W_{\rm B} M_u^{\tau} + b_{\rm B}, \quad C = W_{\rm C} M_u^{\tau} + b_{\rm C},$$
s.t. $||W_{\rm B}||_2 \le 1, \quad ||W_{\rm C}||_2 \le 1,$ (13)

where W_B and W_C are weight matrices, and $\|\cdot\|_2$ is the spectral norm, *i.e.*, the largest singular value of $W_{B/C}$. The spectral norm constraints guarantee Lipschitz continuity, ensuring that B and C remain stable under input perturbations. Overall, DyG-Mamba is robust to input noise, preventing irrelevant samples from being erroneously reinforced. This robustness is guaranteed by Theorem 4.3.

Theorem 4.3. Given $\|\mathbf{W}_B\|_2 \le 1$, $\|\mathbf{W}_C\|_2 \le 1$, $\gamma = \max_i \operatorname{Re}(\lambda_i(\mathbf{A})) < 0$, and T as the total sequence duration, the output perturbation satisfies:

$$\|\Delta\widehat{\boldsymbol{m}}_{k}^{\tau}\| \leq \frac{1}{|\gamma|} \left(1 - e^{\gamma T}\right) \|\Delta \boldsymbol{M}_{u}^{\tau}\|. \tag{14}$$

4.3 DyG-Mamba for Downstream Tasks

For dynamic link prediction, we first process the first-hop interaction sequences of source node u and destination node v through two independent DyG-Mamba models. Based on Eq.(7-10), two models generate sequential output embeddings $\widehat{Z}_{u,\text{out}}^{\tau}$ and $\widehat{Z}_{v,\text{out}}^{\tau}$, respectively. Then, we adopt readout function, *i.e.*, MEAN pooling, to obtain their node embedding, defined as $\widehat{z}_{*}^{\tau} = \text{MEAN}(\widehat{Z}_{*,\text{out}}^{\tau})$ with $* \in \{u,v\}$. Finally, we concatenate two node embedding and adopt an MLP for link prediction:

$$\hat{y} = \text{Signoid}(\text{Linear}(\text{ReLU}(\text{Linear}(\hat{z}_u^{\tau} || \hat{z}_v^{\tau}))). \tag{15}$$

We adopt binary cross-entropy loss for optimization

$$\mathcal{L}_{LP} = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \left[y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i) \right], \tag{16}$$

where $|\mathcal{B}|$ denotes the batch size containing both positive and negative samples, and \hat{y}_i represent the *i*-th ground-truth and predicted label, respectively.

For dynamic node classification, we use one DyG-Mamba model and discard co-occurrence encoding, while keeping other components identical to the link prediction setup.

Computational Complexity. Given batch size b, feature dimension d, and sequence length L. DyG-Mamba achieves linear memory and time complexity of $\mathcal{O}(bdL)$, while DyGFormer has quadratic complexity of $O(bdL^2)$. This highlights the efficiency of DyG-Mamba. Details in Appendix B.

Table 2: **Transductive**: AP for dynamic link prediction with random (*rnd*), historical (*hist*), and inductive (*ind*) negative edge sampling. **bold** and underlined emphasize best and 2nd-best results.

	Datasets	JODIE	DyRep	TGN	CAWN	TGAT	EdgeBank	GraphMixer	TCL	DyGFormer	DyG-Mamba
	Wikipedia	96.50±0.14	94.86±0.06	98.45±0.06	98.76±0.03	96.94±0.06	90.37±0.00	97.25 ± 0.03	96.47±0.16	99.03 ± 0.02	99.06±0.01
	Reddit	98.31±0.14	98.22 ± 0.04	98.63 ± 0.06	99.11 ± 0.01	98.52 ± 0.02	94.86 ± 0.00	97.31 ± 0.01	97.53 ± 0.02	99.22 ± 0.01	99.25 ± 0.00
	MOOC	80.23±2.44	81.97 ± 0.49	89.15 ± 1.60	80.15 ± 0.25	85.84 ± 0.15	57.97 ± 0.00	82.78 ± 0.15	82.38 ± 0.24	87.52 ± 0.49	90.17 ± 0.19
	LastFM	70.85 ± 2.13			86.99 ± 0.06			75.61 ± 0.24	67.27 ± 2.16	93.00 ± 0.12	94.22 ± 0.04
	Enron		82.38 ± 3.36		89.56 ± 0.09			82.25 ± 0.16	79.70 ± 0.71	92.47 ± 0.12	93.22 ± 0.03
	Social Evo.						74.95 ± 0.00		93.13 ± 0.16	94.73 ± 0.01	94.75 ± 0.01
rnd	UCI	89.43±1.09	65.14 ± 2.30		95.18 ± 0.06			93.25 ± 0.57	89.57 ± 1.63	95.79 ± 0.17	96.79 ± 0.08
77166	Can. Parl.		66.54 ± 2.76				64.55 ± 0.00		68.67 ± 2.67	97.36 ± 0.45	98.37 ± 0.07
	US Legis.	75.05 ± 1.52			70.58 ± 0.48			70.74 ± 1.02	69.59 ± 0.48	71.11 ± 0.59	74.11 ± 2.32
	UN Trade		63.21 ± 0.93				60.41 ± 0.00		62.21 ± 0.03	66.46 ± 1.29	68.55 ± 0.16
	UN Vote							52.11 ± 0.16	51.90 ± 0.30	55.55 ± 0.42	65.69 ± 1.10
	Contact						92.58 ± 0.00		92.44 ± 0.12	98.29 ± 0.01	98.37 ± 0.01
	Avg. Rank	6.08	6.00	4.42	8.42	6.33	6.92	4.92	6.58	<u>2.92</u>	2.42
	Wikipedia	83.01±0.66	79.93 ± 0.56	86.86 ± 0.33	71.21 ± 1.67	87.38 ± 0.22	73.35 ± 0.00	90.90 ± 0.10	89.05±0.39	82.23±2.54	82.12 ± 1.22
	Reddit	80.03±0.36	79.83 ± 0.31	81.22 ± 0.61	80.82 ± 0.45	79.55 ± 0.20	73.59 ± 0.00	78.44 ± 0.18	77.14 ± 0.16	81.57 ± 0.67	81.16 ± 0.11
	MOOC	78.94±1.25	75.60 ± 1.12	87.06 ± 1.93	74.05 ± 0.95	82.19 ± 0.62	60.71 ± 0.00	77.77 ± 0.92	77.06 ± 0.41	85.85 ± 0.66	87.33 ± 1.46
	LastFM	74.35±3.81	74.92 ± 2.46	76.87 ± 4.64	69.86 ± 0.43	71.59 ± 0.24	73.03 ± 0.00	72.47 ± 0.49	59.30 ± 2.31	81.57 ± 0.48	84.09 ± 0.44
	Enron	69.85±2.70	71.19 ± 2.76	73.91 ± 1.76	64.73 ± 0.36	64.07 ± 1.05	76.53 ± 0.00	77.98 ± 0.92	70.66 ± 0.39	75.63 ± 0.73	77.41 ± 1.13
	Social Evo.						80.57 ± 0.00		94.74 ± 0.31	97.38 ± 0.14	96.59 ± 0.28
hist	UCI	75.24±5.80	55.10 ± 3.14	80.43 ± 2.12	65.30 ± 0.43	68.27 ± 1.37	65.50 ± 0.00	84.11 ± 1.35	80.25 ± 2.74	82.17 ± 0.82	82.95 ± 2.24
11101	Can. Parl.							74.34 ± 0.87	65.93 ± 3.00	97.00 ± 0.31	97.22 ± 0.29
	US Legis.		86.88 ± 2.25				63.22 ± 0.00		80.53 ± 3.95	85.30 ± 3.88	88.83 ± 0.34
	UN Trade	61.39±1.83					81.32 ± 0.00		55.90 ± 1.17	64.41 ± 1.40	65.19 ± 0.19
	UN Vote							51.20 ± 1.60	52.30 ± 2.35	60.84 ± 1.58	59.51 ± 3.08
	Contact	95.31±2.13					88.81 ± 0.00		93.86 ± 0.21	97.57 ± 0.06	97.80 ± 0.14
	Avg. Rank	6.08	6.00	4.42	8.42	6.33	6.92	4.92	6.58	3.00	2.33
	Wikipedia	75.65±0.79	70.21 ± 1.58	85.62 ± 0.44	74.06 ± 2.62	87.00±0.16	80.63 ± 0.00	88.59 ± 0.17	86.76 ± 0.72	78.29 ± 5.38	84.64±0.77
	Reddit		86.30 ± 0.26				85.48 ± 0.00	85.26 ± 0.11	87.45 ± 0.29	91.11 ± 0.40	91.89 ± 0.42
	MOOC		61.66 ± 0.95					74.27 ± 0.92	74.65 ± 0.54	81.24 ± 0.69	81.15 ± 1.25
	LastFM	62.67±4.49	64.41 ± 2.70	65.95 ± 5.98	67.48 ± 0.77	71.13 ± 0.17	75.49 ± 0.00	68.12 ± 0.33	58.21 ± 0.89	73.97 ± 0.50	74.76 ± 0.40
	Enron		67.79 ± 1.53					75.01 ± 0.79	71.29 ± 0.32	77.41 ± 0.89	79.90 ± 0.90
	Social Evo.						83.69 ± 0.00		94.90 ± 0.36	97.68 ± 0.10	96.91 ± 0.24
ind	UCI						57.43 ± 0.00	80.10 ± 0.51	76.01 ± 1.11	72.25 ± 1.71	73.71 ± 3.88
inci	Can. Parl.	48.42±0.66	58.61 ± 0.86	65.34 ± 2.87	67.75 ± 1.00	68.82 ± 1.21	62.16 ± 0.00	69.48 ± 0.63	65.85 ± 1.75	95.44 ± 0.57	96.58 ± 0.79
	US Legis.	50.27±5.13			65.81 ± 8.52			79.63 ± 0.84	78.15 ± 3.34	81.25±3.62	85.03 ± 0.69
	UN Trade	60.42±1.48			62.54 ± 0.67			60.15 ± 1.29	61.06 ± 1.74	55.79 ± 1.02	61.88 ± 1.46
	UN Vote						66.30 ± 0.00		50.62 ± 0.82	51.91 ± 0.84	57.63 ± 1.15
	Contact	93.43±1.78						90.87 ± 0.35	91.35 ± 0.21	94.75 ± 0.28	94.57 ± 0.22
	Avg. Rank	7.33	7.25	5.17	6.17	5.25	6.75	5.42	5.58	<u>3.75</u>	2.33

5 Experiments

5.1 Experimental Setup

Datasets and Baselines. We evaluate performance on 12 datasets, each split into 70%/15%/15% for training, validation and testing. Details in Appendix D.1. We select nine SOTA baselines for comparison, *e.g.*, four RNN-based methods: JODIE [2], DyRep [29], TGN [5] and CAWN [11], a GNN-based method: TGAT [11], a memory-based method: EdgeBank [30], a MLP-based method: GraphMixer [9], and two Transformer-based methods: TCL [23] and DyGFormer [6].

Implementation Details. For a fair comparison, we use DyGLib [6] to reproduce all baselines via the same training and inference pipeline. We set the same input length for DyGFormer and DyG-Mamba to fairly compare the long-term dynamic graph modeling ability. We train each model for 100 epochs and select the best-performing checkpoint for testing. We repeat each experiment 10 times with different random seeds and report the mean and standard derivation. Details in Appendix D.2.

Evaluation Details. We evaluate baselines in transductive and inductive settings, where the former predicts future links among nodes seen during training, and the latter focus on unseen nodes [6]. For negative sampling, we follow [30] and adopt random (*rnd*), historical (*hist*), and inductive (*ind*) strategies (see Appendix D.3). Metrics are Average Precision (AP) and AUC-ROC.

Table 3: **Inductive**: AP for *dynamic link prediction* with random negative edge sampling strategies. The notations are the same as Table 2.

Datasets	JODIE	DyRep	TGN	CAWN	TGAT	TCL	GraphMixer	DyGFormer	DyG-Mamba
Wikipedia	94.82±0.20	92.43±0.37	97.83±0.04	98.24±0.03	96.22±0.07	96.65±0.02	96.22±0.17	98.59±0.03	98.66±0.02
Reddit	96.50±0.13	96.09 ± 0.11	97.50 ± 0.07	98.62 ± 0.01	97.09 ± 0.04	95.26 ± 0.02	94.09 ± 0.07	98.84 ± 0.02	98.91 ± 0.01
MOOC	79.63±1.92	81.07 ± 0.44	89.04 ± 1.17	81.42 ± 0.24	85.50 ± 0.19	81.41 ± 0.21	80.60 ± 0.22	86.96 ± 0.43	89.98 ± 0.04
LastFM	81.61±3.82	83.02 ± 1.48	81.45 ± 4.29	89.42 ± 0.07	78.63 ± 0.31	82.11 ± 0.42	73.53 ± 1.66	94.23 ± 0.09	95.16 ± 0.05
Enron	80.72±1.39	74.55 ± 3.95	77.94 ± 1.02	86.35 ± 0.51	67.05 ± 1.51	75.88 ± 0.48	76.14 ± 0.79	89.76±0.34	90.97 ± 0.01
Social Evo.	91.96±0.48	90.04 ± 0.47	90.77 ± 0.86	79.94 ± 0.18	91.41 ± 0.16	91.86 ± 0.06	91.55±0.09	93.14 ± 0.04	93.17 ± 0.05
UCI	79.86±1.48	57.48 ± 1.87	88.12 ± 2.05	92.73 ± 0.06	79.54 ± 0.48	91.19 ± 0.42	87.36 ± 2.03	94.54 ± 0.12	93.38 ± 0.21
Can. Parl.	53.92±0.94	54.02 ± 0.76	54.10 ± 0.93	55.80 ± 0.69	55.18 ± 0.79	55.91 ± 0.82	54.30 ± 0.66	87.74 ± 0.71	96.64 ± 0.04
US Legis.	54.93±2.29	57.28 ± 0.71	58.63 ± 0.37	53.17±1.20	51.00 ± 3.11	50.71 ± 0.76	52.59 ± 0.97	54.28 ± 2.87	55.25 ± 4.54
UN Trade	59.65±0.77	57.02 ± 0.69	58.31±3.15	65.24 ± 0.21	61.03 ± 0.18	62.17 ± 0.31	62.21 ± 0.12	64.55 ± 0.62	67.04 ± 0.20
UN Vote	56.64±0.96	54.62 ± 2.22	$58.85{\pm}2.51$	49.94 ± 0.45	52.24 ± 1.46	50.68 ± 0.44	51.60 ± 0.97	55.93 ± 0.39	58.08 ± 0.55
Contact	94.34±1.45	92.18 ± 0.41	93.82 ± 0.99	89.55±0.30	95.87 ± 0.11	90.59 ± 0.05	91.11 ± 0.12	98.03 ± 0.02	98.10 ± 0.02
Avg. Rank	5.83	6.83	4.67	4.92	6.21	6.00	6.71	<u>2.50</u>	1.33

5.2 Effectiveness Evaluation

Table 2 and Table 3 show models' performance in dynamic link prediction under transductive and inductive settings. AUC-ROC score is reported in the Appendix D.4. From these tables, we observe that DyG-Mamba achieves the best performance on most datasets and achieves the best average rank in both AP and AUC-ROC across three negative edge sampling strategies, demonstrating its higher effectiveness and better generalization compared to SOTA baselines. The primary reasons for DyG-Mamba's superior performance can be summarized in three key aspects. (i). DyG-Mamba employs an SSM architecture that effectively handles long-term sequences. In contrast, DyGFormer requires patching under the same input length, which compresses the sequence data and leads to information loss. (ii). DyG-Mamba leverages irregular temporal information to control the compression of historical states, thereby making more efficient use of time information and enhancing model's generalization capability. (iii). DyG-Mamba selectively filters out past noise or irrelevant information, resulting in more robust node embedding and improved prediction accuracy.

To further verify the scalability and efficiency of DyG-Mamba on million-edge temporal graphs, we conducted additional experiments on tgbl-coin-v2 [43], which contains 638K nodes and 22.8M temporal edges. Following the standardized training pipeline and hyperparameter setup in Yu et al. [44], we ensured a fair comparison with the existing baselines. As summarized in Table 4,

Table 4: Scalability on million-edge temporal graphs.

Method	Performance	Running Time	GPU Usage
DyRep TGN GraphMixer DyGFormer	45.20±4.60 58.60±3.70 75.31±0.21 75.17±0.38	49:38:39 38:26:48 11:59:20 45:19:11	48116M 48116M 12204M 41348M
DyG-Mamba	75.17±0.38	12:32:04	18094M

DyG-Mamba not only achieves superior predictive performance but also demonstrates remarkable training efficiency, further validating its scalability on real-world large-scale dynamic graphs.

Scalability of Effectiveness. As shown in Figure 2, to highlight DyG-Mamba's ability to capture long-term temporal dependencies, we compare it to three best-performing baselines. Obviously, DyG-Mamba's performance improves substantially with longer sequences and outperforms baselines even at shorter sequence lengths, showing its effectiveness in modeling long-term dependencies on dynamic graphs.

Ablation Study. In Table 5, we conduct an ablation study on three datasets to evaluate the effectiveness of each component in DyG-Mamba. Specifically, we examine five variants: (i). [w/o timespan] replaces timespan with input sample as control signals, *i.e.*, the same setting with vanilla Mamba with $\Delta t = \text{SiLU}(\text{Linear}(u(t)))$. (ii). [w/o Timeencoding] removes the absolute temporal en-

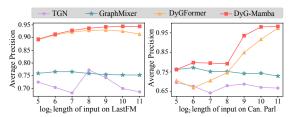


Figure 2: AP score w.r.t. varying sequence lengths.

Table 5: Results (AP) of time information ablations.

Settings	Can. Parl.	Enron	USLegis.
w/o timespan	96.90±0.18	92.14±0.12	72.26 ± 0.76
w/o Time-encoding	97.80±0.43	92.83±0.06	73.33 ± 1.15
w/o Time	96.87±0.18	92.08±0.12	72.19 ± 0.72
w/o Selective	96.32±0.16	91.33±0.10	71.62±0.77
w/o Data-dependent	79.24±0.58	82.25±0.16	70.57±0.84
DyG-Mamba	98.37±0.07	93.22±0.03	74.11±2.23

coding $X_{u,T}^{\tau}$. (iii). [w/o Time] removes both timespan and time-encoding. (iv). [w/o Selective] follows parameter settings of S4 [45], *i.e.*, meaning all parameters are independent of input or timespan. (v). [w/o Data-dependent] changes the parameters B and C from being data-dependent to timespan dependent without spectral norm constraints. We observe that removing any component from DyG-Mamba adversely affects its dynamic graph learning capability. Specifically, [w/o timespan] significantly decreases the performance, as it is crucial for capturing irregular temporal patterns. And [w/o Time-encoding] leads to a slight decline in performance since absolute temporal information is also important. Furthermore, [w/o Selective] also degrades performance since the fixed parameters fail to filter out irrelevant noise. Finally, relying entirely on timespan [w/o Data-dependent] also reduces performance, showing the importance of input-dependent setting for parameters B and C.

Effect of Learnable Δt Function. The design of the learnable function for Δt in Eq. (11) is inspired by the Ebbinghaus forgetting curve, which models memory retention as $R = \exp(-t/S)$, where t is the elapsed time and S a decay constant. We reinterpret this formulation as a timespan-dependent decay coefficient, ensuring that longer intervals induce stronger decay consistent with human memory dynamics. To validate this choice, we compare several monotonic alternatives, including Linear, Logarithmic, Sigmoid, and

Table 6: Ablation on learnable Δt function.

Variants	Can.Parl.	Enron	USLegis.
Linear Logarithmic Sigmoid Exponential $w/o au - t_1$ w/o au mespan	94.64 ± 0.18 97.18 ± 0.12 95.84 ± 0.13 94.68 ± 0.11 97.32 ± 0.20 96.90 ± 0.18	91.13 ± 0.06 92.35 ± 0.05 92.26 ± 0.04 90.84 ± 0.06 92.46 ± 0.16 92.14 ± 0.12	70.28 ± 1.84 72.46 ± 2.34 72.15 ± 2.36 71.23 ± 1.48 72.45 ± 0.84 72.26 ± 0.76
DyG-Mamba	98.37±0.07	93.22±0.03	74.11±2.23

Exponential variants, as well as versions without normalization or timespan inputs. As summarized in Table 6, our formulation consistently achieves the best performance across datasets, demonstrating a balanced and smooth decay behavior. In contrast, Linear and Log variants lack boundedness, Sigmoid saturates early, and the Exp variant over-amplifies long timespans, leading to unstable training.

5.3 Efficiency Evaluation

Given an input length of 256, Figure 3 and Appendix D.4 show the training time per epoch and the size of trainable parameters on Enron data. Obviously, CAWN requires the longest training time and a substantial number of parameters, since it conducts random walks on dynamic graphs to collect time-aware sequences. In contrast, simpler methods, *e.g.*, GraphMixer and JODIE, have fewer parameters, but exhibit a significant performance gap compared to DyG-Former and DyG-Mamba. Overall, DyG-Former and DyG-Mamba.

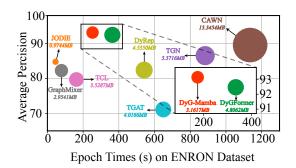


Figure 3: Comparison of efficiency and effectiveness.

Mamba achieves the best performance with a small number of trainable parameters and a moderate training time required per epoch.

Scalability of Efficiency. In Figure 4, to highlight DyG-Mamba's ability to effectively capture long-term temporal dependency on dynamic graphs, we provide a more detailed efficiency comparison between Transformer-based DyG-Former and DyG-Mamba. For a fair comparison, we make the same experimental setting for both frameworks. We observe that, with increasing input sequence length, DyG-Mamba demon-

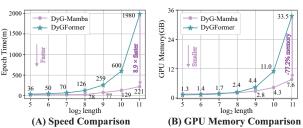


Figure 4: Speed and memory comparison of two layers DyGFormer and DyG-Mamba with varying lengths.

strates a linear growth trend in both runtime and memory consumption, highlighting its efficiency. Specifically, DyG-Mamba is 8.9 times faster than DyGFormer and reduces GPU memory consumption by 77.2% at a sequence length of 2,048. This is because DyG-Mamba only needs a few parameters to compress hidden state and adopts hardware-aware parallel scanning for training.

Table 7: Training convergence comparison on Can. Parl. dataset.

Model	Epoch=20	Epoch=40	Epoch=60	Epoch=80	Epoch=100
Vanilla Mamba	0.2229	0.2026	0.1964	0.1922	0.1914
DyG-Mamba	0.1691	0.1478	0.1480	0.1482	0.1480

5.4 Training Efficiency

To ensure stable and efficient training, DyG-Mamba incorporates several lightweight yet effective designs. The learnable Δt function is implemented as an element-wise exponential decay through a scalar-wise MLP, enabling adaptive modeling of irregular intervals with negligible overhead. A spectral norm constraint regularizes the B and C matrices without adding parameters, preventing instability and ensuring bounded outputs under input noise as supported by Theorem 4.3. Moreover, redefining B and C as input-dependent linear mappings introduces minimal cost while improving the model's adaptability. As shown in Table 7, DyG-Mamba converges smoothly within 100 epochs on the Can.Parl. dataset, confirming its fast and stable optimization behavior across datasets.

5.5 Robustness Evaluation

We conduct a robustness test by randomly inserting 10% to 60% noisy edges with chronological timestamps during the evaluation. In Figure 5, when the proportion of noisy edges increases, DyG-Mamba exhibits only a minor performance decline, indicating stronger noise robustness compared to baselines. We attribute robustness to the review-based selective memory en-

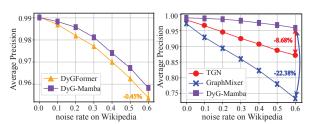


Figure 5: Inserting noisy edges from 10% to 60%.

hancement, which enables the model to identify most relevant information and filter out noise.

5.6 Case Study

In Figure 6, we randomly extract a middle part sample from one long-term input sequence of the Wikipedia dataset, e.g., {1713, 160, 1667, 1667, 1667, 1713, 1667}, to visualize DyG-Mamba's efficiency capability for long-term sequence modeling. Figures 6(A)-(B) show the normalized cosine similarity of node embedding between the last and second-last hidden states. We observe that DyGFormer shows high diagonal similarity and nearly uniform similarity across all neighbors, indicat-

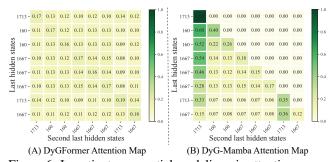


Figure 6: Investigate sequential modeling via attention map between source and destination nodes on link prediction.

ing that it struggles to distinguish important historical information. In contrast, DyG-Mamba assigns greater weights to the historical reappearing destination nodes, better enhancing the node embedding and filtering out irrelevant and noisy historical information.

6 Conclusion

In this work, we propose a novel SSM framework called DyG-Mamba to effectively and efficiently capture long-term temporal dependencies on dynamic graphs. To achieve this goal, we incorporate irregular time spans as controllable signals, thus establishing a strong correlation between dynamic evolution patterns and time information. We also implement a review-based selective memory enhancement to further improve the model's robustness. Experimental evaluations on various downstream tasks show DyG-Mamba's higher performance and better robustness. In the future, we plan to deploy DyG-Mamba in more real-world applications.

References

- [1] Alessio Gravina, Giulio Lovisotto, Claudio Gallicchio, Davide Bacciu, and Claas Grohnfeldt. Long range propagation on continuous-time dynamic graphs. *International Conference on Machine Learning (ICML)*, 2024.
- [2] Srijan Kumar, Xikun Zhang, and Jure Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *International Conference on Knowledge Discovery & Data Mining (KDD)*, page 1269–1278, 2019.
- [3] Lei Bai, Lina Yao, Can Li, Xianzhi Wang, and Can Wang. Adaptive graph convolutional recurrent network for traffic forecasting. In *Neural Information Processing Systems (NeurIPS)*, 2020.
- [4] Le Yu, Zihang Liu, Leilei Sun, Bowen Du, Chuanren Liu, and Weifeng Lv. Continuous-time user preference modelling for temporal sets prediction. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 36(4):1475–1488, 2023.
- [5] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bronstein. Temporal graph networks for deep learning on dynamic graphs. In *ICML Workshop on Graph Representation Learning*, 2020.
- [6] Le Yu, Leilei Sun, Bowen Du, and Weifeng Lv. Towards better dynamic graph learning: New architecture and unified library. In *Conference on Neural Information Processing Systems* (NeurIPS), 2023.
- [7] Yuxia Wu, Yuan Fang, and Lizi Liao. On the feasibility of simple transformer for dynamic graph modeling. In *Proceedings of the ACM on Web Conference (WWW)*, page 870–880, 2024.
- [8] Zihang Jiang, Weihao Yu, Daquan Zhou, Yunpeng Chen, Jiashi Feng, and Shuicheng Yan. Convbert: Improving BERT with span-based dynamic convolution. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [9] 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 *International Conference on Learning Representations (ICLR)*, 2023.
- [10] 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 (ICLR)*, 2024.
- [11] Yanbang Wang, Yen-Yu Chang, Yunyu Liu, Jure Leskovec, and Pan Li. Inductive representation learning in temporal networks via causal anonymous walks. In *International Conference on Learning Representations (ICLR)*, 2021.
- [12] Weilin Cong, Jian Kang, Hanghang Tong, and Mehrdad Mahdavi. On the generalization capability of temporal graph learning algorithms: Theoretical insights and a simpler method. *arXiv preprint*, arXiv:2402.16387, 2024.
- [13] Siwei Zhang, Yun Xiong, Yao Zhang, Yiheng Sun, Xi Chen, Yizhu Jiao, and Yangyong Zhu. Rdgsl: Dynamic graph representation learning with structure learning. In *International Conference on Information and Knowledge Management (CIKM)*, page 3174–3183, 2023.
- [14] ZhengZhao Feng, Rui Wang, TianXing Wang, Mingli Song, Sai Wu, and Shuibing He. A comprehensive survey of dynamic graph neural networks: Models, frameworks, benchmarks, experiments and challenges. *arXiv* preprint, arXiv:2405.00476, 2024.
- [15] Yanping Zheng, Lu Yi, and Zhewei Wei. A survey of dynamic graph neural networks. *Front. Comput. Sci.*, 19(6), 2024.
- [16] Philipp Foth, Lukas Gosch, Simon Geisler, Leo Schwinn, and Stephan Günnemann. Relaxing graph transformers for adversarial attacks. In *ICML 2024 Workshop on Differentiable Almost Everything*, 2024.

- [17] Haonan Yuan, Qingyun Sun, Zhaonan Wang, Xingcheng Fu, Cheng Ji, Yongjian Wang, Bo Jin, and Jianxin Li. Dg-mamba: Robust and efficient dynamic graph structure learning with selective state space models. *Conference on Artificial Intelligence (AAAI)*, 2025.
- [18] Hermann Ebbinghaus. Über das gedächtnis: untersuchungen zur experimentellen psychologie. Duncker & Humblot, 1885.
- [19] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023.
- [20] Bo Ae Chun and Hae Ja Heo. The effect of flipped learning on academic performance as an innovative method for overcoming ebbinghaus' forgetting curve. In *International Conference on Information and Education Technology (IEEE-ICIET)*, page 56–60, 2018.
- [21] Aldo Pareja, Giacomo Domeniconi, Jie Chen, Tengfei Ma, Toyotaro Suzumura, Hiroki Kanezashi, Tim Kaler, Tao B. Schardl, and Charles E. Leiserson. Evolvegen: Evolving graph convolutional networks for dynamic graphs. In *Conference on Artificial Intelligence (AAAI)*, pages 5363–5370, 2020.
- [22] Aravind Sankar, Yanhong Wu, Liang Gou, Wei Zhang, and Hao Yang. Dysat: Deep neural representation learning on dynamic graphs via self-attention networks. In *International Conference on Web Search and Data Mining (WSDM)*, pages 519–527, 2020.
- [23] 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. *arXiv* preprint, arXiv:2105.07944, 2021.
- [24] Dongyuan Li, Satoshi Kosugi, Ying Zhang, Manabu Okumura, Feng Xia, and Renhe Jiang. Revisiting dynamic graph clustering via matrix factorization. In *Proceedings of the ACM on Web Conference (WWW)*, page 1342–1352, 2025.
- [25] Linhao Luo, Gholamreza Haffari, and Shirui Pan. Graph sequential neural ODE process for link prediction on dynamic and sparse graphs. In *International Conference on Web Search and Data Mining (WSDM)*, pages 778–786. ACM, 2023.
- [26] Xiao Luo, Haixin Wang, Zijie Huang, Huiyu Jiang, Abhijeet Sadashiv Gangan, Song Jiang, and Yizhou Sun. CARE: Modeling interacting dynamics under temporal environmental variation. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [27] Ming Jin, Yuan-Fang Li, and Shirui Pan. Neural temporal walks: Motif-aware representation learning on continuous-time dynamic graphs. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2022.
- [28] Zijie Huang, Yizhou Sun, and Wei Wang. Learning continuous system dynamics from irregularly-sampled partial observations. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [29] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, and Hongyuan Zha. Dyrep: Learning representations over dynamic graphs. In *International Conference on Learning Representations* (*ICLR*), 2019.
- [30] Farimah Poursafaei, Andy Huang, Kellin Pelrine, and Reihaneh Rabbany. Towards better evaluation for dynamic link prediction. In *Conference on Neural Information Processing* (NeurIPS), 2022.
- [31] 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 *Conference on Neural Information Processing Systems (NeurIPS)*, pages 572–585, 2021.
- [32] Daniel Y Fu, Tri Dao, Khaled Kamal Saab, Armin W Thomas, Atri Rudra, and Christopher Re. Hungry hungry hippos: Towards language modeling with state space models. In *International Conference on Learning Representations (ICLR)*, 2022.

- [33] Haohao Qu, Liangbo Ning, Rui An, Wenqi Fan, Tyler Derr, Hui Liu, Xin Xu, and Qing Li. A survey of mamba. *arXiv preprint arXiv:2408.01129*, 2024.
- [34] Ali Behrouz and Farnoosh Hashemi. Graph mamba: Towards learning on graphs with state space models. In *Conference on Knowledge Discovery and Data Mining (KDD)*, pages 119–130, 2024.
- [35] Chloe Wang, Oleksii Tsepa, Jun Ma, and Bo Wang. Graph-mamba: Towards long-range graph sequence modeling with selective state spaces. *arXiv preprint*, arXiv:2402.00789, 2024.
- [36] Jintang Li, Ruofan Wu, Xinzhou Jin, Boqun Ma, Liang Chen, and Zibin Zheng. State space models on temporal graphs: A first-principles study. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2024.
- [37] Lincan Li, Hanchen Wang, Wenjie Zhang, and Adelle Coster. Stg-mamba: Spatial-temporal graph learning via selective state space model. *arXiv preprint*, arXiv:2403.12418, 2024.
- [38] Abdulkadir CELIKKANAT, Nikolaos Nakis, and Morten Mørup. Piecewise-velocity model for learning continuous-time dynamic node representations. In *Learning on Graphs Conference* (*LoG*), pages 36–1, 2022.
- [39] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. In *Conference on Neural Information Processing Systems* (NeurIPS), 2020.
- [40] Piotr Woźniak, Edward Gorzelańczyk, and Janusz Murakowski. Two components of long-term memory. Acta neurobiologiae experimentalis, 55(4):301–305, 1995.
- [41] John T Wixted. The psychology and neuroscience of forgetting. Annu. Rev. Psychol., 55(1):235–269, 2004.
- [42] Bo Chang, Minmin Chen, Eldad Haber, and Ed H. Chi. AntisymmetricRNN: A dynamical system view on recurrent neural networks. In *International Conference on Learning Represen*tations (ICLR), 2019.
- [43] Shenyang Huang, Farimah Poursafaei, Jacob Danovitch, Matthias Fey, Weihua Hu, Emanuele Rossi, Jure Leskovec, Michael M. Bronstein, Guillaume Rabusseau, and Reihaneh Rabbany. Temporal graph benchmark for machine learning on temporal graphs. In *Conference on Neural Information Processing Systems (NeurIPS)*, 2023.
- [44] Le Yu. An empirical evaluation of temporal graph benchmark. arXiv preprint arXiv:2307.12510, 2023.
- [45] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations (ICLR)*, 2022.
- [46] Weihao Yu and Xinchao Wang. Mambaout: Do we really need mamba for vision? In *Proceedings of the Computer Vision and Pattern Recognition Conference (CVPR)*, pages 4484–4496, 2025.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We highlight the contributions and scope in abstract and introduction.

Guidelines:

• The answer NA means that the abstract and introduction do not include the claims made in the paper.

- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We provide the limitation section in Appendix A.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We provide proofs of assumption in Appendix A.

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.

• Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide full information for reproducibility in Section 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We provide source code for reproducibility: https://anonymous/DyGMamba.

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).

- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide all the training and test details in Appendix C.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We provide a significance test in Table 2 and report the mean and variance of the experimental results.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: We provide the experimental environments in Appendix D.2.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We conform the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We provide the societal impacts in Appendix E.3.

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We do not contain models or data that have a high risk.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released
 with necessary safeguards to allow for controlled use of the model, for example by
 requiring that users adhere to usage guidelines or restrictions to access the model or
 implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do
 not require this, but we encourage authors to take this into account and make a best
 faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use the open-source toolkit and dataset with the license CC-BY 4.0.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We use the datasets and DyGLib toolkit are open-access. And we anonymize our source code for submission.

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.

- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: the LLM is only used for writing.

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

A Limitations

Although our DyG-Mamba delivers strong performance on dynamic graph modeling, it still has several notable limitations. (i) More Interactions. We focus primarily on edge addition, which is widely studied interaction type in previous research. Extending our framework to other interaction types, such as node addition/deletion, edge deletion, and node/edge feature transformations, remain an avenue for future research. (ii) Larger-scale Datasets. Existing benchmarks are relatively small-scale datasets. And it is unclear whether these findings will generalize to larger or more complex real-world scenarios. (iii) More Domains. Although DyG-Mamba has achieved good results on dynamic network modeling, it is still unknown whether this conclusion can be extended to other domains. Recent insights from Mambaout [46] suggest that Mamba is especially suited for autoregressive and long-sequence tasks. Even though dynamic link prediction and node classification are not strictly auto-regressive, we still obtain state-of-the-art performance, indicating that Mamba's broader effectiveness across diverse tasks warrants further exploration.

B Computational Complexity

In our implementation, we employ two DyG-Mamba layers with batch size b, feature dimension d, expanded state dimension 2d, and SSM dimension d_{ssm} . Note that while Section 4.1 sets the feature dimension of node embeddings as 4d, we use d here for simplicity. On GPUs, high-bandwidth memory (HBM) provides larger capacity, whereas static random-access memory (SRAM) offers higher bandwidth. Building on Mamba, DyG-Mamba first reads $O(bL(2d) + (2d)d_{ssm})$ bytes of (Δ, A, B, C) from slow HBM to fast SRAM. It then derives the discrete \overline{A} , \overline{B} of size $(b, L, 2d, d_{ssm})$ in SRAM, executes the SSM operation in SRAM, and writes the output of size (b, L, 2d) back to HBM. This approach reduces I/O overhead from O(bL(2d)N) to $O(bL(2d) + (2d)d_{ssm})$, yielding a memory complexity of $O(bL(2d) + (2d)d_{ssm})$. Since d_{ssm} is relatively smaller compared to bL, we simplify it as O(bLd). The time complexity to calculate B, C, Δ is $O(3bL(2d)d_{ssm})$, and the SSM process takes $O(bL(2d)d_{ssm})$. Compared to transformer-based methods with quadratic time and memory complexity $O(bL^2d)$, DyG-Mamba scales effectively to large sequence lengths.

C Algorithm Details

Here, we list the detailed workflow of DyG-Mamba in Algorithm.1. We also list the procedures of DyG-Mamba for dynamic link prediction in Algorithm.2 and dynamic node classification in Algorithm.3.

Algorithm 1 Continuous SSM

```
Input: Hidden states Z = \{z_{i,1}, z_{i,2}, \dots, z_{i,L}\}_{i=1}^{B}: (B, L, d), Control signals \Delta t = \{\Delta t_1, \Delta t_2, \dots, \Delta t_L\}: (B, L, d), SSM
dimension d_{ssm}, Expanded dimension: 2d.
// Normalize the input sequence.
Initialize \mathbf{Parameter}_{i}^{A}: (2d,d_{ssm}), \mathbf{Parameter}_{i}^{\Delta}: (d_{ssm})
Initialize \boldsymbol{x}: (B, L, 2d) \leftarrow Linear \boldsymbol{z}(\boldsymbol{Z}), \boldsymbol{z}: (B, L, 2d) \leftarrow Linear \boldsymbol{z}(\boldsymbol{Z}), \Delta t: (B, L, 2d) \leftarrow Linear \boldsymbol{z}(\boldsymbol{\Delta}t),
for o in {forward, backward} do
           \boldsymbol{x}_o': (B, L, 2d) \leftarrow SiLU(Conv1d<sub>o</sub>(\boldsymbol{x}))
            \mathbf{B}_o: (B, L, d_{ssm}) \leftarrow \mathbf{Linear}_o^{\mathbf{B}}(\mathbf{x}_o')
           C_o: (B, L, d_{ssm}) \leftarrow \mathbf{Linear}_o^C(\mathbf{x}_o')
           // Control signal to control the selection of historical information. \Delta_{\circ}\colon (B,L,2d)\leftarrow Using Eq.(11)
           \begin{split} & \bar{\boldsymbol{A}}_o \colon (\mathbf{B}, \mathbf{L}, 2d, d_{ssm}) \leftarrow \boldsymbol{\Delta}_i \otimes \mathbf{Parameter}_o^{\boldsymbol{A}} \\ & \bar{\boldsymbol{B}}_o \colon (\mathbf{B}, \mathbf{L}, 2d, d_{ssm}) \leftarrow \boldsymbol{\Delta}_i \otimes \boldsymbol{B}_o \\ & \boldsymbol{y}_o \colon (\mathbf{B}, \mathbf{L}, 2d) \leftarrow \mathbf{SSMs}(\bar{\boldsymbol{A}}_o, \bar{\boldsymbol{B}}_o, \boldsymbol{C}_o)(\boldsymbol{x}_o') \end{split}
 // Gated y.
{m y}'_{	ext{forward}} \colon (\mathsf{B},\mathsf{L},2d) \leftarrow {m y}_{	ext{forward}} \odot \mathbf{SiLU}({m z})
// Residual connection.
\boldsymbol{y}_{\text{backward}}' \colon (\text{B}, \text{L}, 2d) \leftarrow \boldsymbol{y}_{\text{backward}} \odot \text{SiLU}(\boldsymbol{z})
 \hat{m{Z}}: (B, L, d) \leftarrow Linear ^{m{T}}(m{y}'_{	ext{forward}} + m{y}'_{	ext{backward}}) + m{Z}
 Output: Updated hidden states \hat{Z}
```

Algorithm 2 Dynamic Link Prediction

```
\textbf{Input:} \ \ \text{The dynamic interaction set} \ \mathcal{D} = \{(u_i, v_i, t_i)\}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{readout function Read}(\cdot), \text{output projection } \}_{i=1}^K, \text{Continuous SSM DyG-Mamba}(\cdot), \text{continuous SSM DyG-Mamb
layer \phi(\cdot).
for T Epochs do
                                   for (u_1, v_1, t) \in \mathcal{D} do
                                                                        For a node pair (u_1, v_1). Sampled neighbor sequence: S_1, S_2.
                                                                         for u in \{u_1, v_1\} do
                                                                                                              Node Features: oldsymbol{X}_{u,N}^{	au} \in \mathbb{R}^{|u| 	imes d_V} ,
                                                                                                                Edge Features: oldsymbol{X}_{u,E}^{	au} \in \mathbb{R}^{|u| 	imes d_E},
                                                                                                                Time Features: oldsymbol{X}_{u,T}^{	au} \in \mathbb{R}^{|u| 	imes d_T} ,
                                                                                                                Co-occurance Features: oldsymbol{X}_{u,C}^{	au} \in \mathbb{R}^{|u| 	imes d_C} ,
                                                                                                                Time Span: oldsymbol{\Delta t} = \{\Delta t_1, \overset{\cdot}{\Delta} \overset{\cdot}{t_2}, \ldots, \Delta t_{|u|}\}
                                                                                                              The optimal part of the following part of the part of
                                                                                                                // Continuous SSM encoder.
                                                                                                                \widehat{m{M}}_{u}^{	au} = 	ext{DyG-Mamba}(m{Z}_{u}^{	au}, \Delta t_{u})
                                                                           end
                                                                         \widehat{\boldsymbol{Z}}_{u_1,\text{out}}^{\tau} = \phi(\text{Read}(\widehat{\boldsymbol{M}}_{u_1}^{\tau})), \widehat{\boldsymbol{Z}}_{v_1,\text{out}}^{\tau} = \phi(\text{Read}(\widehat{\boldsymbol{M}}_{v_1}^{\tau})). \text{ // Output.}
                                                                         \hat{y} = \text{Softmax}(\text{Linear}(\text{RELU}(\text{Linear}(\widehat{\boldsymbol{Z}}_{u_1,\text{out}}^{\tau} || \widehat{\boldsymbol{Z}}_{v_1,\text{out}}^{\tau}))).
                                      Compute \mathcal{L}_{\text{Link Prediction}}.
Output: Dynamic link prediction labels.
```

Algorithm 3 Dynamic Node Classification

```
Input: The dynamic interaction set \mathcal{D} = \{(u_i, y_i)\}_{i=1}^N, continuous SSM DyG-Mamba(·), readout function Read(·), output projection layer \phi(\cdot). for T Epochs \mathbf{do}

| for (u,y) \in \mathcal{D} \mathbf{do}
| Sampled neighbor sequence: S = \{(k_1, t_1), (k_2, t_2), \dots, (k_{|u|}, t_{|u|})\}, Node Features: \mathbf{X}_{u,N}^{\tau} \in \mathbb{R}^{|u| \times d} \mathbf{d} \mathbf{v}, Edge Features: \mathbf{X}_{u,T}^{\tau} \in \mathbb{R}^{|u| \times d} \mathbf{d} \mathbf{r}, Time Features: \mathbf{X}_{u,T}^{\tau} \in \mathbb{R}^{|u| \times d} \mathbf{r}, Time Span: \mathbf{\Delta t} = \{\Delta t_1, \Delta t_2, \dots, \Delta t_L\}
| \mathbf{Z}_{u,*}^{\tau} = \mathbf{X}_{u,*}^{\tau} \mathbf{W}_* + \mathbf{b}, where * \in N, E, T.
| \mathbf{Z}_{u}^{\tau} = \mathbf{Z}_{u,N}^{\tau} \| \mathbf{Z}_{u,E}^{\tau} \| \mathbf{Z}_{u,T}^{\tau}
| \widehat{\mathbf{M}}_{u}^{\tau} = \mathrm{DyG-Mamba}(\mathbf{Z}_{u}^{\tau}, \Delta t_{u}),
| \widehat{\mathbf{Z}}_{u,\mathrm{out}}^{\tau} = \phi(\mathrm{Read}(\widehat{\mathbf{M}}_{u}^{\tau})),
| \widehat{y} = \mathrm{Softmax}(\mathrm{Linear}(\mathrm{RELU}(\mathrm{Linear}(\widehat{\mathbf{Z}}_{u,\mathrm{out}}^{\tau}))).
| end | Compute \mathcal{L}_{\mathrm{Node}} Classification labels.
```

Table 8: Statistics of the datasets. N/A denotes that there is no node/edge features. # Node denotes the number of nodes.

Datasets	Domains	#Nodes	#Links	#N&L Feature	Bipartite	Duration	Unique Steps	Time Granularity
Wikipedia	Social	9,227	157,474	N/A & 172	True	1 month	152,757	Unix timestamps
Reddit	Social	10,984	672,447	N/A & 172	True	1 month	669,065	Unix timestamps
MOOC	Interaction	7,144	411,749	N/A & 4	True	17 months	345,600	Unix timestamps
LastFM	Interaction	1,980	1,293,103	N/A & N/A	True	1 month	1,283,614	Unix timestamps
Enron	Social	184	125,235	N/A & N/A	False	3 years	22,632	Unix timestamps
Social Evo.	Proximity	74	2,099,519	N/A & 2	False	8 months	565,932	Unix timestamps
UCI	Social	1,899	59,835	N/A & N/A	False	196 days	58,911	Unix timestamps
Can. Parl.	Politics	734	74,478	N/A & 1	False	14 years	14	years
US Legis.	Politics	225	60,396	N/A & 1	False	12 congresses	12	congresses
UN Trade	Economics	255	507,497	N/A & 1	False	32 years	32	years
UN Vote	Politics	201	1,035,742	N/A & 1	False	72 years	72	years
Contact	Proximity	692	2,426,279	N/A & 1	False	1 month	8,064	5 minutes

D Experimental Details

D.1 Dataset Details

We evaluate our methods on a diverse set of dynamic graph datasets, including twelve publicly available datasets collected by Edgebank [30], which are publicly available¹. We present the statistics of the datasets in Table 8, where #N&L Feature stands for the dimensions of the node and link features. Note that our calculation of the Contact dataset's statistics (694 nodes and 2,426,280 links) slightly differs from the values reported in [30], although both are derived from the same dataset.

Table 9: Configurations

8	
Configuration	Setting
Learning rate	0.0001
Train Epochs	100
Optimizer	Adam
Dimension of time encoding d_T	100
Dimension of co-occurrence d_C	50
Dimension of aligned encoding d	50
Dimension of Δt_i 's encoder $4d$	200
Dimension of output d_{out}	172
Number of Mamba blocks	2
Dimension of SSM d_{ssm}	16
Expanded factor of Mamba	2
Number of Corss-Attention layer	1
5	

Table 10: Sequence Length Settings

Dataset	Sequence Length
Wikipedia	64
Reddit	64
MOOC	256
LastFM	512
Enron	512
Social Evo.	64
UCI	32
Can. Parl.	2048
US Legis.	272
UN Trade	256
UN Vote	128
Contact	32

D.2 Implementation Details

Experiment Environment. We conduct experiments on an Ubuntu 22.04 LTS server equipped with one Intel(R) Core(TM) i9-10900X CPU @ 3.70GHz with 10 physical cores and NVIDIA RTX A6000 GPUs (48GB). The code is written in Python 3.10 and we use PyTorch 2.1.0 on CUDA 11.8 to train the model.

Configuration Details. For all baselines, we follow the configurations as DyGFormer reported [6]. For DyG-Mamba, we list all configurations in Table 9. Then, we perform the grid search to find the optimal sequence length, with a search range spanning from 32 to 2,048 in powers of 2. It is worth noticing that DyG-Mamba can handle nodes with sequence lengths shorter than the defined length. When the sequence length exceeds the specified length, we truncate the sequence and preserve the most recent interactions up to the defined length. Finally, we present the sequence length settings in Table 10. All implementation details could be accessed at the link: https://anonymous.4open/DyGMamba.

¹https://zenodo.org/records/dynamic-graphs

The inconsistency problem between the description of the transductive setting in DyGFormer and the coding in DyGLib. Thanks to other researchers in this community, we observed that CAW-N explicitly avoids including unseen nodes in the validate/test sets. In contrast, DyGLib, TGAT, and TGN adopt a more relaxed version of the transductive setting, where both previously observed and new nodes can appear during evaluation. For a fair comparison, we follow the coding in DyGLib, which uses more relaxed version of the transductive setting. Therefore all baselines in our paper use the same data split with both previously observed and new nodes in transductive setting, to ensure the fairness of experimental comparisons. To compare with CAW-N, we also re-implement CAW-N with the relaxed version of transductive setting.

D.3 Detailed Evaluation Settings

Detailed Settings for Effectiveness Evaluation. To provide a more comprehensive evaluation of baseline performance, and following [30], we adopt three distinct negative edge sampling strategies for the temporal link prediction task. We define the training and test edge sets as E_{train} and E_{test} , respectively. The edges of a given dynamic graph can then be grouped into three categories: (a) edges observed only during training ($E_{\text{train}} \setminus E_{\text{test}}$), (b) edges that appear both in training and test ($E_{\text{train}} \cap E_{\text{test}}$), referred to as *transductive* edges, and (c) edges observed exclusively in the test phase ($E_{\text{test}} \setminus E_{\text{train}}$), regarded as *inductive* edges. Note that inductive negative sampling here is distinct from the usual notion of inductive settings. We elaborate on the three sampling strategies as follows:

- Random Negative Sampling (rnd). Negative edges are sampled at random from all possible node pairs. At each timestep, we retain the timestamps, features, and source nodes of the positive edges, but randomly select their destination nodes from the entire node set.
- Historical Negative Sampling (hist). In historical negative sampling, we focus on edges that were observed at previous time steps but are absent in the current step. This approach assesses whether a method can accurately predict the specific timestamps at which an edge may reappear, rather than simply predicting that it always reoccurs once observed. Formally, for a given time step t, we sample from edges in $(E_{train} \cap \overline{E_t})$. If the number of available historical edges is insufficient to match the number of positive edges, we revert to random sampling for the remainder.
- Inductive Negative Sampling (ind). Whereas historical sampling centers on edges observed during training, inductive negative sampling evaluates whether a model can capture the reoccurrence of edges that first appear only at test time. Once newly appearing edges have been observed in the test, the model is asked to predict if these edges will reoccur in subsequent time steps. Formally, at time t, we sample from $(E_{test} \cap \overline{E_{train}} \cap \overline{E_t})$. If there are not enough such inductive edges to match the number of positive edges, the remaining negative edges are sampled randomly.

Detailed Settings for Efficiency Evaluation. In Figures 3,4, in the evaluation of time and memory consumption, we do not choose the configuration as reported in DyGLib¹, as there is a significant difference in sequence lengths between different baselines. In the LastFM dataset, the reported number of sampled neighbors for DyRep, TGN, and GraphMixer is set to 10, while for DyGFormer and DyG-Mamba, it is 512. This is unfair because the complexity of time and memory is highly dependent on the number of neighbors sampled. Therefore, for a fair comparison in terms of time and memory consumption, we use the same number of sampled neighbors across all models: 32 for UCI, 256 for Enron, 512 for LastFM and 64 for Reddit.

Detailed Settings for Robustness Evaluation. As shown in Figure 5, the purpose of the robustness test is to evaluate the ability to against noise in edges and timestamps. In the training step, we typically train models on a transductive setting with random negative sampling. In the evaluation step, after neighbor sampling, we randomly select $\sigma * L$ positions to insert noise. Specifically, the noise position's node and timestamps are randomly generated. And the noise rate σ is chosen from 0.1 to 0.6.

¹github.com/yule-BUAA/DyGLib

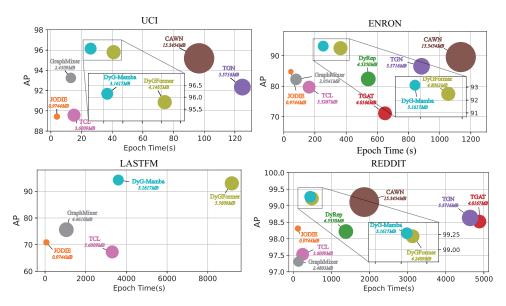


Figure 7: The AP score with different sequence lengths. We use the same sequence length for each model (uci=32, enron=256, lastfm=512, reddit=64), The not appearing model means OOM.

D.4 Additional Experimental Results

Training Time and Parameter Size. Figure 7 shows the additional time and parameter size comparison on the UCI, Enron, LastFM, and Reddit datasets. The models that do not appear in the figure experienced out-of-memory (OOM) errors. We can see that our model achieves the best performance while incurring low time and memory costs.

Toble 11.	Danfarmanasan	drinomio nodo	alassification
Table 11:	Performance on	avnamic node	crassification.

Methods	Wikipedia	Reddit	Avg. Rank
JODIE	88.99±1.05	60.37±2.58	5.00
DyRep	86.39 ± 0.98	63.72 ± 1.32	6.00
TGAT	84.09 ± 1.27	70.04 ± 1.09	5.00
TGN	86.38 ± 2.34	$\overline{63.27\pm0.90}$	7.00
CAWN	84.88 ± 1.33	66.34 ± 1.78	6.00
EdgeBank	N/A	N/A	N/A
TCL	77.83 ± 2.13	68.87 ± 2.15	6.00
GraphMixer	86.80 ± 0.79	64.22 ± 3.32	5.00
DyGFormer	87.44 ± 1.08	68.00 ± 1.74	<u>3.50</u>
DyG-Mamba	88.58 ± 0.92	70.79±1.97	1.50

Performance on Dynamic Node Classification. For dynamic node classification, we estimate the state of a node in a given interaction at a specific time and use *AUC-ROC* as the evaluation metric. Table 11 shows the AUC-ROC results on dynamic node classification. DyG-Mamba achieves SOTA performance on the Reddit dataset and second-best performance on the Wikipedia dataset. In addition, DyG-Mamba achieves the best average rank of 1.5 compared to the second-best DyGFormer with 3.5 AUC-ROC results for all baselines in Table 11.

Robustness Evaluation. We provide additional robustness test results on the UCI and Can. Parl. datasets, as shown in Figure 8. DyG-Mamba consistently exhibits greater robustness compared to DyGFormer and GraphMixer. However, TGN uses a memory bank to store previous node embeddings, making it less sensitive to noise in current neighbors.

Transductive Dynamic Link Prediction. We show the AUC-ROC for transductive dynamic link prediction with three negative sampling strategies in Table 15. Since we cannot reproduce the same performance reported by FreeDyG [10], we directly copy the results from their paper.

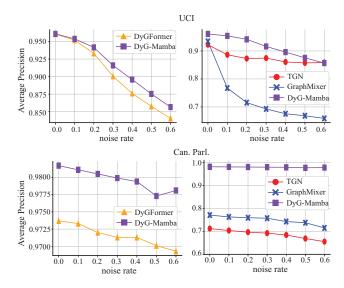


Figure 8: The robustness test on datasets UCI, Can. Parl.

Inductive Dynamic Link Prediction. We present the AP and AUC-ROC for inductive dynamic link prediction with three negative sampling strategies in Table 16 and Table 17.

Comparison of the effectiveness of Co-occurrence. DyGFormer designs co-occurrence module and has conducted ablation studies on it. DyG-Mamba focuses on detailed ablations of the modules we propose. But we also conduct additional ablation to compare the effectiveness of co-occurrence encoding in table 12. From the results, DyG-Mamba still exhibits a strong ability to capture long-term dependencies, outperforming DyGFormer.

Table 12: Comparison of the effectiveness of Co-occurrence Frequency Encoding.

			<u>`</u>
	uci	USLegis	UN Trade
DyG-Mamba	96.79 ± 0.08	74.11 ± 2.32	68.55 ± 0.16
DyG-Mamba w/o Co-occurrence	93.09 ± 1.31	73.53 ± 2.43	66.32 ± 0.55
DyGFormer	95.79 ± 0.17	71.11 ± 0.59	66.46 ± 1.29
DyGFormer w/o Co-occurrence	83.05 ± 0.38	70.59 ± 0.36	61.93 ± 1.79

Ablation Study with Co-occurrence and skip connection. We conduct ablation studies on three datasets. 'w/o Co-occurrence' refers to removing Co-occurrence Frequency Encoding from DyG-Mamba, while 'w/o skip connection' denotes the removal of the skip connection in the SSM.

Table 13: We conduct ablation studies on three datasets. 'w/o Co-occurrence' refers to removing Co-occurrence Frequency Encoding from DyG-Mamba, while 'w/o skip connection' denotes the removal of the skip connection in the SSM.

	UCI	US Legis	UN Trade
DyG-Mamba	96.79 ± 0.08	74.11 ± 2.32	68.55 ± 0.16
w/o Co-occurrence	93.09 ± 1.31	73.53 ± 2.43	66.32 ± 0.55
w/o skip-connection	95.37 ± 0.06	73.64 ± 2.51	67.64 ± 0.25

Performance with different sequence length. As shown in Table 14, DyG-Mamba can achieve better performance with longer input sequences. However, to ensure a fair comparison, we follow the same input sequence length as DyGFormer rather than incorporating additional historical information.

Hyperparameter Sensitivity. DyG-Mamba is insensitive to hyperparameters, see Figure 9. Therefore, the hyperparameters listed in Table 6 are applied consistently across all experiments.

Table 14: AP score across different sequence length. '*' denotes the sequence length used in DyGFormer.

	32	256	512	1024	2048	4096
uci	96.79±0.08*	96.82 ± 0.09	96.95 ± 0.06	97.38 ± 0.05	97.45 ± 0.03	97.47 ± 0.08
USLegis	73.99 ± 1.52	$74.11\pm2.32^*$	74.17 ± 2.14	74.34 ± 2.17	74.47 ± 2.44	74.96 ± 2.04

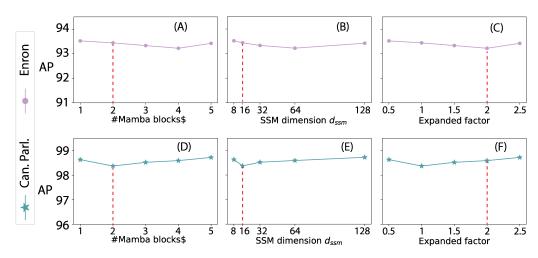


Figure 9: We tune the parameters of DyG-Mamba, including the number of blocks, the SSM dimension, and the expansion factor, on two datasets: Enron (A-C), Can.Parl. (D-F). For each parameter, we adjust its value while keeping the others fixed at DyG-Mamba's default setting (the red line).

Table 15: AUC-ROC for transductive dynamic link prediction with random, historical, and inductive negative sampling strategies.

	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	EdgeBank	TCL	GraphMixer	FreeDyG	DyGFormer	DyG-Mamba
	Wikipedia	96.33 ± 0.07	94.37 ± 0.09	96.67 ± 0.07	98.37 ± 0.07	98.54 ± 0.04	90.78 ± 0.00	95.84 ± 0.18	96.92 ± 0.03	$\textbf{99.41} \pm \textbf{0.01}$	98.91 ± 0.02	98.96 ± 0.00
	Reddit	98.31 ± 0.05	98.17 ± 0.05	98.47 ± 0.02	98.60 ± 0.06	99.01 ± 0.01	95.37 ± 0.00	97.42 ± 0.02	97.17 ± 0.02	$\textbf{99.50} \pm \textbf{0.01}$	99.15 ± 0.01	99.20 ± 0.00
	MOOC										87.91 ± 0.58	
	LastFM										93.05 ± 0.10	
	Enron										93.33 ± 0.13	
											96.30 ± 0.01	
rnd	UCI										94.49 ± 0.26	
	Can. Parl.			75.69 ± 0.78							97.76 ± 0.41	
	US Legis.			75.84 ± 1.99						N/A		78.27 ± 2.80
		69.62 ± 0.44								N/A		72.25 ± 0.07
		68.53 ± 0.95		52.83 ± 1.12 96.95 ± 0.08						N/A N/A	57.12 ± 0.62	
	Contact Avg. Rank	5.33	96.48 ± 0.14 6.75	7.00	97.54 ± 0.35 3.25	5.58	94.34 ± 0.00 8.17	94.15 ± 0.09 8.25	93.94 ± 0.02 6.58	N/A N/A	$\frac{98.53 \pm 0.01}{2.58}$	98.58 ± 0.01 1.50
											78.80 ± 1.95	
	Reddit										80.54 ± 0.29	
	MOOC										87.04 ± 0.35	
	LastFM										$\frac{78.78 \pm 0.35}{10.000}$	
	Enron										76.55 ± 0.52	
	Social Evo.										$\frac{97.28 \pm 0.07}{76.07 \pm 0.24}$	
hist										80.38 ± 0.26 N/A	76.97 ± 0.24 97.61 \pm 0.40	
	US Legis.	62.44 ± 1.11		70.86 ± 0.94 73.47 ± 5.25						N/A N/A	97.61 ± 0.40 90.77 ± 1.96	
	UN Trade			60.37 ± 0.68						N/A	73.86 ± 1.13	
	UN Vote			53.95 ± 3.15						N/A	64.27 ± 1.78	
	Contact			95.39 ± 0.43						N/A	97.17 ± 0.05	
	Avg. Rank	5.00	5.67	7.08	3.92	8.25	6.50	7.25	5.67	N/A	3.33	2.33
											75.09 ± 3.70 86.23 ± 0.51	
	Reddit MOOC										86.23 ± 0.51 80.76 ± 0.76	
	LastFM										$\frac{80.76 \pm 0.76}{69.25 \pm 0.36}$	
	Enron										74.07 ± 0.64	
											97.51 ± 0.04	
ind	UCI										$\frac{97.91 \pm 0.00}{65.96 \pm 1.18}$	
ma	Can. Parl.			72.47 ± 1.18						N/A		96.99 ± 0.65
	US Legis.			71.62 ± 5.42						N/A		90.51 ± 0.05
		66.82 ± 1.27								N/A	62.56 ± 1.51	
	UN Vote			53.04 ± 2.58						N/A	53.37 ± 1.26	
		94.47 ± 1.08								N/A	95.01 ± 0.15	
	Avg. Rank	6.75	7.08	6.00	5.33	5.75	6.08	6.00	5.67	N/A	3.83	2.50
_												

Table 16: AP for inductive dynamic link prediction with random, historical, and inductive negative sampling strategies.

	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	FreeDyG	DyGFormer	DyG-Mamba
rnd	Reddit MOOC LastFM Enron Social Evo. UCI Can. Parl. US Legis.	$\begin{array}{c} 96.50 \pm 0.13 \\ 79.63 \pm 1.92 \\ 81.61 \pm 3.82 \\ 80.72 \pm 1.39 \\ 91.96 \pm 0.48 \\ 79.86 \pm 1.48 \\ 53.92 \pm 0.94 \\ 54.93 \pm 2.29 \\ 59.65 \pm 0.77 \\ 56.64 \pm 0.96 \\ 94.34 \pm 1.45 \end{array}$	$\begin{array}{c} 96.09 \pm 0.11 \\ 81.07 \pm 0.44 \\ 83.02 \pm 1.48 \\ 74.55 \pm 3.95 \\ 90.04 \pm 0.47 \end{array}$	$\begin{array}{c} 97.09 \pm 0.04 \\ 85.50 \pm 0.19 \\ 78.63 \pm 0.31 \\ 67.05 \pm 1.51 \\ 91.41 \pm 0.16 \\ 79.54 \pm 0.48 \\ 55.18 \pm 0.79 \\ 51.00 \pm 3.11 \\ 61.03 \pm 0.18 \\ 52.24 \pm 1.46 \end{array}$	$\begin{array}{c} 97.50 \pm 0.07 \\ 89.04 \pm 1.17 \\ 81.45 \pm 4.29 \\ 77.94 \pm 1.02 \\ 90.77 \pm 0.86 \\ 88.12 \pm 2.05 \\ 54.10 \pm 0.93 \\ 58.63 \pm 0.37 \\ 58.31 \pm 3.15 \\ 58.85 \pm 2.51 \end{array}$	$\begin{array}{c} 98.62 \pm 0.01 \\ 81.42 \pm 0.24 \\ 89.42 \pm 0.07 \\ 86.35 \pm 0.51 \\ 79.94 \pm 0.18 \\ 92.73 \pm 0.06 \\ 55.80 \pm 0.69 \\ 53.17 \pm 1.20 \\ \underline{65.24 \pm 0.21} \\ 49.94 \pm 0.45 \end{array}$	$\begin{array}{c} 94.09 \pm 0.07 \\ 80.60 \pm 0.22 \\ 73.53 \pm 1.66 \\ 76.14 \pm 0.79 \\ 91.55 \pm 0.09 \\ 87.36 \pm 2.03 \\ 54.30 \pm 0.66 \\ 52.59 \pm 0.97 \\ 62.21 \pm 0.12 \\ 51.60 \pm 0.97 \end{array}$	$\begin{array}{c} 95.26\pm0.02\\ 81.41\pm0.21\\ 82.11\pm0.42\\ 75.88\pm0.48\\ 91.86\pm0.06\\ 91.19\pm0.42\\ 55.91\pm0.82\\ 50.71\pm0.76\\ 62.17\pm0.31\\ 50.68\pm0.44 \end{array}$	$\begin{array}{c} 98.91 \pm 0.01 \\ 87.75 \pm 0.62 \\ 94.89 \pm 0.01 \\ 89.69 \pm 0.17 \\ \textbf{94.76} \pm \textbf{0.05} \end{array}$	54.28 ± 2.87 64.55 ± 0.62 55.93 ± 0.39	98.91 ± 0.01 89.98 ± 0.46 95.16 ± 0.05 90.97 ± 0.01 93.17 ± 0.05
hist	Reddit MOOC LastFM Enron Social Evo. UCI Can. Parl. US Legis.	$\begin{array}{c} 62.34 \pm 0.54 \\ 63.22 \pm 1.55 \\ 70.39 \pm 4.31 \\ 65.86 \pm 3.71 \\ 88.51 \pm 0.87 \\ 63.11 \pm 2.27 \\ 52.60 \pm 0.88 \\ 52.94 \pm 2.11 \\ 55.46 \pm 1.19 \\ 61.04 \pm 1.30 \\ \hline 90.42 \pm 2.34 \end{array}$	$\begin{array}{c} 61.60 \pm 0.72 \\ 62.93 \pm 1.24 \\ 71.45 \pm 1.76 \\ 62.08 \pm 2.27 \\ 88.72 \pm 1.10 \end{array}$	$\begin{array}{c} \overline{63.47\pm0.36} \\ 76.73\pm0.29 \\ 76.27\pm0.25 \\ 61.40\pm1.31 \\ 93.97\pm0.54 \\ 70.52\pm0.93 \\ 56.72\pm0.47 \\ 51.83\pm3.95 \\ 55.28\pm0.71 \\ 53.05\pm3.10 \end{array}$	$\begin{array}{c} 64.85\pm0.85\\ 77.07\pm3.41\\ 66.65\pm6.11\\ 62.91\pm1.16\\ 90.66\pm1.62\\ 70.78\pm0.78\\ 54.42\pm0.77\\ \hline 61.18\pm1.10\\ \hline 52.80\pm3.19\\ \hline \textbf{63.74}\pm3.00 \end{array}$	$\begin{array}{c} 63.67 \pm 0.41 \\ 74.68 \pm 0.68 \\ 71.33 \pm 0.47 \\ 60.70 \pm 0.36 \\ 79.83 \pm 0.38 \\ 64.54 \pm 0.47 \\ 57.14 \pm 0.07 \\ 55.56 \pm 1.71 \\ 55.00 \pm 0.38 \\ 47.98 \pm 0.84 \end{array}$	$\begin{array}{c} 60.83 \pm 0.25 \\ 74.27 \pm 0.53 \\ 65.78 \pm 0.65 \\ 67.11 \pm 0.62 \\ 94.10 \pm 0.31 \\ 76.71 \pm 1.00 \\ 55.71 \pm 0.74 \\ 53.87 \pm 1.41 \\ \underline{55.76 \pm 1.03} \\ \overline{54.19 \pm 2.17} \end{array}$	$\begin{array}{c} 64.50 \pm 0.26 \\ 74.00 \pm 0.97 \\ 76.42 \pm 0.22 \\ 72.37 \pm 1.37 \\ 94.01 \pm 0.47 \\ \underline{81.66 \pm 0.49} \\ 55.84 \pm 0.73 \\ 52.03 \pm 1.02 \\ 54.94 \pm 0.97 \\ 48.09 \pm 0.43 \end{array}$	$\begin{array}{c} 66.02 \pm 0.41 \\ \underline{81.63 \pm 0.33} \\ \overline{77.28 \pm 0.21} \\ \underline{73.01 \pm 0.88} \\ \underline{96.69 \pm 0.14} \end{array}$	$\begin{array}{c} 56.31 \pm 3.46 \\ 53.20 \pm 1.07 \\ 52.63 \pm 1.26 \end{array}$	$\begin{array}{c} \textbf{66.74} \pm \textbf{0.13} \\ \textbf{81.64} \pm \textbf{0.67} \\ \textbf{79.22} \pm \textbf{0.33} \\ \textbf{75.12} \pm \textbf{1.43} \\ \textbf{95.45} \pm \textbf{0.30} \end{array}$
ind	Reddit MOOC LastFM Enron Social Evo. UCI Can. Parl. US Legis.	$\begin{array}{c} 62.32 \pm 0.54 \\ 63.22 \pm 1.55 \\ 70.39 \pm 4.31 \\ 65.86 \pm 3.71 \\ 88.51 \pm 0.87 \\ 63.16 \pm 2.27 \\ 52.58 \pm 0.86 \\ 52.94 \pm 2.11 \\ 55.43 \pm 1.20 \\ 61.17 \pm 1.33 \end{array}$	$\begin{array}{c} 61.58 \pm 0.72 \\ 62.92 \pm 1.24 \\ 71.45 \pm 1.75 \\ 62.08 \pm 2.27 \\ 88.72 \pm 1.10 \end{array}$	$\begin{array}{c} 63.40 \pm 0.36 \\ 76.72 \pm 0.30 \\ 76.28 \pm 0.25 \\ 61.40 \pm 1.30 \\ 93.97 \pm 0.54 \\ 70.49 \pm 0.93 \\ 56.46 \pm 0.50 \\ 51.83 \pm 3.95 \\ 55.58 \pm 0.68 \\ 53.08 \pm 3.10 \end{array}$	$\begin{array}{c} 64.84\pm0.84\\ 77.07\pm3.40\\ 69.46\pm4.65\\ 62.90\pm1.16\\ 90.65\pm1.62\\ 70.73\pm0.79\\ 54.18\pm0.73\\ \hline 61.18\pm1.10\\ \hline 52.80\pm3.24\\ \hline \textbf{63.71}\pm2.97 \end{array}$	$\begin{array}{c} 63.65 \pm 0.41 \\ 74.69 \pm 0.68 \\ 71.33 \pm 0.47 \\ 60.72 \pm 0.36 \\ 79.83 \pm 0.39 \\ 64.54 \pm 0.47 \\ 57.06 \pm 0.08 \\ 55.56 \pm 1.71 \\ 54.97 \pm 0.38 \\ 48.01 \pm 0.82 \end{array}$	$\begin{array}{c} 60.81 \pm 0.26 \\ 74.28 \pm 0.53 \\ 65.78 \pm 0.65 \\ 67.11 \pm 0.62 \\ 94.10 \pm 0.32 \\ 76.65 \pm 0.99 \\ 55.46 \pm 0.69 \\ 53.87 \pm 1.41 \\ \underline{55.66 \pm 0.98} \\ \overline{54.13 \pm 2.16} \end{array}$	$\begin{array}{c} 64.49 \pm 0.25 \\ 73.99 \pm 0.97 \\ 76.42 \pm 0.22 \\ 72.37 \pm 1.38 \\ 94.01 \pm 0.47 \\ \hline 81.64 \pm 0.49 \\ \hline 55.76 \pm 0.65 \\ 52.03 \pm 1.02 \\ 54.88 \pm 1.01 \\ 48.10 \pm 0.40 \end{array}$	$\begin{array}{c} {\bf 64.98 \pm 0.20} \\ {\bf 81.41 \pm 0.31} \\ {\bf 77.01 \pm 0.43} \\ {\bf 72.85 \pm 0.81} \\ {\bf 96.91 \pm 0.12} \end{array}$	$\overline{56.31 \pm 3.46}$ 52.56 ± 1.70 52.61 ± 1.25	$\begin{array}{c} \textbf{66.74} \pm \textbf{0.13} \\ \textbf{81.64} \pm \textbf{0.67} \\ \textbf{79.22} \pm \textbf{0.33} \\ \textbf{75.12} \pm \textbf{1.43} \\ \textbf{95.45} \pm \textbf{0.30} \end{array}$

Table 17: AUC-ROC for inductive dynamic link prediction with random, historical, and inductive negative sampling strategies.

	Datasets	JODIE	DyRep	TGAT	TGN	CAWN	TCL	GraphMixer	FreeDyG	DyGFormer	DyG-Mamba
_	Wikipedia	194 33 + 0 27	91.49 ± 0.45	95 90 + 0 09	97.72 + 0.03	08 03 + 0 04	95.57 ± 0.20	96.30 ± 0.04	99.01 ± 0.02	98.48 ± 0.03	98 55 ± 0.01
	Reddit									98.71 ± 0.01	
	MOOC									87.62 ± 0.51	
	LastFM									94.08 ± 0.08	
	Enron									90.69 ± 0.26	
	Social Evo.									95.29 ± 0.03	
rnd	UCI	78.80 ± 0.94	58.08 ± 1.81	77.64 ± 0.38	86.68 ± 2.29	90.40 ± 0.11	84.49 ± 1.82	89.30 ± 0.57	$\textbf{93.01} \pm \textbf{0.08}$	92.63 ± 0.13	92.05 ± 0.23
	Can. Parl.	53.81 ± 1.14	55.27 ± 0.49	56.51 ± 0.75	55.86 ± 0.75	58.83 ± 1.13	55.83 ± 1.07	58.32 ± 1.08	N/A	89.33 ± 0.48	$\textbf{97.11} \pm \textbf{0.09}$
	US Legis.	58.12 ± 2.35	61.07 ± 0.56	48.27 ± 3.50	$\textbf{62.38} \pm \textbf{0.48}$	51.49 ± 1.13	50.43 ± 1.48	47.20 ± 0.89	N/A	53.21 ± 3.04	52.73 ± 1.24
	UN Trade		58.82 ± 0.98						N/A		69.37 ± 0.06
	UN Vote		55.13 ± 3.46						N/A		60.03 ± 0.02
	Contact		91.89 ± 0.38						N/A		98.33 ± 0.01
	Avg. Rank	5.67	6.83	6.25	4.25	5.17	6.92	6.08	N/A	<u>2.25</u>	1.67
	Wikipedia	$ 61.86 \pm 0.53 $	57.54 ± 1.09	78.38 ± 0.20	75.75 ± 0.29	62.04 ± 0.65	79.79 ± 0.96	$\textbf{82.87} \pm \textbf{0.21}$	82.08 ± 0.32	68.33 ± 2.82	69.73 ± 0.48
	Reddit	61.69 ± 0.39	60.45 ± 0.37	64.43 ± 0.27	64.55 ± 0.50	64.94 ± 0.21	61.43 ± 0.26	64.27 ± 0.13	66.79 ± 0.31	64.81 ± 0.25	$\textbf{67.82} \pm \textbf{0.30}$
	MOOC									80.77 ± 0.63	
	LastFM									70.73 ± 0.37	
	Enron									65.78 ± 0.42	
											96.03 ± 0.24
hist	UCI									65.55 ± 1.01	
			52.38 ± 0.46						N/A		91.54 ± 0.39
			67.94 ± 0.98						N/A		56.15 ± 0.15
	UN Trade		57.90 ± 1.33						N/A		62.81 ± 0.21
	UN Vote		63.98 ± 2.12						N/A		62.69 ± 1.23
	Contact		88.88 ± 0.68						N/A		94.18 ± 0.36
	Avg. Rank	5.92	7.00	5.33	5.17	6.25	5.08	4.58	N/A	<u>3.50</u>	2.17
										68.33 ± 2.82	
	Reddit									67.82 ± 0.30	
	MOOC									80.77 ± 0.63	
	LastFM									70.73 ± 0.37	
	Enron										$\textbf{75.35} \pm \textbf{1.06}$
										$\textbf{96.91} \pm \textbf{0.09}$	
ind	UCI									65.58 ± 1.00	
	Can. Parl.		52.35 ± 0.52						N/A		91.54 ± 0.39
	US Legis.		67.94 ± 0.98						N/A N/A		56.15 ± 0.15 62.81 ± 0.21
	UN Trade UN Vote		57.87 ± 1.36 64.10 ± 2.10						N/A N/A		62.81 ± 0.21 62.69 ± 1.23
	Contact		88.87 ± 0.67						N/A N/A		62.69 ± 1.23 94.18 ± 0.36
	Avg. Rank	5.92	6.92	5.25	88.85 ± 1.39 5.17	6.33	5.08	90.04 ± 0.29 4.67	N/A N/A	$\frac{94.14 \pm 0.26}{3.42}$	94.18 ± 0.36 2.25
	Avg. Kank	3.92	0.92	3.23	5.17	0.33	5.08	4.07	IV/A	<u>3.42</u>	4,45

E Proofs for Theorems

E.1 Proofs of Theorem 1

Proof. Consider the following linear time-invariant (LTI) system on the interval $[t_{k-1}, t_k]$, where $t_k - t_{k-1} = \Delta t_{k,i}$ for the *i*-th coordinate:

$$\frac{d\boldsymbol{h}(s)}{ds} = \boldsymbol{A}_k \boldsymbol{h}(s) + \boldsymbol{B}_k \boldsymbol{m}_k, \quad \boldsymbol{h}(t_{k-1}) = \boldsymbol{h}_{k-1},$$

where $\boldsymbol{A}_k = \operatorname{diag}(\lambda_1,\dots,\lambda_n)$ and $\boldsymbol{B}_k = \begin{bmatrix} \boldsymbol{B}_{k,1},\dots,\boldsymbol{B}_{k,n} \end{bmatrix}^{\top}$ (also arranged diagonally for each coordinate). Since \boldsymbol{A}_k is diagonal with $\operatorname{Re}(\lambda_i) < 0$, the fundamental matrix solution for this system (i.e., the matrix exponential) is also diagonal and can be written as

$$e^{\mathbf{A}_k \Delta t_k} = \operatorname{diag}(e^{\lambda_1 \Delta t_{k,1}}, \dots, e^{\lambda_n \Delta t_{k,n}}).$$

Step 1: Expressing \overline{A}_k .

By definition of the matrix exponential for a diagonal matrix:

$$\overline{A}_k = e^{A_k \Delta t_k} = \operatorname{diag}(e^{\lambda_1 \Delta t_{k,1}}, \dots, e^{\lambda_n \Delta t_{k,n}}).$$

Step 2: Expressing \overline{B}_k .

We compute the convolution term associated with the inhomogeneous part $B_k m_k$. For a diagonal system, the solution for each coordinate i is:

$$\boldsymbol{h}_{k,i} = e^{\lambda_i \Delta t_{k,i}} \, \boldsymbol{h}_{k-1,i} + \int_0^{\Delta t_{k,i}} e^{\lambda_i (\Delta t_{k,i} - \tau)} \, \boldsymbol{B}_{k,i} \, \boldsymbol{m}_{k,i} \, d\tau.$$

Since $m_{k,i}$ is constant w.r.t. τ and $B_{k,i}$ is also constant in this interval, we can pull them out of the integral:

$$m{h}_{k,i} = e^{\lambda_i \Delta t_{k,i}} \, m{h}_{k-1,i} \, + \, m{B}_{k,i} \, m{m}_{k,i} \, \int_0^{\Delta t_{k,i}} e^{\lambda_i (\Delta t_{k,i} - au)} \, d au.$$

Evaluating the integral:

$$\int_{0}^{\Delta t_{k,i}} e^{\lambda_{i}(\Delta t_{k,i}-\tau)} d\tau = \int_{0}^{\Delta t_{k,i}} e^{\lambda_{i}(\Delta t_{k,i}-u)} du = \left[-\frac{1}{\lambda_{i}} e^{\lambda_{i}(\Delta t_{k,i}-u)}\right]_{0}^{\Delta t_{k,i}} = \frac{1}{\lambda_{i}} (e^{\lambda_{i}\Delta t_{k,i}}-1).$$

Hence,

$$\boldsymbol{h}_{k,i} = e^{\lambda_i \Delta t_{k,i}} \, \boldsymbol{h}_{k-1,i} + \frac{1}{\lambda_i} \left(e^{\lambda_i \Delta t_{k,i}} - 1 \right) \boldsymbol{B}_{k,i} \, \boldsymbol{m}_{k,i}.$$

This leads to

$$\overline{\boldsymbol{B}}_{k} = \operatorname{diag}\left(\lambda_{1}^{-1}\left(e^{\lambda_{1}\Delta t_{k,1}}-1\right)\boldsymbol{B}_{k,1},\ldots,\lambda_{n}^{-1}\left(e^{\lambda_{n}\Delta t_{k,n}}-1\right)\boldsymbol{B}_{k,n}\right),\,$$

since each diagonal entry of \overline{B}_k matches the integral factor $\lambda_i^{-1}(e^{\lambda_i \Delta t_{k,i}}-1) B_{k,i}$.

Step 3: Final Form of $h_{k,i}$.

Combining the above results yields the stated coordinate-wise update for $h_{k,i}$:

$$m{h}_{k,i} = e^{\lambda_i \Delta t_{k,i}} \, m{h}_{k-1,i} \, + \, \lambda_i^{-1} \left(e^{\lambda_i \Delta t_{k,i}} - 1 \right) m{B}_{k,i} \, m{m}_{k,i}.$$

This confirms the forms of both \overline{A}_k and \overline{B}_k as well as the final expression for each coordinate $h_{k,i}$.

E.2 Proofs of Theorem 2

Proof. Since DyG-Mamba has three core parameters, *i.e.*, Δ , B and C that control its effectiveness. We define Theorem 2 to explain the main effects of the parameters B and C.

We first copy Eq.(8) from the paper, SSM-based node representation learning process, as follows

$$\boldsymbol{h}_{k} = \bar{\boldsymbol{A}}_{k} \boldsymbol{h}_{k-1} + \bar{\boldsymbol{B}}_{k} \boldsymbol{u}_{k}, \quad \hat{\boldsymbol{z}}_{k}^{\tau} = \bar{\boldsymbol{C}}_{k} \boldsymbol{h}_{k}, \tag{17}$$

 h_k in Eq.(17) can be further decomposed as follows:

$$h_k = \prod_{i=0}^{k-2} \bar{A}_{k-i} \bar{B}_1 u_1 + \dots + \prod_{i=0}^{k-j} \bar{A}_{k-i} \bar{B}_{j-1} u_{j-1} + \dots + \bar{B}_k u_k,$$
(18)

Then, considering ${\bf A}$ is one fixed parameter and $\bar{{\bf A}}_k = \exp(\Delta t_k {\bf A}), \, \widehat{z}_k^{\tau}$ in Eq.(8) can be formulated as:

$$\hat{z}_{k}^{T} = \bar{C}_{k} \prod_{i=0}^{k-2} \bar{A}_{k-i} \bar{B}_{1} u_{1} + \dots + \bar{C}_{k} \prod_{i=0}^{k-j} \bar{A}_{k-i} \bar{B}_{j-1} u_{j-1} + \dots + \bar{C}_{k} \bar{B}_{k} u_{k},
= e^{(\sum_{i=0}^{k-2} \Delta t_{k-i} A)} \bar{C}_{k} \bar{B}_{1} u_{1} + \dots + e^{(\sum_{i=0}^{k-j-1} \Delta t_{k-i} A)} \bar{C}_{k} \bar{B}_{j} u_{j} + \dots + \bar{C}_{k} \bar{B}_{k} u_{k},$$
(19)

where u_k is the k-th input of M_u^{τ} , i.e., $u_k = M_u^{\tau}[k,:]$, and parameters $B = \operatorname{Linear}_B(M_u^{\tau})$ and $C = \operatorname{Linear}_C(M_u^{\tau})$. Thus, \bar{C}_k can be considered as one query of k-th input u_k and \bar{B}_j can be considered as the key of j-th input u_j . Then, $\bar{C}_k\bar{B}_j$ can measure the similarity between u_k and u_j , similar to the self-attention mechanism. According to the above-mentioned proof, we can conclude that parameters B and C can measure the similarity between current input to the previous ones and selectively copy the previous input.

E.3 Proofs of Theorem 3

Proof. Step 1: Parameter perturbation bounds. Under spectral normalization constraints $(\|W_B\|_2 \le 1, \|W_C\|_2 \le 1)$, the parameter perturbations induced by input noise ΔM_u^{τ} satisfy:

$$\|\Delta B\| \le \|W_B\|_2 \|\Delta M_u^{\tau}\| \le \|\Delta M_u^{\tau}\|, \quad \|\Delta C\| \le \|W_C\|_2 \|\Delta M_u^{\tau}\| \le \|\Delta M_u^{\tau}\|. \tag{20}$$

This follows directly from the Lipschitz continuity enforced by spectral normalization, where the spectral norm $\|\mathbf{W}\|_2$ precisely defines the maximum amplification factor of the linear transformation.

Step 2: Perturbation propagation analysis. From the state update decomposition in Theorem 4.2, the output perturbation $\Delta \widehat{m}_k^{\tau}$ can be expressed as:

$$\|\Delta\widehat{\boldsymbol{m}}_{k}^{\tau}\| \leq \sum_{j=1}^{k} \left(\|\Delta\overline{\boldsymbol{C}}_{k}\| \|\overline{\boldsymbol{B}}_{j}\| + \|\overline{\boldsymbol{C}}_{k}\| \|\Delta\overline{\boldsymbol{B}}_{j}\| \right) \|\boldsymbol{m}_{j}\| \prod_{i=0}^{k-j-1} \|e^{\Delta t_{k-i}\boldsymbol{A}_{k-i}}\|.$$
 (21)

Using $\|\overline{B}_j\|$, $\|\overline{C}_k\| \leq 1$ (spectral normalization) and $\|\Delta \overline{B}_j\|$, $\|\Delta \overline{C}_k\| \leq \|\Delta M_u^{\tau}\|$ (Step 1), we derive:

$$\|\Delta \widehat{\boldsymbol{m}}_{k}^{\tau}\| \leq \|\Delta \boldsymbol{M}_{u}^{\tau}\| \sum_{j=1}^{k} e^{\gamma(t_{k} - t_{j})}. \tag{22}$$

Here, the exponential term $\prod_{i=0}^{k-j-1} \|e^{\Delta t_{k-i} A_{k-i}}\| \le e^{\gamma(t_k-t_j)}$ arises from the eigenvalue constraint $\gamma = \max_i \operatorname{Re}(\lambda_i(A)) < 0$.

Step 3: Stability via exponential decay. For $\gamma < 0$, the summation over exponentially decaying terms can be approximated as:

$$\sum_{j=1}^{k} e^{\gamma(t_k - t_j)} \approx \int_0^T e^{\gamma(T - t)} dt = \frac{1}{|\gamma|} \left(1 - e^{\gamma T} \right), \tag{23}$$

where T is the total sequence duration. This yields the final perturbation bound:

$$\|\Delta \widehat{\boldsymbol{m}}_{k}^{\tau}\| \le \kappa \|\Delta \boldsymbol{M}_{u}^{\tau}\|, \quad \text{where } \kappa = \frac{1}{|\gamma|} (1 - e^{\gamma T}).$$
 (24)

Impact Statement

This paper aims to advance the field of Machine Learning by introducing a Mamba-based framework for continuous-time dynamic graph modeling. We examine its potential impacts from two primary perspectives as follows. (i) dynamic graph modeling. With the rapid expansion of social and economic networks, dynamic graph modeling has emerged as a prominent research topic in the machine learning community. In contrast to existing methods, our approach introduces a novel Mamba-based framework, which incorporates irregular timespans as control signals for continuous SSMs. This design enhances the model's ability to effectively and efficiently capture long-term temporal dependency on dynamic graphs. (ii) Time Information Effects. Despite significant advances in dynamic graph modeling, there is still a lack of theoretical foundations regarding the influence of temporal information on the evolution of dynamic graphs. Our work highlights the significant potential of the timespan-based memory forgetting mechanism to deepen the theoretical understanding of time in this domain. Overall, we do not foresee any direct negative societal implications resulting from this research. Although more advanced modeling capabilities can be applied across a range of domains, we believe that this methodology itself does not introduce any ethical or societal concerns beyond those typically associated with improvements in machine learning.