# RETHINKING EVALUATION FOR TEMPORAL LINK PREDICTION THROUGH COUNTERFACTUAL ANALYSIS

Anonymous authors

Paper under double-blind review

# ABSTRACT

In response to critiques of existing evaluation methods for temporal link prediction (TLP) models, we propose a novel approach to verify if these models truly capture temporal patterns in the data. Our method involves a sanity check formulated as a counterfactual question: "What if a TLP model is tested on a temporally distorted version of the data instead of the real data?" Ideally, a TLP model that effectively learns temporal patterns should perform worse on temporally distorted data compared to real data. We provide an in-depth analysis of this hypothesis and introduce two data distortion techniques to assess well-known TLP models. Our contributions are threefold: (1) We introduce two simple techniques to distort temporal patterns within a graph, generating temporally distorted test splits of well-known datasets for sanity checks. These distortion methods are applicable to any temporal graph dataset. (2) We perform counterfactual analysis on six TLP models JODIE, TGAT, TGN, CAWN, GraphMixer, and DyGFormer to evaluate their capability in capturing temporal patterns across different datasets. (3) We introduce two metrics - average time difference (ATD) and average count difference (ACD) – to provide a comprehensive measure of a model's predictive performance.

025 026

004

006 007 008

009 010

011

012

013

014

015

016

017

018

019

021

# 1 INTRODUCTION

027 028

In static graphs, link prediction refers to the task of predicting whether an edge exists between two
nodes after having observed other edges in the graph. Temporal link prediction (TLP) is a dynamic
extension of link prediction wherein the task is to predict whether a link (edge) exists between any
two nodes in the future based on the historical observations (Qin and Yeung, 2023). The predictive
capability of TLP models make them useful in applications pertaining to dynamic graphs, such as
product recommendations (Qin et al., 2024; Fan et al., 2021), social network content or account
recommendation (Fan et al., 2019; Daud et al., 2020), fraud detection in financial networks (Kim
et al., 2024), and resource allocation, to name a few.

In the TLP literature (Kumar et al., 2019; Trivedi et al., 2019; Xu et al., 2020; Rossi et al., 2020; Wang et al., 2020; Cong et al., 2023; Yu et al., 2023), the TLP task is treated as a binary classification problem where the query

040 041

 $q_1$ : "Does an edge exist between the nodes u and v at time t?"

042 is processed by a model and then compared with the ground truth following which metrics such as 043 area under the receiver operating characteristic curve (AU-ROC), and average precision (AP) are 044 reported. The ground truth consists of positive samples, and a fixed number of random negative samples. There are a couple of issues in the binary classification approach. Firstly, the timestamps in the query are restricted to the timestamps present in the ground truth, which makes the evaluation 046 biased and does not test the model's performance in the continuous time range. Secondly, checking 047 for the existence of an edge at a specific timestamp is an ill-posed question, and instead the existence 048 of an edge should be queried within a finite time-interval. Lastly, the negative edge sampling strategy, 049 and the number of negative samples per positive sample impact the performance metrics as seen in 050 EXH (Poursafaei and Rabbany, 2023). 051

- 052 Alternatively, in a rank-based approach, the query is formulated as:
- 053

 $q_2$ : "Which nodes are likely to have an edge with node u at time t?"

054 In this case, the model returns an ordered list of nodes arranged from most likely to least likely. 055 Then, the rank of the ground truth edge is returned if a match is found, and if not, a high number is 056 reported. For all the edges in the test data, metrics such as Mean Average Rank (MAR) or Mean 057 Reciprocal Rank (MRR) can be reported to assess the performance of the model (Huang et al., 2024). 058 While the rank-based metrics are more intuitive than AU-ROC and AP, the issues regarding binary classification mentioned above still remain unaddressed. To give a true picture of the predictive power of the TLP models, a penalty term should be introduced to account for the nodes that are incorrectly 060 estimated to form an edge with node u at time t. 061

062 In a recent work, Poursafaei et al. (2022) highlighted that the state-of-the-art (SoTA) performance 063 of some TLP models on the standard benchmark datasets is near-perfect. This is counterintuitive 064 because TLP is a challenging task, even more challenging than link prediction of static graphs, due to the additional degree of freedom in the data induced by the temporal dimension. The flaw in the 065 evaluation method is attributed to the limited negative sampling strategy, and the authors propose a 066 new negative edge sampling strategy which results in a different ranking of the baselines. 067

068 Inspired by the critique of the evaluation method, we propose a method to conduct sanity check of 069 the TLP models to determine if they truly capture the temporal patterns in the data. The sanity check is formulated as the counterfactual question (Pearl, 2009):

- 071
- 072

078

079

081

082

084

085

090

092

093 094

095

107

073

"What if a TLP model which is trained on a temporal graph is tested on temporally distorted version of the data instead of the real data?"

074 Ideally, a TLP model which is capable of learning the temporal patterns should perform worse on 075 temporally distorted data compared to the real data. We conduct an in-depth analysis of this argument 076 and introduce various data distortion techniques to assess well-known TLP models. 077

**Contributions** The contributions of our work can be summarised as follows:

- We introduce simple **techniques** to distort the temporal patterns within a graph. These techniques are then used to generate temporally distorted version of the test split of some famous datasets which can be used for **sanity check**. Moreover, the distortion methods can be applied to any temporal graph dataset [Link to code repository].
- We perform **counterfactual analysis** on TLP models such as JODIE (Kumar et al., 2019), TGAT (Xu et al., 2020), TGN (Rossi et al., 2020), CAWN (Wang et al., 2020), GraphMixer (Cong et al., 2023), and DyGFormer (Yu et al., 2023) to check whether they are capable of capturing the temporal patters within various datasets.
  - We propose two **metrics**: average time difference (ATD), and average count difference (ACD) to measure the performance of TLP models. These metrics can provide a holistic picture of a model's predictive performance.
    - Lastly, we propose an alternative evaluation strategy for TLP through which the existing pitfalls of binary classification and ranking methods can be avoided.

**PRELIMINARIES** 2

098 In TLP literature, continuous-time temporal graphs with *instantaneous edges* are often considered, 099 where edges represent interaction events between two nodes at a specific point in time. Alternatively, 100 temporal graphs can be defined with edges that appear at a certain time and either persist for a duration (Celikkanat et al., 2024) or accumulate indefinitely. In this work, we focus on the instantaneous 101 edge temporal graph, also known as interaction graphs (Qin et al., 2024) or unevenly sampled edge 102 sequence (Qin and Yeung, 2023). 103

**Definition 2.1.** A temporal graph with  $m \in \mathbb{N}$  instantaneous edges formed between nodes in  $\mathcal{U}$  and 104 105  $\mathcal{V}$  is defined as  $\mathcal{G} = (\mathcal{U}, \mathcal{V}, \mathcal{E})$ , where  $\mathcal{E} \triangleq \{(u_i, v_i, t_i) : i \in [m], u_i \in \mathcal{U}, v_i \in \mathcal{V}, t_i \in \mathbb{R}\}$  denotes the set of edges. The tuple (u, v, t) is referred to as an edge event. 106

While the definition caters to bipartite structure, with  $\mathcal{U} = \mathcal{V}$ , it can also represent general graphs.

<sup>096</sup> 2.1 **DEFINITIONS** 097

108 **Definition 2.2.** The occurrences of a particular edge (u, v) in  $\mathcal{E}$  is denoted as  $\mathcal{E}_{(u,v)}$  and defined as 109  $\mathcal{E}_{(u,v)} \triangleq \{(u,v,t) : (u,v,t) \in \mathcal{E}\}.$ 110

**Definition 2.3.** The slice of edges in  $\mathcal{E}$  with timestamps in the range  $(t_1, t_2)$  is denoted as  $\mathcal{E}(t_1, t_2)$ , 111 and defined as  $\mathcal{E}(t_1, t_2) \triangleq \{(u, v, t) : (u, v, t) \in \mathcal{E}, t \in (t_1, t_2)\}.$ 112

**Definition 2.4.** The timestamps in  $\mathcal{E}$  consisting of  $m \in \mathbb{N}$  edges can be extracted through a function 113  $\mathscr{T}: (\mathcal{U} \times \mathcal{V} \times \mathbb{R})^m \to \mathbb{R}^m \text{ as } \mathscr{T}(\mathcal{E}) \triangleq \{t: (u, v, t) \in \mathcal{E}\}.$ 114

#### 116 2.2 POINT PROCESSES

115

117

123 124 125

129 130 131

132

133

138 139

140

141 142

143

144 145 146

147

151

Perry and Wolfe (2013) modelled the interaction events of a directed edge (u, v) as an inhomogeneous 118 Poisson point process. In a recent work on continuous-time representation learning on temporal 119 graphs, Modell et al. (2024) followed suit, and assumed  $\mathcal{E}_{(u,v)}$  to be sampled from an independent 120 inhomogeneous Poisson point process with intensity  $\lambda_{(u,v)}(t)$ . The number of edge events (u, v)121 between timestamps  $t_1$  and  $t_2$  follow a Poisson distribution with rate  $\int_{t_1}^{t_2} \lambda_{(u,v)}(t) dt$ , i.e., 122

$$\mathcal{E}_{(u,v)}(t_1, t_2) | \sim \operatorname{Poisson}\left(\int_{t_1}^{t_2} \lambda_{(u,v)}(t) \, dt\right). \tag{1}$$

To connect the present to the past, Du et al. (2016) view the intensity function  $\lambda_{(u,v)}^{\star}(t)$  as a nonlinear 126 function of the sample history, where  $\star$  indicates that the function is conditioned on the history. The 127 conditional density function for edge (u, v) is written as 128

$$p_{(u,v)}^{\star}(t) = \lambda_{(u,v)}^{\star}(t) \exp\left(-\int_{t'}^{t} \lambda_{(u,v)}^{\star}(\tau) \, d\tau\right),\tag{2}$$

where t' < t is the last time when edge (u, v) was observed. The goal is to find the parameters  $\lambda^{\star}_{(u,v)}(t) : 0 < t \leq T$  which can describe the observation  $\mathcal{E}_{(u,v)}$ . This is done by minimizing the negative log likelihood (NLL) at the timestamps of edge occurrence (Shchur et al., 2021):

$$\min_{\lambda_{(u,v)}^{\star}(t): 0 < t \le T} - \sum_{t \in \mathscr{T}\left(\mathcal{E}_{(u,v)}\right)} \log\left(\lambda_{(u,v)}^{\star}(t)\right) + \int_{0}^{T} \lambda_{(u,v)}^{\star}(\tau) \, d\tau, \quad T = \max \mathscr{T}\left(\mathcal{E}_{(u,v)}\right). \tag{3}$$

Shchur et al. (2021) summarize the operation of a neural temporal point process as follows:

- The edge events in  $\{(u, v, t_i) : i \in [m]\}$  are represented as feature vectors  $\boldsymbol{x}_i = f_{\boldsymbol{c}}(u, v, t_i)$ ,
- The historical feature vectors are encoded into a state vector  $h_i = f_h(x_1, \cdots x_{i-1})$ ,
- The distribution of  $t_i$  conditioned on the past is simply conditioned on  $h_i$ .

The functions  $f_{\mathfrak{e}}$  and  $f_{\mathfrak{h}}$ , as well as the conditioning on  $h_i$  are implemented using neural networks.

#### 3 **COUNTERFACTUAL ANALYSIS**

A temporal graph is characterized by (1) the *order* in which the edges appear, (2) the *frequency* 148 with which edges appear over time, and (3) the time gap between any two edge events. In this 149 work, we refer to these characteristics as **temporal patterns**. Furthermore, if temporal patterns 150 observed in the past enable predictions of future temporal patterns that outperform naïve estimates on a specific performance metric, then the temporal data is considered **learnable**. This does not require 152 the temporal pattern to remain consistent over time; rather, it suggests that future changes can be 153 estimated from past observations. 154

**Experiment Setup** A model f is trained on a temporal graph  $\mathcal{E}_{train}$  and tested on  $\mathcal{E}_{test}$  through the 156 binary classification approach resulting in a performance metric such as AP. The train and test data 157 are chronologically split from the same temporal graph which is assumed to be generated through a 158 common causal mechanism, i.e.,  $\mathcal{E}_{\text{train}} = \mathcal{E}(0, \tau_0)$ , and  $\mathcal{E}_{\text{test}} = \mathcal{E}(\tau_0, T)$ . In light of the experimental 159 setup, we ask the following question: 160

- Would the model f which is trained on  $\mathcal{E}_{\text{train}}$  perform well if tested on a distorted 161 version of  $\mathcal{E}_{test}$  instead of  $\mathcal{E}_{test}$ ?
  - 3

162 To formalise the question in the counterfactual framework (Pearl, 2019), we consider the following 163 statements: 164

- x': The model f is tested on  $\mathcal{E}_{\text{test}}$
- y': The performance metric is  $\alpha$ 166

167

168 169

170

171

172

173

179

181

183

185

186 187

188

191

197

201

203

207

- x: The model f is tested on a *temporally distorted* version of  $\mathcal{E}_{\text{test}}$
- y: The performance metric is less than  $\alpha$

Additionally,  $y_x$  is read as y when x. The counterfactual question is framed as  $P(y_x \mid x', y')$ , i.e.,

The probability that the performance metric would be less than  $\alpha$  had the test data been a temporally distorted version of  $\mathcal{E}_{test}$ , given the performance metric was observed to be at least  $\alpha$  when the model was tested on  $\mathcal{E}_{test}$ .

174 To answer the question above, we design the *intervention* as graphically depicted in Fig. 1. The TLP 175 model f is trained on the data  $\mathcal{E}_{train}$ . The true test data  $\mathcal{E}_{test}$  is temporally distorted through some 176 function  $\mathfrak{D}(\cdot)$  resulting in  $\mathcal{E}' = \mathfrak{D}(\mathcal{E}_{test})$ . Finally, we test the model f on the true data  $\mathcal{E}_{test}$  and the 177 temporally distorted data  $\mathcal{E}'$  and compare the metrics which may result in either of the two scenarios 178 shown in the figure based on which we can comment on the effectiveness of f.



Figure 1: The intervention setup to verify the counterfactual question above.

189 **Motivation** To motivate the counterfactual analysis, we present a simplified example of *classifica*-190 tion on binary sequences, viewed as a discretized version of the temporal graph described in Def. 2.1. Let the set of all binary sequences of length  $m \in \mathbb{N}$  be denoted by  $\mathbb{B}_m = \{0, 1\}^m$ , and let  $\boldsymbol{b} \in \mathbb{B}_m$ 192 be a binary sequence representing the true data. Moreover, we consider a model f whose output 193 is  $\hat{b} \in \mathbb{B}_m$ . The performance metric achieved by  $\hat{b}$  on ground truth b is denoted as  $\phi(\hat{b}, b)$ . Next, 194 let  $b' \in \mathbb{B}_m \setminus \{b\}$  denote a distorted version of b. Building on the above setup, the counterfactual 195 question is  $P(\phi(\mathbf{b}, \mathbf{b}') < \phi(\mathbf{b}, \mathbf{b}))$ , i.e., the probability that the model performs relatively worse on 196 the distorted sequence. Next, we find the **conditions** on model output  $\hat{b}$  and distorted sequence b' such that  $P(\phi(\hat{\boldsymbol{b}}, \boldsymbol{b}') < \phi(\hat{\boldsymbol{b}}, \boldsymbol{b})) = 1$ .

Figure 2: Scatter plot showing the normalised Hamming dis-200 tance between **b** and **b'** on the x-axis and the performance (AP) of the classifier on the distorted sequence b' on the y-axis. The normalised Hamming distance serves as a distortion metric 202 for binary sequences. Each point corresponds to a random  $\boldsymbol{b} \in \mathbb{B}_m$ , and  $\boldsymbol{b}$  such that  $\phi(\boldsymbol{b}, \boldsymbol{b}) \geq \alpha$ . In the figure, m = 16, 204 and  $\alpha = 0.9$ . We observe that for  $\|\boldsymbol{b} - \boldsymbol{b}'\|_1 > \beta$ , all the 205 206 points lie below  $\alpha$ , i.e.,  $P(\phi(\mathbf{b}, \mathbf{b}') < \alpha) = 1$ .



208 Let  $\mathfrak{M}$  be a **causal model** which generates **b** succeeding another binary sequence  $b_0$ . The causal model 209 can produce multiple sequences succeeding  $b_0$ , i.e.,  $b \sim \mathfrak{M}(b_0)$ . We assume that  $\mathfrak{M}(b_0) \in \Omega \subset \mathbb{B}_m$ 210 with  $|\Omega| \ll 2^m$ . If the model f produces output  $b \in \Omega$  after being trained on  $b_0$ , we can say that it 211 has learnt the causal mechanism underlying  $\mathfrak{M}$ , which gives us the lower bound on the performance 212 metric as  $\alpha_0 = \inf_{\hat{b}, b \in \Omega} \phi(b, b)$ . Then, based on Fig. 2, we make the following Assumption: 213

**Assumption 3.1.** For some  $\alpha \leq \alpha_0$ , there exists  $\beta \in (0,1)$  such that a model which admits a 214 performance metric of  $\alpha$  on the true sequence **b**, reports a performance metric lower than  $\alpha$  for all 215 distorted samples b' with distortion greater than  $\beta$ .

228

229 230

231

232

233

234

235

246 247

248 249

250

251

257

258

259

260 261 262

216

217

218

Going back to the experiment setup in Fig. 1, we conjecture the following: **Conjecture 3.1.**  $P(y_x \mid x', y') \neq 1 \implies$  model f is not capable of discerning the temporal patterns distorted through  $\mathfrak{D}$ . We now present the logic behind this conjecture. Please consider the following statements,  $s_1$ : The model f is capable of discerning temporal patterns in  $(\mathcal{E}_{\text{train}}, \mathcal{E}_{\text{test}})$  $s_2$ : The function  $\mathfrak{D}$  generates temporally distorted test data  $\mathcal{E}' = \mathfrak{D}(\mathcal{E}_{test})$  $s_3$ : The data ( $\mathcal{E}_{train}, \mathcal{E}_{test}$ ) is learnable  $s_4$ : The performance metric reported by the model f on true test data  $\mathcal{E}_{test}$  is always higher than that reported on the distorted test data  $\mathcal{E}'$ , i.e.,  $P(y_x \mid x', y') = 1$ . We start with  $s_1 \wedge s_2 \wedge s_3 \implies s_4$ . Assuming that the data is learnable, i.e.,  $s_3 = 1$ , we get  $s_1 \wedge s_2 \implies s_4$ . Through contraposition, we arrive at  $\neg s_4 \implies \neg s_1 \vee \neg s_2$ , where  $\neg s_4 \equiv P(y_x \mid x', y') \neq 1.$ Further, we impose that  $\mathfrak{D}$  satisfies Assumption 3.1, i.e.,  $\neg s_2 = 0$ , allowing us to conclude  $\neg s_4 \implies$  $\neg s_1$  which reads as the conjecture above. In Fig. 3 we present an example depicting what temporal distortion means using point processes. In the remainder of this section, we define *temporal distortion metrics* and then discuss two *temporal* 

236 Figure 3: Let  $\mathcal{E}_{train} \cup \mathcal{E}_{test}$  be sampled from a 237 point process with intensity  $\lambda^{\star}(t), t \in [0, T]$ . 238 We generate  $\mathcal{E}'$  from another point process with 239 intensity  $\lambda'(t), t \in [\tau_0, T]$ . We depict the in-240 tensity functions as two sinusoidal waves with 241 different frequency and phase. If a model f242 learns this intensity function by observing  $\mathcal{E}_{train}$ , and then generates samples for prediction, they 243 would be more similar to  $\mathcal{E}_{test}$  than  $\mathcal{E}'$ , resulting 244 in a lower performance on the distorted test set. 245

distortion techniques to distort the temporal graphs.



3.1 TEMPORAL DISTORTION METRICS

Let  $\mathcal{E}$  be a temporal graph sampled from a temporal point process with intensity  $\lambda^{\star}(t)$  for  $t \in [0, T]$ . Let  $\mathcal{E}'$  be data sampled from another point process with intensity  $\lambda'(t)$  for  $t \in [0, T]$ .

**Definition 3.1.** The temporal graph  $\mathcal{E}'$  is  $\delta$ -temporally distorted w.r.t.  $\mathcal{E}$  if for some  $\delta > 0$ ,

$$\sum_{(u,v)\in\mathcal{U}\times\mathcal{V}}\frac{1}{T}\int_0^T |\lambda^*_{(u,v)}(t) - \lambda'_{(u,v)}(t)|\,dt > \delta.$$
(4)

In practice, we do not have access to the true intensity functions, so we have to compare the realisations instead. Let  $\mathcal{E}$  and  $\mathcal{E}'$  be two temporal graphs, then we measure the difference in their temporal patterns through the metrics defined below.

**Definition 3.2.** The average time difference (ATD) between  $\mathcal{E}$  and  $\mathcal{E}'$  is defined as:

$$\mathsf{ATD}(\mathcal{E}, \mathcal{E}') \triangleq \frac{1}{T|\mathcal{E}|} \sum_{(u, v, t) \in \mathcal{E}} \min_{t' \in \mathscr{T}\left(\mathcal{E}'_{(u, v)}\right) \cup \{T\}} |t - t'|, \quad T = \max \mathscr{T}\left(\mathcal{E}\right) - \min \mathscr{T}\left(\mathcal{E}\right).$$
(5)

In ATD, we measure the time difference between the edge event  $(u, v, t) \in \mathcal{E}$  and the closest  $(u, v, t') \in \mathcal{E}'$ , reporting the average over all the edge events in  $\mathcal{E}$ . In Fig. 4 we show two temporal graphs as *impulse trains* with each impulse color coded to represent the edge of the sample 3-node graph. Through ATD we can measure the overall difference in the occurrence of an edge event. However, ATD fails to capture the difference in the frequency with which an edge occurs in the two temporal graphs  $\mathcal{E}$  and  $\mathcal{E}'$ . Therefore, we define average count difference ACD to measure the difference in the frequency with which edges occur in the temporal graph.

 **Definition 3.3.** The average count difference (ACD) between  $\mathcal{E}$  and  $\mathcal{E}'$  is defined as:

$$\mathsf{ACD}(\mathcal{E}, \mathcal{E}') \triangleq \frac{1}{|\mathcal{E}|} \sum_{(u,v,t) \in \mathcal{E}} \left| |\mathcal{E}_{(u,v)}(t - \bar{\tau}, t + \bar{\tau})| - |\mathcal{E}'_{(u,v)}(t - \bar{\tau}, t + \bar{\tau})| \right|, \quad \bar{\tau} \in \mathbb{R}^+.$$
(6)

For each edge event  $(u, v, t) \in \mathcal{E}$ , we count the number of occurrences of (u, v) in the time range  $(t - \overline{\tau}, t + \overline{\tau})$  in both  $\mathcal{E}$  and  $\mathcal{E}'$  and measure the count difference. In Fig. 4 we depict the time interval as a light blue box centred around each edge event in  $\mathcal{E}$ . For  $\overline{\tau} \to 0$ , the search becomes restricted to an infinitesimal time interval, with  $\mathsf{ACD}(\mathcal{E}, \mathcal{E}')$  approaching  $1 - \frac{1}{|\mathcal{E}|} \sum_{(u,v,t) \in \mathcal{E}} \mathbb{I}\{(u, v, t) \in \mathcal{E}'\}$ .



Figure 4: Comparing two temporal graphs by measuring the time difference, and the count difference within intervals of duration  $2\bar{\tau}$  centred around the edge event.

Table 1: Joint interpretation of the distortion metrics.

	ATD↓	ATD↑
ACD↓	similar	the edge events are shifted near the extremities of the $2\bar{\tau}$ interval.
ACD ↑	the edge event is duplicated multiple times in the vicinity of the original edge event.	the edge events are either duplicated or reduced, and also shifted away from the original edge interval.

### 3.2 TEMPORAL DISTORTION TECHNIQUES

Now that we are equipped with metrics to measure the difference between two temporal graphs, we device distortion functions  $\mathfrak{D}(\cdot)$  which can enable us to investigate the counterfactual question posed earlier. We propose two distortion techniques  $\mathfrak{D}_{INTENSE}(\cdot, K)$  which creates K time-perturbed copies of each edge events, and  $\mathfrak{D}_{SHUFFLE}(\cdot)$  wherein the timestamps of different edge events are shuffled.

**INTENSE** Let the real temporal graph data be denoted by  $\mathcal{E} = \bigcup_{(u,v) \in \mathcal{U} \times \mathcal{V}} \mathcal{E}_{(u,v)}$ , and the distorted version be denoted by  $\mathcal{E}' = \bigcup_{(u,v) \in \mathcal{U} \times \mathcal{V}} \mathcal{E}'_{(u,v)}$ . Then, for each edge event (u, v, t) in the real data  $\mathcal{E}$ , we create K edge events  $(u, v, t + \tau)$  with  $\tau$  sampled uniformly from  $(-\bar{\tau}, \bar{\tau})$  for some  $\bar{\tau} \in \mathbb{R}^{+1}$ . Alternatively, if it is known that  $\mathcal{E}_{(u,v)}$  is sampled from a point process with intensity  $\lambda'_{(u,v)}(t)$ , then we can generate  $\mathcal{E}'_{(u,v)}$  by sampling from another point process with intensity  $\lambda'_{(u,v)}(t)$ , such that

$$\lambda'_{(u,v)}(t) = K\lambda^{\star}_{(u,v)}(t), \quad \forall (u,v) \in \mathcal{U} \times \mathcal{V}.$$
<sup>(7)</sup>

SHUFFLE For any two edge events  $(u, v, t), (u', v', t') \in \mathcal{E}$ , we shuffle the timestamps in the distorted version, i.e.  $(u, v, t'), (u', v', t) \in \mathcal{E}'$ . The shuffling process is also called label permutation (Chatterjee, 2018). In terms of the point process, we can explain shuffling as follows. If  $\mathcal{E}_{(u,v)}$  is known to be sampled from a point process with intensity  $\lambda_{(u,v)}^*(t)$ , then  $\mathcal{E}'_{(u,v)}$  can be generated by sampling from an inhomogeneous Poisson point process with intensity  $\lambda'_{(u,v)}(t)$ , where

 $\left(\int_{-\infty}^{T} \chi_{\star}(t) dt\right) \sum_{\lambda} \chi_{\star}(t)$ 

322  
323 
$$\lambda'_{(u,v)}(t) = \frac{\left(\int_{0}^{t} \lambda_{(u,v)}(t) \, dt\right) \sum_{(u',v') \in \mathcal{U} \times \mathcal{V}} \lambda_{(u',v')}(t)}{\sum_{(u',v') \in \mathcal{U} \times \mathcal{V}} \int_{0}^{T} \lambda_{(u',v')}^{\star}(t) \, dt}.$$
(8)

324	Algorithm 1 $\mathfrak{D}_{INTENSE}$	Algorithm 2 $\mathfrak{D}_{SHUFFLE}$
325	<b>Input</b> $\mathcal{E}$ $K \subset \mathbb{N}$ $\overline{\tau} \subset \mathbb{R}^+$	Innut <i>E</i>
326	Output $\mathcal{E}'$	Output $\mathcal{E}'$
327	1: $\mathcal{E}' = \emptyset$	1: $\mathcal{E}' = \varnothing$
328	2: for $(u, v, t) \in \mathcal{E}$ do	2: $\mathcal{T} \leftarrow \mathscr{T}(\mathcal{E})$
329	3: for $k \in [K]$ do	3: for $(u, v, t) \in \mathcal{E}$ do
330	4: $\tau \sim Uniform(-\bar{\tau}, \bar{\tau})$	4: $ au \sim \mathcal{T}$
331	5: $\mathcal{E}' \leftarrow \mathcal{E}' \cup \{(u, v, t + \tau)\}$	5: $\mathcal{E}' \leftarrow \mathcal{E}' \cup \{(u, v, \tau)\}$
332	6: end for	6: $\mathcal{T} \leftarrow \mathcal{T} \setminus \{\tau\}$
333	7: end for	7: end for
004		

The operations of  $\mathfrak{D}_{INTENSE}$  and  $\mathfrak{D}_{SHUFFLE}$  are described in Algorithms 1 and 2, respectively. Moreover, a visual example is provided in Fig. 5. The computational complexity of  $\mathfrak{D}_{INTENSE}(\cdot, K)$  is  $\mathcal{O}(K|\mathcal{E}|)$ , and of  $\mathfrak{D}_{SHUFFLE}(\cdot)$  is  $\mathcal{O}(|\mathcal{E}|)$ . In short, both distortion techniques are linear in the number of edge events in  $\mathcal{E}$ .

Figure 5: Visual representation of INTENSE and SHUFFLE distortions. In  $\mathfrak{D}_{INTENSE}(\mathcal{E}, 5)$ , 5 edge events are created in the vicinity of the true edge event in  $\mathcal{E}$ . This increases the frequency with which edges appear in an interval, thereby distorting the temporal pattern. In  $\mathfrak{D}_{SHUFFLE}(\mathcal{E})$ , as the name suggests, the order in which the edges appear is shuffled and thus the temporal pattern is distorted, the edges now appear where they should not be.



We use  $\mathfrak{D}_{INTENSE}(\cdot, 5)$  and  $\mathfrak{D}_{SHUFFLE}(\cdot)$  to create 10 temporally distorted samples of the test splits of each dataset. In Table 2, we present the ATD, and ACD by comparing the distorted samples with the original test data of different datasets. The metrics ATD and ACD should be considered in conjunction to measure the dissimilarity of two temporal graphs. For each real test data, we create 10 distorted samples and report the mean and 95% confidence interval of the metrics to ensure statistical reliability.

Table 2: Distortion measures on different datasets.

		wikipedia	reddit	uci	lastfm	mooc
INTENSE	ATD ACD	$\begin{array}{c} 6.9e\text{-}6 \pm 2e\text{-}8 \\ 4.479 \pm 1.9e\text{-}3 \end{array}$	$\begin{array}{c} 1.6\text{e-}6\pm2\text{e-}9\\ 4.112\pm3.9\text{e-}4\end{array}$	$\begin{array}{c} 1.6\text{e-5} \pm \text{1.2e-7} \\ 7.214 \pm \text{1.2e-2} \end{array}$	$\begin{array}{c} 8.6e\text{-}7 \pm 9.4\text{e}\text{-}10 \\ 4.046 \pm 1.8\text{e}\text{-}4 \end{array}$	$\begin{array}{c} 2.5e\text{-}6 \pm \text{7e-9} \\ 4.627 \pm 1.4e\text{-}3 \end{array}$
SHUFFLE	ATD ACD	$\begin{array}{c} 0.078 \pm \text{5.7e-4} \\ 1.093 \pm \text{3.4e-4} \end{array}$	$\begin{array}{c} 0.099 \pm {\scriptstyle 3e-4} \\ 1.033 \pm {\scriptstyle 8e-5} \end{array}$	$\begin{array}{c} 0.132 \pm \text{8.4e-4} \\ 1.877 \pm \text{3.3e-3} \end{array}$	$\begin{array}{c} 0.0800 \pm {\scriptstyle 1.7e\text{-}4} \\ 1.0011 \pm {\scriptstyle 1.4e\text{-}4} \end{array}$	$\begin{array}{c} 0.1906 \pm \text{6.7e-4} \\ 1.1896 \pm \text{8.9e-5} \end{array}$

# 4 Results

We evaluate the performance of the following TLP models in light of Proposition 3.1: JODIE (Kumar et al., 2019), TGAT (Xu et al., 2020), TGN (Rossi et al., 2020), CAWN (Wang et al., 2020), GraphMixer (Cong et al., 2023), DyGFormer (Yu et al., 2023)

The models are evaluated under two **settings**: *transductive* and *inductive*. In transductive TLP, the nodes u, v in the positive sample  $(u, v, t) \in \mathcal{E}_{test}$  were observed during training. In contrast, in inductive TLP, at least one node in u, v is novel, and was not observed during training.

In Table 3, we have arranged the datasets in increasing order of their size (more details can be found
 in Appendix A.1). We notice that all the TLP models pass the counterfactual test for SHUFFLE
 distortion on the smallest dataset: uci, and some of them {TGAT, GraphMixer, DyGFormer}

377

1

335

336

337

338

339 340

341

342

343

345

347

348 349

350

351

352

353

354

355 356

357

364 365

366 367

368

369

For the experiments we set $\bar{\tau}$ –	_	$\max \mathscr{T}(\mathcal{E}) {-} \min \mathscr{T}(\mathcal{E})$	
For the experiments we set $7 =$	-		•

pass for SHUFFLE on the second-smallest dataset wikipedia, and only GraphMixer and TGN pass on reddit. Surprisingly, JODIE passes on INTENSE distortion for two of the largest datasets lasstfm and mooc. And overall, none of the TLP models pass the counterfactual test on the INTENSE distortions. This allows us to conclude the followg: (1) The TLP models are able to discern the temporarl order of edge occurrence, however this capability worsens for larger datasets, and (2) the TLP models do not keep count of the frequency with which edges appear over time.

Table 3: Performance (AP) of the models JODIE, TGAT, TGN, CAWN, GraphMixer, and DyGFormer on five datasets, and their temporally distorted versions denoted as INTENSE, and SHUFFLE. For each metric, we report the mean, and the 95% confidence interval (CI) as mean ± CI. We have marked the metrics in blue for distortions that showed that a model was incapable of learning on a certain dataset as per Conjecture 3.1, and orange otherwise.

	AP	uci	wikipedia	reddit	lastfm	mooc
	transductive	$0.8726 \pm 5e^{3}$	$0.9137 \pm 58.3$	$0.9654 \pm 58.3$	0.7036 + 28.3	$0.8068 \pm 6.4$
	INTENSE	$0.9129 \pm 5e^{-3}$	$0.9078 \pm 363$	$0.9567 \pm 363$	$0.7090 \pm 2e^{-9}$ $0.7090 \pm 3e^{-4}$	0.7556 + 4e-4
DIE	SHUFFLE	$0.8509 \pm 3e-3$	$0.8962 \pm 4e-2$	$0.9613 \pm 4e-2$	$0.7036 \pm 1e-3$	$0.8072 \pm 5e-4$
О́Г	inductive	$0.7310 \pm 2\text{e-}2$	$0.8970 \pm 5\text{e-}3$	$0.9138 \pm \texttt{2e-2}$	$0.8431 \pm 4\text{e-}3$	$0.7931 \pm 1\text{e-}3$
	INTENSE	$0.8332 \pm 8\text{e-}3$	$0.8972 \pm 1e-2$	$0.9308 \pm 4\text{e-}2$	$0.8361 \pm 5e-4$	$0.7658 \pm 3e-4$
	SHUFFLE	$0.6994 \pm 8\text{e-}3$	$0.9078~\pm 2\text{e-}2$	$0.9251 \pm \text{6e-3}$	$0.8431 \pm 2e-3$	$0.7931 \pm 7e-4$
	transductive	$0.7694 \pm \textrm{7e-3}$	$0.9528 \pm \texttt{2e-3}$	$0.9818 \pm \text{6e-4}$	$0.7309 \pm 3\text{e-}4$	$0.8458 \pm \text{5e-4}$
	INTENSE	$0.8637 \pm 2e-2$	$0.9691 \pm 2e-3$	$0.9825 \pm 6e-4$	$0.9840 \pm 1e$ -4	$0.9610 \pm 1e-4$
GAJ	SHUFFLE	$0.7336 \pm 2e-2$	$0.9532 \pm 5e-3$	$0.9826 \pm 6e-3$	$0.7308 \pm 3e-4$	$0.8458 \pm 4e-4$
Ĥ	inductive	$0.7008 \pm 1\text{e-}2$	$0.9401~_{\text{2e-3}}$	$0.9658 \pm 1\text{e-}3$	$0.7817 \pm \mathtt{2e-4}$	$0.8430 \pm \textbf{3e-4}$
	Intense	$0.8095 \pm 2e-2$	$0.9621 \pm 2e-3$	$0.9676 \pm 1e-3$	$0.9841 \pm 1e-4$	$0.9621 \pm 1e-4$
	SHUFFLE	$0.6324 \pm 1e-2$	$0.9304 \pm 7e-3$	$0.9664 \pm 3e-3$	$0.7817 \pm 2e-4$	$0.8430 \pm 3e-4$
	transductive	$0.7975 \pm 1\text{e-}2$	$0.9472~\pm~\text{le-3}$	$0.9578 \pm 1\text{e-}3$	$0.7764 \pm 5\text{e-}3$	$0.8855 \pm 4\text{e-}3$
	INTENSE	$0.9709 \pm 3e-3$	$0.9911 \pm 6e-4$	$0.9744 \pm 2e-3$	$0.9916 \pm 1e-5$	$0.9629 \pm 6e-4$
lGN	SHUFFLE	$0.6520 \pm 4e-2$	$0.8487 \pm 3e-2$	$0.9563 \pm 2e-3$	$0.7764 \pm 9e-4$	$0.8848 \pm 1e-3$
	inductive	$0.7948 \pm \text{6e-3}$	$0.9463 \pm \text{1e-3}$	$0.9346 \pm 1\text{e-}3$	$0.8336 \pm 4\text{e-}3$	$0.8873 \pm \text{1e-3}$
	Intense	$0.9650 \pm 2e-3$	$0.9908 \pm 6e-4$	$0.9645 \pm 3e-3$	$0.9927 \pm 2e-5$	$0.9641 \pm 2e-4$
	SHUFFLE	$0.6193 \pm 9e-3$	$0.8376 \pm 3e-2$	$0.9299 \pm 3e-3$	$0.8337 \pm 6e-3$	$0.8872 \pm 2e-3$
	transductive	$0.9397 \pm \text{8e-4}$	$0.9901 \pm 1\text{e-4}$	$0.9884 \pm 3\text{e-}3$	$0.8755 \pm 3\text{e-}4$	$0.8667 \pm \mathtt{2e-4}$
5	INTENSE	$0.9889 \pm 7e-4$	$0.9975 \pm 8e-5$	$0.9942 \pm 7e-5$	$0.9879 \pm 2e-4$	$0.9719 \pm 1e-4$
AW	SHUFFLE	$0.8866 \pm 2e-3$	$0.9887 \pm 3e-4$	$0.9880 \pm 2e-3$	$0.8755 \pm 3e-4$	$0.8666 \pm 3e-4$
U	inductive	$0.9273~\pm\textrm{2e-3}$	$0.9896 \pm 4\text{e-}4$	$0.9859 \pm 3\text{e-}3$	$0.9031 \pm 5\text{e-}4$	$0.8543 \pm \text{4e-4}$
	INTENSE	$0.9857 \pm 2e-3$	$0.9971 \pm 1e-5$	$0.9938 \pm 8e-5$	$0.9889 \pm 3e-4$	$0.9731 \pm 2e-4$
	SHUFFLE	$0.8783 \pm 3e-2$	$0.9896 \pm 6e-3$	$0.9851 \pm 1e-3$	$0.9030 \pm 5e-4$	$0.8541 \pm 4e-4$
Ч	transductive	$0.9323~\pm\textrm{2e-3}$	$0.9690 \pm 4\text{e-}4$	$0.9738 \pm 3\text{e-}4$	$0.7630 \pm \text{1e-4}$	$0.8233 \pm \textbf{3e-4}$
ĽX.	INTENSE	$0.9923 \pm 6e-4$	$0.9966 \pm 2e-4$	$0.9965 \pm 1e-4$	$0.9858 \pm 1e-4$	$0.9537 \pm 1e-4$
МЧ	SHUFFLE	$0.8553 \pm 3e-3$	$0.9096 \pm 1e-3$	$0.9725 \pm 2e-4$	$0.7630 \pm 1e-4$	$0.8230 \pm 2e-4$
apl	inductive	$0.9133 \pm 1\text{e-}3$	$0.9639 \pm 1\text{e-}4$	$0.9517 \pm \text{8e-4}$	$0.8261 \pm 3\text{e-}4$	$0.8077~\pm\textrm{2e-4}$
ц	INTENSE	$0.9771 \pm 5e-4$	$0.9939 \pm 1e-4$	$0.9937 \pm 2e-4$	$0.9867 \pm 1e-4$	$0.9555 \pm 1e-4$
	SHUFFLE	$0.7945 \pm 3e-4$	$0.8900 \pm 2e-3$	$0.9477 \pm 7e-4$	$0.8261 \pm 3e-4$	$0.8072 \pm 3e-4$
ч	transductive	$0.9596 \pm 3\text{e-}4$	$0.9901~\pm \texttt{2e-4}$	$0.9921 \pm 1\text{e-}4$	$0.9096 \pm 1\text{e-}4$	$0.8622~^{\text{2e-4}}$
me.	ÍNTENSE	$0.9938 \pm 1e-4$	$0.9983 \pm 1e-4$	$0.9984 \pm 1e-4$	$0.9912 \pm 1e-4$	$0.9709 \pm 1e-4$
р	SHUFFLE	$0.9515 \pm 1e-3$	$0.9892 \pm 1e-4$	$0.9924 \pm 1e-4$	$0.9096 \pm 2e-4$	$0.8620 \pm 4e-4$
γGF	inductive	$0.9437 \pm 1\text{e-}4$	$0.9854 \pm 5\text{e-}4$	$0.9880 \pm 3\text{e-}4$	$0.9293 \pm 1\text{e-}4$	$0.8509 \pm \textbf{3e-4}$
Ď	INTENSE	$0.9854 \pm 1e-4$	$0.9965 \pm 4e-5$	$0.9973 \pm 1e-5$	$0.9918 \pm 2e-4$	$0.9723 \pm 1e-4$
	SHUFFLE	$0.9291 \pm 4e-4$	$0.9833 \pm 3e-4$	$0.9878 \pm 3e-4$	$0.9293 \pm 1e-4$	$0.8506 \pm 5e-4$

# 432 5 DISCUSSION

Some of the TLP models used in this work such as GraphMixer, and DyGFormer are considered the SoTA on most datasets, with near-perfect performance. However, as we showed earlier, a higher metric alone is not indicative of good performance without sanity checks. The counterfactual question helps make the evaluation more explainable, as models that perform worse on temporally distorted data with high ATD and ACD can claim superiority over models that do not. An ideal TLP model should be able to capture the difference in the count of edge events, their order, and the temporal shifts in the edge events.

To reiterate, if the performance of the model on the temporally distorted test data is similar or better
 than the performance on the original test data, then it implies one the following:

444

446

447 448

- 444 445
- the model has not made use of the temporal information in the training set,
- there is no useful temporal information in the dataset,
- the temporal distortion is weak.

In the absence of a guarantee that the dataset has useful temporal information that can aid prediction, we can compare different models by comparing the performance gaps.

449 450

451 **Future Work** Moving away from the binary classification approach to assess the performance of 452 temporal link prediction, future research should explore a generative approach where after observing 453 a temporal graph from time  $t \in (0, \tau_0)$ , the model can generate a temporal graph in  $t \in (\tau_0, T)$ . This 454 generated temporal graph can then be compared with the ground truth to measure similarity and 455 assess the performance of the model. The proposed metrics ATD and ACD can be used to measure the difference in the timestamps, as well as the occurrence frequency of the edges. Furthermore, the 456 same model architecture used by the TLP models discussed in this paper, can be further improved 457 by devising new training objectives that incorporate counterfactual analysis through the distortions 458 SHUFFLE and INTENSE. 459

460
 461
 462
 462
 463
 464
 464
 465
 466
 466
 466
 467
 468
 468
 469
 469
 469
 469
 469
 460
 460
 461
 462
 463
 464
 464
 465
 466
 466
 466
 467
 468
 468
 469
 469
 469
 469
 460
 460
 461
 462
 463
 464
 464
 465
 466
 466
 466
 467
 468
 468
 469
 469
 469
 469
 460
 460
 461
 462
 463
 464
 464
 464
 465
 466
 466
 466
 466
 466
 467
 468
 468
 469
 469
 469
 469
 460
 461
 462
 463
 464
 464
 464
 465
 466
 466
 466
 466
 467
 468
 468
 469
 469
 469
 469
 469
 460
 461
 462
 463
 464
 464
 464
 464
 464
 464
 465
 466
 466
 466
 466
 466
 466
 466
 467
 468
 468
 468
 469
 469

465 466 467

468

469

470

471

477

478

479

# References

- A. Celikkanat, N. Nakis, and M. Mørup. Continuous-time graph representation with sequential survival process. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 11177–11185, 2024.
- S. Chatterjee. Learning and memorization. In *International conference on machine learning*, pages 755–763. PMLR, 2018.
- W. Cong, S. Zhang, J. Kang, B. Yuan, H. Wu, X. Zhou, H. Tong, and M. Mahdavi. Do We Really Need Complicated Model Architectures For Temporal Networks? In *The Eleventh International Conference on Learning Representations*, Sept. 2023.
  - N. N. Daud, S. H. Ab Hamid, M. Saadoon, F. Sahran, and N. B. Anuar. Applications of link prediction in social networks: A review. *Journal of Network and Computer Applications*, 166:102716, 2020.
- N. Du, H. Dai, R. Trivedi, U. Upadhyay, M. Gomez-Rodriguez, and L. Song. Recurrent marked temporal point processes: Embedding event history to vector. In *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pages 1555–1564, 2016.
- 485 W. Fan, Y. Ma, Q. Li, Y. He, E. Zhao, J. Tang, and D. Yin. Graph neural networks for social recommendation. In *The world wide web conference*, pages 417–426, 2019.

490

491

492

495

496

497

498 499

500

501

505

507

508

509

510

511 512

513

514

520

521

522 523

524

530

531

- Z. Fan, Z. Liu, J. Zhang, Y. Xiong, L. Zheng, and P. S. Yu. Continuous-time sequential recommendation with temporal graph collaborative transformer. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pages 433–442, 2021.
  - S. Huang, F. Poursafaei, J. Danovitch, M. Fey, W. Hu, E. Rossi, J. Leskovec, M. Bronstein, G. Rabusseau, and R. Rabbany. Temporal graph benchmark for machine learning on temporal graphs. *Advances in Neural Information Processing Systems*, 36, 2024.
- Y. Kim, Y. Lee, M. Choe, S. Oh, and Y. Lee. Temporal graph networks for graph anomaly detection
   in financial networks. *arXiv preprint arXiv:2404.00060*, 2024.
  - S. Kumar, X. Zhang, and J. Leskovec. Predicting dynamic embedding trajectory in temporal interaction networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1269–1278, 2019.
  - A. Modell, I. Gallagher, E. Ceccherini, N. Whiteley, and P. Rubin-Delanchy. Intensity profile projection: A framework for continuous-time representation learning for dynamic networks. *Advances in Neural Information Processing Systems*, 36, 2024.
- P. Panzarasa, T. Opsahl, and K. M. Carley. Patterns and dynamics of users' behavior and interaction: Network analysis of an online community. *Journal of the American Society for Information Science and Technology*, 60(5):911–932, 2009.
- J. Pearl. *Causality*. Cambridge university press, 2009.
  - J. Pearl. The seven tools of causal inference, with reflections on machine learning. *Communications* of the ACM, 62(3):54–60, 2019.
  - J. W. Pennebaker, M. E. Francis, and R. J. Booth. Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*, 71(2001):2001, 2001.
  - P. O. Perry and P. J. Wolfe. Point process modelling for directed interaction networks. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 75(5):821–849, 2013.
- F. Poursafaei and R. Rabbany. Exhaustive Evaluation of Dynamic Link Prediction. In 2023 IEEE International Conference on Data Mining Workshops (ICDMW), pages 1121–1130, Shanghai, China, Dec. 2023. IEEE. ISBN 9798350381641. doi: 10.1109/ICDMW60847.2023.00147.
- F. Poursafaei, S. Huang, K. Pelrine, and R. Rabbany. Towards better evaluation for dynamic link
  prediction. *Advances in Neural Information Processing Systems*, 35:32928–32941, 2022.
  - M. Qin and D.-Y. Yeung. Temporal Link Prediction: A Unified Framework, Taxonomy, and Review, June 2023.
  - Y. Qin, W. Ju, H. Wu, X. Luo, and M. Zhang. Learning graph ode for continuous-time sequential recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 2024.
- E. Rossi, B. Chamberlain, F. Frasca, D. Eynard, F. Monti, and M. Bronstein. Temporal graph networks for deep learning on dynamic graphs. *arXiv preprint arXiv:2006.10637*, 2020.
- 528 O. Shchur, A. C. Türkmen, T. Januschowski, and S. Günnemann. Neural temporal point processes: A review. *arXiv preprint arXiv:2104.03528*, 2021.
  - R. Trivedi, M. Farajtabar, P. Biswal, and H. Zha. Dyrep: Learning representations over dynamic graphs. In *International conference on learning representations*, 2019.
- Y. Wang, Y.-Y. Chang, Y. Liu, J. Leskovec, and P. Li. Inductive representation learning in temporal networks via causal anonymous walks. In *International Conference on Learning Representations*, 2020.
- D. Xu, C. Ruan, E. Korpeoglu, S. Kumar, and K. Achan. Inductive representation learning on temporal graphs. In *International Conference on Learning Representations*, 2020.
- 539 L. Yu, L. Sun, B. Du, and W. Lv. Towards better dynamic graph learning: New architecture and unified library. *Advances in Neural Information Processing Systems*, 36:67686–67700, 2023.

540 541	А	DATASE	ets & Moe	DELS								
542 543	A.1	TEMPOR	IPORAL GRAPH DATASETS									
544	We use the following datasets <sup>2</sup> to perform counterfactual analysis <sup>3</sup> :											
545												
546		• wikipedia (Kumar et al., 2019) describes a dynamic graph of interaction between the aditors and Wikipedia pages over a span of any month. The antrice consist of the user ID										
547		nage I	nage ID, and timestamp. The edge features are LIWC-feature vectors (Pennebaker et al									
548		2001)	2001) of the edit text. The edge feature dimension is 172.									
549		• modd	• us delit (Kumon at al. 2010) describes a kinemitic interaction and hater and									
550		and si	• reduce (Kulliar et al., 2019) describes a objartite interaction graph between the Users and subreddits. The interaction event is recorded with the IDs of the user subreddit and									
551		timest	timestamp. Similar to wikipedia, the post content is converted into a LIWC-feature									
552		vector of dimension 172 which serves as the edge feature.										
553		• uci (	Panzarasa et	al 2009) is a dy	znamic graph des	cribing message-ex	change among					
555		the stu	idents at Univ	versity of Californ	ia at Irvine (UCI)	from April to Octo	ber 2004. The					
556	interaction event consists of the user IDs, and timestamp.											
557		• last	fm (Kumar ei	t al., 2019) is also	a bipartite graph	depicting the interac	ctions between					
558	1000 users and 1000 most listened songs over a span of one month.											
559		• mooc	(Kumar et al	2019) as the nam	e suggests is a stu	dent interaction netw	ork enrolled in					
560		the sa	me online cou	rse.			ork enroned in					
561												
562				Table 4: The sca	le of different data	isets.						
563												
564			Dataset	Total nodes $(10^3)$	Total Edges (10 <sup>3</sup> )	Unique Edges (10 <sup>3</sup> )						
565			uci	1.89	59.84	20.29						
566			wikipedia	9.23	157.47	18.25						
567			reddit	10.98	672.45	78.52						
568			mooc	7.14	411.75	178.44						
569												
570												
571	A.2	TEMPOR	ral Link Pr	EDICTION MODE	LS							
572	War	naka usa o	f the followin	a models <sup>4</sup> to test t	ha countarfactual	framawark						
574	we i	liake use 0		g moders to test t		Inamework.						
575		• JODI	E (Kumar et	al., 2019) uses a	recurrent neural	network (RNN) to	generate node					
576		embec	ldings for eacl	n interaction event	. The future embed	lding of a node is est	imated through					
577		a nove	el projection o	perator which is t	urn in used to prec	lict future edge even	its.					
578		• TGAT	(Xu et al., 20	20) relies on self-	attention mechanis	sm to generate node	embeddings to					
579		captur	e the tempora	l evolution of the	graph structure.							
580		• TGN (	Rossi et al., 2	020) combine mer	nory modules with	graph-based operate	ors to create an					
581		encod	er-decoder pa	ir capable of creat	ing temporal node	embeddings.						
582		• CAWN	(Wang et al.	. 2020) propose a	novel strategy b	ased on the law of	triadic closure.					
583		where	temporal wa	lks retrieve the dy	namic graph mot	ifs without explicitly	y counting and					
584		selecti	ng the motifs.	The node IDs are	replaced with the	hitting counts to faci	ilitate inductive					
585		infere	nce.									
586		• Grap	hMixer(Co	ng et al., 2023) use	e a simple architec	ture where the encod	ler and decoder					
587		are de	signed using	multi-layer percep	trons (MLPs).							
588		• DyGF	ormer (Yue	t al., 2023) use a t	ransformer to lear	n from nodes' first-h	op interactions					
589		and re	port SoTA res	sults on most of th	e datasets.		•					
590												
591			1 1 .	1.10 1	1 / 1	12707						
592	2'.	i ne datasets	can be downlo	aded from https://ze	enodo.org/records//2	213/96						

<sup>&</sup>lt;sup>3</sup>The datasets are chronologically split in the ratio 0.7 : 0.15 : 0.15 into train, validation, and test sets. <sup>4</sup>The optimal hyper-parameters reported by the models are used.