# A Closer Look at In-Context Learning for Temporal Knowledge Graph Forecasting

Anonymous ACL submission

### Abstract

While temporal knowledge graph forecasting (TKGF) approaches have traditionally relied heavily on complex graph neural network architectures, recent advances in large language models, specifically in-context learning (ICL), have presented promising out-of-the-box alternatives. While previous works have shown the 007 potential of using ICL, its limitations and generalization capabilities for TKGF are underexplored. In this study, we conduct a comparative analysis of complexity (e.g., number 011 of hops) and sparsity (e.g., relation frequency) confounders between ICL and supervised mod-013 els using two annotated TKGF benchmarks. Our experimental results showcase that while ICL performs on par or outperforms supervised models in lower complexity scenarios, its effec-017 tiveness diminishes in more complex settings (e.g., multi-step, more number of hops, etc.), where supervised models are superior.

#### 1 Introduction

024

027

Knowledge graphs (KGs) are commonly used structures that store relational information as a graph (Bollacker et al., 2008; Vrandečić and Krötzsch, 2014). While using KGs for keeping static facts is common, they are unsuitable for holding complex dynamic (*i.e.*, temporal) information. Temporal knowledge graphs (TKGs) are extensions of KGs that enable the storage of such information (Leetaru and Schrodt, 2013; García-Durán et al., 2018). Consequently, TKGs allow practitioners to do various predictive tasks on complex temporal data. One critical task that has been empowered by TKGs is temporal knowledge graph forecasting (TKGF) (Gastinger et al., 2023), where the objective is to predict future facts from a set of prior facts before a specific time in a TKG. A hypothetical real-world example of TKGF is to answer the question, "Who is USA going to Meet in June 2025?" based on previous political



Figure 1: **Example of Graph to Text Prompt Conver sion for ICL.** The given task is to predict which country is gonna meet the USA during the G7 summit, based on the previous interactions between the countries.

events. This scenario can be represented by the query quadruple q = (USA, Meets, ?, June 2025) and the time-constrained TKG  $\mathcal{G}_t = \{(USA, Meets, UK, June 2024), (USA, Attends, G7, June 2025), (USA, Meets, Germany, June 2025), ... \}.$ 

041

042

044

046

051

054

Recent studies have demonstrated large language models' (LLMs) effectiveness as general estimators across various function classes (Garg et al., 2022; Mirchandani et al., 2023). Consequently, these advancements have sparked interest in employing LLMs for temporal knowledge graph forecasting (TKGF). Specifically, LLMs have shown remarkable potential for TKGF, surpassing state-ofthe-art supervised models in some scenarios using

	(	Complexity		Spars	sity
Temporal Rule	# Unique Entities	# Unique Relations	# Hops	Relation Frequency	Time Interval
$(E_1, express intent to meet^{-1}, E_2, T_1) \Rightarrow (E_1, \underbrace{share information}_R, E_2, T_2)$	2	2	1	$f_R$	$T_2 - T_1$
$(E_1, provide \ military \ aid, E_2, T_1) \land (E_2, intend \ to \ protect^{-1}, E_3, T_2) \Rightarrow (E_1, \underbrace{provide \ military \ aid}_R, E_3, T_3)$	3	2	2	$f_R$	$T_3 - T_1$
$(E_1, riot, E_2, T_1) \land (E_2, make statement, E_1, T_2) \land (E_1, riot, E_2, T_3) \Rightarrow (E_1, \underbrace{demonstrate \ or \ rally}_R, E_2, T_4)$	2	3	3	$f_R$	$T_4 - T_1$

Table 1: Confounder values examples. The samples are taken from Liu et al. (2022) with some small modifications. Note that  $f_R$  refers to the frequency of relation R among all quadruples in the dataset.

in-context-learning (ICL) (Lee et al., 2023) (see Figure 1 for an example). Methods such as ICL present a cheap, fast, and ready-to-use alternative to traditional methods, many of which use computationally heavy graph neural network (GNN) architectures. However, despite all their benefits, the broad applications of such solutions for forecasting problems and LLMs' "grey box" nature (*e.g.*, opaque reasoning process, unpredictability across different temporal patterns) raise concerns regarding their limitations and generalizability.

In this study, we provide insights into the effect of various confounders - arising from relational and temporal patterns - on the effectiveness of ICL for TKGF. To this end, first, we utilize a state-of-the-art rule-based model to generate reasoning rules from well-known TKG benchmarks, ICEWS14 and ICEWS18 (García-Durán et al., 2018). Then, based on the generated rules, we create two labeled datasets containing confounder annotations for the test sets. Finally, we use these datasets to compare ICL-based models to stateof-the-art supervised models in single-step and multi-step settings (Gastinger et al., 2023) across complexity (e.g., number of unique entities), and sparsity (e.g., relation frequency) confounders (see Table 1 for more thorough examples). Our experimental results on the annotated datasets show that (1) ICL-based models outperform supervised models in scenarios with lower complexity, such as annotated samples with 1-hop patterns in singlestep settings or samples involving only one unique relation, and (2) increasing the complexity of the patterns results in ICL-based models to underperform massively compared to the supervised models. This phenomenon is particularly evident in multistep settings, where ICL-based models lag behind

supervised models in all scenarios.

#### 2 Background and Related Work

**Formal Definition of TKGF.** Formally, a TKG  $\mathcal{G} = (\mathcal{Q}, \mathcal{E}, \mathcal{R}, \mathcal{T})$  comprises a set of quadruples  $\mathcal{Q}$  in the form (s, r, o, t), where s and o are entities within  $\mathcal{E}$ , r is a relation within  $\mathcal{R}$ , and t is a timestamp from  $\mathcal{T}$ . The TKG forecasting task aims to predict a missing entity in future quadruples, either as (s, r, ?, t) for tail prediction or (?, r, o, t) for head prediction, using historical data from the graph. This process involves scoring all entities so that the true entity receives the highest ranking.

092

093

094

097

099

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

**Supervised Models.** Recent supervised models primarily utilize embedding-based GNNs to enhance their structural and sequential learning capabilities. Specifically, they have used autoregressive architectures to aggregate information both globally and locally in RE-Net (Jin et al., 2020), combined convolutional and recurrent architectures for modeling temporal sequences in RE-GCN (Li et al., 2021), introduced neural ordinary differential equations to model temporal sequences in TANGO (Han et al., 2021), and extended convolutional architectures to learn evolutionary patterns in CEN (Li et al., 2022). In parallel to these models, other approaches have been introduced in prior works, such as using a copy-mechanism in CyGNet (Zhu et al., 2021), leveraging reinforcement learning on temporal paths in TiTer (Sun et al., 2021), and learning temporal logic rules via temporal random walks in TLogic (Liu et al., 2022).

ICL-based Models. Recent advances in LLMs have drastically improved their capabilities, leading to emergent behaviors such as ICL. ICL allows LLMs to perform tasks conditioned solely on

087

880

Dataset	C	$ \mathcal{D} $	# of Fac	Time	
Dataset	$ \mathcal{L} $	$ \mathcal{K} $	Train/Valid/Test	Annotated	Granularity
ICEWS14 ICEWS18	6,869 23,033	230 256	75k/8.5k/7.3k 373k/46k/50k	11,625 65,003	1 day 1 day

Table 2: **Dataset statistics.** Each dataset consists of historical facts divided into three subsets based on time.

the provided context without parameter optimization. Utilizing ICL, Lee et al. (2023) introduced the first LLM-based TKGF model, which showed performance on par with state-of-the-art supervised models without any training. Moreover, Xia et al. (2024) introduced an improved historical fact retriever and an alignment training procedure, posting better performances than the state-of-the-art supervised models<sup>1</sup>. While these ICL-based models have shown interesting achievements toward the TKGF task, we still lack a proper understanding of their limitations, a gap we aim to bridge.

#### **3** Experimental Setup

#### 3.1 Datasets

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

162

163

164

165

Our experiments focused on two prominent TKGF datasets: ICEWS14 (García-Durán et al., 2018) and ICEWS18 (Jin et al., 2020) (see Table 2). We specifically chose these datasets because 1) they are commonly used by almost all the prior works in the literature and 2) they pose a much more significant challenge to the forecasting models compared to other existing datasets such as WIKI (Leblay and Chekol, 2018) and YAGO (Rebele et al., 2016). Moreover, to keep our results consistent and comparable to previous works, we use the same splits as Gastinger et al. (2023).

#### 3.2 Weak Supervision

One of the challenges we faced in our experiments was the absence of annotations for different confounders in the existing datasets. To overcome this issue, we employ weak supervision (Voskarides et al., 2018; Zhang et al., 2024; Tong et al., 2024) using TLogic (Liu et al., 2022), a state-of-the-art rule-learning-based TKG model, to annotate test samples with temporal multi-hop patterns. To this end, first, we ran the rule-learning part of TLogic on the combination of all quadruples from the train, valid, and test sets with the number of hops  $\in \{1, 2, 3\}$ . Then, we annotated each test sample using the matching pattern with the highest score<sup>2</sup>, if such a rule existed. Finally, for the annotated test quadruples, we extract various confounders from their associated patterns, including the number of unique entities and relations, the pattern's length denoted as "hop", the relation frequency of the test query, and the time interval (see Table 2 for annotation statistics).

#### 3.3 Models

For our ICL-based baseline, we utilize the model as described by Lee et al. (2023), which employs gpt-neox-20b (Black et al., 2022). This method is an inference-time approach that demonstrates performance comparable to supervised models. Moreover, for the TKG baselines we used state-of-the-art models with the hyperparameters and implementation as provided by (Gastinger et al., 2023): RE-Net (Jin et al., 2020), RE-GCN (Li et al., 2021), TANGO (Han et al., 2021), CyGNet (Zhu et al., 2021), and CEN (Li et al., 2022).

#### **3.4 Implementation Details**

We retain the top 100 entities with the highest scores (or the highest log probability) to evaluate each prediction. This protocol is done due to a limitation of the ICL-based model preventing it from predicting entities that do not appear in its context, which at most contains 100 historical facts, bounded by the context length of the underlying model (*i.e.*, gpt-neox-20b). Moreover, this protocol allows us to evaluate and fairly compare the ICL-based and supervised models across our experiments. As for our metrics, we report the Hits@ $\{1,3\}$  based on the list of retained entities for each prediction. All baseline models report both head and tail prediction performance by generating a head query (?, r, o, t) and a tail query (s, r, ?, t)for each test quadruple (s, r, o, t), following standard practices in the literature. We report the average head and tail prediction performances. The codebase uses PyTorch (Paszke et al., 2019) and Huggingface (Wolf et al., 2020) libraries.

## 4 Experiments

Table 3 presents our experimental results on bothsingle-step (top) and multi-step (bottom) queries,grouped by the *number of hops* as the confounder.We can observe that the ICL-based models only

205

206

207

208

210

211

166

<sup>&</sup>lt;sup>1</sup>The implementation has not been made public.

<sup>&</sup>lt;sup>2</sup>For each matched reasoning path, TLogic combines rule confidence and temporal recency scores into one score.

				ICEV	WS14					ICEV	WS18		
Single-step	Train		H@1			H@3			H@1			H@3	
		1-hop	2-hop	3-hop	1-hop	2-hop	3-hop	1-hop	2-hop	3-hop	1-hop	2-hop	3-hop
RE-GCN	1	42.6	15.2	38.7	63.6	32.2	56.1	34.5	19.5	31.9	54.7	35.5	50.2
TANGO	1	36.4	12.0	36.2	54.5	24.8	50.2	29.7	16.3	28.3	48.8	31.0	45.5
CEN	1	43.3	15.2	39.0	63.2	30.0	56.2	33.9	18.9	31.1	54.0	34.3	49.1
Average		40.8	14.1	38.0	60.4	29.0	54.2	32.7	18.3	30.4	52.5	33.6	48.3
Median		42.6	15.2	38.7	63.2	30.0	56.1	33.9	18.9	31.1	54.0	34.3	49.1
gpt-neox-20b-entity	X	46.4	10.6	37.7	65.8	17.8	54.0	29.8	11.2	27.9	48.1	22.4	43.6
$\Delta$ Average		5.6	-3.5	-0.2	5.4	-11.2	-0.2	-2.9	-7.1	-2.5	-4.4	-11.2	-4.6
$\Delta$ Median		3.7	-4.6	-0.1	2.7	-12.2	-2.1	-4.1	-7.8	-3.2	-5.9	-11.9	-5.5
gpt-neox-20b-pair	X	43.6	6.8	37.9	58.3	9.8	51.4	31.0	10.9	30.1	48.7	17.9	47.1
$\Delta$ Average		2.9	-7.4	-0.1	-2.2	-19.2	-2.7	-1.7	-7.3	-0.3	-3.8	-15.7	-1.2
$\Delta$ Median		1.0	-8.5	-0.8	-4.9	-20.2	-4.6	-2.9	-8.0	-0.4	-5.3	-16.4	-2.0
				ICEV	WS14					ICEV	WS18		
Multi-step	Train		H@1	ICEV	WS14	H@3			H@1	ICEV	WS18	H@3	
Multi-step	Train	1-hop	H@1 2-hop	ICEV 3-hop	WS14	H@3 2-hop	3-hop	1-hop	H@1 2-hop	ICEV 3-hop	WS18	H@3 2-hop	3-hop
Multi-step	Train √	1-hop 37.3	H@1 2-hop 13.3	ICEN 3-hop 36.0	WS14 1-hop 54.1	H@3 2-hop 25.9	3-hop 51.3	1-hop 28.8	H@1 2-hop 16.0	ICEV 3-hop 27.8	WS18 1-hop 48.0	H@3 2-hop 31.4	3-hop 45.0
Multi-step RE-NET RE-GCN	Train ✓ ✓	1-hop 37.3 36.6	H@1 2-hop 13.3 15.7	3-hop 36.0 34.9	WS14 1-hop 54.1 55.4	H@3 2-hop 25.9 30.0	3-hop 51.3 49.0	1-hop 28.8 29.5	H@1 2-hop 16.0 18.2	ICEV 3-hop 27.8 28.9	WS18 1-hop 48.0 48.3	H@3 2-hop 31.4 33.0	3-hop 45.0 45.8
Multi-step RE-NET RE-GCN CyGNet	Train ✓ ✓ ✓	1-hop 37.3 36.6 35.5	H@1 2-hop 13.3 15.7 11.9	ICEN 3-hop 36.0 34.9 34.5	WS14 1-hop 54.1 55.4 53.6	H@3 2-hop 25.9 30.0 26.0	3-hop 51.3 49.0 49.9	1-hop 28.8 29.5 25.5	H@1 2-hop 16.0 18.2 13.4	ICEN 3-hop 27.8 28.9 26.1	WS18 1-hop 48.0 48.3 44.9	H@3 2-hop 31.4 33.0 28.3	3-hop 45.0 45.8 44.1
Multi-step RE-NET RE-GCN CyGNet Average	Train	1-hop 37.3 36.6 35.5 36.4	H@1 2-hop 13.3 15.7 11.9 13.6	ICEN 3-hop 36.0 34.9 34.5 35.1	WS14 1-hop 54.1 55.4 53.6 54.3	H@3 2-hop 25.9 30.0 26.0 27.3	3-hop 51.3 49.0 49.9 50.1	1-hop 28.8 29.5 25.5 27.9	H@1 2-hop 16.0 18.2 13.4 15.9	ICEN 3-hop 27.8 28.9 26.1 27.6	WS18 1-hop 48.0 48.3 44.9 47.1	H@3 2-hop 31.4 33.0 28.3 30.9	3-hop 45.0 45.8 44.1 45.0
Multi-step RE-NET RE-GCN CyGNet Average Median	Train ✓ ✓ ✓	1-hop 37.3 36.6 35.5 36.4 36.6	H@1 2-hop 13.3 15.7 11.9 13.6 13.3	<b>ICEN</b> 3-hop 36.0 34.9 34.5 35.1 34.9	WS14 1-hop 54.1 55.4 53.6 54.3 54.1 54.1	H@3 2-hop 25.9 30.0 26.0 27.3 26.0	3-hop 51.3 49.0 49.9 50.1 49.9	1-hop 28.8 29.5 25.5 27.9 28.8	H@1 2-hop 16.0 18.2 13.4 15.9 16.0	ICEV 3-hop 27.8 28.9 26.1 27.6 27.8	WS18 1-hop 48.0 48.3 44.9 47.1 48.0	H@3 2-hop 31.4 33.0 28.3 30.9 31.4	3-hop 45.0 45.8 44.1 45.0 45.0
RE-NET         RE-GCN         CyGNet         Average         Median         gpt-neox-20b-entity	Train	1-hop 37.3 36.6 35.5 36.4 36.6 34.3	H@1 2-hop 13.3 15.7 11.9 13.6 13.3 8.7	<b>ICEN</b> 3-hop 36.0 34.9 34.5 35.1 34.9 32.1	WS14 1-hop 54.1 55.4 53.6 54.3 54.1 49.6	H@3 2-hop 25.9 30.0 26.0 27.3 26.0 16.9	3-hop 51.3 49.0 49.9 50.1 49.9 44.6	1-hop 28.8 29.5 25.5 27.9 28.8 19.7	H@1 2-hop 16.0 18.2 13.4 15.9 16.0 8.9	<b>ICE</b> 3-hop 27.8 28.9 26.1 27.6 27.8 19.7	WS18 1-hop 48.0 48.3 44.9 47.1 48.0 31.3	H@3 2-hop 31.4 33.0 28.3 30.9 31.4 17.8	3-hop 45.0 45.8 44.1 45.0 45.0 30.7
RE-NET         RE-GCN         CyGNet         Average         Median         gpt-neox-20b-entity         △ Average	Train	1-hop           37.3           36.6           35.5           36.4           36.6           34.3           -2.1	H@1 2-hop 13.3 15.7 11.9 13.6 13.3 8.7 -4.9	3-hop 36.0 34.9 34.5 35.1 34.9 32.1 -3.0	WS14 1-hop 54.1 55.4 53.6 54.3 54.1 49.6 -4.8	H@3 2-hop 25.9 30.0 26.0 27.3 26.0 16.9 -10.4	3-hop 51.3 49.0 49.9 50.1 49.9 44.6 -5.4	1-hop           28.8           29.5           25.5           27.9           28.8           19.7           -8.2	H@1 2-hop 16.0 18.2 13.4 15.9 16.0 8.9 -7.0	ICE 3-hop 27.8 28.9 26.1 27.6 27.8 19.7 -7.9	WS18 1-hop 48.0 48.3 44.9 47.1 48.0 31.3 -15.8	H@3 2-hop 31.4 33.0 28.3 30.9 31.4 17.8 -13.1	3-hop 45.0 45.8 44.1 45.0 45.0 30.7 -14.2
RE-NET         RE-GCN         CyGNet         Average         Median         gpt-neox-20b-entity $\Delta$ Average $\Delta$ Median	Train ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓	1-hop           37.3           36.6           35.5           36.4           36.6           34.3           -2.1           -2.3	H@1 2-hop 13.3 15.7 11.9 13.6 13.3 8.7 -4.9 -4.6	ICEV 3-hop 36.0 34.9 34.5 35.1 34.9 32.1 -3.0 -2.8	WS14 1-hop 54.1 55.4 53.6 54.3 54.1 49.6 -4.8 -4.5	H@3 2-hop 25.9 30.0 26.0 27.3 26.0 16.9 -10.4 -9.1	3-hop 51.3 49.0 49.9 50.1 49.9 44.6 -5.4 -5.3	1-hop           28.8           29.5           25.5           27.9           28.8           19.7           -8.2           -9.0	H@1 2-hop 16.0 18.2 13.4 15.9 16.0 8.9 -7.0 -7.2	ICEV 3-hop 27.8 28.9 26.1 27.6 27.8 19.7 -7.9 -8.1	WS18 1-hop 48.0 48.3 44.9 47.1 48.0 31.3 -15.8 -16.7	H@3 2-hop 31.4 33.0 28.3 30.9 31.4 17.8 -13.1 -13.6	3-hop 45.0 45.8 44.1 45.0 45.0 30.7 -14.2 -14.3
RE-NET         RE-GCN         CyGNet         Average         Median         gpt-neox-20b-entity $\Delta$ Average $\Delta$ Median         gpt-neox-20b-pair	Train	1-hop           37.3           36.6           35.5           36.4           36.6           34.3           -2.1           -2.3           30.9	H@1 2-hop 13.3 15.7 11.9 13.6 13.3 8.7 -4.9 -4.6 6.5	ICEV 3-hop 36.0 34.9 34.5 35.1 34.9 32.1 -3.0 -2.8 32.6	WS14 1-hop 54.1 55.4 53.6 54.3 54.1 49.6 -4.8 -4.5 43.7	H@3 2-hop 25.9 30.0 26.0 27.3 26.0 16.9 -10.4 -9.1 8.9	3-hop 51.3 49.0 49.9 50.1 49.9 44.6 -5.4 -5.3 43.4	1-hop           28.8           29.5           25.5           27.9           28.8           19.7           -8.2           -9.0           23.7	H@1 2-hop 16.0 18.2 13.4 15.9 16.0 8.9 -7.0 -7.2 8.7	ICEV 3-hop 27.8 28.9 26.1 27.6 27.8 19.7 -7.9 -8.1 25.6	WS18 1-hop 48.0 48.3 44.9 47.1 48.0 31.3 -15.8 -16.7 37.9	H@3 2-hop 31.4 33.0 28.3 30.9 31.4 17.8 -13.1 -13.6 14.4	3-hop 45.0 45.8 44.1 45.0 45.0 30.7 -14.2 -14.3 38.5
Multi-step         RE-NET         RE-GCN         CyGNet         Average         Median         gpt-neox-20b-entity         Δ Average         Δ Median         gpt-neox-20b-pair         Δ Average	Train	1-hop           37.3           36.6           35.5           36.4           36.6           34.3           -2.1           -2.3           30.9           -5.5	H@1 2-hop 13.3 15.7 11.9 13.6 13.3 8.7 -4.9 -4.6 6.5 -7.1	ICEV 3-hop 36.0 34.9 34.5 35.1 34.9 32.1 -3.0 -2.8 32.6 -2.5	WS14 1-hop 54.1 53.6 54.3 54.3 54.1 49.6 -4.8 -4.8 -4.5 43.7 -10.6	H@3 2-hop 25.9 30.0 26.0 27.3 26.0 16.9 -10.4 -9.1 8.9 -18.3	3-hop 51.3 49.0 50.1 49.9 50.1 49.9 44.6 -5.4 -5.3 43.4 -6.7	1-hop 28.8 29.5 25.5 27.9 28.8 19.7 -8.2 -9.0 23.7 -4.2	H@1 2-hop 16.0 18.2 13.4 15.9 16.0 8.9 -7.0 -7.2 8.7 -7.2	ICEV 3-hop 27.8 28.9 26.1 27.6 27.8 19.7 -7.9 -8.1 25.6 -2.0	WS18 1-hop 48.0 48.3 44.9 47.1 48.0 31.3 -15.8 -16.7 37.9 -9.2	H@3 2-hop 31.4 33.0 28.3 30.9 31.4 17.8 -13.1 -13.6 14.4 -16.4	3-hop 45.0 45.8 44.1 45.0 45.0 30.7 -14.2 -14.3 38.5 -6.5

Table 3: **Performance (Hits@K)** comparison between supervised models and ICL for single-step (top) and multistep (bottom) prediction, grouped by the **number of hops** as the confounder. The first group consists of supervised models, whereas the second group consists of ICL models, *i.e.*, GPT-NeoX. The green and red colors represent where ICL is outperforming and underperforming the average performance of the supervised models.

212 perform better with 1-hop queries in the ICEWS14 dataset. Moreover, as the number of hops, an in-213 dicator of the pattern complexity, increases, super-214 215 vised models outperform ICL-based models. Inter-216 estingly, this decline in performance is not monotonic in terms of complexity, making it even more 217 challenging to predict the potential pitfalls. For 218 example, LLMs' worst performance in ICEWS14 219 occurs in 2-hop queries, while the performance on 3-hop queries stays competitive. Moreover, we 221 observe the same trend when analyzing other confounders related to pattern complexity. For example, ICL-based models outperform the supervised models in patterns involving two unique entities 225 on ICEWS14. However, as the number of unique entities increases, the performance of ICL-based models declines (see Table 4 in Appendix A). Similarly, this trend is evident when the samples are grouped by number of unique relations (see Table 5 in Appendix A). When the samples are grouped by 231 *relation frequency*, the ICL-based models perform on par or moderately outperform the supervised 233 models only in the ICEWS14 single-step setting. 234 In all other cases, the supervised models outperform the ICL-based models. However, the upward

trend in Figure 2 (Appendix A) indicates that as relation frequency increases, the performance gap between the ICL-based and supervised models decreases. Moreover, when the samples are grouped by *time interval* (see Figure 3 in Appendix A), the supervised models consistently outperform the ICLbased models. We observe that ICL-based models perform worse in the multi-step setup across all confounders than their counterpart average supervised models. Finally, the performance gap is wider on the ICEWS18 (compared to ICEWS14), which could be attributed to it being more challenging. 237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

## 5 Conclusion

In this paper, we presented an in-depth analysis of the effect of various confounders on the predictive power of ICL-based and supervised models for TKGF. Specifically, we created two annotated benchmarks for testing models across varied complexities and sparsity levels. Our experimental results indicate that while ICL is effective in lowcomplexity scenarios, its performance rapidly deteriorates as the complexity of the patterns increases. These findings highlight the need for more granular evaluation and testing of LLMs for TKGF. 261

277

278

279

281

282

284

289

290

291

296

297

301

302

307

308

310

311

312

313

## Limitations

262 Our work is the first step toward a more granular evaluation of TKGF. As such, expanding the 263 presented findings with more annotated datasets, identifying additional confounders, and evaluat-265 ing a broader range of supervised and LLM-based 267 models should be explored in future works. With the growing utilization of LLMs, comprehensive benchmarks allow us to make more grounded comparisons across models rather than being misled by potential spurious biases. Moreover, we observed 271 fluctuations in performance gaps with increased 272 complexities in different confounders. This phe-273 nomenon makes the performance comparison more uncertain, which should be further investigated in future works.

## References

- Sidney Black, Stella Biderman, Eric Hallahan, Quentin Anthony, Leo Gao, Laurence Golding, Horace He, Connor Leahy, Kyle McDonell, Jason Phang, Michael Pieler, Usvsn Sai Prashanth, Shivanshu Purohit, Laria Reynolds, Jonathan Tow, Ben Wang, and Samuel Weinbach. 2022. GPT-NeoX-20B: An opensource autoregressive language model. In Proceedings of BigScience Episode #5 – Workshop on Challenges & Perspectives in Creating Large Language Models, pages 95–136, virtual+Dublin. Association for Computational Linguistics.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In Proceedings of the 2008 ACM SIGMOD International Conference on Management of Data, SIGMOD '08, page 1247–1250, New York, NY, USA. Association for Computing Machinery.
- Alberto García-Durán, Sebastijan Dumančić, and Mathias Niepert. 2018. Learning sequence encoders for temporal knowledge graph completion. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 4816–4821, Brussels, Belgium. Association for Computational Linguistics.
- Shivam Garg, Dimitris Tsipras, Percy S Liang, and Gregory Valiant. 2022. What can transformers learn incontext? a case study of simple function classes. In Advances in Neural Information Processing Systems, volume 35, pages 30583–30598. Curran Associates, Inc.
- Julia Gastinger, Timo Sztyler, Lokesh Sharma, Anett Schuelke, and Heiner Stuckenschmidt. 2023. Comparing apples and oranges? on the evaluation of methods for temporal knowledge graph forecasting. In Machine Learning and Knowledge Discovery in

*Databases: Research Track*, pages 533–549, Cham. Springer Nature Switzerland.

314

315

316

317

318

319

321

322

323

324

325

327

329

330

331

332

333

334

335

336

337

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

360

361

362

363

364

365

366

367

368

369

- Zhen Han, Zifeng Ding, Yunpu Ma, Yujia Gu, and Volker Tresp. 2021. Learning neural ordinary equations for forecasting future links on temporal knowledge graphs. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8352–8364, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Woojeong Jin, Meng Qu, Xisen Jin, and Xiang Ren. 2020. Recurrent event network: Autoregressive structure inferenceover temporal knowledge graphs. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 6669–6683, Online. Association for Computational Linguistics.
- Julien Leblay and Melisachew Wudage Chekol. 2018. Deriving validity time in knowledge graph. In *Companion Proceedings of the The Web Conference 2018*, WWW '18, page 1771–1776, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.
- Dong-Ho Lee, Kian Ahrabian, Woojeong Jin, Fred Morstatter, and Jay Pujara. 2023. Temporal knowledge graph forecasting without knowledge using incontext learning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 544–557, Singapore. Association for Computational Linguistics.
- Kalev Leetaru and Philip A Schrodt. 2013. Gdelt: Global data on events, location, and tone, 1979–2012. In *ISA annual convention*, volume 2, pages 1–49. Citeseer.
- Zixuan Li, Saiping Guan, Xiaolong Jin, Weihua Peng, Yajuan Lyu, Yong Zhu, Long Bai, Wei Li, Jiafeng Guo, and Xueqi Cheng. 2022. Complex evolutional pattern learning for temporal knowledge graph reasoning. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics* (Volume 2: Short Papers), pages 290–296, Dublin, Ireland. Association for Computational Linguistics.
- Zixuan Li, Xiaolong Jin, Wei Li, Saiping Guan, Jiafeng Guo, Huawei Shen, Yuanzhuo Wang, and Xueqi Cheng. 2021. Temporal knowledge graph reasoning based on evolutional representation learning. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '21, page 408–417, New York, NY, USA. Association for Computing Machinery.
- Yushan Liu, Yunpu Ma, Marcel Hildebrandt, Mitchell Joblin, and Volker Tresp. 2022. Tlogic: Temporal logical rules for explainable link forecasting on temporal knowledge graphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(4):4120– 4127.

371

Suvir Mirchandani, Fei Xia, Pete Florence, brian ichter,

Danny Driess, Montserrat Gonzalez Arenas, Kan-

ishka Rao, Dorsa Sadigh, and Andy Zeng. 2023.

Large language models as general pattern machines.

In 7th Annual Conference on Robot Learning.

Adam Paszke, Sam Gross, Francisco Massa, Adam

Lerer, James Bradbury, Gregory Chanan, Trevor

Killeen, Zeming Lin, Natalia Gimelshein, Luca

Antiga, Alban Desmaison, Andreas Köpf, Edward

Yang, Zachary DeVito, Martin Raison, Alykhan Te-

jani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang,

Junjie Bai, and Soumith Chintala. 2019. Pytorch: An

imperative style, high-performance deep learning li-

brary. In Advances in Neural Information Processing

Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, De-

cember 8-14, 2019, Vancouver, BC, Canada, pages

Thomas Rebele, Fabian Suchanek, Johannes Hoffart,

Joanna Biega, Erdal Kuzey, and Gerhard Weikum.

2016. Yago: A multilingual knowledge base from wikipedia, wordnet, and geonames. In *The Semantic* 

Web – ISWC 2016: 15th International Semantic Web Conference, Kobe, Japan, October 17–21, 2016, Proceedings, Part II, page 177–185, Berlin, Heidelberg.

Haohai Sun, Jialun Zhong, Yunpu Ma, Zhen Han, and Kun He. 2021. TimeTraveler: Reinforcement learning for temporal knowledge graph forecasting. In Proceedings of the 2021 Conference on Empirical

Methods in Natural Language Processing, pages

8306-8319, Online and Punta Cana, Dominican Re-

public. Association for Computational Linguistics.

Yongqi Tong, Sizhe Wang, Dawei Li, Yifan Wang,

Simeng Han, Zi Lin, Chengsong Huang, Jiaxin

Huang, and Jingbo Shang. 2024. Optimizing lan-

guage model's reasoning abilities with weak supervi-

Nikos Voskarides, Edgar Meij, Ridho Reinanda, Abhinav Khaitan, Miles Osborne, Giorgio Stefanoni,

Prabhanjan Kambadur, and Maarten de Rijke. 2018.

Weakly-supervised contextualization of knowledge graph facts. In *The 41st International ACM SIGIR* 

Conference on Research & Development in Informa-

tion Retrieval, SIGIR '18, page 765–774, New York,

NY, USA. Association for Computing Machinery.

Denny Vrandečić and Markus Krötzsch. 2014. Wiki-

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien

Chaumond, Clement Delangue, Anthony Moi, Pier-

ric Cistac, Tim Rault, Remi Louf, Morgan Funtow-

icz, Joe Davison, Sam Shleifer, Patrick von Platen,

Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame,

Quentin Lhoest, and Alexander Rush. 2020. Trans-

ACM, 57(10):78-85.

data: a free collaborative knowledgebase. Commun.

sion. arXiv preprint arXiv:2405.04086.

8024-8035.

Springer-Verlag.

- 37
- 37

37

376 377

- 3
- 31
- 3
- 3
- 3

3

- 3
- 3
- 3
- 3
- 3
- 39

399 400

401 402

403

404 405 406

407 408

409 410 411

412 413

414 415

- 416
- 417

418 419

420 421

422 423 424

427formers: State-of-the-art natural language processing.428In Proceedings of the 2020 Conference on Empirical

Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics. 429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

- Yuwei Xia, Ding Wang, Qiang Liu, Liang Wang, Shu Wu, and Xiaoyu Zhang. 2024. Enhancing temporal knowledge graph forecasting with large language models via chain-of-history reasoning. *arXiv preprint arXiv:2402.14382*.
- Tianyi Zhang, Linrong Cai, Jeffrey Li, Nicholas Roberts, Neel Guha, and Frederic Sala. 2024. Stronger than you think: Benchmarking weak supervision on realistic tasks. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*
- Cunchao Zhu, Muhao Chen, Changjun Fan, Guangquan Cheng, and Yan Zhang. 2021. Learning from history: Modeling temporal knowledge graphs with sequential copy-generation networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(5):4732–4740.

## A Full Experimental Results

				ICE	WS14		ICEWS18								
Single-step	Train		H@1	H@1		H@3			H@1			H@3			
		2	3	4	2	3	4	2	3	4	2	3	4		
RE-GCN	1	45.2	36.8	10.6	66.0	55.2	23.2	35.0	32.9	16.0	55.2	51.7	30.2		
TANGO	1	39.4	34.2	8.3	57.2	49.0	18.0	30.2	29.3	12.8	49.5	47.0	25.5		
CEN	1	46.1	36.9	9.8	65.8	54.6	23.0	34.4	32.3	14.7	54.6	50.9	28.3		
Average		43.6	36.0	9.6	63.0	52.9	21.4	33.2	31.5	14.5	53.1	49.9	28.0		
gpt-neox-20b-entity	×	49.2	35.3	7.8	68.1	50.8	16.8	30.3	27.3	12.9	48.6	43.7	22.2		
$\Delta$ Average		5.6	-0.7	-1.7	5.1	-2.1	-4.6	-2.9	-4.2	-1.6	-4.5	-6.1	-5.7		
$\Delta$ Median		3.9	-1.5	-2.0	2.3	-3.8	-6.3	-4.1	-5.0	-1.8	-6.0	-7.1	-6.0		
gpt-neox-20b-pair	×	41.6	34.5	6.8	61.3	48.3	9.9	31.6	30.4	12.0	49.4	46.9	19.3		
$\Delta$ Average		2.1	-1.4	-2.8	-1.6	-4.7	-11.5	-1.7	-1.1	-2.5	-3.8	-3.0	-8.7		
$\Delta$ Median		0.4	-2.2	-3.1	-4.4	-6.3	-13.1	-2.9	-1.9	-2.7	-5.3	-4.0	-9.0		
				ICEV	WS14					ICE	WS18				
Multi-step	Train		H@1	ICEV	WS14	H@3			H@1	ICE	WS18	H@3			
Multi-step	Train	2	H@1 3	<b>ICEN</b>	<b>WS14</b>	H@3 3	4	2	H@1 3	<b>ICE</b>	WS18	H@3 3	4		
Multi-step	Train	2 40.0	H@1 3 34.1	<b>ICEN</b> 4 9.9	WS14 2 57.3	<b>H@3</b> 3 49.2	4 20.2	2 29.3	H@1 3 28.7	<b>ICE</b> 4 13.0	<b>WS18</b> 2 48.7	H@3 3 46.6	4		
Multi-step RE-NET RE-GCN	Train ✓ ✓	2 40.0 38.2	H@1 3 34.1 35.0	<b>ICEN</b> 4 9.9 10.6	<b>WS14</b> 2 57.3 56.4	H@3 3 49.2 50.5	4 20.2 20.8	29.3 30.1	H@1 3 28.7 30.1	ICE 4 13.0 13.7	WS18 2 48.7 49.0	H@3 3 46.6 47.5	4 25.6 26.7		
RE-NET RE-GCN CyGNet	Train ✓ ✓ ✓	2 40.0 38.2 37.9	H@1 3 34.1 35.0 33.3	4 9.9 10.6 8.3	<b>WS14</b> 2 57.3 56.4 56.2	H@3 3 49.2 50.5 49.6	4 20.2 20.8 17.0	2 29.3 30.1 26.0	H@1 3 28.7 30.1 26.7	<b>ICE</b> 4 13.0 13.7 11.4	WS18 2 48.7 49.0 45.4	H@3 3 46.6 47.5 45.6	4 25.6 26.7 23.7		
Multi-step RE-NET RE-GCN CyGNet Average	Train ✓ ✓ ✓	40.0 38.2 37.9 38.7	H@1 3 34.1 35.0 33.3 34.1	4 9.9 10.6 8.3 9.6	<b>WS14</b> 2 57.3 56.4 56.2 56.6	H@3 3 49.2 50.5 49.6 49.8	4 20.2 20.8 17.0 19.3	2 29.3 30.1 26.0 28.5	H@1 3 28.7 30.1 26.7 28.5	<b>ICE</b> 4 13.0 13.7 11.4 12.7	WS18 2 48.7 49.0 45.4 47.7	H@3 3 46.6 47.5 45.6 46.6	4 25.6 26.7 23.7 25.3		
Multi-step RE-NET RE-GCN CyGNet Average gpt-neox-20b-entity	Train	2 40.0 38.2 37.9 38.7 37.1	H@1 3 34.1 35.0 33.3 34.1 29.4	4 9.9 10.6 8.3 9.6 6.4	WS14 2 57.3 56.4 56.2 56.6 52.5	H@3 3 49.2 50.5 49.6 49.8 41.9	4 20.2 20.8 17.0 19.3 13.6	29.3 30.1 26.0 28.5 20.2	H@1 3 28.7 30.1 26.7 28.5 19.7	<b>ICE</b> 4 13.0 13.7 11.4 12.7 8.2	WS18 2 48.7 49.0 45.4 47.7 31.9	H@3 3 46.6 47.5 45.6 46.6 31.2	4 25.6 26.7 23.7 25.3 15.6		
Multi-step RE-NET RE-GCN CyGNet Average gpt-neox-20b-entity $\Delta$ Average	Train	2 40.0 38.2 37.9 38.7 37.1 -1.5	H@1 3 34.1 35.0 33.3 34.1 29.4 -4.7	4 9.9 10.6 8.3 9.6 6.4 -3.2	WS14 2 57.3 56.4 56.2 56.6 52.5 -4.1	H@3 3 49.2 50.5 49.6 49.8 41.9 -7.8	4 20.2 20.8 17.0 19.3 13.6 -5.7	2 29.3 30.1 26.0 28.5 20.2 -8.2	H@1 3 28.7 30.1 26.7 28.5 19.7 -8.7	4 13.0 13.7 11.4 12.7 8.2 -4.5	WS18 2 48.7 49.0 45.4 47.7 31.9 -15.8	H@3 3 46.6 47.5 45.6 46.6 31.2 -15.4	4 25.6 26.7 23.7 25.3 15.6 -9.7		
Multi-step RE-NET RE-GCN CyGNet Average gpt-neox-20b-entity $\Delta$ Average $\Delta$ Median	Train	2 40.0 38.2 37.9 38.7 37.1 -1.5 -1.1	H@1 3 34.1 35.0 33.3 34.1 29.4 -4.7 -4.7	4 9.9 10.6 8.3 9.6 6.4 -3.2 -3.5	WS14 2 57.3 56.4 56.2 56.6 52.5 -4.1 -3.9	H@3 3 49.2 50.5 49.6 49.8 41.9 -7.8 -7.7	4 20.2 20.8 17.0 19.3 13.6 -5.7 -6.6	2 29.3 30.1 26.0 28.5 20.2 -8.2 -9.1	H@1 3 28.7 30.1 26.7 28.5 19.7 -8.7 -8.9	ICE           4           13.0           13.7           11.4           12.7           8.2           -4.5           -4.8	WS18 2 48.7 49.0 45.4 47.7 31.9 -15.8 -16.8	H@3 3 46.6 47.5 45.6 46.6 31.2 -15.4 -15.4	4 25.6 26.7 23.7 25.3 15.6 -9.7 -10.0		
Multi-step RE-NET RE-GCN CyGNet Average gpt-neox-20b-entity <u>A Average</u> <u>A Median</u> gpt-neox-20b-pair	Train	2 40.0 38.2 37.9 38.7 37.1 -1.5 -1.1 34.1	H@1           3           34.1           35.0           33.3           34.1           29.4           -4.7           -4.7           30.2	ICEN           4           9.9           10.6           8.3           9.6           6.4           -3.2           -3.5           5.8	WS14 2 57.3 56.4 56.2 56.6 52.5 -4.1 -3.9 46.7	H@3           3           49.2           50.5           49.6           49.8           41.9           -7.8           -7.7           41.4	4 20.2 20.8 17.0 19.3 13.6 -5.7 -6.6 8.4	2 29.3 30.1 26.0 28.5 20.2 -8.2 -9.1 24.4	H@1 3 28.7 30.1 26.7 28.5 19.7 -8.7 -8.9 25.4	ICE           4           13.0           13.7           11.4           12.7           8.2           -4.5           -4.8           8.8	WS18 2 48.7 49.0 45.4 47.7 31.9 -15.8 -16.8 38.7	H@3 3 46.6 47.5 45.6 46.6 31.2 -15.4 -15.4 38.4	4 25.6 26.7 23.7 25.3 15.6 -9.7 -10.0 14.4		
Multi-step RE-NET RE-GCN CyGNet Average gpt-neox-20b-entity <u>A Average</u> <u>A Median</u> gpt-neox-20b-pair <u>A Average</u>	Train ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓ ✓	2 40.0 38.2 37.9 38.7 37.1 -1.5 -1.1 34.1 -4.6	H@1           3           34.1           35.0           33.3           34.1           29.4           -4.7           -4.7           -4.7           -4.7           -4.7           -4.7	ICEN           4           9.9           10.6           8.3           9.6           6.4           -3.2           -3.5           5.8           -3.8	WS14 2 57.3 56.4 56.2 56.6 52.5 -4.1 -3.9 46.7 -9.9	H@3 3 49.2 50.5 49.6 49.8 41.9 -7.8 -7.7 41.4 -8.4	4 20.2 20.8 17.0 19.3 13.6 -5.7 -6.6 8.4 -10.9	2 29.3 30.1 26.0 28.5 20.2 -8.2 -9.1 24.4 -4.0	H@1 3 28.7 30.1 26.7 28.5 19.7 -8.7 -8.9 25.4 -3.1	4 13.0 13.7 11.4 12.7 8.2 -4.5 -4.8 8.8 -3.9	WS18 2 48.7 49.0 45.4 47.7 31.9 -15.8 -16.8 38.7 -9.1	H@3 3 46.6 47.5 45.6 46.6 31.2 -15.4 -15.4 -15.4 38.4 -8.2	4 25.6 26.7 23.7 25.3 15.6 -9.7 -10.0 14.4 -10.9		

Table 4: **Performance** (**Hits**@**K**) comparison between supervised models and ICL for single-step (top) and multistep (bottom) prediction, grouped by the number of **number of unique entities** as confounder. The first group consists of supervised models, whereas the second group consists of ICL models, *i.e.*, GPT-NeoX with a history length of 100. The green and red colors represent where LLM is outperforming and underperforming the average performance of the supervised models.



Figure 2: **Hits@1 difference** between the average performance of ICL and the average performance of supervised models, grouped by the **relation frequency** confounder, for single-step (top) and multi-step (bottom) prediction.

<i>.</i>		ICEWS14									ICEWS18								
Single-step	Train		Н	@1			Н	@3			H	@1		H@3					
		1	2	3	4	1	2	3	4	1	2	3	4	1	2	3	4		
RE-GCN TANGO CEN Average Median gpt-neox-20b-entity <b>A Average</b>	√ √ √ ×	49.5 45.5 50.1 48.4 49.5 57.0 8.6	35.7 29.2 36.6 33.8 35.7 36.3 2.4	34.4 31.8 34.5 33.5 34.4 35.1 1.6	40.3 38.0 40.3 39.6 40.3 35.8 -3.7	71.7 64.6 70.2 68.9 70.2 77.0 8.1	54.8 45.4 54.9 51.7 54.8 53.1 1.4	53.0 46.9 52.4 50.8 52.4 50.0 -0.8	55.3 49.7 56.3 53.8 55.3 51.2 -2.6	34.6 30.9 33.8 33.1 33.8 30.1 -3.0	31.9 26.9 31.3 30.0 31.3 28.2 -1.8	30.8 27.6 30.1 29.5 30.1 24.2 -5.2	28.4 24.8 27.9 27.0 27.9 22.4 -4.6	55.0 50.8 53.9 53.3 53.9 48.2 -5.1	51.2 45.1 50.4 48.9 50.4 45.6 -3.3	48.6 44.1 47.7 46.8 47.7 39.1 -7.7	45.6 41.1 44.9 43.9 44.9 35.7 -8.2		
$\begin{tabular}{lllllllllllllllllllllllllllllllllll$	X	7.4 58.5 10.1 9.0	0.5 29.5 -4.3 -6.2	0.8 34.4 0.9 0.1	-4.5 32.2 -7.4 -8.1	6.7 82.0 13.1 11.8	-1.7 40.8 -10.9 -14.0	-2.4 48.0 -2.8 -4.4	-4.1 41.4 -12.4 -13.9	-3.7 35.8 2.7 2.0	-3.1 27.9 -2.1 -3.3	-5.8 26.8 -2.7 -3.3	-5.5 22.6 -4.4 -5.2	-5.8 56.6 3.3 2.6	-4.7 44.1 -4.8 -6.3	-8.6 41.1 -5.7 -6.6	-9.3 33.4 -10.5 -11.5		
Multi-step	Train		H	91	ICE	WS14	H@	93			H@	1	ICE	WS18	H	@3			
Multi-step	Train	1	<b>H</b> @	@1 3	<b>ICE</b>	WS14	<b>H</b> @	<b>3</b> 3	4	1	H@	2 <b>1</b> 3	<b>ICE</b>	WS18	<b>H</b> (	@ <b>3</b> 3	4		
Multi-step RE-NET RE-GCN CyGNet Average Median	Train ✓ ✓ ✓	1 45.5 42.5 46.6 44.8 45.5	H0 2 30.8 31.2 27.1 29.7 30.8	<b>3</b> <b>3</b> <b>3</b> <b>1</b> .7 <b>3</b> <b>1</b> .7 <b>3</b> <b>1</b> .4 <b>3</b> <b>1</b> .6 <b>3</b> <b>1</b> .7	ICE 4 37.0 35.8 34.2 35.7 35.8	WS14 1 62.0 63.6 66.2 63.9 63.6	<b>H</b> @ 2 46.6 47.6 43.3 45.8 46.6	<b>3</b> <b>47.6</b> <b>46.1</b> <b>46.9</b> <b>46.9</b> <b>46.9</b> <b>46.9</b> <b>46.9</b>	4 51.7 48.5 49.4 49.9 49.4	1 30.4 31.3 28.7 30.1 30.4	H@ 2 26.1 26.9 22.7 25.2 26.1	21 3 26.8 28.6 25.2 26.8 26.8 26.8	ICE 4 24.7 25.0 21.6 23.7 24.7	WS18 1 50.8 50.6 49.0 50.1 50.6	<b>H</b> (0) 2 44.0 44.8 41.1 43.3 44.0	<b>3</b> <b>44.1</b> <b>45.0</b> <b>42.8</b> <b>43.9</b> <b>44.1</b>	4 40.8 41.0 37.8 39.9 40.8		
Multi-step RE-NET RE-GCN CyGNet Average Median gpt-neox-20b-entity $\Delta$ Average $\Delta$ Median	Train ✓ ✓ ✓ ×	1           45.5           42.5           46.6           44.8           45.5           41.5           -3.3           -3.9	H0 2 30.8 31.2 27.1 29.7 30.8 28.0 -1.7 -2.8	<b>2</b> 1 <b>3</b> <b>3</b> 1.7 <b>3</b> 1.7 <b>3</b> 1.4 <b>3</b> 1.6 <b>3</b> 1.7 <b>2</b> 8.7 <b>-</b> 2.9 <b>-</b> 3.0	ICE 4 37.0 35.8 34.2 35.7 35.8 31.1 -4.6 -4.7	WS14 1 62.0 63.6 66.2 63.9 63.6 56.9 -7.0 -6.7	H@ 2 46.6 47.6 43.3 45.8 46.6 41.4 -4.4 -5.2	<b>3</b> <b>47.6</b> <b>46.1</b> <b>46.9</b> <b>46.9</b> <b>46.9</b> <b>46.9</b> <b>46.9</b> <b>41.3</b> <b>-5.6</b> <b>-5.6</b>	4 51.7 48.5 49.4 49.9 49.4 44.0 -5.8 -5.4	30.4           31.3           28.7           30.1           30.4           21.9           -8.2           -8.5	H@           2           26.1           26.9           22.7           25.2           26.1           18.1           -7.2           -8.0	3           26.8           28.6           25.2           26.8           17.9           -9.0           -8.9	ICE 4 24.7 25.0 21.6 23.7 24.7 15.5 -8.2 -9.1	WS18 	H0 2 44.0 44.8 41.1 43.3 44.0 29.9 -13.4 -14.1	23         3         44.1         45.0         42.8         43.9         44.1         28.4         -15.5         -15.6	4 40.8 41.0 37.8 39.9 40.8 24.6 -15.3 -16.2		

Table 5: **Performance** (**Hits**@**K**) comparison between supervised models and ICL for single-step (top) and multistep (bottom) prediction, grouped by the number of **number of unique relations** as confounder. The first group consists of supervised models, whereas the second group consists of ICL models, *i.e.*, GPT-NeoX with a history length of 100. The green and red colors represent where LLM is outperforming and underperforming the average performance of the supervised models.



Figure 3: **Hits@1 difference** between the average performance of ICL and the average performance of supervised models, grouped by the **time interval** confounder, for single-step (top) and multi-step (bottom) prediction.