
Provably expressive temporal graph networks

Anonymous Author(s)

Affiliation

email

Abstract

1 Temporal graph networks (TGNs) have gained prominence as models for embed-
2 ding dynamic interactions, but little is known about their theoretical underpinnings.
3 We establish fundamental results about the representational power and limits of the
4 two main categories of TGNs: those that aggregate temporal walks (WA-TGNs),
5 and those that augment local message passing with recurrent memory modules
6 (MP-TGNs). Specifically, novel constructions reveal the inadequacy of MP-TGNs
7 and WA-TGNs, proving that neither category subsumes the other. We extend the
8 1-WL (Weisfeiler-Leman) test to temporal graphs, and show that the most powerful
9 MP-TGNs should use injective updates, as in this case they become as expressive
10 as the temporal WL. Also, we show that sufficiently deep MP-TGNs cannot benefit
11 from memory, and MP/WA-TGNs fail to compute graph properties such as girth.

12 These theoretical insights lead us to introduce PINT — a novel architecture that
13 leverages injective temporal message passing and relative positional features. Im-
14 portantly, PINT is provably more expressive than both MP-TGNs and WA-TGNs.
15 Our experiments demonstrate that PINT significantly outperforms existing TGNs
16 on several real-world benchmarks.

17 1 Introduction

18 Graph neural networks (GNNs) [11, 30, 38] have recently led to breakthroughs in many applications
19 [7, 28, 31] by resorting to message passing between neighboring nodes in input graphs. While
20 message passing imposes an important inductive bias, it does not account for the dynamic nature of
21 interactions in time-evolving graphs arising from many real-world domains such as social networks
22 and bioinformatics [16, 39]. In several scenarios, these temporal graphs are only given as a sequence
23 of timestamped events. Recently, temporal graph nets (TGNs) [16, 27, 32, 37, 41] have emerged as a
24 prominent learning framework for temporal graphs and have become particularly popular due to their
25 outstanding predictive performance. Aiming at capturing meaningful structural and temporal patterns,
26 TGNs combine a variety of building blocks, such as self-attention [33, 34], time encoders [15, 40],
27 recurrent models [5, 13], and message passing [10].

28 Unraveling the learning capabilities of (temporal) graph networks is imperative to understanding
29 their strengths and pitfalls, and designing better, more nuanced models that are both theoretically
30 well-grounded and practically efficacious. For instance, the enhanced expressivity of higher-order
31 GNNs has roots in the inadequacy of standard message-passing GNNs to separate graphs that are
32 indistinguishable by the Weisfeiler-Leman isomorphism test, known as 1-WL test or color refinement
33 algorithm [21, 22, 29, 36, 42]. Similarly, many other notable advances on GNNs were made possible
34 by untangling their ability to generalize [9, 17, 35], extrapolate [44], compute graph properties [4, 6, 9],
35 and express Boolean classifiers [1]; by uncovering their connections to distributed algorithms [19, 29],
36 graph kernels [8], dynamic programming [43], diffusion processes [3], graphical models [45], and
37 combinatorial optimization [2]; and by analyzing their discriminative power [20, 23]. In stark contrast,
38 the theoretical foundations of TGNs remain largely unexplored. For instance, unresolved questions
39 include: How does the expressive power of existing TGNs compare? When do TGNs fail? Can we
40 improve the expressiveness of TGNs? What are the limits on the power of TGNs?

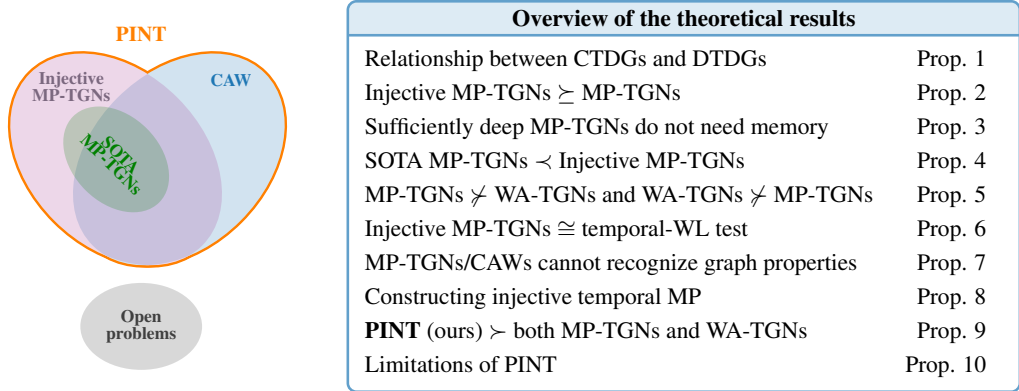


Figure 1: Schematic diagram and summary of our contributions.

41 We establish a series of results to address these fundamental questions. We begin by showing that
 42 discrete-time dynamic graphs (DTDGs) can always be converted to continuous-time analogues
 43 (CTDGs) without loss of information, so we can focus on analyzing the ability of TGNs to distinguish
 44 nodes/links of CTDGs. We consider a general framework for message-passing TGNs (MP-TGNs)
 45 [27] that subsumes a wide variety of methods [e.g., 16, 32, 41]. We prove that equipping MP-TGNs
 46 with injective aggregation and update functions leads to the class of most expressive *anonymous*
 47 *MP-TGNs* (i.e., those that do not leverage node ids). Extending the color-refinement algorithm to
 48 temporal settings, we show that these most powerful MP-TGNs are as expressive as the temporal
 49 WL method. Notably, existing MP-TGNs do not enforce injectivity. We also delineate the role of
 50 memory in MP-TGNs: nodes in a network with only a few layers of message passing fail to aggregate
 51 information from a sufficiently wide receptive field (i.e., from distant nodes), so memory serves
 52 to offset this highly local view with additional global information. In contrast, sufficiently deep
 53 architectures obviate the need for memory modules.

54 Different from MP-TGNs, walk-aggregating TGNs (WA-TGNs) such as CAW [37] obtain represen-
 55 tations from anonymized temporal walks. We provide constructions that expose shortcomings of
 56 each framework, establishing that WA-TGNs can distinguish links in cases where MP-TGNs fail and
 57 vice-versa. Consequently, neither class is more expressive than the other. Additionally, we show that
 58 MP-TGNs and CAWs cannot decide temporal graph properties such as diameter, girth, or number of
 59 cycles. Strikingly, our analysis unravels the subtle relationship between the walk computations in
 60 CAWs and the MP steps in MP-TGNs.

61 Equipped with these theoretical insights, we propose PINT (short for *position-encoding injective*
 62 *temporal graph net*), founded on a new temporal layer that leverages the strengths of both MP-TGNs
 63 and WA-TGNs. Like the most expressive MP-TGNs, PINT defines injective message passing and
 64 update steps. PINT also augments memory states with novel relative positional features, and these
 65 features can replicate all the discriminative benefits available to WA-TGNs. Interestingly, the time
 66 complexity of computing our positional features is less severe than the sampling overhead in CAW,
 67 thus PINT can often be trained faster than CAW. Importantly, we establish that PINT is provably
 68 more expressive than CAW as well as MP-TGNs.

69 **Our contributions** are three-fold:

- 70 • a rigorous theoretical foundation for TGNs is laid - elucidating the role of memory, benefits
 71 of injective message passing, limits of existing TGN models, temporal extension of the
 72 1-WL test and its implications, impossibility results about temporal graph properties, and
 73 the relationship between main classes of TGNs - as summarized in Figure 1;
- 74 • explicit injective temporal functions are introduced, and a novel method for temporal graphs
 75 is proposed that is provably more expressive than state-of-the-art TGNs;
- 76 • extensive empirical investigations underscore practical benefits of this work. The proposed
 77 method is either competitive or significantly better than existing models on several real
 78 benchmarks for dynamic link prediction, in transductive as well as inductive settings.

79 **2 Preliminaries**

80 We denote a *static graph* G as a tuple $(V, E, \mathcal{X}, \mathcal{E})$, where $V = \{1, 2, \dots, n\}$ denotes the set of
 81 nodes and $E \subseteq V \times V$ the set of edges. Each node $u \in V$ has a feature vector $x_u \in \mathcal{X}$ and each
 82 edge $(u, v) \in E$ has a feature vector $e_{uv} \in \mathcal{E}$, where \mathcal{X} and \mathcal{E} are countable sets of features.

83 **Dynamic graphs** can be roughly split according to their discrete- or continuous-time nature [14].
 84 A *discrete-time dynamic graph* (DTDG) is of a sequence of graph snapshots (G_1, G_2, \dots) usually
 85 sampled at regular intervals, each snapshot being a static graph $G_t = (V_t, E_t, \mathcal{X}_t, \mathcal{E}_t)$.

86 A *continuous-time dynamic graph* (CTDG) evolves with node- and edge-level *events*, such as addition
 87 and deletion. We represent a CTDG as a sequence of time-stamped graphs $(G(t_0), G(t_1), \dots)$ such
 88 that $t_k < t_{k+1}$, and $G(t_{k+1})$ results from updating $G(t_k)$ with all events at time t_{k+1} . We assume
 89 no event occurs between t_k and t_{k+1} . We denote an interaction (i.e., edge addition event) between
 90 nodes u and v at time t as a tuple (u, v, t) associated with a feature vector $e_{uv}(t)$. Unless otherwise
 91 stated, interactions correspond to undirected edges, i.e., (u, v, t) is a shorthand for $(\{u, v\}, t)$.

92 Noting the CTDGs allow for finer (irregular) temporal resolution, we now formalize the intuition that
 93 DTDGs can be reduced to, and thus analyzed as CTDGs, but the converse does not always hold.

94 **Proposition 1** (Relationship between DTDG and CTDG). *For any DTDG we can build a CTDG with*
 95 *the same sets of node and edge features that contains the same information, i.e., we can reconstruct*
 96 *the original DTDG from the converted CTDG. The converse holds if the CTDG timestamps form a*
 97 *subset of a uniformly spaced countable set.*

98 Following the usual practice [16, 37, 41], we focus on CTDGs with edge addition events (see
 99 Appendix E for a discussion on deletion). Thus, we can represent temporal graphs as sets $\mathcal{G}(t) =$
 100 $\{(u_k, v_k, t_k) \mid t_k < t\}$. We also assume each distinct node v in $\mathcal{G}(t)$ has an initial feature vector x_v .

101 **Message-passing temporal graph nets (MP-TGNs).** Rossi et al. [27] introduced MP-TGN as
 102 a general representation learning framework for temporal graphs. The goal is to encode the graph
 103 dynamics into node embeddings, capturing information that is relevant for the task at hand. To
 104 achieve this, MP-TGNs rely on three main ingredients: memory, aggregation, and update. Memory
 105 comprises a set of vectors that summarizes the history of each node, and is updated using a recurrent
 106 model whenever an event occurs. The aggregation and update components resemble those in message-
 107 passing GNNs, where the embedding of each node is refined using messages from its neighbors.

108 We define the *temporal neighborhood* of node v at time t as $\mathcal{N}(v, t) = \{(u, e_{uv}(t'), t') \mid \exists (u, v, t') \in$
 109 $\mathcal{G}(t)\}$, i.e., the set of neighbor/feature/timestamp triplets from all interactions of node v prior to t .

110 MP-TGNs compute the temporal representation $h_v^{(\ell)}(t)$ of v at layer ℓ by recursively applying

$$\tilde{h}_v^{(\ell)}(t) = \text{AGG}^{(\ell)}(\{\{h_u^{(\ell-1)}(t), t - t', e\} \mid (u, e, t') \in \mathcal{N}(v, t)\}\}) \quad (1)$$

$$h_v^{(\ell)}(t) = \text{UPDATE}^{(\ell)}\left(h_v^{(\ell-1)}(t), \tilde{h}_v^{(\ell)}(t)\right), \quad (2)$$

111 where $\{\{\cdot\}\}$ denotes multisets, $h_v^{(0)}(t) = s_v(t)$ is the *state* of v at time t , and $\text{AGG}^{(\ell)}$ and $\text{UPDATE}^{(\ell)}$
 112 are arbitrary parameterized functions. The memory block updates the states as events occur. Let
 113 $\mathcal{J}(v, t)$ be the set of events involving v at time t . The state of v is updated due to $\mathcal{J}(v, t)$ as

$$m_v(t) = \text{MSGAGG}(\{\{s_v(t^-), s_u(t^-), t - t_v, e_{vu}(t)\} \mid (v, u, t) \in \mathcal{J}(v, t)\}\}) \quad (3)$$

$$s_v(t) = \text{MEMORY}(s_v(t^-), m_v(t)), \quad (4)$$

114 where $s_v(0) = x_v$ (initial node features), $s_v(t^-)$ denotes the state of v right before time t , and t_v
 115 denotes the time of the last update to v . MSGAGG combines information from simultaneous events
 116 involving node v and MEMORY usually implements a GRU [5]. Notably, some MP-TGNs do not use
 117 memory, or equivalently, they employ *identity memory*, i.e., $s_v(t) = x_v$. See Appendix A for details.

118 **Causal Anonymous Walks (CAWs).** Wang et al. [37] proposed CAW as an effective approach for
 119 link prediction on temporal graphs. To predict if an event (u, v, t) occurs, CAW first obtains sets S_u
 120 and S_v of temporal walks starting at nodes u and v . An L -length *temporal walk* is represented as $W =$
 121 $((w_1, t_1), (w_2, t_2), \dots, (w_L, t_L))$, with $t_1 > t_2 > \dots > t_L$ and $(w_{i-1}, w_i, t_i) \in \mathcal{G}(t)$ for all i . Then,
 122 CAW anonymizes walks replacing each node w with a set $I_{\text{CAW}}(w; S_u, S_v) = \{g(w; S_u), g(w; S_v)\}$
 123 of two feature vectors. The ℓ -th entry of $g(w; S_u)$ stores how many times w appears at the ℓ -th
 124 position in a walk of S_u , i.e. $g(w, S_u)[\ell] = |\{W \in S_u : (w, t_\ell) = W_\ell\}|$ where W_ℓ is ℓ -th pair of W .

125 To encode a walk W with respect to the sets S_u and S_v , CAW applies $\text{ENC}(W; S_u, S_v) =$
 126 $\text{RNN}([f_1(I_{\text{CAW}}(w_i; S_u, S_v)) \| f_2(t_{i-1} - t_i)]_{i=1}^L)$ where f_1 is a permutation-invariant function, f_2 is
 127 a time encoder, and $t_0 = t$. Finally, CAW combines the embeddings of each walk in $S_u \cup S_v$ using
 128 mean-pooling or self-attention to obtain the representation for the event (u, v, t) .

129 In practice, TGNs often rely on sampling schemes for computational reasons. However we are
 130 concerned with the expressiveness of TGNs, so our analysis assumes complete structural information,
 131 i.e., S_u is the set of all temporal walks from u and MP-TGNs combine information from all neighbors.

132 3 The representational power and limits of TGNs

133 We now study the expressiveness of TGNs on node/edge-level prediction. We also establish connec-
 134 tions to a variant of the WL test and show limits of specific TGN models. Proofs are in Appendix B.

135 3.1 Distinguishing nodes with MP-TGNs

136 We analyze MP-TGNs w.r.t. their ability to map different nodes to different locations in the embedding
 137 space. In particular, we say that an L -layer MP-TGN distinguishes two nodes u, v of a temporal
 138 graph at time t , if the last layer embeddings of u and v are different, i.e., $h_u^{(L)}(t) \neq h_v^{(L)}(t)$.

139 We can describe the MP computations of a node v at time t via its *temporal computation tree* (TCT)
 140 $T_v(t)$. $T_v(t)$ has v as its root and height equal to the number of TGN layers L . We will keep the
 141 dependence on depth L implicit for notational simplicity. For each element $(u, e, t') \in \mathcal{N}(v, t)$
 142 associated with v , we have a node, say i , in the next layer of the TCT linked to the root by an edge
 143 annotated with (e, t') . The remaining TCT layers are built recursively using the same mechanism.
 144 We denote by $\#_v^t$ the (possibly many-to-one) operator that maps nodes in $T_v(t)$ back to nodes in $\mathcal{G}(t)$,
 145 e.g., $\#_v^t i = u$. Each node i in $T_v(t)$ has a state vector $s_i = s_{\#_v^t i}^t$. To get the embedding of the root
 146 v , information is propagated bottom-up, i.e., starting from the leaves all the way up to the root —
 147 each node aggregates the message from the layer below and updates its representation along the way.
 148 Whenever clear from context, we denote $\#_v^t$ simply as $\#$ for a cleaner notation.

149 We study the expressive power of MP-TGNs through the lens of functions on multisets adapted to
 150 temporal settings, i.e., comprising triplets of node states, edge features, and timestamps. Intuitively,
 151 injective functions ‘preserve’ the information as it is propagated, so should be essential for maximally
 152 expressive MP-TGNs. We formalize this idea in Lemma 1 and Proposition 2 via Definition 1.

153 **Definition 1** (Isomorphic TCTs). *Two TCTs $T_z(t)$ and $T_{z'}(t)$ at time t are isomorphic if there is a*
 154 *bijection $f : V(T_z(t)) \rightarrow V(T_{z'}(t))$ between the nodes of the trees such that the following holds:*

$$155 (u, v, t') \in E(T_z(t)) \iff (f(u), f(v), t') \in E(T_{z'}(t))$$

$$156 \forall (u, v, t') \in E(T_z(t)) : e_{uv}(t') = e_{f(u)f(v)}(t') \text{ and } \forall u \in V(T_z(t)) : s_u = s_{f(u)} \text{ and } k_u = k_{f(u)}$$

157 Here, k_u denotes the level (depth) of node u in the tree. The root node has level 0, and for a node u
 158 with level k_u , the children of u have level $k_u + 1$.

159 **Lemma 1.** *If an MP-TGN Q with L layers distinguishes two nodes u, v of a dynamic graph $\mathcal{G}(t)$,*
 160 *then the L -depth TCTs $T_u(t)$ and $T_v(t)$ are not isomorphic.*

161 For non-isomorphic TCTs, Proposition 2 shows that improving MP-TGNs with *injective* message
 162 passing layers suffices to achieve node distinguishability, extending results from static GNNs [42].

163 **Proposition 2** (Most expressive MP-TGNs). *If the L -depth TCTs of two nodes u, v of a temporal*
 164 *graph $\mathcal{G}(t)$ at time t are not isomorphic, then there exists an MP-TGN Q with L layers and injective*
 165 *aggregation and update functions at each layer that is able to distinguish nodes u and v .*

166 So far, we have considered TCTs with general memory modules, i.e., nodes are annotated with
 167 memory states. However, an important question remains: *How does the expressive power of MP-*
 168 *TGNs change as a function of the memory?* Our next result - Proposition 3 - shows that adding
 169 GRU-based memory does not increase the expressiveness of suitably deep MP-TGNs.

170 **Proposition 3** (The role of memory). *Let $\mathcal{Q}_L^{[M]}$ denote the class of MP-TGNs with recurrent memory*
 171 *and L layers. Similarly, we denote by \mathcal{Q}_L the family of memoryless MP-TGNs with L layers. Let Δ*
 172 *be the temporal diameter of $\mathcal{G}(t)$ (see Definition B2). Then, it holds that:*

- 173 1. *If $L < \Delta$: $\mathcal{Q}_L^{[M]}$ is strictly more powerful than \mathcal{Q}_L in distinguishing nodes of $\mathcal{G}(t)$;*
- 174 2. *For any L : $\mathcal{Q}_{L+\Delta}$ is at least as powerful as $\mathcal{Q}_L^{[M]}$ in distinguishing nodes of $\mathcal{G}(t)$.*

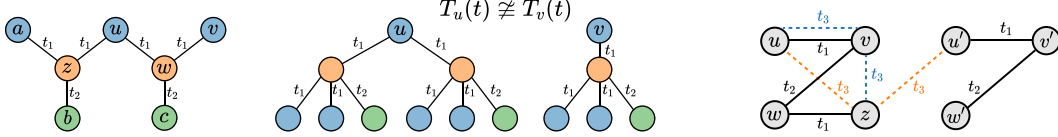


Figure 2: **Limitations of TGNs.** [Left] Temporal graph with nodes u, v that TGN-Att/TGAT cannot distinguish. Colors are node features, edge features are identical, and $t_3 > t_2 > t_1$. [Center] TCTs of u and v are non-isomorphic. However, the attention layers of TGAT/TGN-Att compute weighted averages over a same multiset of values, returning identical messages for u and v . [Right] MP-TGNs fail to distinguish the events (u, v, t_3) and (v, z, t_3) as TCTs of z and u are isomorphic. Meanwhile, CAW cannot separate (u, z, t_3) and (u', z, t_3) : the 3-depth TCTs of u and u' are not isomorphic, but the temporal walks from u and u' have length 1, keeping CAW from capturing structural differences.

175 The MP-TGN framework is rather general and subsumes many modern methods for temporal graphs
 176 [e.g., 16, 32, 41]. We now analyze the theoretical limitations of two concrete instances of MP-TGNs:
 177 TGAT [41] and TGN-Att [27]. Remarkably, these models are among the best-performing MP-TGNs.
 178 Nonetheless, we can show that there are nodes of very simple temporal graphs that TGAT and
 179 TGN-Att cannot distinguish (see Figure 2). We formalize this in Proposition 4 by establishing that
 180 there are cases in which TGNs with injective layers can succeed, but TGAT and TGN-Att cannot.

181 **Proposition 4** (Limitations of TGAT/TGN-Att). *There exist temporal graphs containing nodes u, v*
 182 *that have non-isomorphic TCTs, yet no TGAT nor TGN-Att with mean message aggregator (i.e., using*
 183 *MEAN as MSGAGG) can distinguish u and v .*

184 This limitation stems from the fact that the attention mechanism employed by TGAT and TGN-Att is
 185 proportion invariant [26]. The memory module of TGN-Att cannot counteract this limitation due to
 186 its mean-based message aggregation scheme. We provide more details in Appendix B.6.

187 3.2 Predicting temporal links

188 Models for dynamic graphs are usually trained and evaluated on temporal link prediction [18], which
 189 consists in predicting whether an event would occur at a given time. To predict an event between
 190 nodes u and v at t , MP-TGNs combine the node embeddings $h_u^{(L)}(t)$ and $h_v^{(L)}(t)$, and evaluate the
 191 resulting vector using an MLP. On the other hand, CAW is originally designed for link prediction
 192 tasks and directly computes edge embeddings, bypassing the computation of node representations.

193 We can extend the notion of node distinguishability to edges/events. We say that a model distinguishes
 194 two synchronous events $\gamma = (u, v, t)$ and $\gamma' = (u', v', t)$ of a temporal graph if it assigns different
 195 edge embeddings $h_\gamma \neq h_{\gamma'}$ for γ and γ' . Proposition 5 asserts that CAWs are not strictly more
 196 expressive than MP-TGNs, and vice-versa. Intuitively, CAW’s advantage over MP-TGNs lies in its
 197 ability to exploit node identities and capture correlation between walks. However, CAW imposes
 198 temporal constraints on random walks, i.e., walks have timestamps in decreasing order, which can
 199 limit its ability to distinguish events. Figure 2(Right) sketches constructions for Proposition 5.

200 **Proposition 5** (Limitations of MP-TGNs and CAW). *There exist distinct synchronous events of a*
 201 *temporal graph that CAW can distinguish but MP-TGNs with injective layers cannot, and vice-versa.*

202 3.3 Connections with the WL test

203 The Weisfeiler-Leman test (1-WL) has been used as a key tool to analyze the expressive power of
 204 GNNs. We now study the power of MP-TGNs under a temporally-extended version of 1-WL, and
 205 prove negative results regarding whether MP-TGN’s can recognize properties of temporal graphs.

206 **Temporal WL test.** We can extend the WL test for temporal settings in a straightforward manner
 207 by exploiting the equivalence between temporal graphs and multi-graphs with timestamped edges
 208 [24]. In particular, the temporal variant of 1-WL assigns colors for all nodes in an input dynamic
 209 graph $\mathcal{G}(t)$ by applying the following iterative procedure:

210 *Initialization:* The colors of all nodes in $\mathcal{G}(t)$ are initialized using the initial node features: $\forall v \in$
 211 $V(\mathcal{G}(t)), c^0(v) = x_v$. If node features are not available, all nodes receive identical colors;

212 *Refinement:* At step ℓ , the colors of all nodes are refined using a hash (injective) function: for all
 213 $v \in V(\mathcal{G}(t))$, we apply $c^{\ell+1}(v) = \text{HASH}(c^\ell(v), \{(c^\ell(u), e_{uv}(t'), t') : (u, v, t') \in \mathcal{G}(t)\})$;

214 *Termination:* The test is carried out for two temporal graphs at time t in parallel and stops when
 215 the multisets of corresponding colors diverge, returning non-isomorphic. If the algorithm
 216 runs until the number of different colors stops increasing, the test is deemed inconclusive.

217 We note that the temporal WL test trivially reduces to the standard 1-WL test if all timestamps and
 218 edge features are identical. The resemblance between MP-TGNs and GNNs and their corresponding
 219 WL tests suggests that the power of MP-TGNs is bounded by the temporal WL test. Proposition 6
 220 conveys that MP-TGNs with injective layers are as powerful as the temporal WL test.

221 **Proposition 6.** *Assume finite spaces of initial node features \mathcal{X} , edge features \mathcal{E} , and timestamps \mathcal{T} .
 222 Let the number of events of any temporal graph be bounded by a fixed constant. Then, there is an
 223 MP-TGN with suitable parameters using injective aggregation/update functions that outputs different
 224 representations for two temporal graphs if and only if the temporal-WL test outputs ‘non-isomorphic’.*

225 A natural consequence of the limited power of MP-TGNs is that even the most powerful MP-TGNs
 226 fail to distinguish relevant graph properties, and the same applies to CAWs (see Proposition 7).

227 **Proposition 7.** *There exist non-isomorphic temporal graphs differing in properties such as diameter,
 228 girth, and total number of cycles, which cannot be differentiated by MP-TGNs and CAWs.*

229 Figure 3 provides a construction for Proposition 7.
 230 The temporal graphs $\mathcal{G}(t)$ and $\mathcal{G}'(t)$ differ in diameter
 231 (∞ vs. 3), girth (3 vs. 6), and number of cycles (2
 232 vs. 1). By inspecting the TCTs, one can observe that,
 233 for any node in $\mathcal{G}(t)$, there is a corresponding one
 234 in $\mathcal{G}'(t)$ whose TCTs are isomorphic, e.g., $T_{u_1}(t) \cong$
 235 $T_{u'_1}(t)$ for $t > t_3$. As a result, the multisets of node
 236 embeddings for these temporal graphs are identical.
 237 We provide more details and a construction - where CAW fails to decide properties - in the Appendix.

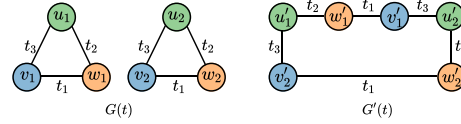


Figure 3: Examples of temporal graphs for which MP-TGNs cannot distinguish the diameter, girth, and number of cycles.

238 4 Position-encoding injective temporal graph net

239 We now leverage insights from our analysis in Section 3 to build more powerful TGNs. First, we
 240 discuss how to build injective aggregation and update functions in the temporal setting. Second, we
 241 propose an efficient scheme to compute positional features based on counts from TCTs. In addition,
 242 we show that the proposed method, called *position-encoding injective temporal graph net* (PINT), is
 243 more powerful than both WA-TGNs and MP-TGNs in distinguishing events in temporal graphs.

244 **Injective temporal aggregation.** An important design principle in TGNs is to prioritize (give higher
 245 importance to) events based on recency [37, 41]. Proposition 8 introduces an injective aggregation
 246 scheme that captures this principle employing a linearly exponential time decay.

247 **Proposition 8** (Injective function on temporal neighborhood). *Let \mathcal{X} and \mathcal{E} be countable, and \mathcal{T}
 248 countable and bounded. There exists a function f and scalars α and β such that $\sum_i f(x_i, e_i) \alpha^{-\beta t_i}$
 249 is unique on any multiset $M = \{(x_i, e_i, t_i)\} \subseteq \mathcal{X} \times \mathcal{E} \times \mathcal{T}$ with $|M| < N$, where N is a constant.*

250 Leveraging Proposition 8 and the approximation capabilities of multi-layer perceptrons (MLPs), we
 251 propose *position-encoding injective temporal graph net* (PINT). In particular, PINT computes the
 252 embedding of node v at time t and layer ℓ using the following message passing steps:

$$\tilde{h}_v^{(\ell)}(t) = \sum_{(u, e, t') \in \mathcal{N}(v, t)} \text{MLP}_{\text{agg}}^{(\ell)} \left(h_u^{(\ell-1)}(t) \parallel e \right) \alpha^{-\beta(t-t')} \quad (5)$$

$$h_v^{(\ell)}(t) = \text{MLP}_{\text{upd}}^{(\ell)} \left(h_v^{(\ell-1)}(t) \parallel \tilde{h}_v^{(\ell)}(t) \right) \quad (6)$$

253 where \parallel denotes concatenation, $h_v^{(0)} = s_v(t)$, α and β are scalar (hyper-)parameters, and $\text{MLP}_{\text{agg}}^{(\ell)}$
 254 and $\text{MLP}_{\text{upd}}^{(\ell)}$ denote the nonlinear transformations of the aggregation and update steps, respectively.

255 We note that to guarantee that the MLPs in PINT implement injective aggregation/update, we must
 256 further assume that the edge and node features (states) take values from some finite support. In
 257 addition, we highlight that there may exist many other ways to achieve injective temporal MP - we
 258 have presented a solution that captures the ‘recency’ inductive bias of real-world temporal networks.

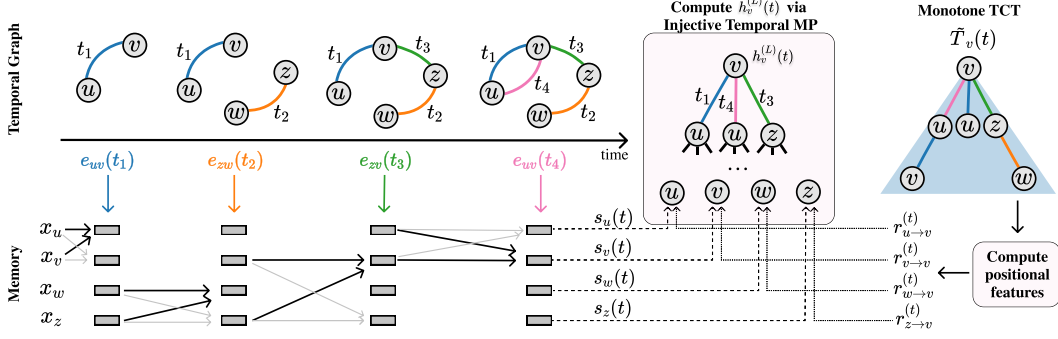


Figure 4: **PINT**. Following the MP-TGN protocol, PINT updates memory states as events unroll. Meanwhile, we use Eqs. (7-11) to update positional features. To extract the embedding for node v , we build its TCT, annotate nodes with memory + positional features, and run (injective) MP.

259 **Relative positional features.** To boost the power of PINT, we propose augmenting memory states
 260 with *relative* positional features. These features count how many temporal walks of a given length
 261 exist between two nodes, or equivalently, how many times nodes appear at different levels of TCTs.

262 Formally, let P be the $d \times d$ matrix obtained by padding a $(d - 1)$ -dimensional identity matrix
 263 with zeros on its top row and its rightmost column. Also, let $r_{j \rightarrow u}^{(t)} \in \mathbb{N}^d$ denote the positional
 264 feature vector of node j relative to u 's TCT at time t . For each event (u, v, t) , with u and
 265 v not participating in other events at t , we recursively update the positional feature vectors as

$$\mathcal{V}_i^{(0)} = \{i\} \quad \forall i \quad (7) \quad \mathcal{V}_u^{(t)} = \mathcal{V}_v^{(t)} = \mathcal{V}_v^{(t^-)} \cup \mathcal{V}_u^{(t^-)} \quad (9)$$

$$266 \quad r_{i \rightarrow j}^{(0)} = \begin{cases} [1, 0, \dots, 0]^\top & \text{if } i = j \\ [0, 0, \dots, 0]^\top & \text{if } i \neq j \end{cases} \quad (8) \quad r_{i \rightarrow v}^{(t)} = P r_{i \rightarrow u}^{(t^-)} + r_{i \rightarrow v}^{(t^-)} \quad \forall i \in \mathcal{V}_u^{(t^-)} \quad (10)$$

$$r_{j \rightarrow u}^{(t)} = P r_{j \rightarrow v}^{(t^-)} + r_{j \rightarrow u}^{(t^-)} \quad \forall j \in \mathcal{V}_v^{(t^-)} \quad (11)$$

267 where we use t^- to denote values “right before” t . The set \mathcal{V}_i keeps track of the nodes for which
 268 we need to update positional features when i participates in an interaction. For simplicity, we have
 269 assumed that there are no other events involving u or v at time t . Appendix B.10 provides equations
 270 for the general case where nodes can participate in multiple events at the same timestamp.

271 The value $r_{i \rightarrow v}^{(t)}[k]$ (the k -th component of $r_{i \rightarrow v}^{(t)}$) corresponds to how many different ways we can get
 272 from v to i in k steps through temporal walks. Additionally, we provide in Lemma 2 an interpretation
 273 of relative positional features in terms of the so-called monotone TCTs (Definition 2).

274 **Definition 2.** The monotone TCT of a node u at time t , denoted by $\tilde{T}_u(t)$, is the maximal subtree of
 275 the TCT of u s.t. for any path $p = (u, t_1, u_1, t_2, u_2, \dots)$ from the root u to leaf nodes of $\tilde{T}_u(t)$ time
 276 monotonically decreases with respect to time, i.e., we have that $t_1 > t_2 > \dots$.

277 **Lemma 2.** For any pair of nodes i, u of a temporal graph $\mathcal{G}(t)$, the k -th component of the positional
 278 feature vector $r_{i \rightarrow u}^{(t)}$ stores the number of times i appears at the k -th layer of the monotone TCT of u .

279 **Edge and node embeddings.** To obtain the embedding h_γ for an event $\gamma = (u, v, t)$, an L -layer
 280 PINT computes embeddings for node u and v using L steps of temporal message passing. However,
 281 when computing the embedding $h_u^L(t)$ of u , we concatenate node states $s_j(t)$ with the positional
 282 features $r_{j \rightarrow u}^{(t)}$ and $r_{j \rightarrow v}^{(t)}$ for all node j in the L -hop temporal neighborhood of u . We apply the same
 283 procedure to obtain $h_v^L(t)$, and then combine $h_v^L(t)$ and $h_u^L(t)$ using a readout function.

284 Similarly, to compute representations for node-level prediction, for each node j in the L -hop neigh-
 285 borhood of u , we concatenate node states $s_j(t)$ with features $r_{j \rightarrow u}^{(t)}$. Then, we use our injective MP to
 286 combine the information stored in u and its neighboring nodes. Figure 4 illustrates the process.

287 Notably, Proposition 9 states that PINT is strictly more powerful than existing TGNs. In fact, the
 288 relative positional features mimic the discriminative power of WA-TGNs, while eliminate their
 289 temporal monotonicity constraints. Additionally, PINT can implement injective temporal message
 290 passing, achieving maximally-expressive MP-TGNs.

291 **Proposition 9** (Expressiveness of PINT: link prediction). *PINT (with relative positional features) is*
 292 *strictly more powerful than both MP-TGNs and CAW in distinguishing events in temporal graphs.*

293 **When does PINT fail?** Naturally, whenever the TCTs (annotated with
 294 positional features) for the endpoints of two edges (u, v, t) and (u', v', t)
 295 are pairwise isomorphic, PINT returns the same edge embedding and is
 296 not able to differentiate the events. Figure 5 shows an example in which
 297 this happens — we assume that all node/edge features are identical. Due
 298 to graph symmetries, u and z occur the same number of times in each level
 299 of v ’s monotone TCT. Also, the sets of temporal walks starting at u and z
 300 are identical if we swap the labels of these nodes. Importantly, CAWs and
 301 MP-TGNs also fail here, as stated in Proposition 9.

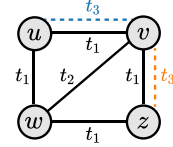


Figure 5: PINT cannot distinguish the events (u, v, t_3) and (v, z, t_3) .

302 **Proposition 10** (Limitations of PINT). *There are synchronous events of temporal graphs that PINT*
 303 *cannot distinguish (as seen in Figure 5).*

304 **Implementation and computational cost.** The online updates for PINT’s positional features have
 305 complexity $\mathcal{O}(d|\mathcal{V}_u^{(t^-)}| + d|\mathcal{V}_v^{(t^-)}|)$. Similarly to CAW’s sampling procedure, our online update
 306 is a sequential process better done in CPUs. However, while CAW may require significant CPU-
 307 GPU memory exchange — proportional to both the number of walks and their depth —, we only
 308 communicate the positional features. We can also speed-up the training of PINT by pre-computing
 309 the positional features for each batch, avoiding redundant computations at each epoch. Apart from
 310 positional features, the computational cost of PINT is similar to that of TGN-Att. Following standard
 311 MP-TGN procedure, we control the branching factor of TCTs using neighborhood sampling.

312 Note that the positional features monotonically increase with time, which is undesirable for practical
 313 generalization purposes. Since our theoretical results hold for any fixed t , this issue can be solved
 314 by dividing the positional features by a time-dependent normalization factor. Nonetheless, we have
 315 found that employing L_1 -normalization leads to good empirical results for all evaluated datasets.

316 5 Experiments

317 We now assess the performance of PINT on several popular and large-scale benchmarks for TGNs.
 318 We have implemented experiments using PyTorch [25] and code is available as additional material.

319 **Tasks and datasets.** We evaluate PINT on dynamic link prediction, closely following the evaluation
 320 setup employed by Rossi et al. [27] and Xu et al. [41]. We use six popular benchmark datasets:
 321 Reddit, Wikipedia, Twitter, UCI, Enron, and LastFM [16, 27, 37, 41]. Notably, UCI, Enron, and
 322 LastFM are non-attributed networks, i.e., they do not contain feature vectors associated with the
 323 events. Node features are absent in all datasets, thus following previous works we set them to vectors
 324 of zeros [27, 41]. Since Twitter is not publicly available, we follow the guidelines by Rossi et al. [27]
 325 to create our version. We provide more details regarding datasets in the supplementary material.

326 **Baselines.** We compare PINT against five prominent TGNs: Jodie [16], DyRep [32], TGAT [41],
 327 TGN-Att [27], and CAW [37]. For completeness, we also report results using two static GNNs: GAT
 328 [34] and GraphSage [12]. Since we adopt the same setup as TGN-Att, we use their table numbers
 329 for all baselines but CAW on Wikipedia and Reddit. The remaining results were obtained using the
 330 implementations and guidelines available from the official repositories. As an ablation study, we also
 331 include a version of PINT without relative positional features in the comparison. We provide detailed
 332 information about hyperparameters and the training of each model in the supplementary material.

333 **Experimental setup.** We follow Xu et al. [41] and use a 70%-15%-15% (train-val-test) temporal
 334 split for all datasets. We adopt average precision (AP) as the performance metric. We also analyze
 335 separately predictions involving only nodes seen during training (transductive), and those involving
 336 novel nodes (inductive). We report mean and standard deviation of the AP over ten runs. For further
 337 details, see Appendix D. We provide additional results in the supplementary material.

338 **Results.** Table 1 shows that PINT is the best-performing method on five out of six datasets for the
 339 transductive setting. Notably, the performance gap between PINT and TGN-Att amounts to over
 340 15% AP on UCI. The gap is also relatively high compared to CAW on LastFM, Enron, and UCI;
 341 with CAW being the best model only on Enron. We also observe that many models achieve relatively
 342 high AP on the attributed networks (Reddit, Wikipedia, and Twitter). This aligns well with findings

Table 1: **Average Precision (AP)** results for link prediction. We denote the best-performing model (highest mean AP) in **blue**. In 5 out of 6 datasets, PINT achieves the highest AP in the transductive setting. For the inductive case, PINT outperforms previous MP-TGNs and competes with CAW. We also show the performance of PINT with and without relative positional features. For all datasets, adopting positional features leads to significant performance gains.

	Model	Reddit	Wikipedia	Twitter	UCI	Enron	LastFM
Transductive	GAT	97.33 ± 0.2	94.73 ± 0.2	-	-	-	-
	GraphSAGE	97.65 ± 0.2	93.56 ± 0.3	-	-	-	-
	Jodie	97.11 ± 0.3	94.62 ± 0.5	98.23 ± 0.1	86.73 ± 1.0	77.31 ± 4.2	69.32 ± 1.0
	DyRep	97.98 ± 0.1	94.59 ± 0.2	98.48 ± 0.1	54.60 ± 3.1	77.68 ± 1.6	69.24 ± 1.4
	TGAT	98.12 ± 0.2	95.34 ± 0.1	98.70 ± 0.1	77.51 ± 0.7	68.02 ± 0.1	54.77 ± 0.4
	TGN-Att	98.70 ± 0.1	98.46 ± 0.1	98.00 ± 0.1	80.40 ± 1.4	79.91 ± 1.3	80.69 ± 0.2
	CAW	98.39 ± 0.1	98.63 ± 0.1	98.72 ± 0.1	92.16 ± 0.1	92.09 ± 0.7	81.29 ± 0.1
	PINT (w/o pos. feat.)	98.62 ± .04	98.43 ± .04	98.53 ± 0.1	92.68 ± 0.5	83.06 ± 2.1	81.35 ± 1.6
	PINT	99.03 ± .01	98.78 ± 0.1	99.35 ± .01	96.01 ± 0.1	88.71 ± 1.3	88.06 ± 0.7
	Inductive	GAT	95.37 ± 0.3	91.27 ± 0.4	-	-	-
GraphSAGE		96.27 ± 0.2	91.09 ± 0.3	-	-	-	-
Jodie		94.36 ± 1.1	93.11 ± 0.4	96.06 ± 0.1	75.26 ± 1.7	76.48 ± 3.5	80.32 ± 1.4
DyRep		95.68 ± 0.2	92.05 ± 0.3	96.33 ± 0.2	50.96 ± 1.9	66.97 ± 3.8	82.03 ± 0.6
TGAT		96.62 ± 0.3	93.99 ± 0.3	96.33 ± 0.1	70.54 ± 0.5	63.70 ± 0.2	56.76 ± 0.9
TGN-Att		97.55 ± 0.1	97.81 ± 0.1	95.76 ± 0.1	74.70 ± 0.9	78.96 ± 0.5	84.66 ± 0.1
CAW		97.81 ± 0.1	98.52 ± 0.1	98.54 ± 0.4	92.56 ± 0.1	91.74 ± 1.7	85.67 ± 0.5
PINT (w/o pos. feat.)		97.22 ± 0.2	97.81 ± 0.1	96.10 ± 0.1	90.25 ± 0.3	75.99 ± 2.3	88.44 ± 1.1
PINT		98.25 ± .04	98.38 ± .04	98.20 ± .03	93.97 ± 0.1	81.05 ± 2.4	91.76 ± 0.7

from [37], where TGN-Att was shown to have competitive performance against CAW on Wikipedia and Reddit. The performance of GAT and TGAT (static GNNs) on Reddit and Wikipedia reinforces the hypothesis that the edge features add significantly to the discriminative power. On the other hand, PINT and CAW, which leverage relative identities, show superior performance relative to other methods when only time and degree information is available, i.e., on unattributed networks (UCI, Enron, and LastFM). Table 1 also shows the effect of using relative positional features. While including these features boosts PINT’s performance systematically, our ablation study shows that PINT w/o positional features still outperforms other MP-TGNs on unattributed networks. In the inductive case, we observe a similar behavior: PINT is consistently the best MP-TGN, and is better than CAW on 3/6 datasets. Overall, PINT (w/ positional features) also yields the lowest standard deviations. This suggests that positional encodings might be a useful inductive bias for TGNs.

Time comparison. Figure 6 compares the training times of PINT against other TGNs. For fairness, we use the same architecture (number of layers & neighbors) for all MP-TGNs: i.e., the best-performing PINT. For CAW, we use the one that yielded results in Table 1. As expected, TGAT is the fastest model. Note that the average time/epoch of PINT gets amortized since positional features are pre-computed. Without these features, PINT’s runtime closely matches TGN-Att. When trained for over 25 epochs, PINT runs considerably faster than CAW. More details and results are provided in the Appendix.

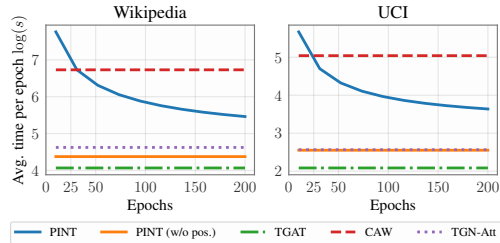


Figure 6: Time comparison: PINT versus TGNs (in log-scale). The cost of pre-computing positional features is quickly diluted as the number of epochs increases.

6 Conclusion

We laid a rigorous theoretical foundation for TGNs, including the role of memory modules, relationship between classes of TGNs, and failure cases for MP-TGNs. Together, our theoretical results shed light on the representational capabilities of TGNs, and connections with their static counterparts. We also introduced a novel TGN method, provably more expressive than existing TGNs.

Key practical takeaways from this work: (a) temporal models should be designed to have injective update rules and to exploit both neighborhood and walk aggregation, and (b) deep architectures can likely be made more compute-friendly as the role of memory gets diminished with depth, provably.

375 References

- 376 [1] P. Barceló, E. V. Kostylev, M. Monet, J. Pérez, J. L. Reutter, and J.-P. Silva. The logical expressiveness of
377 graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- 378 [2] Q. Cappart, D. Chételat, E. B. Khalil, A. Lodi, C. Morris, and P. Velickovic. Combinatorial optimization
379 and reasoning with graph neural networks. In *International Joint Conference on Artificial Intelligence*
380 *(IJCAI)*, 2021.
- 381 [3] B. Chamberlain, J. Rowbottom, M. Gorinova, M. M. Bronstein, S. Webb, and E. Rossi. GRAND: graph
382 neural diffusion. In *International Conference on Machine Learning (ICML)*, 2021.
- 383 [4] M. Chen, Z. Wei, Z. Huang, B. Ding, and Y. Li. Simple and deep graph convolutional networks. In
384 *International Conference on Machine Learning (ICML)*, 2020.
- 385 [5] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio.
386 Learning phrase representations using RNN encoder–decoder for statistical machine translation. In
387 *Empirical Methods in Natural Language Processing (EMNLP)*, 2014.
- 388 [6] N. Dehmamy, A.-L. Barabási, and R. Yu. Understanding the representation power of graph neural networks
389 in learning graph topology. In *Advances in neural information processing systems (NeurIPS)*, 2019.
- 390 [7] A. Darrow-Pinion, J. She, D. Wong, O. Lange, T. Hester, L. Perez, M. Nunkesser, S. Lee, X. Guo,
391 B. Wiltshire, P. W. Battaglia, V. Gupta, A. Li, Z. Xu, A. Sanchez-Gonzalez, Y. Li, and P. Velickovic. Eta
392 prediction with graph neural networks in google maps. In *Conference on Information and Knowledge*
393 *Management (CIKM)*, 2021.
- 394 [8] Simon S. Du, Kangcheng Hou, Ruslan Salakhutdinov, Barnabás Póczos, Ruosong Wang, and Keyulu Xu.
395 Graph neural tangent kernel: Fusing graph neural networks with graph kernels. In *Advances in Neural*
396 *Information Processing Systems (NeurIPS)*, 2019.
- 397 [9] V. Garg, S. Jegelka, and T. Jaakkola. Generalization and representational limits of graph neural networks.
398 In *International Conference on Machine Learning (ICML)*, 2020.
- 399 [10] J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum
400 chemistry. In *International Conference on Machine Learning (ICML)*, 2017.
- 401 [11] M. Gori, G. Monfardini, and F. Scarselli. A new model for learning in graph domains. In *IEEE International*
402 *Joint Conference on Neural Networks (IJCNN)*, 2005.
- 403 [12] W. Hamilton, Z. Ying, and J. Leskovec. Inductive representation learning on large graphs. In *Advances in*
404 *Neural Information Processing Systems (NeurIPS)*, 2017.
- 405 [13] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- 406 [14] S. Kazemi, R. Goel, K. Jain, I. Kobyzev, A. Sethi, P. Forsyth, and P. Poupart. Representation learning for
407 dynamic graphs: A survey. *Journal of Machine Learning Research*, 21(70):1–73, 2020.
- 408 [15] S. M. Kazemi, R. Goel, S. Eghbali, J. Ramanan, J. Sahota, S. Thakur, S. Wu, C. Smyth, P. Poupart, and
409 M. Brubaker. Time2vec: Learning a vector representation of time. *ArXiv: 1907.05321*, 2019.
- 410 [16] S. Kumar, X. Zhang, and J. Leskovec. Predicting dynamic embedding trajectory in temporal interaction
411 networks. In *International Conference on Knowledge Discovery & Data Mining (KDD)*, 2019.
- 412 [17] Renjie Liao, Raquel Urtasun, and Richard Zemel. A PAC-bayesian approach to generalization bounds for
413 graph neural networks. In *International Conference on Learning Representations (ICLR)*, 2021.
- 414 [18] D. Liben-Nowell and J. Kleinberg. The link prediction problem for social networks. *Journal of the*
415 *American Society for Information Science and Technology*, 58(7):1019–1031, 2007.
- 416 [19] A. Loukas. What graph neural networks cannot learn: depth vs width. In *International Conference on*
417 *Learning Representations (ICLR)*, 2020.
- 418 [20] Andreas Loukas. How hard is to distinguish graphs with graph neural networks? In *Advances in Neural*
419 *Information Processing Systems (NeurIPS)*, 2020.
- 420 [21] H. Maron, H. Ben-Hamu, H. Serviansky, and Y. Lipman. Provably powerful graph networks. In *Advances*
421 *in Neural Information Processing Systems (NeurIPS)*, 2019.

- 422 [22] C. Morris, M. Ritzert, M. Fey, W. L. Hamilton, J. E. Lenssen, G. Rattan, and M. Grohe. Weisfeiler and
423 leman go neural: Higher-order graph neural networks. In *AAAI Conference on Artificial Intelligence*
424 (*AAAI*), 2019.
- 425 [23] H. Nguyen and T. Maehara. Graph homomorphism convolution. In *International Conference on Machine*
426 *Learning (ICML)*, 2020.
- 427 [24] F. Orsini, P. Frasconi, and L. D. Raedt. Graph invariant kernels. In *International Joint Conference on*
428 *Artificial Intelligence (IJCAI)*, 2015.
- 429 [25] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and
430 A. Lerer. Automatic differentiation in pytorch. In *Advances in Neural Information Processing Systems*
431 (*NeurIPS - Workshop*), 2017.
- 432 [26] Jorge Pérez, Javier Marinković, and Pablo Barceló. On the turing completeness of modern neural network
433 architectures. In *International Conference on Learning Representations (ICLR)*, 2019.
- 434 [27] Emanuele Rossi, Ben Chamberlain, Fabrizio Frasca, Davide Eynard, Federico Monti, and Michael Bron-
435 stein. Temporal graph networks for deep learning on dynamic graphs. In *ICML 2020 Workshop on Graph*
436 *Representation Learning*, 2020.
- 437 [28] A. Sanchez-Gonzalez, J. Godwin, T. Pfaff, R. Ying, J. Leskovec, and P. Battaglia. Learning to simulate
438 complex physics with graph networks. In *International Conference on Machine Learning (ICML)*, 2020.
- 439 [29] R. Sato, M. Yamada, and H. Kashima. Approximation ratios of graph neural networks for combinatorial
440 problems. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- 441 [30] F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model.
442 *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- 443 [31] J. M. Stokes, K. Yang, K. Swanson, W. Jin, A. Cubillos-Ruiz, N. M. Donghia, C. R. MacNair, S. French,
444 L. A. Carfrae, Z. Bloom-Ackermann, V. M. Tran, A. Chiappino-Pepe, A. H. Badran, I. W. Andrews, E. J.
445 Chory, G. M. Church, E. D. Brown, T. S. Jaakkola, R. Barzilay, and J. J. Collins. A deep learning approach
446 to antibiotic discovery. *Cell*, 180(4):688 – 702, 2020.
- 447 [32] R. Trivedi, M. Farajtabar, P. Biswal, and H. Zha. DyRep: Learning representations over dynamic graphs.
448 In *International Conference on Learning Representations (ICLR)*, 2019.
- 449 [33] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin.
450 Attention is all you need. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- 451 [34] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio.
452 Graph Attention Networks. In *International Conference on Learning Representations (ICLR)*, 2018.
- 453 [35] S. Verma and Z.-L. Zhang. Stability and generalization of graph convolutional neural networks. In
454 *International Conference on Knowledge Discovery & Data Mining (KDD)*, 2019.
- 455 [36] C. Vignac, A. Loukas, and P. Frossard. Building powerful and equivariant graph neural networks with
456 structural message-passing. In *Neural Information Processing Systems (NeurIPS)*, 2020.
- 457 [37] Y. Wang, Y. Chang, Y. Liu, J. Leskovec, and P. Li. Inductive representation learning in temporal networks
458 via causal anonymous walks. In *International Conference on Learning Representations (ICLR)*, 2021.
- 459 [38] Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu. A comprehensive survey on graph neural
460 networks. *IEEE Transactions on Neural Networks and Learning Systems*, pages 1–21, 2020.
- 461 [39] D. Xu, W. Cheng, D. Luo, Y. Gu, X. Liu, J. Ni, B. Zong, H. Chen, and X. Zhang. Adaptive neural network
462 for node classification in dynamic networks. In *IEEE International Conference on Data Mining (ICDM)*,
463 2019.
- 464 [40] D. Xu, C. Ruan, E. Korpeoglu, S. Kumar, and K. Achan. Self-attention with functional time representation
465 learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2019.
- 466 [41] D. Xu, C. Ruan, E. Korpeoglu, S. Kumar, and K. Achan. Inductive representation learning on temporal
467 graphs. In *International Conference on Learning Representations (ICLR)*, 2020.
- 468 [42] K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In *International*
469 *Conference on Learning Representations (ICLR)*, 2019.

- 470 [43] K. Xu, J. Li, M. Zhang, S. S. Du, K.-I. Kawarabayashi, and S. Jegelka. What can neural networks reason
471 about? In *International Conference on Learning Representations (ICLR)*, 2020.
- 472 [44] K. Xu, M. Zhang, J. Li, S. S. Du, K.-I. Kawarabayashi, and S. Jegelka. How neural networks extrapolate:
473 From feedforward to graph neural networks. In *International Conference on Learning Representations*
474 *(ICLR)*, 2021.
- 475 [45] Z. Zhang, F. Wu, and W. S. Lee. Factor graph neural networks. In *Advances in Neural Information*
476 *Processing Systems (NeurIPS)*, 2020.

Provably expressive temporal graph networks (Supplementary material)

Anonymous Author(s)

Affiliation

email

477 A Further details on temporal graph networks

478 In this section we present more details about the models TGAT, TGN-Att, and CAW.

479 A.1 Temporal graph attention (TGAT)

480 Temporal graph attention networks [41] combine time encoders [15] and self-attention [33]. In
481 particular, the time encoder ϕ is given by

$$\phi(t - t') = [\cos(\omega_1(t - t') + b_1), \dots, \cos(\omega_d(t - t') + b_d)], \quad (\text{S1})$$

482 where ω_i 's and b_i 's are learned scalar parameters. The time embeddings are concatenated to the
483 edge features before being fed into a typical self-attention layer, where the query q is a function of a
484 reference node v , and both values V and keys K depend on v 's temporal neighbors. Formally, TGAT
485 first computes a matrix $C_v^{(\ell)}(t)$ whose u -th row is $c_{vu}^{(\ell)}(t) = [h_u^{(\ell-1)}(t) \parallel \phi(t - t_{uv}) \parallel e_{uv}]$ for all
486 $(u, e_{uv}, t_{uv}) \in \mathcal{N}(v, t)$. Then, the output $\tilde{h}_v^{(\ell)}(t)$ of the AGG $^{(\ell)}$ function is given by

$$q = [h_v^{(\ell-1)}(t) \parallel \phi(0)]W_q^{(\ell)} \quad K = C_v^{(\ell)}(t)W_K^{(\ell)} \quad V = C_v^{(\ell)}(t)W_V^{(\ell)} \quad (\text{S2})$$

$$\tilde{h}_v^{(\ell)}(t) = \text{softmax}(qK^\top)V \quad (\text{S3})$$

487 where $W_q^{(\ell)}$, $W_K^{(\ell)}$, and $W_V^{(\ell)}$ are model parameters. Regarding the UPDATE function, TGAT applies
488 a multilayer perceptron, i.e., $h_v^{(\ell)}(t) = \text{MLP}^{(\ell)}(h_v^{(\ell-1)}(t) \parallel \tilde{h}_v^{(\ell)}(t))$.

489 A.2 Temporal graph networks with attention (TGN-Att)

490 We now discuss details regarding the MP-TGN framework omitted from the main paper for simplicity.

491 For the sake of generality, Rossi et al. [27] present a formulation for MP-TGNs that can handle
492 node-level events, e.g., node feature updates. These events lead to *i*) updating node memory states,
493 and *ii*) using the time-evolving node features as additional inputs for the message passing functions.
494 Nonetheless, to the best of our knowledge, all relevant CTDG benchmarks comprise only edge
495 events. Therefore, for ease of presentation, we omit node events and temporal node features from our
496 treatment. In Appendix E, we discuss how to handle node-level events.

497 Equation 1, Equation 2, and Equation 3 might make the reader believe that memory states are updated
498 before prediction, i.e., computing node embeddings (message passing) using the features of the edge
499 we want to predict. However, this would incur information leakage. On the other hand, not doing so
500 prevents us from propagating gradients through the memory modules. To get around this problem,
501 Rossi et al. [27] propose updating the memory with messages coming from previous batches, and
502 then predicting the interactions.

503 **Revisiting memory with message aggregators.** To speed up computations, MP-TGNs employ a
504 form of batch learning where events/messages in a same batch are aggregated. In our analysis, we
505 assume that two events belong to the same batch only if they occur at the same timestamp. Importantly,
506 message aggregators allow removing ambiguity in the way the memory of a node participating in
507 multiple events (at the same timestamp) is updated — two events involving a given node i at the same
508 time could lead to different ways of updating the state of i .

509 Suppose the event $\gamma = (i, u, t)$ occurs. MP-TGNs proceed by computing a message function MSG_e
 510 for each endpoint of γ , i.e.,

$$\begin{aligned} m_{i,u}(t) &= \text{MSG}_e(s_i(t^-), s_u(t^-), t - t_i, e_{iu}(t)) \\ m_{u,i}(t) &= \text{MSG}_e(s_u(t^-), s_i(t^-), t - t_u, e_{iu}(t)) \end{aligned}$$

511 Following the original formulation, we assume an identity message function — simply the concatenation
 512 of the inputs, i.e., $\text{MSG}_e(s_i(t^-), s_u(t^-), t - t_i, e_{iu}(t)) = [s_i(t^-), s_u(t^-), t - t_i, e_{iu}(t)]$.

513 Now, suppose two events (i, u, t) and (i, v, t) happen. MP-TGNs aggregate the messages from these
 514 events using a function MSGAGG to obtain a single message for i :

$$\bar{m}_i(t) = \text{MSGAGG}(m_{i,u}(t), m_{i,v}(t))$$

515 Rossi et al. [27] propose non-learnable message aggregators, such as the mean aggregator (average
 516 all messages for a given node), that we denote as MEANAGG and adopt throughout our
 517 analysis. As an example, under events (i, u, t) and (i, v, t) , the aggregated message for i is
 518 $\bar{m}_i(t) = 0.5([s_i(t^-), s_u(t^-), t - t_i, e_{iu}(t)] + [s_i(t^-), s_v(t^-), t - t_i, e_{iv}(t)])$.

519 The memory update of our query node i is given by

$$s_i(t) = \text{MEMORY}(s_i(t^-), \bar{m}_i(t)).$$

520 Finally, we note that TGAT does not have a memory module. TGN-Att consists of the model resulting
 521 from augmenting TGAT with a GRU-based memory.

522 A.3 Causal anonymous walks (CAW)

523 We now provide details regarding how CAW obtains edge embeddings for a query event $\gamma = (u, v, t)$.

524 A temporal walk is represented as $W = ((w_1, t_1), (w_2, t_2), \dots, (w_L, t_L))$, with $t_1 > t_2 > \dots > t_L$
 525 and $(w_{i-1}, w_i, t_i) \in \mathcal{G}(t)$ for all i . We denote by $S_u(t)$ the set of maximal temporal walks starting at
 526 u of size at most L obtained from the temporal graph at time t . Following the original paper, we drop
 527 the time dependence henceforth.

528 A given walk W gets anonymized through replacing each element w_i belonging to W by a 2-element
 529 set of vectors $I_{\text{CAW}}(w_i; S_u, S_v)$ accounting for how many times w_i appears at each position of walks
 530 in S_u and S_v . These vectors are denoted by $g(w_i, S_u)$ and $g(w_i, S_v)$. The walk is encoded using a
 531 RNN:

$$\text{ENC}(W; S_u, S_v) = \text{RNN}([f_1(I_{\text{CAW}}(w_i; S_u, S_v)) \| f_2(t_i - t_{i-1})]_{i=1}^L),$$

532 where f_1 is

$$f_1(I_{\text{CAW}}(w_i; S_u, S_v)) = \text{MLP}(g(w_i, S_u)) + \text{MLP}(g(w_i, S_v)).$$

533 We note that the MLPs share parameters. The function f_2 is given by

$$f_2(t) = [\cos(\omega_i t), \sin(\omega_i t), \dots, \cos(\omega_d t), \sin(\omega_d t)]$$

534 where ω_i 's are learned parameters.

535 To compute the embedding h_γ for (u, v, t) , CAW considers two readout functions: mean and self-
 536 attention. Finally, the final link prediction is obtained from a 2-layer MLP over h_γ .

537 B Proofs

538 B.1 Further definitions and Lemmata

539 **Definition B1** (Monotone walk.). An N -length monotone walk in a temporal graph $\mathcal{G}(t)$ is a sequence
 540 $(w_1, t_1, w_2, t_2, \dots, w_{N+1})$ such that $t_i > t_{i+1}$ and $(w_i, w_{i+1}, t_i) \in \mathcal{G}(t)$ for all i .

541 **Definition B2** (Temporal diameter.). We say the temporal diameter of a graph $\mathcal{G}(t)$ is Δ if the longest
 542 monotone walk in $\mathcal{G}(t)$ has length (i.e. number of edges) exactly Δ .

543 **Lemma B1.** If the TCTs of two nodes are isomorphic, then their monotone TCTs (Definition 2) are
 544 also isomorphic, i.e., $T_u(t) \cong T_v(t) \Rightarrow \tilde{T}_u(t) \cong \tilde{T}_v(t)$ for two nodes u and v of a dynamic graph.

545 *Proof.* Since $T_u(t) \cong T_v(t)$, we have that

$$p = (u_0, t_1, u_1, t_2, u_2, \dots) \text{ from } T_u(t) \iff p' = (f(u_0), t_1, f(u_1), t_2, f(u_2), \dots) \text{ from } T_v(t),$$

$$\text{with } s_{u_i} = s_{f(u_i)} \text{ and } e_{u_i u_{i+1}}(t_{i+1}) = e_{f(u_i) f(u_{i+1})}(t_{i+1}) \text{ and } k_{u_i} = k_{f(u_i)} \quad \forall i$$

546 where $f : V(T_u(t)) \rightarrow V(T_v(t))$ is a bijection.

547 Assume that $\tilde{T}_u(t) \not\cong \tilde{T}_v(t)$. Then, either there exists a path $p_s = (u'_0, t'_1, u'_1, t'_2, \dots)$ in $\tilde{T}_u(t)$, such
 548 that $t'_{k+1} < t'_k$ for all k (i.e., a monotone walk), with no corresponding one in $\tilde{T}_v(t)$ or vice-versa.
 549 Without loss of generality, let us consider the former case.

550 We can construct the path p'_s in $T_v(t)$ by applying f in all elements of p_s , i.e., $p'_s =$
 551 $(f(u'_0), t'_1, f(u'_1), t'_2, \dots)$. Note that p'_s is a monotone walk in $T_v(t)$. Since $\tilde{T}_v(t)$ is the maximal
 552 monotone subtree of $T_v(t)$, it must contain p'_s , leading to contradiction. \square

553 **Lemma B2.** Let $\mathcal{G}(t)$ and $\mathcal{G}'(t)$ be any two non-isomorphic temporal graphs. If an MP-TGN obtains
 554 different multisets of node embeddings for $\mathcal{G}(t)$ and $\mathcal{G}'(t)$. Then, the temporal WL test decides $\mathcal{G}(t)$
 555 and $\mathcal{G}'(t)$ are not isomorphic.

556 *Proof.* Recall Proposition 3 shows that if an MP-TGN with memory is able to distinguish two nodes,
 557 then there is a memoryless MP-TGN with Δ (temporal diameter) additional layers that does the
 558 same. Thus, it suffices to show that if the multisets of colors from temporal WL for $\mathcal{G}(t)$ and $\mathcal{G}'(t)$
 559 after ℓ iterations are identical, then the multisets of embeddings from the memoryless MP-TGN
 560 are also identical, i.e., if $\{\{c^\ell(u)\}_{u \in V(\mathcal{G}(t))}\} = \{\{c^\ell(u')\}_{u' \in V(\mathcal{G}'(t))}\}$, then $\{\{h_u^{(\ell)}(t)\}_{u \in V(\mathcal{G}(t))}\} =$
 561 $\{\{h_{u'}^{(\ell)}(t)\}_{u' \in V(\mathcal{G}'(t))}\}$. To do so, we repurpose the proof of Lemma 2 in [42].

562 More broadly, we show that for any two nodes of a temporal graph $\mathcal{G}(t)$, if the temporal WL returns
 563 $c^\ell(u) = c^\ell(v)$, we have that corresponding embeddings from MP-TGN without memory are identical
 564 $h_u^\ell(t) = h_v^\ell(t)$. We proceed with a proof by induction.

565 [*Base case*] For $\ell = 0$, the proposition trivially holds as the temporal WL has the initial node features
 566 as colors, and memoryless MP-TGNs have these features as embeddings.

567 [*Induction step*] Assume the proposition holds for iteration ℓ . Thus, for any two nodes u, v , if
 568 $c^{\ell+1}(u) = c^{\ell+1}(v)$, we have

$$(c^\ell(u), \{\{c^\ell(i), e_{iu}(t'), t') : (u, i, t') \in \mathcal{G}(t)\}\}) = (c^\ell(v), \{\{c^\ell(j), e_{jv}(t'), t') : (v, j, t') \in \mathcal{G}(t)\}\})$$

569 and, by the induction hypothesis, we know

$$(h_u^{(\ell)}(t), \{\{h_i^{(\ell)}(t), e_{iu}(t'), t') : (u, i, t') \in \mathcal{G}(t)\}\}) =$$

$$(h_v^{(\ell)}(t), \{\{h_j^{(\ell)}(t), e_{jv}(t'), t') : (v, j, t') \in \mathcal{G}(t)\}\})$$

570 We also note that this last identity also implies

$$(h_u^{(\ell)}(t), \{\{h_i^{(\ell)}(t), t - t', e \mid (i, e, t') \in \mathcal{N}(u, t)\}\}) =$$

$$(h_v^{(\ell)}(t), \{\{h_j^{(\ell)}(t), t - t', e \mid (j, e, t') \in \mathcal{N}(v, t)\}\})$$

571 since there exists an event $(u, i, t') \in \mathcal{G}(t)$ with feature $e_{ui}(t') = e$ iff there is an element $(i, e, t') \in$
 572 $\mathcal{N}(u, t)$.

573 As a result, the inputs of the MP-TGN’s aggregation and update functions are identical, which leads
574 to identical outputs $h_u^{(\ell+1)}(t) = h_v^{(\ell+1)}(t)$. Therefore, if the temporal WL test obtains identical
575 multisets of colors for two temporal graphs after ℓ steps, the multisets of embeddings at layer ℓ for
576 these graphs are also identical. \square

577 **Lemma B3** (Lemma 5 in [42]). *Assume \mathcal{X} is countable. There exists a function $f : \mathcal{X} \rightarrow \mathbb{R}^n$ so that*
578 *$h(X) = \sum_{x \in X} f(x)$ is unique for each multiset $X \subset \mathcal{X}$ of bounded size. Moreover, any multiset*
579 *function g can be decomposed as $g(X) = \varphi(\sum_{x \in X} f(x))$ for some function φ .*

580 B.2 Proof of Proposition 1: Relationship between DTDGs and CTDGs

581 *Proof.* We prove the two statements in Proposition 1 separately. In the following, we treat CTDGs as
582 sets of events up to a given timestamp.

583 **Statement 1:** For any DTDG we can build a CTDG that contains the same information.

584 A DTDG consists of a sequence of graphs with no temporal information. We can model this using the
585 CTDG formalism by setting a fixed time difference δ between consecutive elements $G(t_i), G(t_{i+1})$
586 of the CTDG, i.e., $t_{i+1} - t_i = \delta$ for all $i \geq 0$.

587 Consider a DTDG given by the sequence (G_1, G_2, \dots) . To build the equivalent CTDG, we define
588 $S(G_i)$ as the set of edge events corresponding to G_i , i.e., $S(G_i) = \{(u, v, i\delta) : (u, v) \in E(G_i)\}$.
589 We also make the edge features of these events match those in the DTDG, i.e., $e_{uv}(i\delta) = e_{uv} \in \mathcal{E}_i$.
590 To account for node features, for all $u \in V(G_i)$, we create an event $(u, \diamond, i\delta)$ between u and a
591 dummy node \diamond , with feature $e_{u\diamond}(i\delta) = x_u \in \mathcal{X}_i$. Let $C(G_i)$ denote the set comprising these
592 node-level events. Then, we can construct the CTDG $G(t_i) = \cup_{j=1}^i S(G_j) \cup C(G_j)$ for $i = 1, \dots$
593 Reconstructing the DTDG (G_1, G_2, \dots) is trivial. To build G_i , it suffices to select all events at time
594 $i\delta$ in the CTDG. Events involving \diamond determine node features and the remaining ones constitute edges
595 in the DTDG.

596 **Statement 2:** The converse holds if the CTDG timestamps form a subset of some uniformly spaced
597 countable set.

598 We say that a countable set $A \subset \mathbb{R}$ is uniformly spaced if there exists some $\delta \in \mathbb{R}$ such that
599 $a_{i+1} - a_i = \delta$ for all i where (a_1, a_2, \dots) is the ordered sequence formed from elements a_r of A ,
600 i.e., $a_1 < a_2 < \dots < a_i < a_{i+1}, \dots$

601 Note that DTDGs are naturally represented by a set of uniformly spaced timestamps. This is because
602 DTDGs correspond to sequences that do not contain any time information. Let us denote the set of
603 CTDG timestamps $T \subseteq \mathcal{T}$ such that \mathcal{T} is countable and uniformly spaced. Our idea is to construct a
604 DTDG sequence with timestamps that coincide with the elements in \mathcal{T} . Then, since $T \subseteq \mathcal{T}$, we do
605 not lose any information pertaining to events occurring at timestamps given by T . Without loss of
606 generality, in the following we assume that the elements of T and \mathcal{T} are arranged in their increasing
607 order respectively, i.e., $t_i < t_{i+1}$ for all i , and $\tau_k < \tau_{k+1}$ for all k .

608 Consider a CTDG $(G(t_1), G(t_2), \dots)$ such that $G(t_i) = \{(u, v, t) : t \in T \text{ and } t \leq t_i\}$ for $t_i \in T$.
609 Also, let us denote $H(t_i) = \{(u, v, t) \in G(t_i) : t = t_i\}$ the set of events at time $t_i \in T$. We can
610 build a corresponding DTDG (G_1, G_2, \dots) such that for all $\tau_k \in \mathcal{T}$ the k -th snapshot G_k is

$$V(G_k) = \begin{cases} \{u : (u, \cdot, \tau_k) \in H(\tau_k)\}, & \text{if } \tau_k \in T; \\ \emptyset, & \text{otherwise.} \end{cases}$$

$$E(G_k) = \begin{cases} \{(u, v) : (u, v, \tau_k) \in H(\tau_k)\}, & \text{if } \tau_k \in T; \\ \emptyset, & \text{otherwise.} \end{cases}$$

611 To recover the original CTDG, we can adapt the reconstruction procedure we used in the previous
612 part of the proof. We define

$$\tilde{I} = \{(i, k) \in \mathbb{N} \times \mathbb{N} : \tau_k = t_i \text{ for } t_i \in T \text{ and } \tau_k \in \mathcal{T}\}. \quad (\text{S4})$$

613 Note that we can treat \tilde{I} as a map by defining $\tilde{I}(i) = k$ if and only if $(i, k) \in \tilde{I}$. To recover the
614 original CTDG, we first create the set of events $S(G_k) = \{(u, v, k\delta) : (u, v) \in E(G_k)\}$. Then, we
615 build $G(t_i) = \cup_{j: j \leq \tilde{I}(i)} S(G_j)$ for $t_i \in T$. \square

616 **B.3 Proof of Lemma 1**

617 *Proof.* Here we show that if two nodes u and v have isomorphic (L -depth) TCTs, then MP-TGNs
 618 (with L -layers) compute identical embeddings for u and v . Formally, let $T_{u,\ell}(t)$ denote the TCT of u
 619 with ℓ layers. We want to show that $T_{u,\ell}(t) \cong T_{v,\ell}(t) \Rightarrow h_u^{(\ell)}(t) = h_v^{(\ell)}(t)$. We employ a proof by
 620 induction on ℓ . Since there is no ambiguity, we drop the dependence on time in the following.

621 [*Base case*] Consider $\ell = 1$. By the isomorphism assumption $T_{u,1} \cong T_{v,1}$, $h_u^{(0)} = s_u = s_v = h_v^{(0)}$ —
 622 roots of both trees have the same states. Also, for any children i of u in $T_{u,1}$ there is a corresponding
 623 one (with annotated edges) $f(i)$ in $T_{v,1}$ with $s_i = s_{f(i)}$. Recall that the ℓ -th layer aggregation
 624 function $\text{AGG}^{(\ell)}(\cdot)$ acts on multisets of tiplets of previous-layer embeddings, edge features and
 625 timestamps of temporal neighbors (see Equation 1). Since the temporal neighbors of u correspond to
 626 its children in $T_{u,1}$, then the output of the aggregation function for u and v are identical: $\tilde{h}_u^{(1)} = \tilde{h}_v^{(1)}$.
 627 In addition, since the initial embeddings of u and v are also equal (i.e., $h_u^{(0)} = h_v^{(0)}$), we can ensure
 628 that the update function returns $h_u^{(1)} = h_v^{(1)}$.

629 [*Induction step*] Assuming that $T_{u,\ell-1} \cong T_{v,\ell-1} \Rightarrow h_u^{(\ell-1)} = h_v^{(\ell-1)}$ for any pair of nodes u and
 630 v , we will show that $T_{u,\ell} \cong T_{v,\ell} \Rightarrow h_u^{(\ell)} = h_v^{(\ell)}$. For any children i of u , let us define the subtree
 631 of $T_{u,\ell}$ rooted at i by T_i . We know that T_i has depth $\ell - 1$, and since $T_{u,\ell} \cong T_{v,\ell}$, there exists a
 632 corresponding subtree of T_v (of depth $\ell - 1$) rooted at $f(i)$ such that $T_i \cong T_{f(i)}$. Using the induction
 633 hypothesis, we obtain that the multisets of embeddings from the children i of u and children $f(i)$ of
 634 v are identical. Note that if two ℓ -depth TCTs are isomorphic, they are also isomorphic up to depth
 635 $\ell - 1$, i.e., $T_{u,\ell} \cong T_{v,\ell}$ implies $T_{u,\ell-1} \cong T_{v,\ell-1}$ and, consequently, $h_u^{(\ell-1)} = h_v^{(\ell-1)}$ (by induction
 636 hypothesis). Thus, the input of the aggregation and update functions are identical and they compute
 637 the same embeddings for u and v . \square

638 **B.4 Proof of Proposition 2: Most expressive MP-TGNs**

639 *Proof.* Consider MP-TGNs with parameter values that make $\text{AGG}^{(\ell)}(\cdot)$ and $\text{UPDATE}^{(\ell)}(\cdot)$ injective
 640 functions on multisets of triples of hidden representations, edge features and timestamps. The
 641 existence of these parameters is guaranteed by the fact that, at any given time t , the space of observed
 642 node states (and hidden embeddings from temporal neighbors), edge features and timestamps is finite
 643 (see Lemma B3).

644 Again, let $T_{u,\ell}(t)$ denote the TCT of u with ℓ layers. We want to prove that, under the injectivity
 645 assumption, if $T_{u,\ell}(t) \not\cong T_{v,\ell}(t)$, then $h_u^{(\ell)}(t) \neq h_v^{(\ell)}(t)$ for any two nodes u and v . In the following,
 646 we simplify notation by removing the dependence on time. We proceed with proof by induction on
 647 the TCT's depth ℓ . Also, keep in mind that $\varphi_\ell = \text{UPDATE}^{(\ell)} \circ \text{AGG}^{(\ell)}$ is injective for any ℓ .

648 [*Base case*] For $\ell = 1$, if $T_{u,1} \not\cong T_{v,1}$ then the root node states are different (i.e., $s_u \neq s_v$) or the
 649 multiset of states/edge features/ timestamps triples from u and v 's children are different. In both
 650 cases, the inputs of φ_ℓ are different and it therefore outputs different embeddings for u and v .

651 [*Induction step*] The inductive hypothesis is $T_{u,\ell-1} \not\cong T_{v,\ell-1} \Rightarrow h_u^{(\ell-1)} \neq h_v^{(\ell-1)}$ for any pair of
 652 nodes u and v . If $T_{u,\ell} \not\cong T_{v,\ell}$, at least one of the following holds: i) the states of u and v are different,
 653 ii) the multisets of edges (edge features/ timestamps) with endpoints in u and v are different, or iii)
 654 there is no pair-wise isomorphism between the TCTs rooted at u and v 's children. In the first two
 655 cases, φ_ℓ trivially outputs different embeddings for u and v . We are left with the case in which only
 656 the latter occurs. Using our inductive hypothesis, the lack of a (isomorphism ensuring) bijection
 657 between the TCTs rooted at u and v 's children implies there is also no bijection between their multiset
 658 of embeddings. In turn, this guarantees that φ will output different embeddings for u and v . \square

659 **B.5 Proof of Proposition 3: The role of memory**

660 *Proof.* We prove the two parts of the proposition separately. In the following proofs, we rely on the
 661 concept of monotone TCTs (see Definition 2).

662 **Statement 1:** If $L < \Delta$: $\mathcal{Q}_L^{[M]}$ is strictly stronger than \mathcal{Q}_L .

663 We know that the family of L -layer MP-TGNs with memory comprises the family of L -layer MP-
664 TGNs without memory (we can assume identity memory). Therefore, $\mathcal{Q}_L^{[M]}$ is at least as powerful
665 as \mathcal{Q}_L . To show that $\mathcal{Q}_L^{[M]}$ is strictly stronger (more powerful) than \mathcal{Q}_L , when $L < \Delta$, it suffices
666 to create an example for which memory can help distinguish a pair of nodes. We provide a trivial
667 example in Figure S1 for $L = 1$. Note that the 1-depth TCTs of u and v are isomorphic when no
668 memory is used. However, when equipped with memory, the interaction (b, c, t_1) affects the states of
669 v and c , making the 1-depth TCTs of u and v (at time $t > t_2$) no longer isomorphic.

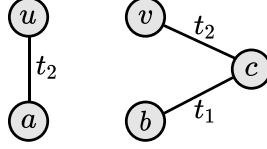


Figure S1: Temporal graph where all initial node features and edge features are identical. We assume that $t_2 > t_1$.

670 **Statement 2:** For any L : $\mathcal{Q}_{L+\Delta}$ is at least as powerful as $\mathcal{Q}_L^{[M]}$.

671 It suffices to show if $\mathcal{Q}_{L+\Delta}$ cannot distinguish a pair of nodes u and v , $\mathcal{Q}_L^{[M]}$ cannot distinguish them
672 too. Let $T_{u,L}^M(t)$ and $T_{u,L}(t)$ denote the L -depth TCTs of u with and without memory respectively.
673 Using Lemma 1, this is equivalent to showing that $T_{u,L+\Delta}(t) \cong T_{v,L+\Delta}(t) \Rightarrow T_{u,L}^M(t) \cong T_{v,L}^M(t)$,
674 since no MP-TGN can separate nodes associated with isomorphic TCTs. In the following, when we
675 omit the number of layers from TCTs, we assume TCTs of arbitrary depth.

676 **Step 1:** Characterizing the dependence of memory on initial states and events in the dynamic graph.

677 We now show that the memory for a node u , after processing all events with timestamp $\leq t_n$, depends
678 on the initial states of a set of nodes \mathcal{V}_u^n , and a set of events annotated with their respective timestamps
679 and features \mathcal{B}_u^n . If at time t_n no event involves a node z , we set $\mathcal{B}_z^n = \mathcal{B}_z^{n-1}$ and $\mathcal{V}_z^n = \mathcal{V}_z^{n-1}$. We
680 also initialize $\mathcal{B}_u^0 = \emptyset$ and $\mathcal{V}_u^0 = \{u\}$ for all nodes u . We proceed with a proof by induction on the
681 number of observed timestamps n .

682 [*Base case*] Let $\mathcal{I}_1(u) = \{v : (u, v, t_1) \in \mathcal{G}(t_1^+)\}$ be the set of nodes interacting with u at time
683 t_1 , where $\mathcal{G}(t_1^+)$ is the temporal graph right after t_1 . Similarly, let $\mathcal{J}_1(u) = \{(u, \cdot, t_1) \in \mathcal{G}(t_1^+)\}$
684 be the set of events involving u at time t_1 . Recall that before t_1 , all memory states equal initial
685 node features (i.e., $s_u(t_1^-) = s_u(0)$). Then, the updated memory (see Equation 3) for u depends on
686 $\mathcal{V}_u^1 = \mathcal{V}_u^0 \cup_{v \in \mathcal{I}(u)} \mathcal{V}_v^0$, $\mathcal{B}_u^1 = \mathcal{B}_u^0 \cup \mathcal{J}_1(u)$.

687 [*Induction step*] Assume that for timestamp t_{n-1} the proposition holds. We now show that it holds for
688 t_n . Since the proposition holds for $n-1$ timestamps, we know that the memory of any w that interacts
689 with u in t_n , i.e. $w \in \mathcal{I}_n(u)$, depends on \mathcal{V}_w^{n-1} and \mathcal{B}_w^{n-1} , and the memory of u so far depends on
690 \mathcal{V}_u^{n-1} and \mathcal{B}_u^{n-1} . Then, the updated memory for u depends on $\mathcal{V}_u^n = \mathcal{V}_u^{n-1} \cup_{w \in \mathcal{I}(u)} \{w, u\} \cup \mathcal{V}_w^{n-1}$
691 and $\mathcal{B}_u^n = \mathcal{B}_u^{n-1} \cup \mathcal{J}_n(u) \cup_{w \in \mathcal{I}(u)} \mathcal{B}_w^{n-1}$.

692 **Step 2:** $(z, w, t_{zw}) \in \mathcal{B}_u^n$ if and only if there is a path $(u_k, t_k = t_{zw}, u_{k+1})$ in $\tilde{T}_u(t_n^+)$ — the
693 monotone TCT of u (see Definition 2) after processing events with timestamp $\leq t_n$ — with either
694 $\#u_k = z, \#u_{k+1} = w$ or $\#u_k = w, \#u_{k+1} = z$.

695 [*Forward direction*] An event (z, w, t_{zw}) with $t_{zw} \leq t_n$ will be in \mathcal{B}_u^n only if $z = u$ or $w = u$, or
696 if there is a subset of events $\{(u, \#u_1, t_1), (\#u_1, \#u_2, t_2), \dots, (\#u_k, \#u_{k+1}, t_{zw})\}$ with $\#u_k = z$ and
697 $\#u_{k+1} = w$ such that $t_n \geq t_1 > \dots > t_{zw}$. In either case, this will lead to root-to-leaf path in $\tilde{T}_u(t_n^+)$
698 passing through (u_k, t_{zw}, u_{k+1}) . This subset of events can be easily obtained by backtracking edges
699 that caused unions/updates in the procedure from **Step 1**.

700 [*Backward direction*] Assume there is a subpath $p = (u_k, t_k = t_{zw}, u_{k+1}) \in \tilde{T}_u(t_n^+)$ with $\#u_k = z$
701 and $\#u_{k+1} = w$ such that $(z, w, t_{zw}) \notin \mathcal{B}_u^n$. Since we can obtain p from $\tilde{T}_u(t_n^+)$, we know that the
702 sequence of events $r = ((u, \#u_1, t_1), \dots, (\#u_{k-2}, \#u_{k-1} = z, t_{k-1}), (z, w, t_k = t_{zw}))$ happened and
703 that $t_i > t_{i+1} \forall i$. However, since $(z, w, t_{zw}) \notin \mathcal{B}_u^n$, there must be no monotone walk starting from
704 u going through the edge (z, w, t_{zw}) to arrive at w , which is exactly what r characterizes. Thus, we
705 reach contradiction.

706 Note that the nodes in \mathcal{V}_u^n are simply the nodes that have an endpoint in the events \mathcal{B}_u^n , and therefore
 707 are also nodes in $\tilde{T}_u(t_n^+)$ and vice-versa.

708 **Step 3:** For any node u , there is a bijection that maps $(\mathcal{V}_u^n, \mathcal{B}_u^n)$ to $\tilde{T}_u(t_n^+)$.

709 First, we note that $(\mathcal{V}_u^n, \mathcal{B}_u^n)$ depends on a subset of all events, which we represent as $\mathcal{G}' \subseteq \mathcal{G}(t_n^+)$.
 710 Since \mathcal{B}_u^n contains all events in \mathcal{G}' and $(\mathcal{V}_u^n, \mathcal{B}_u^n)$ can be uniquely constructed from \mathcal{G}' , then there is a
 711 bijection g that maps from \mathcal{G}' to $(\mathcal{V}_u^n, \mathcal{B}_u^n)$.

712 Similarly, $\tilde{T}_u(t_n^+)$ also depends on a subset of events which we denote by $\mathcal{G}'' \subseteq \mathcal{G}(t_n^+)$. We note that
 713 the unique events in $\tilde{T}_u(t_n^+)$ correspond to \mathcal{G}'' , and we can uniquely build the tree $\tilde{T}_u(t_n^+)$ from \mathcal{G}'' .
 714 This implies that there is a bijection h that maps from \mathcal{G}'' to $\tilde{T}_u(t_n^+)$.

715 Previously, we have shown that all events in \mathcal{B}_u^n are also in $\tilde{T}_u(t_n^+)$ and vice-versa. This implies that
 716 both sets depend on the same events, and thus on the same subset of all events, i.e., $\mathcal{G}' = \mathcal{G}'' = \mathcal{G}_S$.
 717 Since there is a bijection g between \mathcal{G}_S and $(\mathcal{V}_u^n, \mathcal{B}_u^n)$, and a bijection h between \mathcal{G}_S and $\tilde{T}_u(t_n^+)$,
 718 there exists a bijection f between $(\mathcal{V}_u^n, \mathcal{B}_u^n)$ and $\tilde{T}_u(t_n^+)$.

719 **Step 4:** If $T_{u,L+\Delta}(t^+) \cong T_{v,L+\Delta}(t^+)$, then $T_{u,L}^M(t^+) \cong T_{v,L}^M(t^+)$.

720 To simplify notation, we omit here the dependence on time.

721 Any node $w \in T_{u,L}^M$ also appears in $T_{u,L+\Delta}$ at the same level. The subtree of $T_{u,L+\Delta}$ rooted at w ,
 722 denoted here by T'_w , has depth at least $k \geq \Delta$. Note that T'_w corresponds to the k -depth TCT of $\sharp w$.
 723 Since the depth of T'_w is at least Δ , we know that the $\tilde{T}'_w \cong \tilde{T}_{\sharp w}$ — i.e., imposing time-constraints to
 724 T'_w results in the monotone TCT of node $\sharp w$. Also, because the memory of $\sharp w$ depends on $\tilde{T}_{\sharp w}$, T'_w
 725 comprises the information used to compute the memory state of $\sharp w$. Note that this applies to any w in
 726 $T_{u,L}^M$; thus, $T_{u,L+\Delta}$ contains all we need to compute the states of any node of the dynamic graph that
 727 appears in $T_{u,L}^M$. The same argument applies to $T_{v,L+\Delta}$ and $T_{v,L}^M$. Finally, since $T_{u,L}$ can be uniquely
 728 computed from $T_{u,L+\Delta}$, and $T_{v,L}^M$ from $T_{v,L+\Delta}$, if $T_{u,L+\Delta} \cong T_{v,L+\Delta}$, then $T_{u,L}^M \cong T_{v,L}^M$. \square

729 B.6 Proof of Proposition 4: Limitations of TGAT and TGN-Att

730 *Proof.* In this proof, we first provide an example of a dynamic graph where the TCTs of two
 731 nodes u and v are not isomorphic. Then, we show that we can not find a TGAT model such that
 732 $h_u^{(L)}(t) \neq h_v^{(L)}(t)$, i.e., TGAT does not distinguish u and v . Next, we show that even if we consider
 733 TGATs with memory (TGN-Att), it is still not possible to distinguish nodes u and v in our example.

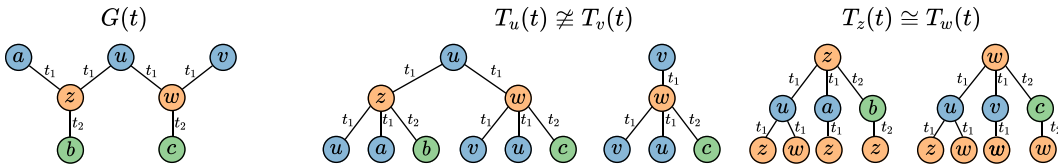


Figure S2: (Leftmost) Example of a temporal graph for which TGN-Att and TGAT cannot distinguish nodes u and v even though their TCTs are non-isomorphic. Colors denote node features and all edge features are identical, and $t_2 > t_1$ (and $t > t_2$). (Right) The 2-depth TCTs of nodes u, v, z and w . The TCTs of u and v are non-isomorphic whereas the TCTs of z and w are isomorphic.

734 Figure S2(leftmost) provides a temporal graph where all edge events have the same edge features.
 735 Colors denote node features. As we can observe, the TCTs of nodes u and v are not isomorphic. In
 736 the following, we consider node distinguishability at time $t > t_2$.

737 **Statement 1:** TGAT cannot distinguish the nodes u and v in our example.

738 **Step 1:** For any TGAT with ℓ layers, we have that $h_w^{(\ell)}(t) = h_z^{(\ell)}(t)$.

739 We note that the ℓ -layer TCTs of nodes w and z are isomorphic, for any ℓ . To see this, one can consider
 740 the symmetry around node u that allows us to define a node permutation function (bijection) f given
 741 by $f(z) = w, f(w) = z, f(a) = v, f(u) = u, f(b) = c, f(c) = b, f(v) = a$. Figure S2(right)
 742 provides an illustration of the 2-depth TCTs of z and w at time $t > t_2$.

743 By Lemma 1, if the ℓ -layer TCTs of two nodes z and w are isomorphic, then no ℓ -layer MP-TGN can
 744 distinguish them. Thus, we conclude that $h_w^{(\ell)}(t) = h_z^{(\ell)}(t)$ for any TGAT with arbitrary number of
 745 layers ℓ .

746 **Step 2:** There is no TGAT such that $h_v^{(\ell)}(t) \neq h_u^{(\ell)}(t)$.

747 To compute $h_v^{(\ell)}(t)$, TGAT aggregates the messages of v 's temporal neighbors at layer $\ell - 1$, and
 748 then combines $h_v^{(\ell-1)}(t)$ with the aggregated message $\tilde{h}_v^{(\ell-1)}(t)$ to obtain $h_v^{(\ell)}(t)$.

749 Note that $\mathcal{N}(u, t) = \{(z, e, t_1), (w, e, t_1)\}$ and $\mathcal{N}(v, t) = \{(w, e, t_1)\}$, where e denotes an edge
 750 feature vector. Also, we have previously shown that $h_w^{(\ell-1)}(t) = h_z^{(\ell-1)}(t)$.

751 Using the TGAT aggregation layer (Equation S2), the query vectors of u and v are $q_u =$
 752 $[h_u^{(\ell-1)}(t) \parallel \phi(0)]W_q^{(\ell)}$ and $q_v = [h_v^{(\ell-1)}(t) \parallel \phi(0)]W_q^{(\ell)}$, respectively.

753 Since all events have the common edge features e , the matrices $C_u^{(\ell)}$ and $C_v^{(\ell)}$ share the same vector in
 754 their rows. The single-row matrix $C_v^{(\ell)}$ is given by $C_v^{(\ell)} = [h_w^{(\ell-1)}(t) \parallel \phi(t - t_1) \parallel e]$, while the two-row
 755 matrix $C_u^{(\ell)} = \begin{bmatrix} [h_w^{(\ell-1)}(t) \parallel \phi(t - t_1) \parallel e]; [h_z^{(\ell-1)}(t) \parallel \phi(t - t_1) \parallel e] \end{bmatrix}$, with $h_w^{(\ell-1)}(t) = h_z^{(\ell-1)}(t)$. We
 756 can express $C_u^{(\ell)} = [1, 1]^\top r$ and $C_v^{(\ell)} = r$, where r denotes the row vector $r = [h_z^{(\ell-1)}(t) \parallel \phi(t -$
 757 $t_1) \parallel e]$.

758 Using the key and value matrices of node v , i.e., $K_v = C_v^{(\ell)}W_K^{(\ell)}$ and $V_v = C_v^{(\ell)}W_V^{(\ell)}$, we have that

$$\begin{aligned} \tilde{h}_v^{(\ell)}(t) &= \text{softmax}(q_v K_v^\top) V_v \\ &= \underbrace{\text{softmax}(q_v K_v^\top)}_{=1} r W_V^{(\ell)} && \text{[softmax of a single element is 1]} \\ &= r W_V^{(\ell)} \\ &= \underbrace{\text{softmax}(q_u K_u^\top) [1, 1]^\top}_{=1} r W_V^{(\ell)} = \tilde{h}_u^{(\ell)}(t) && \text{[softmax outputs a convex combination]} \end{aligned}$$

759 We have shown that the aggregated messages of nodes u and v are the same at any layer ℓ . We note
 760 that the initial embeddings are also identical $h_v^{(0)}(t) = h_u^{(0)}(t)$ as u and v have the same color. Recall
 761 that the update step is $h_v^{(\ell)}(t) = \text{MLP}(h_v^{(\ell-1)}(t), \tilde{h}_v^{(\ell)}(t))$. Therefore, if the initial embeddings are
 762 identical, and the aggregated messages at each layer are also identical, we have that $h_u^{(\ell)}(t) = h_v^{(\ell)}(t)$
 763 for any ℓ .

764 **Statement 2:** TGN-Att cannot distinguish the nodes u and v in our example.

765 We now show that adding a memory module to TGAT produces node states such that $s_u(t) = s_v(t) =$
 766 $s_a(t)$, $s_z(t) = s_w(t)$, and $s_b(t) = s_c(t)$. If that is the case, then these node states could be treated
 767 as node features in a equivalent TGAT model of our example in Figure S2, proving that there is no
 768 TGN-Att such that $h_v^{(\ell)}(t) \neq h_u^{(\ell)}(t)$. In the following, we consider TGN-Att with average message
 769 aggregators (see Appendix A).

770 We begin by showing that $s_a(t) = s_u(t) = s_v(t)$ after memory updates. We note
 771 that the message node a receives is $[e \parallel t_1 \parallel s_z(t_1^-)]$. The message node u receives is
 772 $\text{MEANAGG}([e \parallel t_1 \parallel s_w(t_1^-)], [e \parallel t_1 \parallel s_z(t_1^-)])$, but since $s_w(t_1^-) = s_z(t_1^-)$, both messages are the same,
 773 and the average aggregator outputs $[e \parallel t_1 \parallel s_z(t_1^-)]$. Finally, the message that node v receives is
 774 $[e \parallel t_1 \parallel s_w(t_1^-)] = [e \parallel t_1 \parallel s_z(t_1^-)]$. Since all three nodes receive the same message and have the same
 775 initial features, their updated memory states are identical.

776 Now we show that $s_z(t) = s_w(t)$, for $t_1 \leq t < t_2$. Note that the message that node z receives is
 777 $\text{MEANAGG}([e \parallel t_1 \parallel s_a(t_1^-)], [e \parallel t_1 \parallel s_u(t_1^-)]) = [e \parallel t_1 \parallel s_u(t_1^-)]$, with $s_u(t_1^-) = s_a(t_1^-)$. The message
 778 that node w receives is $\text{MEANAGG}([e \parallel t_1 \parallel s_u(t_1^-)], [e \parallel t_1 \parallel s_v(t_1^-)]) = [e \parallel t_1 \parallel s_u(t_1^-)]$. Again, since
 779 the initial features and the messages received by each node are equal, $s_z(t) = s_w(t)$ for $t_1 \leq t < t_2$.

780 We can then use this to show that $s_z(t) = s_w(t)$ for $t \geq t_2$. Note that at time t_2 , the message
 781 that nodes z and w receive are $[e \parallel t_2 - t_1 \parallel s_b(t_2^-)]$ and $[e \parallel t_2 - t_1 \parallel s_c(t_2^-)]$, respectively. Also, note

782 that $s_b(t_2^-) = s_c(t_2^-) = s_b(0) = s_c(0)$ as the states of b and c are only updated at t_2 . Because
 783 the received messages and the previous states (before t_2) of z and w are identical, we have that
 784 $s_z(t) = s_w(t)$ for $t \geq t_2$.

785 Finally, we show that $s_b(t) = s_c(t)$. Using that $s_z(t_2^-) = s_w(t_2^-)$ in conjunction with the fact that
 786 node b receives message $[[e||t_2 - t_1||s_z(t_2^-)]]$, and node c receives $[e||t_2 - t_1||s_w(t_2^-)]$, we obtain
 787 $s_b(t) = s_c(t)$ since initial memory states and messages that the nodes received are the same. \square

788 B.7 Proof of Proposition 5: Limitations of MP-TGN and CAWs

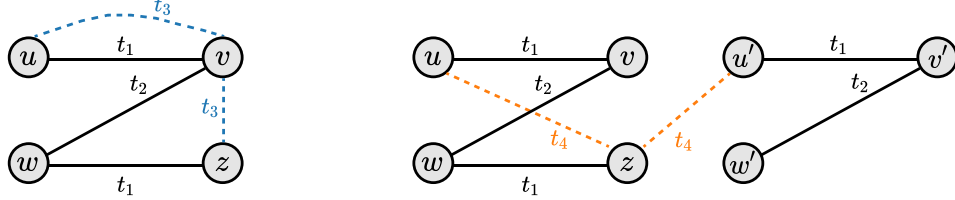


Figure S3: **(Left)** Example of a temporal graph for which CAW can distinguish the events (u, v, t_3) and (z, v, t_3) but MP-TGNs cannot. We assume that all edge and node features are identical, and $t_{k+1} > t_k$ for all k . **(Right)** Example for which MP-TGNs can distinguish (u, z, t_4) and (u', z, t_4) but CAW cannot.

789 *Proof.* Using the example in Figure S3(Left), we adapt a construction by Wang et al. [37] to show
 790 that CAW can separate events that MP-TGNs adopting node embedding concatenation cannot. We
 791 first note that the TCTs of u and z are isomorphic. Thus, since v is a common endpoint in (u, v, t_3)
 792 and (z, v, t_3) , no MP-TGN can distinguish these two events. Nonetheless, CAW obtains the following
 793 anonymized walks for the event (u, v, t_3) :

$$\begin{aligned} & \underbrace{\{[1, 0, 0], [0, 1, 0]\}}_{I_{CAW}(u; S_u, S_v)} \xrightarrow{t_1} \underbrace{\{[0, 1, 0], [2, 0, 0]\}}_{I_{CAW}(v; S_u, S_v)} \\ & \underbrace{\{[0, 1, 0], [2, 0, 0]\}}_{I_{CAW}(v; S_u, S_v)} \xrightarrow{t_1} \underbrace{\{[1, 0, 0], [0, 1, 0]\}}_{I_{CAW}(u; S_u, S_v)} \\ & \underbrace{\{[0, 1, 0], [2, 0, 0]\}}_{I_{CAW}(v; S_u, S_v)} \xrightarrow{t_2} \underbrace{\{[0, 0, 0], [0, 1, 0]\}}_{I_{CAW}(w; S_u, S_v)} \xrightarrow{t_1} \underbrace{\{[0, 0, 0], [0, 0, 1]\}}_{I_{CAW}(z; S_u, S_v)} \end{aligned}$$

794 and the walks associated with (z, v, t_3) are (here we omit underbraces for readability):

$$\begin{aligned} & \{[1, 0, 0], [0, 0, 1]\} \xrightarrow{t_1} \{[0, 1, 0], [0, 1, 0]\} \\ & \{[0, 0, 0], [2, 0, 0]\} \xrightarrow{t_1} \{[0, 0, 0], [0, 1, 0]\} \\ & \{[0, 0, 0], [2, 0, 0]\} \xrightarrow{t_2} \{[0, 1, 0], [0, 1, 0]\} \xrightarrow{t_1} \{[1, 0, 0], [0, 0, 1]\} \end{aligned}$$

795 In this example, assume that MLPs used to encode each walk correspond to identity mappings. Then,
 796 the sum of the elements in each set is injective since each element of the sets in the anonymized
 797 walks are one-hot vectors. We note that, in this example, we can simply choose a RNN that sums the
 798 vectors in each sequence (walks), and then apply a mean readout layer (or pooling aggregator) to
 799 obtain distinct representations for (u, v, t_3) and (z, v, t_3) .

800 We now use the example in Figure S3(Right) to show that MP-TGNs can separate events that CAW
 801 cannot. To see why MP-TGNs can separate the events (u, z, t_4) and (u', z, t_4) , it suffices to observe
 802 that the 4-depth TCTs of u and u' are non-isomorphic. Thus, a MP-TGN with injective layers could
 803 distinguish such events. Now, let us take a look at the anonymized walks for (u, z, t_4) :

$$\begin{aligned} & \underbrace{\{[1, 0], [0, 0]\}}_{I_{CAW}(u; S_u, S_z)} \xrightarrow{t_1} \underbrace{\{[0, 1], [0, 0]\}}_{I_{CAW}(v; S_u, S_z)} \\ & \underbrace{\{[0, 0], [1, 0]\}}_{I_{CAW}(z; S_u, S_z)} \xrightarrow{t_1} \underbrace{\{[0, 0], [0, 1]\}}_{I_{CAW}(w; S_u, S_z)} \end{aligned}$$

804 and for (u', z, t_4) :

$$\begin{aligned} \underbrace{\{[1, 0], [0, 0]\}}_{I_{CAW}(u'; S_{u'}, S_z)} &\xrightarrow{t_1} \underbrace{\{[0, 1], [0, 0]\}}_{I_{CAW}(v'; S_{u'}, S_z)} \\ \underbrace{\{[0, 0], [1, 0]\}}_{I_{CAW}(z; S_{u'}, S_z)} &\xrightarrow{t_1} \underbrace{\{[0, 0], [0, 1]\}}_{I_{CAW}(w; S_{u'}, S_z)} \end{aligned}$$

805 Since the sets of walks are identical, they must have the same embedding. Therefore, there is no
806 CAW model that can separate these two events. \square

807 **B.8 Proof of Proposition 6: Injective MP-TGNs and the temporal WL test**

808 We want to prove that injective MP-TGNs can separate two temporal graphs if and only if the
809 temporal WL does the same. Our proof comprises two parts. We first show that if an MP-TGN
810 produces different multisets of embeddings for two non-isomorphic temporal graphs $\mathcal{G}(t)$ and $\mathcal{G}'(t)$,
811 then the temporal WL decides these graphs are not isomorphic. Then, we prove that, if the temporal
812 WL decides $\mathcal{G}(t)$ and $\mathcal{G}'(t)$ are non-isomorphic, there is an injective MP-TGN (i.e., with injective
813 message-passing layers) that outputs distinct multisets of embeddings.

814 **Statement 1:** Temporal WL is at least as powerful as MP-TGNs.

815 See Lemma B2 for proof.

816 **Statement 2:** Injective MP-TGN is at least as powerful as temporal WL.

817 *Proof.* To prove this, we can repurpose the proof of Theorem 3 in [42]. In particular, we assume MP-
818 TGNs that meet the injective requirements of Proposition 2, i.e., MP-TGNs that implement injective
819 aggregate and update functions on multisets of hidden representations from temporal neighbors.
820 Following their footprints, we prove that there is a injection φ to the set of embeddings of all nodes
821 in a temporal graph from their respective colors in the temporal WL test. We do so via induction
822 on the number of layers ℓ . To achieve our purpose, we can assume identity memory without loss of
823 generality.

824 The base case ($\ell = 0$) is straightforward since the temporal WL test initializes colors with node
825 features. We now focus on the inductive step. Suppose the proposition holds for $\ell - 1$. Note that our
826 update function:

$$h_v^{(\ell)}(t) = \text{UPDATE}^{(\ell)} \left(h_v^{(\ell-1)}(t), \text{AGG}^{(\ell)}(\{(h_u^{(\ell-1)}(t), t - t', e) \mid (u, e, t') \in \mathcal{N}(v, t)\}) \right)$$

827 can be rewritten using φ as a function of node colors:

$$h_v^{(\ell)}(t) = \text{UPDATE}^{(\ell)} \left(\varphi(c^{\ell-1}(v)), \text{AGG}^{(\ell)}(\{(\varphi(c^{\ell-1}(u)), t - t', e) \mid (u, e, t') \in \mathcal{N}(v, t)\}) \right).$$

828 Note that the composition of injective functions is also injective. In addition, time-shifting operations
829 are also injective. Thus we can construct an injection ψ such that:

$$\begin{aligned} h_v^{(\ell)}(t) &= \psi \left(c^{\ell-1}(v), \{(c^{\ell-1}(u), t', e) \mid (u, e, t') \in \mathcal{N}(v, t)\} \right) \\ &= \psi \left(c^{\ell-1}(v), \{(c^{\ell-1}(u), t', e_{uv}(t')) \mid (v, u, t') \in \mathcal{G}(t)\} \right) \end{aligned}$$

830 since there exists an element $(u, e, t') \in \mathcal{N}(v, t)$ if and only if there is an event $(u, v, t') \in \mathcal{G}(t)$ with
831 feature $e_{uv}(t') = e$.

832 Then, we can write:

$$\begin{aligned} h_v^{(\ell)}(t) &= \psi \circ \text{HASH}^{-1} \circ \text{HASH} \left(c^{\ell-1}(v), \{(c^{\ell-1}(u), t', e_{uv}(t')) \mid (u, v, t') \in \mathcal{G}(t)\} \right) \\ &= \psi \circ \text{HASH}^{-1}(c^{(\ell)}(v)) \end{aligned}$$

833 Note $\varphi = \psi \circ \text{HASH}^{-1}$ is injective since it is a composition of two injective functions. We then
834 conclude that if the temporal WL test outputs different multisets of colors, then a suitable MP-TGN
835 outputs different multisets of embeddings. \square

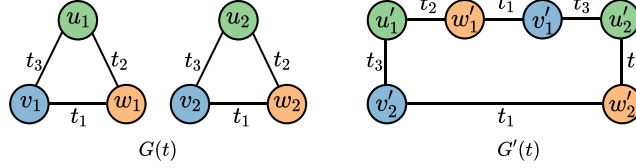


Figure S4: Examples of temporal graphs for which MP-TGNs cannot distinguish the diameter, girth, and number of cycles. For any node in $\mathcal{G}(t)$ (e.g., u_1), there is a corresponding one in $\mathcal{G}'(t)$ (u'_1) whose TCTs are isomorphic.

836 **B.9 Proof of Proposition 7: MP-TGNs and CAWs fail to decide some graph properties**

837 **Statement 1:** MP-TGNs fail to decide some graph properties.

838 *Proof.* Adapting a construction by Garg et al. [9], we provide in Figure S4 an example that demon-
 839 strates Proposition 7. Colors denote node features, and all edge features are identical. The temporal
 840 graphs $\mathcal{G}(t)$ and $\mathcal{G}'(t)$ are non-isomorphic and differ in properties such as diameter (∞ for $\mathcal{G}(t)$ and
 841 3 for $\mathcal{G}'(t)$), girth (3 for $\mathcal{G}(t)$ and 6 for $\mathcal{G}'(t)$), and number of cycles (2 for $\mathcal{G}(t)$ and 1 for $\mathcal{G}'(t)$). In
 842 spite of that, for $t > t_3$, the set of embeddings of nodes in $\mathcal{G}(t)$ is the same as that of nodes in $\mathcal{G}'(t)$
 843 and, therefore, MP-TGNs cannot decide these properties. In particular, by constructing the TCTs of
 844 all nodes at time $t > t_3$, we observe that the TCTs of the pairs (u_1, u'_1) , (u_2, u'_2) , (v_1, v'_1) , (v_2, v'_2) ,
 845 (w_1, w'_1) , (w_2, w'_2) are isomorphic and, therefore, they can not be distinguished (Lemma 1). \square

846 **Statement 2:** CAWs fail to decide some graph properties.

847 Since CAW does not provide a recipe to obtain graph-level embeddings, we first define such a
 848 procedure. Let $\mathcal{G}(t)$ be a temporal graph given as a set of events. We sequentially compute event
 849 embeddings h_γ for each event $\gamma = (u, v, t') \in \mathcal{G}(t)$ respecting the temporal order (two or more
 850 events at the same timestamp are computed in parallel). We then apply a readout layer to the set of
 851 event embeddings to obtain a graph-level representation. We provide a proof assuming this procedure.

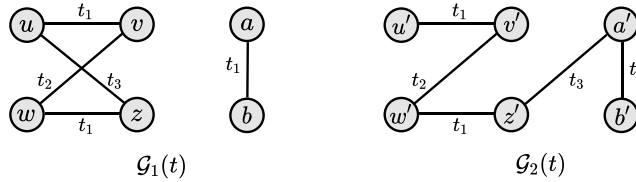


Figure S5: Examples of temporal graphs with different diameters, girths, and numbers of cycles. CAWs fail to distinguish $\mathcal{G}_1(t)$ and $\mathcal{G}_2(t)$.

852 *Proof.* We can adapt our construction in Figure 2 [rightmost] to extend Proposition 7 to CAW. The
 853 idea consists of creating two temporal graphs with different diameters, girths, and numbers of cycles
 854 that comprise events that CAW cannot separate — Figure S5 provides one such construction. In
 855 particular, CAW obtains identical embeddings for (u, z, t_3) and (a', z', t_3) (as shown in Proposition 5).
 856 The remaining events are the same up to node re-labelling and thus also lead to identical embeddings.
 857 Therefore, CAW cannot distinguish $\mathcal{G}_1(t)$ and $\mathcal{G}_2(t)$ although they clearly differ in diameter, girth,
 858 and number of cycles. \square

859 **B.10 Proof of Lemma 2**

860 We now show that the k -th component of the relative positional features $r_{u \rightarrow v}^{(t)}$ corresponds to the
 861 number of occurrences of u at the k -th layer of the monotone TCT of v , and this is valid for all pairs
 862 of nodes u and v of the dynamic graph. We proceed with a proof by induction.

863 [*Base case*] Let us consider $t = 0$, i.e., no events have occurred. By definition, $r_{u \rightarrow v}^{(0)}$ is the zero vector
 864 if $u \neq v$, indicating that node u does not belong to the TCT of v . If $u = v$, then $r_{u \rightarrow u}^{(0)} = [1, 0, \dots, 0]$
 865 corresponds to count 1 for the root of the TCT of u . Thus, for $t = 0$, the proposition holds.

866 [Induction step] Assume that the proposition holds
 867 for all nodes and any time instant before t . We will
 868 show that after the event $\gamma = (u, v, t)$ at time t , the
 869 proposition remains true.

870 Note that the event γ only impacts the monotone
 871 TCTs of u and v . The reason is that the TCTs of
 872 all other nodes have timestamps lower than t , which
 873 prevents the event γ from belonging to any path (with
 874 decreasing timestamps) from the root.

875 Without loss of generality, let us now consider the
 876 impact of γ on the monotone TCT of v . Figure S6
 877 shows how the TCT of v changes after γ , i.e., how it
 878 goes from $\tilde{T}_v(t^-)$ to $\tilde{T}_v(t)$. In particular, the process
 879 attaches the TCT of u to the root node v . Under this
 880 change, we need to update the counts of all node i
 881 in $\tilde{T}_u(t^-)$ regarding how many times it appears in
 882 $\tilde{T}_v(t)$. We do so by adding the counts in $\tilde{T}_u(t^-)$ (i.e.,
 883 $r_{i \rightarrow u}^{(t^-)}$) to $\tilde{T}_v(t^-)$ (i.e., $r_{i \rightarrow v}^{(t^-)}$), accounting for the 1-
 884 layer mismatch, since $\tilde{T}_u(t^-)$ is attached to the first layer. This can be easily achieved with the shift
 885 matrix $P = \begin{bmatrix} 0 & 0 \\ I_{d-1} & 0 \end{bmatrix}$ applied to the counts of any node i in $\tilde{T}_u(t^-)$, i.e.,

$$r_{i \rightarrow v}^{(t)} = P r_{i \rightarrow u}^{(t^-)} + r_{i \rightarrow v}^{(t^-)} \quad \forall i \in \mathcal{V}_u^{(t^-)},$$

886 where $\mathcal{V}_u^{(t^-)}$ comprises the nodes of the original graph that belong to $\tilde{T}_u^{(t^-)}$.

887 Similarly, the event γ also affects the counts of nodes in the TCT of v w.r.t. the TCT of u . To account
 888 for that change, we follow the same procedure and update $r_{j \rightarrow u}^{(t)} = P r_{j \rightarrow v}^{(t^-)} + r_{j \rightarrow u}^{(t^-)}$, $\forall j \in \mathcal{V}_v^{(t^-)}$.

889 **Handling multiple events at the same time.** We now consider the setting where a given node v
 890 interacts with multiple nodes u_1, u_2, \dots, u_J at time t . We can extend the computation of positional
 891 features to this setting in a straightforward manner by noting that each event leads to an independent
 892 branch in the TCT of v . Therefore, the update of the positional features with respect to v is given by

$$r_{i \rightarrow v}^{(t)} = P \sum_{j=1}^J r_{i \rightarrow u_j}^{(t^-)} + r_{i \rightarrow v}^{(t^-)} \quad \forall i \in \bigcup_{j=1}^J \mathcal{V}_{u_j}^{(t^-)}$$

$$\mathcal{V}_v^{(t)} = \mathcal{V}_v^{(t^-)} \bigcup_{j=1}^J \mathcal{V}_{u_j}^{(t^-)}.$$

893 We note that the updates of the positional features of u_1, \dots, u_J remain untouched if they do not
 894 interact with other nodes at time t .

895 B.11 Proof of Proposition 8: Injective function on temporal neighborhood

896 *Proof.* To capture the intuition behind the proof, first consider a multiset M such that $|M| < 4$. We
 897 can assign a unique number $\psi(m) \in \{1, 2, 3, 4\}$ to any distinct element $m \in M$. Also, the function
 898 $h(m) = 10^{-\psi(m)}$ denotes the decimal expansion of $\psi(m)$ and corresponds to reserving one decimal
 899 place for each unique element $m \in M$. Since there are less than 10 elements in the multiset, note
 900 that $\sum_m h(m)$ is unique for any multiset M .

901 To prove the proposition, we also leverage the well-known fact that the Cartesian product of two
 902 countable sets is countable — the Cantor’s (bijective) pairing function $z : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$, with
 903 $z(n_1, n_2) = \frac{(n_1+n_2)(n_1+n_2+1)}{2} + n_2$, provides a proof for that.

904 Here, we consider multisets $M = \{(x_i, e_i, t_i)\}$ whose tuples take values on the Cartesian product of
 905 the countable sets \mathcal{X} , \mathcal{E} , and \mathcal{T} — the latter is also assumed to be bounded. In addition, we assume
 906 the lengths of all multisets are bounded by N , i.e., $|M| < N$ for all M . Since \mathcal{X} and \mathcal{E} are countable,

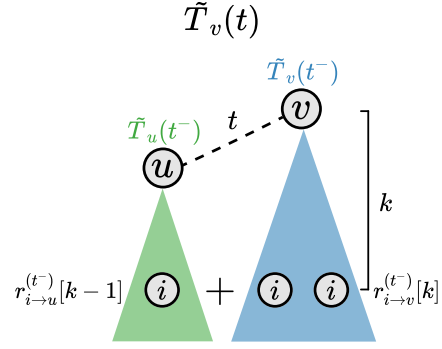


Figure S6: Illustration of how the monotone TCT of v changes after an event between u and v at time t . This allows us to see how to update the positional features of any node i of the dynamic graph that belongs to $\tilde{T}_u(t^-)$ relative to v .

907 there exists an enumeration function $\psi : \mathcal{X} \times \mathcal{E} \rightarrow \mathbb{N}$ for all M . Without loss of generality, we
 908 assume $\mathcal{T} = \{1, 2, t_{\max}\}$. We want to show that exists a function of the form $\sum_i 10^{-k\psi(x_i, e_i)} \alpha^{-\beta t_i}$
 909 that is unique on any multiset M .

910 Our idea is to separate a range of k decimal slots for each unique element (x_i, e_i, \cdot) in the multiset.
 911 Each such a range has to accommodate at least t_{\max} decimal slots (one for each value of t_i). Finally,
 912 we need to make sure we can add up to N values at each decimal slot.

913 Formally, we map each tuple (x_i, e_i, \cdot) to one of k decimal slots starting from $10^{-k\psi(x_i, e_i)}$. In
 914 particular, for each element $(x_i, e_i, t_i = j)$ we add one unit at the j -th decimal slot after $10^{-k\psi(x_i, e_i)}$.
 915 Also, to ensure the counts for (x_i, e_i, j) and $(x_i, e_i, l \neq j)$ do not overlap, we set $\beta = \lceil \log_{10} N \rceil$
 916 since no tuple can repeat more than N times. We use $\alpha = 10$ as we shift decimals. Finally, to
 917 guarantee that each range encompasses t_{\max} slots of β decimals, we set $k = \beta(t_{\max} + 1)$. Therefore,
 918 the function

$$\sum_i 10^{-k\psi(x_i, e_i)} \alpha^{-\beta t_i}$$

919 is unique on any multiset M . We note that, without loss of generality, one could choose a different
 920 basis (other than 10). \square

921 B.12 Proof of Proposition 9: Expressiveness of PINT: link prediction

922 *Proof.* We now show that PINT (with relative positional features) is strictly more powerful than
 923 MP-TGN and CAW in distinguishing edges of temporal graphs. Leveraging Proposition 5, it suffices
 924 to show that PINT is at least as powerful as both CAW and MP-TGN.

925 **Statement 1:** PINT is at least as powerful as MP-TGNs.

926 Since PINT is a generalization of MP-TGNs with injective aggregation/update layers, it derives that
 927 it is at least as powerful as MP-TGNs. We can set the model's parameters associated with positional
 928 features to zero and obtain an equivalent MP-TGN.

929 **Statement 2:** PINT is at least as powerful as CAW.

930 We wish to show that for any pair of events that PINT cannot distinguish, CAW also cannot distinguish
 931 it. Let us consider the events (u, v, t) and (u', v', t) of a temporal graph $\mathcal{G}(t)$. Formally, we want to
 932 prove that if $\{\{T_u(t), T_v(t)\}\} = \{\{T_{u'}(t), T_{v'}(t)\}\}$ (i.e., the multisets contain TCTs that are pairwise
 933 isomorphic), then $\{\{\text{ENC}(W; S_u, S_v)\}\}_{W \in \{S_u \cup S_v\}} = \{\{\text{ENC}(W'; S_{u'}, S_{v'})\}\}_{W' \in \{S_{u'} \cup S_{v'}\}}$, where
 934 ENC denotes the walk-encoding function of CAW. Importantly, for the sake of this proof, we assume
 935 that all TCTs here are augmented with positional features, characterizing edge embeddings obtained
 936 from PINT.

937 Without loss of generality, we can assume that $T_u(t) \cong T_{u'}(t)$ and $T_v(t) \cong T_{v'}(t)$. By Lemma B1,
 938 we know that the corresponding monotone TCTs are also isomorphic: $\tilde{T}_u(t) \cong \tilde{T}_{u'}(t)$, and $\tilde{T}_v(t) \cong$
 939 $\tilde{T}_{v'}(t)$ with associated bijections $f_1 : V(\tilde{T}_u(t)) \rightarrow V(\tilde{T}_{u'}(t))$ and $f_2 : V(\tilde{T}_v(t)) \rightarrow V(\tilde{T}_{v'}(t))$.

940 We can construct a tree T_{uv} by attaching $\tilde{T}_u(t)$ and $\tilde{T}_v(t)$ to a (virtual) root node uv — without loss
 941 of generality, u and v are the left-hand and right-hand child of uv , respectively. We can follow the
 942 same procedure and create the tree $T_{u'v'}$ by attaching the TCTs $\tilde{T}_{u'}(t)$ and $\tilde{T}_{v'}(t)$ to a root node
 943 $u'v'$. Since the left-hand and right-hand subtrees of T_{uv} and $T_{u'v'}$ are isomorphic, then T_{uv} and
 944 $T_{u'v'}$ are also isomorphic. Let $f : V(T_{uv}) \rightarrow V(T_{u'v'})$ denote the bijection associated with the
 945 augmented trees. We also assume f is constructed by preserving the bijections f_1 and f_2 defined
 946 between the original monotone TCTs: this ensures that f does not map any node in $V(\tilde{T}_u(t))$ to a
 947 node in $V(\tilde{T}_{v'}(t))$, for instance. We have that

$$[r_{\#i \rightarrow u}^{(t)} \| r_{\#i \rightarrow v}^{(t)}] = [r_{\#f(i) \rightarrow u'}^{(t)} \| r_{\#f(i) \rightarrow v'}^{(t)}] \quad \forall i \in V(T_{uv}) \setminus \{uv\}$$

948 Note that we use the function $\#$ (that maps nodes in the TCT to nodes in the dynamic graph) here
 949 because the positional feature vectors are defined for nodes in the dynamic graph.

950 To guarantee that two encoded walks are identical $\text{ENC}(W; S_u, S_v) = \text{ENC}(W'; S_{u'}, S_{v'})$, it
 951 suffices to show that the anonymized walks are equal. Thus, we turn our problem into show-
 952 ing that for any walk $W = (w_0, t_0, w_1, t_1, \dots)$ in $S_u \cup S_v$, there exists a corresponding one

953 $W' = (w'_0, t_0, w'_1, t_1, \dots)$ in $S_{u'} \cup S_{v'}$ such that $I_{\text{CAW}}(w_i; S_u, S_v) = I_{\text{CAW}}(w'_i; S_{u'}, S_{v'})$ for all
 954 i . Recall that $I_{\text{CAW}}(w_i; S_u, S_v) = \{g(w_i; S_u), g(w_i; S_v)\}$, where $g(w_i; S_u)$ is a vector whose k -
 955 component stores how many times w_i appears in a walk from S_u at position k .

956 A key observation is that there is an equivalence between deanonymized root-leaf paths in T_{uv} and
 957 walks in $S_u \cup S_v$ (disregarding the virtual root node). By deanonymized, we mean paths where node
 958 identities (in the temporal graph) are revealed by applying the function \sharp . Using this equivalence, it
 959 suffices to show that

$$g(\sharp i; S_u) = g(\sharp f(i); S_{u'}) \text{ and } g(\sharp i; S_v) = g(\sharp f(i); S_{v'}) \quad \forall i \in V(T_{uv}) \setminus \{uv\}$$

960 Suppose there is an $i \in V(T_{uv}) \setminus \{uv\}$ such that $g(\sharp i; S_u) \neq g(\sharp f(i); S_{u'})$. Without loss of generality,
 961 suppose this holds for the ℓ -th entry of the vectors.

962 We know there are exactly $r_{a \rightarrow u}^{(t)}[\ell]$ nodes at the ℓ -th level of $\tilde{T}_u(t)$ that are associated with $a = \sharp i \in$
 963 $V(\mathcal{G}(t))$. We denote by Ψ the set comprising such nodes. It also follows that computing $g(\sharp i; S_u)[\ell]$
 964 is the same as summing up the amount of leaves of each subtree of $\tilde{T}_u(t)$ rooted at $\psi \in \Psi$, which we
 965 denote as $\mathcal{L}(\psi; \tilde{T}_u(t))$, i.e.,

$$g(\sharp i; S_u)[\ell] = \sum_{\psi \in \Psi} \mathcal{L}(\psi; \tilde{T}_u(t)).$$

966 Since we assume $g(\sharp i; S_u)[\ell] \neq g(\sharp f(i); S_{u'})[\ell]$, then it holds that

$$g(\sharp i; S_u)[\ell] \neq g(\sharp f(i); S_{u'})[\ell] \Rightarrow \sum_{\psi \in \Psi} \mathcal{L}(\psi; \tilde{T}_u(t)) \neq \sum_{\psi \in \Psi} \mathcal{L}(f(\psi); \tilde{T}_{u'}(t)) \quad (\text{S5})$$

967 Note that the subtree of \tilde{T}_u rooted at ψ should be isomorphic to the subtree of $\tilde{T}_{u'}$ rooted at $f(\psi)$,
 968 and therefore have the same number of leaves. However, the RHS of Equation S5 above implies there
 969 is a $\psi \in \Psi$ for which $\mathcal{L}(\psi; \tilde{T}_u) \neq \mathcal{L}(f(\psi); \tilde{T}_{u'})$, reaching a contradiction. The same argument can
 970 be applied to v and v' to prove that $g(\sharp i; S_v) = g(\sharp f(i); S_{v'})$.

971 □

972 C Additional related works

973 **Structural features for static GNNs.** Using structural features to enhance the power of GNNs is
 974 an active research topic. Bouritsas et al. [47] improved GNN expressivity by incorporating counts
 975 of local structures in the message-passing procedure, e.g, the number triangles a node appears on.
 976 These counts depend on identifying subgraph isomorphisms and, naturally, can become intractable
 977 depending on the chosen substructure. Li et al. [51] proposed increasing the power of GNNs using
 978 distance encodings, i.e., augmenting original node features with distance-based ones. In particular,
 979 they compute the distance between a node set whose representation is to be learned and each node
 980 in the graph. To alleviate the cost of distance encoding, an alternative is to learn absolute position
 981 encoding schemes [50, 59, 60], that try to summarize the role each node plays in the overall graph
 982 topology. We note that another class of methods uses random features to boost the power of GNNs
 983 [46, 56]. However, these models are referred to be hard to converge and to obtain noisy predictions
 984 [60].

985 The most trivial difference between these approaches and PINT is that our relative positional features
 986 account for temporal information. On a deeper level, our features summarize the role each node plays
 987 in each other’s monotone TCT instead of measuring, e.g., pair-wise distances in the original graph or
 988 counting substructures. Also, our scheme leverages the temporal aspect to achieve computational
 989 tractability, updating features incrementally as events unroll. Finally, while some works proposing
 990 structural features for static GNNs present marginal gains [60], PINT exhibits significant performance
 991 gains in real-world temporal link prediction tasks.

992 **Other models for temporal graphs.** Representation learning for dynamic graphs is a broad and
 993 diverse field. In fact, strategies to cope with the challenge of modeling dynamic graphs can come
 994 in many flavors, including simple aggregation schemes [18], walk-aggregating methods [52], and
 995 combinations of sequence models with GNNs [54, 57]. For instance, Seo et al. [57] used a spectral
 996 graph convolutional network [48] to encode graph snapshots followed by a graph-level LSTM [13].

997 Manessi et al. [53] followed a similar approach but employed a node-level LSTM, with parameters
998 shared across the nodes. Sankar et al. [55] proposed a fully attentive model based on graph attention
999 networks [34]. Pareja et al. [54] applied a recurrent neural net to dynamically update the parameters
1000 of a GCN. Gao and Ribeiro [49] compared the expressive power of two classes of models for discrete
1001 dynamic graphs: time-and-graph and time-then-graph. The former represents the standard approach
1002 of interleaving GNNs and sequence (e.g., RNN) models. In the latter class, the models first capture
1003 node and edge dynamics using RNNs, and are then feed into graph neural networks. The authors
1004 showed that time-then-graph has expressive advantage over time-and-graph approaches under mild
1005 assumptions. For an in-depth review of representation learning for dynamic graphs, we refer to the
1006 survey works by Kazemi et al. [14] and Skarding et al. [58].

1007 While most of the early works focused on discrete-time dynamic graphs, we have recently witnessed
1008 a rise in interest in models for event-based temporal graphs (i.e., CTDGs). The reason is that
1009 models for DTDGs may fail to leverage fine-grained temporal and structural information that can be
1010 crucial in many applications. In addition, it is hard to specify meaningful time intervals for different
1011 tasks. Thus, modern methods for temporal graphs explicitly incorporate timestamp information
1012 into sequence/graph models, achieving significant performance gains over approaches for DTDGs
1013 [37]. Appendix A provides a more detailed presentation of CAW, TGN-Att, and TGAT, which are
1014 among the best performing models for link prediction on temporal graphs. Besides these methods,
1015 JODIE [16] applies two RNNs (for the source and target nodes of an event) with a time-dependent
1016 embedding projection to learn node representations of item-user interaction networks. Trivedi et al.
1017 [32] employed RNNs with a temporally attentive module to update node representations. APAN
1018 [61] consists of a memory-based TGN that uses attention mechanism to update memory states using
1019 multi-hop temporal neighborhood information.

1020 D Datasets and implementation details

1021 D.1 Datasets

1022 In our empirical evaluation, we have considered six datasets for dynamic link prediction: Reddit¹,
1023 Wikipedia², UCI³, LastFM⁴, Enron⁵, and Twitter. Reddit is a network of posts made by users on
1024 subreddits, considering the 1,000 most active subreddits and the 10,000 most active users. Wikipedia
1025 comprises edits made on the 1,000 most edited Wikipedia pages by editors with at least 5 edits. Both
1026 Reddit and Wikipedia networks include links collected over one month, and text is used as edge
1027 features, providing informative context. The LastFM dataset is a network of interactions between
1028 user and the songs they listened to. UCI comprises students’ posts to a forum at the University
1029 of California Irvine. Enron contains a collection of email events between employees of the Enron
1030 Corporation, before its bankruptcy. The Twitter dataset is a non-bipartite net where nodes are users
1031 and interactions are retweets. Since Twitter is not publicly available, we build our own version by
1032 following the guidelines by Rossi et al. [27]. We use the data available from the 2021 Twitter RecSys
1033 Challenge and select 10,000 nodes and their associated interactions based on node participation:
1034 number of interactions the node participates in. We also apply multilingual BERT to obtain text
1035 representations of retweets (edge features).

1036 Table S1 reports statistics of the datasets such as number of temporal nodes and links, and the
1037 dimensionality of the edge features. We note that UCI, Enron, and LastFM represent non-attributed
1038 networks and therefore do not contain feature vectors associated with the events. Also, the node
1039 features for all datasets are vectors of zeros [41].

Table S1: Summary statistics of the datasets.

Dataset	#Nodes	#Events	#Edge feat.	Bipartite?
Reddit	10,984 (10,000 / 984)	672,447	172	Yes
Wikipedia	9,227 (8,227 / 1,000)	157,474	172	Yes
Twitter	8,925	406,564	768	No
UCI	1,899	59,835	-	No
Enron	184	125,235	-	No
LastFM	1,980 (980 / 1,000)	1,293,103	-	Yes

1040 D.2 Implementation details

1041 We train all models in link prediction tasks in a self-supervised approach. During training, we
1042 generate negative samples: for each actual event (z, w, t) (class 1), we create a fake one (z, w', t)
1043 (class 0) where w' is uniformly sampled from the set of nodes, and both events have the same edge
1044 feature vector.

1045 To ensure a fair comparison, we mainly rely on the original repositories and guidelines. For instance,
1046 regarding the training of MP-TGNs (including PINT), we mostly follow the setup and choices in the
1047 implementation available in [27]. In particular, we apply the Adam optimizer with learning rate 10^{-4}
1048 during 50 epochs with early stopping if there is no improvement in validation AP for 5 epochs. In
1049 addition, we use batch size 200 for all methods. We report statistics (mean and standard deviation) of
1050 the performance metric (AP) over ten runs.

1051 **MP-TGNs.** For TGN-Att, we follow Rossi et al. [27] and sample either ten or twenty temporal
1052 neighbors with memory dimensionality equal to 172, node embedding dimension equal to 100, and
1053 two attention heads. We use a memory unit implemented as a GRU, and update the state of each node
1054 based on only its most recent message. For TGAT, we use twenty temporal neighbors and two layers.

1055 **CAW.** We conduct model selection using grid search over: i) time decay $\alpha \in$
1056 $\{0.01, 0.1, 0.25, 0.5, 1.0, 2.0, 4.0, 10.0, 100.0\} \times 10^{-6}$, ii) number of walks $M \in \{1, 2, 3, 4, 5\}$;
1057 and iii) walk length $L \in \{32, 64, 128\}$. The best combination of hyperparameters is shown in

¹<http://snap.stanford.edu/jodie/reddit.csv>

²<http://snap.stanford.edu/jodie/wikipedia.csv>

³<http://konect.cc/networks/opsahl-ucforum/>

⁴<http://snap.stanford.edu/jodie/lastfm.csv>

⁵<https://www.cs.cmu.edu/~./enron/>

1058 Table S2. The remaining training choices follows the default values from the original implementation.
 1059 Importantly, we note that TGN-Att’s original evaluation setup is different from CAW’s. Thus, we
 1060 adapted CAW’s original repo to reflect these differences and ensure a valid comparison.

Table S2: Optimal hyperparameters for CAW.

Dataset	Time decay α	#Walks	Walk length
Reddit	10^{-8}	32	3
Wikipedia	4×10^{-6}	64	4
Twitter	10^{-6}	64	3
UCI	10^{-5}	64	2
Enron	10^{-6}	64	5
LastFM	10^{-6}	64	2

1061 **PINT.** We use $\alpha = 2$ (in the exponential aggregation function), and experiment with learned and
 1062 fixed β . We apply a relu function to avoid negative values of β , which could lead to unstable training.
 1063 We do grid search as the follow: when learning beta, we consider initial values for $\beta \in \{0.1, 0.5\}$;
 1064 for the fixed case (`requires_grad=False`), we evaluate $\beta \in \{0.001, 0.0001, 0.0001\}$ and always
 1065 apply memory as in the original implementation of TGN-Att. We consider number message passing
 1066 layers ℓ in $\{1, 2\}$. Also, we apply neighborhood sampling with the number of neighbors in $\{10, 20\}$,
 1067 and update the state of a node based on its most recent message. We then carry out model selection
 1068 based on AP values obtained during validation. Overall, the models with fixed β led to better results.
 1069 Table S3 reports the optimal hyperparameters for PINT found via automatic model selection.

1070 In all experiments, we use relative positional features with $d = 4$ dimensions. For computational
 1071 efficiency, we update the relative positional features only after processing a batch, factoring in all
 1072 events from that batch. Note that this prevents information linkage as these positional features take
 1073 effect after prediction. In addition, since temporal events repeat (in the same order) at each epoch, we
 1074 also speed up PINT’s training procedure by precomputing and saving the positional features for each
 1075 batch. To save up space, we store the positional features as sparse matrices.

Table S3: Optimal hyperparameters for PINT.

Dataset	β	#Neighbors	#Layers
Reddit	10^{-5}	10	2
Wikipedia	10^{-4}	10	2
Twitter	10^{-5}	20	2
UCI	10^{-5}	10	2
Enron	10^{-5}	20	2
LastFM	10^{-4}	10	2

1076 **Hardware.** For all experiments, we use Tesla V100 GPU cards and consider a memory budget of
 1077 32GB of RAM. We omit further details to preserve anonymity during the review process.

1078 E Deletion and node-level events

1079 Rossi et al. [27] propose handling edge deletions by simply updating memory states of the edge’s
 1080 endpoints, as if we were dealing with a usual edge addition. However, it is not discussed whether
 1081 the edge in question should be excluded from the event list or if we should just add a novel event
 1082 with edge features that characterize deletion. If we choose the former, we may be unable to recover
 1083 the memory state of a node from its monotone TCT and the original node features. Removing an
 1084 edge from the event list also affects the computation of node embeddings. Therefore, we advise
 1085 practitioners to do the latter when using PINT. It is worth mentioning that the vast majority of models
 1086 for temporal interaction prediction do not consider the possibility of deletion events.

1087 Regarding node-level events, PINT can accommodate node addition by simply creating novel memory
 1088 states. To deal with node feature updates, we can create an edge event with both endpoints on that
 1089 node, inducing a self-loop in the dynamic graph. Also, we can combine (e.g., concatenate) the

1090 temporal features in message-passage operations, similarly to the general formulation of the MP-TGN
 1091 framework [27]. Finally, we can deal with the removal of a node v by following our previous (edge
 1092 deletion) procedure to delete all edges with endpoints in v .

1093 F Additional experiments

1094 **Time comparison.** Figure S7 compares the time per epoch for PINT and for the prior art (CAW,
 1095 TGN-Att, and TGAT) in the Enron and LastFM datasets. Following the trend in Figure 6, Figure S7
 1096 further supports that PINT is generally slower than other MP-TGNs but, after a few training epochs,
 1097 is orders of magnitude faster than CAW. In the case of Enron, the time CAW takes to complete an
 1098 epoch is much higher than the time we need to preprocess PINT’s positional features.

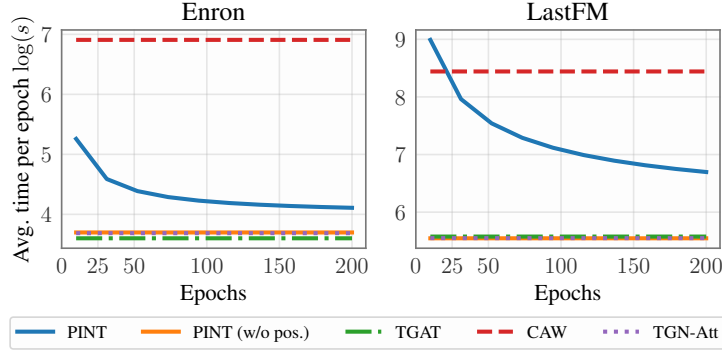


Figure S7: Time comparison: PINT versus TGNs (in log-scale) on Enron and LastFM.

1099 **Dimensionality of positional features.** We now assess the performance of PINT as a function of
 1100 the dimensionality d of the relative positional features. Figure S8 shows the performance of PINT
 1101 for $d \in \{4, 10, 15, 20\}$ on UCI and Enron. We report mean and standard deviation of AP on test set
 1102 obtained from five independent runs. In all experiments, we re-use the optimal hyper-parameters
 1103 found with $d = 4$ (the configuration shown in the main paper). Increasing the dimensionality of the
 1104 positional features leads to performance gains on both datasets. Notably, we obtain a significant boost
 1105 for Enron with $d = 10$: $92.696 \pm 0.095\%$ AP in the transductive setting and $88.342 \pm 0.294\%$ in the
 1106 inductive case. Thus, PINT becomes the best-performing model on Enron (transductive). On the UCI
 1107 dataset, higher d results in more significant gains in the inductive case. On UCI, with $d = 20$, we
 1108 obtain 96.36 ± 0.067 (transductive) and 94.77 ± 0.12 (inductive).

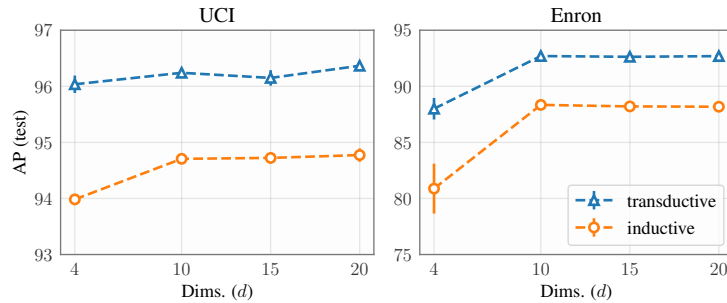


Figure S8: PINT: AP (mean and std) as a function of the dimensionality of the positional features.

1109 **Incorporating relative positional encoding into MP-TGNs.** We can also use PINT’s relative
 1110 positional encoding to boost the performance of MP-TGNs. As a proof of concept, we have imple-
 1111 mented TGN-Att with PE. We use the same model selection procedure as TGN-Att in Table 1, and
 1112 incorporate $d = 4$ -dimensional positional features.

1113 Table S4 shows the performance of TGN-Att+PE on three unattributed link prediction benchmarks
 1114 (UCI, Enron, and LastFM). Notably, TGN-Att receives a significant boost from our PE. However,
 1115 PINT still beats TGN-Att+PE on 5 out of 6 cases. The values for TGN-Att+PE reflect the outcome of
 1116 5 repetitions.

Table S4: Link prediction: TGN-Att + relative positional encoding (AP).

	Transductive			Inductive		
	UCI	Enron	LastFM	UCI	Enron	LastFM
TGN-Att	80.40 ± 1.4	79.91 ± 1.3	80.69 ± 0.2	74.70 ± 0.9	78.96 ± 0.5	84.66 ± 0.1
TGN-Att + PE	95.64 ± 0.06	85.04 ± 2.54	89.41 ± 0.91	92.82 ± 0.4	76.27 ± 3.37	91.63 ± 0.31
PINT	96.01 ± 0.1	88.71 ± 1.3	88.06 ± 0.7	93.97 ± 0.1	81.05 ± 2.4	91.76 ± 0.7

1117 **Experiments on node classification.** For completeness, we also evaluate PINT on node-level
1118 tasks (Wikipedia and Reddit). We follow closely the experimental setup in Rossi et al. [27] and
1119 compare against the baselines therein. Table S5 shows that PINT ranks first on Reddit and second on
1120 Wikipedia. The values for PINT reflect the outcome of 5 repetitions.

Table S5: Results for node classification (AUC).

	Wikipedia	Reddit
CTDNE	75.89 ± 0.5	59.43 ± 0.6
JODIE	84.84 ± 1.2	61.83 ± 2.7
TGAT	83.69 ± 0.7	65.56 ± 0.7
DyRep	84.59 ± 2.2	62.91 ± 2.4
TGN-Att	87.81 ± 0.3	67.06 ± 0.9
PINT	87.59 ± 0.6	67.31 ± 0.2

1121 G Societal and broader impact

1122 Temporal graph networks have shown remarkable performance in relevant domains such as social
1123 networks, e-commerce, and education. In this paper, we establish fundamental results that delineate
1124 the representational power of TGNs. We expect that our findings will help declutter the literature and
1125 serve as a seed for future developments. Moreover, our analysis culminates with PINT, a method
1126 that is provably more powerful than the prior art and shows superior predictive performance in a
1127 series of benchmarks. We believe that PINT (and its underlying concepts) will help engineers and
1128 researchers build better recommendation engines, improving the quality of systems that permeate
1129 our lives. Furthermore, we do not foresee any negative societal impact stemming directly from our
1130 developments.

1131 Supplementary References

- 1132 [46] R. Abboud, I. I. Ceylan, M. Grohe, and T. Lukasiewicz. The surprising power of graph neural networks
1133 with random node initialization. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2021.
- 1134 [47] G. Bouritsas, F. Frasca, S. Zafeiriou, and M. M. Bronstein. Improving graph neural network expressivity
1135 via subgraph isomorphism counting. In *Arxiv e-prints*, 2020.
- 1136 [48] M. Defferrard, X. Bresson, and P. Vandergheynst. Convolutional neural networks on graphs with fast
1137 localized spectral filtering. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.
- 1138 [49] J. Gao and B. Ribeiro. On the equivalence between temporal and static graph representations for observa-
1139 tional predictions. *ArXiv*, 2103.07016, 2021.
- 1140 [50] D. Kreuzer, D. Beaini, W. L. Hamilton, V. Letourneau, and P. Tossou. Rethinking graph transformers with
1141 spectral attention. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2021.
- 1142 [51] P. Li, Y. Wang, H. Wang, and J. Leskovec. Distance encoding: Design provably more powerful neural
1143 networks for graph representation learning. In *Advances in Neural Information Processing Systems*
1144 (*NeurIPS*), 2020.
- 1145 [52] S. Mahdavi, S. Khoshraftar, and A. An. dynnode2vec: Scalable dynamic network embedding. In
1146 *International Conference on Big Data*, 2018.
- 1147 [53] Franco Manessi, Alessandro Rozza, and Mario Manzo. Dynamic graph convolutional networks. *Pattern*
1148 *Recognition*, 97, 2020.

- 1149 [54] A. Pareja, G. Domeniconi, J. Chen, T. Ma, H. Kanezashi T. Suzumura, T. Kaler, T. B. Schardl, and C. E.
1150 Leiserson. EvolveGCN: Evolving graph convolutional networks for dynamic graphs. In *AAAI Conference*
1151 *on Artificial Intelligence (AAAI)*, 2020.
- 1152 [55] A. Sankar, Y. Wu, L. Gou, W. Zhang, and H. Yang. Dysat: Deep neural representation learning on dynamic
1153 graphs via self-attention networks. In *International Conference on Web Search and Data Mining (WSDM)*,
1154 2020.
- 1155 [56] R. Sato, M. Yamada, and H. Kashima. Random features strengthen graph neural networks. In *SIAM*
1156 *International Conference on Data Mining (SDM)*, 2021.
- 1157 [57] Y. Seo, M. Defferrard, P. Vandergheynst, and X. Bresson. Structured sequence modeling with graph
1158 convolutional recurrent networks. In *International Conference on Neural Information Processing (ICONIP)*,
1159 2018.
- 1160 [58] Joakim Skarding, Bogdan Gabrys, and Katarzyna Musial. Foundations and modeling of dynamic networks
1161 using dynamic graph neural networks: A survey. *IEEE Access*, 9:79143–79168, 2021.
- 1162 [59] B. Srinivasan and B. Ribeiro. On the equivalence between positional node embeddings and structural graph
1163 representations. In *International Conference on Learning Representations (ICLR)*, 2020.
- 1164 [60] H. Wang, H. Yin, M. Zhang, and P. Li. Equivariant and stable positional encoding for more powerful graph
1165 neural networks. In *International Conference on Learning Representations (ICLR)*, 2022.
- 1166 [61] X. Wang, D. Lyu, M. Li, Y. Xia, Q. Yang, X. Wang, X. Wang, P. Cui, Y. Yang, B. Sun, and Z. Guo. APAN:
1167 Asynchronous propagation attention network for real-time temporal graph embedding. *International*
1168 *Conference on Management of Data*, 2021.