# Software Engineering Event Modeling using Relative Time in Temporal Knowledge Graphs

Kian Ahrabian<sup>1</sup> Daniel Tarlow<sup>21</sup> Hehuimin Cheng<sup>1</sup> Jin L.C. Guo<sup>1</sup>

#### Abstract

We present a multi-relational temporal Knowledge Graph based on the daily interactions between artifacts in GitHub, one of the largest social coding platforms. Such representation enables posing many user-activity and project management questions as link prediction and time queries over the knowledge graph. In particular, we introduce two new datasets for i) interpolated time-conditioned link prediction and ii) extrapolated time-conditioned link/time prediction queries, each with distinguished properties. Our experiments on these datasets highlight the potential of adapting knowledge graphs to answer broad software engineering questions. Meanwhile, it also reveals the unsatisfactory performance of existing temporal models on extrapolated queries and time prediction queries in general. To overcome these shortcomings, we introduce an extension to current temporal models using relative temporal information with regards to past events.

# 1. Introduction

Hosting over 100 million repositories, GitHub (GH) is one of the biggest social coding platforms (GitHub, 2018). Over the past decade, the available artifacts hosted on GH have become one of the most important resources for software engineering (SE) researchers to study various aspects of programming, software development, the characteristics of open source users and ecosystem (Cosentino et al., 2017). Example questions of interest include when an issue will be closed (Kikas et al., 2016; Rees-Jones et al., 2017), how likely a pull request will be merged and when (Gousios & Zaidman, 2014; Soares et al., 2015), and who should review a pull request (Yu et al., 2016; Hannebauer et al., 2016). Our aim in this work is to connect the above SE research questions to the literature on learning knowledge graph (KG) embeddings, with a particular emphasis on temporal KGs due to the importance of the temporal component in the above questions. Methods for time prediction and timeconditioned link prediction in KGs (Kazemi et al., 2020, Section 5.1-5.2) are generally based on point process models (Trivedi et al., 2017; 2019; Knyazev et al., 2019) or adaptations of KG embeddings that additionally use time to compute scores (Dasgupta et al., 2018; Leblay & Chekol, 2018; García-Durán et al., 2018; Goel et al., 2019). While point processes are elegant, they are more challenging to work with and require strong assumptions on the underlying intensity functions (see e.g., Trivedi et al., 2019, Equation 1 & Section 4). Thus we focus on KG embedding-based methods, in particular starting with Diachronic Embeddings (DE-X) (Goel et al., 2019) for time-varying embeddings and RotatE (Sun et al., 2019) for scoring facts.

Our contributions are the following:

- Collecting two new temporal KG datasets<sup>1</sup> from GH public events that allow casting the above SE questions as time-conditioned prediction queries (Section 2).
- Benchmarking existing temporal KG embedding methods on the new datasets (Section 4).
- Based on the observation that existing temporal KG embedding models do not well capture patterns in relative time that are important to SE applications, particularly in extrapolated settings, e.g., *How long will it take from pull request being opened to being closed?*, we propose a new relative temporal KG embedding inspired by the use of *relative time* in attention-based neural networks like Music Transformer (Huang et al., 2019) and Transformer-XL (Dai et al., 2019) (Section 3).

In total, our work brings together SE research questions and temporal KG models by introducing new datasets, benchmarking existing recent methods on the datasets, and suggesting a new direction of leveraging Transformer-style relative time modeling in KG embeddings.

<sup>&</sup>lt;sup>1</sup>School of Computer Science, McGill University, Montreal, Canada <sup>2</sup>Google Research, Brain Team, Montreal, Canada. Correspondence to: Kian Ahrabian <kian.ahrabian@mail.mcgill.ca>.

Proceedings of the  $37^{th}$  International Conference on Machine Learning, Vienna, Austria, PMLR 108, 2020. Copyright 2020 by the author(s).

<sup>&</sup>lt;sup>1</sup>https://zenodo.org/record/3928580

DATASET	$ \mathbf{V} $	$ \mathbf{E} $	$ \mathbf{R} $	$ \mathbf{T} $	$\mathbf{D}_{\mathbf{MAX}}$	$\mathbf{D}_{\mathbf{MED}}$
GITHUB-SE 1Y	125,455,982	249,375,075	19	365	3,519,105	2
GITHUB-SE 1Y-REPO	133,335	285,788	18	365	73,345	1
GITHUB-SE 1M	13,690,824	23,124,510	19	31	1,324,179	2
GITHUB-SE 1M-NODE	139,804	293,014	14	31	639	2
GITHUB (ORIGINAL) (TRIVEDI ET AL., 2019)	12,328	771,214	3	366	-	-
GITHUB (SUBNETWORK) (KNYAZEV ET AL., 2019)	284	20,726	8	366	4,790	53.5
ICEWS14 (GARCÍA-DURÁN ET AL., 2018)	7,128	90,730	230	365	6,083	3
ICEWS05-15 (GARCÍA-DURÁN ET AL., 2018)	10,488	461,329	251	4017	52,890	5
YAGO15K (GARCÍA-DURÁN ET AL., 2018)	15,403	138,056	34	198	6,611	5
WIKIDATA (LEBLAY & CHEKOL, 2018)	11,153	150,079	96	328	586	5
GDELT (LEETARU & SCHRODT, 2013)	500	3,419,607	20	366	53,857	10,336

Table 1. Characteristics comparison of the introduced datasets to existing temporal KG datasets.

### 2. Dataset

**Creation** To create the dataset, we retrieved from GH Archive<sup>2</sup> all of the raw public events in GH in 2019. The knowledge graph was then constructed by tuples, each of which represents an individual event containing temporal information based on its type and a predefined set of extraction rules. The properties of the constructed KG is shown in the first row of Table 1, referred to as GITHUB-SE 1Y.

Due to the substantial size of the GITHUB-SE 1Y and unrealistic computational demands of training KG embedding models on this KG, we sampled the GITHUB-SE 1Y using two distinct strategies, described in more details in Appendix A.1. The first strategy aims to retain maximum temporal information about particular SE projects. To achieve this, first, an induced sub-graph containing all related nodes was extracted for each node with type Repository. Then, for each sub-graph a popularity score was calculated as  $P(G) = W_1 \times S_G + W_2 \times T_G$  where  $S_G$  is the size of the graph,  $T_G$  is the time-span of the graph, and  $W_1, W_2 \in \mathbb{R}^+$ are weight values. Finally, from the top three ranked repositories, we selected the Visual Studio Code repository to extract a one-year slice as it exercised more functionalities related to the target entities in this work, i.e. issues and pull requests. We name this dataset GITHUB-SE 1Y-REPO due to its repository-centric characteristics.

The second strategy aims at preserving the most informative nodes regardless of their type. We used a variation of Snowball Sampling (Goodman, 1961) on all the events in December 2019. This sampled dataset, i.e. GITHUB-SE 1M-NODE, captures events across various projects and therefore, can be used to answer queries such as which project does a user start contributing at a certain time.

**Characteristics** In Table 1, we compare the variations of GITHUB-SE KG proposed in this work with commonly

used datasets in the literature. Even the sampled down versions of our datasets are considerably larger in terms number of nodes. They have much higher edge to node ratios which translates into sparsity in graphs, but this sparsity level is close to what appears in GitHub as a whole. Additionally, similar to relations, each node in our datasets is also typed.

Trivedi et al. (2019) also collects a temporal KG dataset from GitHub. However, this dataset is exclusively focused on the social aspects of GitHub, discarding repositories and only including user-user interactions, and it does not appear to be publicly available beyond raw data and a small subnetwork extracted in a follow-up work (Knyazev et al., 2019). To differentiate our datasets, which focus on the SE aspects of GitHub, we append – SE to the dataset names.

The distinguishing characteristics of the proposed datasets, i.e. size, sparsity, node-typing, diversity, focus on SE aspects, and temporal nature, introduce a variety of engineering and theoretical challenges that make these datasets a suitable choice for exploring and exposing the limitations of temporal knowledge graph embedding models.

## 3. Method

### 3.1. Existing KG embedding models

We first examine the performance of the state-of-the-art KG embedding models on the GITHUB-SE 1M-NODE and GITHUB-SE 1Y-REPO datasets. We select RotatE (Sun et al., 2019) for the static settings considering its ability to infer *Symmetry*, *Antisymmetry*, *Inversion*, and *Composition* relational patterns. Moreover, we use DE-X (Goel et al., 2019) for the dynamic setting due to its superior performance on existing benchmarks and the fact that for any static model X there exists an equivalent DE-X model ensuring the ability to learn aforementioned patterns.

Notationally, a KG G is a set of tuples of the form x = (s, r, o, t) respectively representing the subject, relation, ob-

<sup>&</sup>lt;sup>2</sup>https://www.gharchive.org

ject, and timestamp. The diachronic embeddings are defined as

$$D(e,t) = E(s) \oplus (E_A(e) + sin(t * E_F(e) + E_\phi(e)))$$

where  $E, E_A, E_F, E_{\phi}$  are embedding lookup tables and the last three respectively represent *Amplitude*, *Frequency*, and *Phase* of a sinusoid. Similar to Goel et al. (2019), whenever t consists of multiple set of numbers rather than one, e.g. year, month, and day, for each set of numbers a separate D is defined, and the values are summed up. Subsequently, the scoring function of the DE-RotatE model is defined as

$$score(s, r, o, t) = D(s, t) \circ E_R(r) - D(o, t)$$

where  $E_R$  is an embedding lookup table for relations.

#### 3.2. Relative Temporal Context

The idea of using relative temporal information has been successfully employed in natural language processing (Vaswani et al., 2017; Dai et al., 2019) and music generation (Huang et al., 2019). These models formalize the intuition that temporal spacing between events is more central than the absolute time at which an event happened. We believe this framing is also appropriate for SE applications of temporal knowledge graphs: to predict if a pull request is closed at time t, it is more important to know how long it has been since the pull request was opened than it is to know t.

A challenge is that there are a lot of events, and we do not want to hard-code which durations are relevant. Instead, we would like the model to learn which temporal durations are important for scoring a temporal fact. As the number of related facts to an entity could be as high as a few thousand, we propose to pick a fixed-number of facts as temporal context to provide as an input to the models.

Let  $\mathcal{H}(e, r, t) = \{t' \mid (s', r', o', t') \in G \land (t' < t) \land (r' = r) \land (s' = e \lor o' = e)\}$  be the set of times associated with facts involving entity e and relation r occurring before time t, and let  $\delta(e, r, t) = t - \max \mathcal{H}(e, r, t)$  be the relative time since a fact involving e and relation r has occurred. Hence, an entity's *relative temporal context* at query time  $t_q$  is  $\Delta(e, t_q) = [\delta(e, 1, t_q), \dots, \delta(e, |R|, t_q)]^\top \in \mathbb{R}^{|R|}$ .

#### 3.3. Relative Time DE-RotatE (RT-DE-ROTATE)

We now turn attention to using the relative temporal context  $\Delta(e, t)$  as an input to temporal KG embeddings. Our inspiration is the *Transformer* encoder, which has emerged as a successful substitute to more traditional Recurrent Neural Network approaches used for sequential (Vaswani et al., 2017; Dai et al., 2019; Huang et al., 2019) and structural (Parmar et al., 2018) tasks. The core idea is to employ a variation of attention mechanism called *Self-Attention* that assigns importance scores to the elements of the same sequence. Unlike recurrence mechanism, the positional information is injected to the Transformer styled encoders by 1) adding sine/cosine functions of different frequencies to the input (Vaswani et al., 2017; Dai et al., 2019), or 2) directly infusing relative distance information to attention computation in form of a matrix addition (Shaw et al., 2018; Huang et al., 2019). Vaswani et al. (2017) introduced a positional encoding scheme in form of sinusoidal vectors defined as  $\rho(i) =$  $[\rho_1(i), \ldots, \rho_d(i)]$ , where  $\rho_j(i) = \sin(i/10000^{\frac{\lfloor j/2 \rfloor}{d}})$  if j is even and  $\cos(i/10000^{\frac{\lfloor j/2 \rfloor}{d}})$  if j is odd, i is the absolute position, and d is the embedding dimension. In the followup *Transformer-XL* model, Dai et al. (2019) introduce a reparameterization of the relative attention where the attention score between a query element at position i and a key element at position j is defined as

$$A_{i,j}^{rel} = \underbrace{E(i)^\top W_Q^\top W_{K,E} E(j)}_{(a)} + \underbrace{E(i)^\top W_Q^\top W_{K,\rho} \rho(i-j)}_{(b)} + \underbrace{u^\top W_{K,E} E(j)}_{(c)} + \underbrace{v^\top W_{K,\rho} \rho(i-j)}_{(d)}$$

where  $E(i), E(j) \in \mathbb{R}^{d \times 1}$  are the *i* and *j* element embeddings,  $u, v \in \mathbb{R}^{d \times 1}, W_Q, W_{K,E}, W_{K,\rho}$  are trainable  $d \times d$  matrices, and i - j is the relative position between *i* and *j*.

The main difference in our setting is that the above models compute a score based on a single relative time i - j, while our relative temporal context  $\Delta$  contains |R| relative times for each entity. Our approach is to score a tuple (s, r, o, t)based on the information available at query time  $t_q$ .<sup>3</sup> For each entity  $e \in (s, o)$  we define a positional embeddings matrix P of relative times between t and the events in its relative temporal context  $\Delta(e, t_q)$  as

$$P(e,t,t_q) = \begin{bmatrix} \rho(t-t_q+\delta(e,1,t_q))\\ \rho(t-t_q+\delta(e,2,t_q))\\ \vdots\\ \rho(t-t_q+\delta(e,|R|,t_q)) \end{bmatrix} \in \mathbb{R}^{|R| \times d}.$$

Intuitively, these relative times encode "if the event happened at time t, how long would it have been since the events in the relative time context." A learned, relation-specific row vector  $W_P(r) \in \mathbb{R}^{1 \times |R|}$  for  $r = 1, \ldots, |R|$  chooses which rows of P are important, and then  $\gamma(r, e, t, t_q) =$  $W_P(r)P(e, t, t_q) \in \mathbb{R}^{1 \times d}$ , abbreviated  $\gamma(r, e, t)$ , embeds the relative temporal context of e, replacing  $W_{K,\rho}\rho(i-j)$ :

$$A_x^{rel} = \underbrace{D(s,t)W(r)D(o,t)^{\top}}_{(a)} + \underbrace{\gamma(r,s,t)W_P\gamma(r,o,t)^{\top}}_{(d)} + \underbrace{E(s)W_E\gamma(r,o,t)^{\top}}_{(b)} + \underbrace{\gamma(r,s,t)W_E^{\top}E(o)^{\top}}_{(c)}$$

<sup>3</sup>During training,  $t_q$  is the t of the positive sample, and during evaluation,  $t_q$  is set to the maximum timestamp in the training set.

DATASET	Туре	Model	HITS@1	HITS@3	HITS@10	$\mathbf{MR}$	MRR
		ROTATE	47.58	76.66	88.95	807.40	0.6328
	INTERPOLATED	DE-ROTATE	47.98	76.92	88.87	779.50	0.6349
GITHUB-SE		RT-DE-ROTATE (OURS)	49.70	78.67	90.48	773.90	0.6522
1M-NODE		ROTATE	25.40	49.02	57.54	4762.87	0.3797
	EXTRAPOLATED	DE-ROTATE	26.28	48.53	57.33	4840.16	0.3838
		RT-DE-ROTATE (OURS)	26.50	49.54	57.94	4891.81	0.3888
		ROTATE	44.05	57.14	80.95	18.54	0.5460
	INTERPOLATED	DE-ROTATE	42.17	53.88	76.88	24.67	0.5233
GITHUB-SE 1Y-REPO		RT-DE-ROTATE (OURS)	48.93	60.96	78.32	14.47	0.5815
		ROTATE	2.11	4.82	9.71	1917.03	0.0464
	EXTRAPOLATED	DE-ROTATE	1.77	4.08	9.10	1961.75	0.0402
		RT-DE-ROTATE (OURS)	38.25	40.08	64.06	1195.02	0.4345

Table 2. Performance comparison on time-conditioned Link Prediction. Results within the 95% confidence interval of the best are bolded.

where W(r) is a relation-specific weight matrix and  $W_E$ and  $W_P$  are tuple-agnostic weight matrices; however, this formulation is suitable for bilinear models. Hence, we derive a translational variation for the DE-RotatE model as

$$A_x^{rel} = \underbrace{\|D(s,t) \circ E_R(r) - D(o,t)\|}_{(a)} + \underbrace{\|E(s)W_E - \gamma(r,o,t)\|}_{(b)} + \underbrace{\|\gamma(r,s,t) - W_E E(o)\|}_{(c)}$$

where W(r) is replaced by an embedding lookup table  $E_R(r)$ . Intuitively, under this formulation (a) capture entities compatibility and (b) and (c) capture entity-specific temporal context compatibility. In comparison, the existing models only include term (a) discarding terms (b) and (c).

# 4. Experiments

**Datasets:** We use a 90%-5%-5% events split for constructing the train, validation, and test sets. For the interpolated queries the split was done randomly, whereas we split the data using event timestamps for the extrapolated queries. Table 3 in the Appendix presents details of the splits.

**Queries:** For time-conditioned link prediction, we selected events related to the resolution of Github issues and pull requests due to their direct impact on software development and maintenance practice. Particularly, we used "Who will close issue X at time T?" and 'Who will close pull-request X at time T?" for evaluation. For time prediction, we used the analogous time queries of the aforementioned queries for evaluation, i.e. "When will issue X be closed by user Y?" and 'When will pull-request X be closed by user Y?".

**Evaluation and Results:** We calculated the standard metrics to evaluate the model performance on the test set. For

the extrapolated time-conditioned link prediction queries, after using the validation set for hyperparameter tuning, we retrained the selected models using both training and validation sets for evaluation. We also report the model performance without retraining in the Appendix Table 5.

In Table 2 we compare the model performance on the timeconditioned link prediction queries. On the GITHUB-SE 1M-NODE queries, our model slightly outperforms existing models in some cases, but the difference is statistically insignificant in others. On the GITHUB-SE 1Y-REPO, on the other hand, our RT-DE-ROTATE model shows a significant performance boost, particularly on the extrapolated time-conditioned link prediction queries, indicating the importance of using relative time as temporal context.

For the extrapolated time prediction queries on GITHUB-SE 1Y-REPO dataset, our model performed slightly better on HITS@1, HITS@3, and Mean Reciprocal Rank (MRR) than the other existing models while marginally surpassing the random baseline on all metrics. These results, detailed in the Appendix Table 8, stress the necessity of having further studies on extrapolated time prediction queries.

### 5. Conclusion

In this work, we bridge between the SE domain questions and the literature on KG embedding models by introducing two new datasets based on the daily interactions in the GH platform and casting those questions as queries on an appropriate KG. Furthermore, we introduce RT-X, a novel extension to existing KG embedding models that make use of relative time with regards to past relevant events. Our experiments highlight shortcomings of existing temporal KG embedding models, notably on extrapolated time-conditioned link prediction, and exhibit the advantage of leveraging relative time as introduced in the RT-X model. In total, this work highlights new opportunities for improving temporal KG embedding models on time prediction queries.

# Acknowledgements

We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC). This research was enabled in part by support provided by Google, Calcul Québec, and Compute Canada. We thank Daniel Johnson for helpful comments.

# References

- Cosentino, V., Izquierdo, J. L. C., and Cabot, J. A systematic mapping study of software development with github. *IEEE Access*, 5:7173–7192, 2017.
- Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q. V., and Salakhutdinov, R. Transformer-xl: Attentive language models beyond a fixed-length context. arXiv preprint arXiv:1901.02860, 2019.
- Dasgupta, S. S., Ray, S. N., and Talukdar, P. Hyte: Hyperplane-based temporally aware knowledge graph embedding. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pp. 2001–2011, 2018.
- García-Durán, A., Dumančić, S., and Niepert, M. Learning sequence encoders for temporal knowledge graph completion. *arXiv preprint arXiv:1809.03202*, 2018.
- GitHub. Thank you for 100 million repositories. https: //github.blog/2018-11-08-100m-repos/, 2018. (Accessed on 03/01/2020).
- Goel, R., Kazemi, S. M., Brubaker, M., and Poupart, P. Diachronic embedding for temporal knowledge graph completion. *arXiv preprint arXiv:1907.03143*, 2019.
- Goodman, L. A. Snowball sampling. The annals of mathematical statistics, pp. 148–170, 1961.
- Gousios, G. and Zaidman, A. A dataset for pull-based development research. In *Proceedings of the 11th Working Conference on Mining Software Repositories*, pp. 368– 371, 2014.
- Hannebauer, C., Patalas, M., Stünkel, S., and Gruhn, V. Automatically recommending code reviewers based on their expertise: An empirical comparison. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*, pp. 99–110, 2016.
- Huang, C.-Z. A., Vaswani, A., Uszkoreit, J., Simon, I., Hawthorne, C., Shazeer, N., Dai, A. M., Hoffman, M. D., Dinculescu, M., and Eck, D. Music transformer. In *International Conference on Learning Representations*, 2019.

- Kazemi, S. M., Goel, R., Jain, K., Kobyzev, I., Sethi, A., Forsyth, P., and Poupart, P. Representation learning for dynamic graphs: A survey. *Journal of Machine Learning Research*, 21(70):1–73, 2020.
- Kikas, R., Dumas, M., and Pfahl, D. Using dynamic and contextual features to predict issue lifetime in github projects. In 2016 IEEE/ACM 13th Working Conference on Mining Software Repositories (MSR), pp. 291–302. IEEE, 2016.
- Knyazev, B., Augusta, C., and Taylor, G. W. Learning temporal attention in dynamic graphs with bilinear interactions. arXiv preprint arXiv:1909.10367, 2019.
- Lacroix, T., Usunier, N., and Obozinski, G. Canonical tensor decomposition for knowledge base completion. *arXiv preprint arXiv:1806.07297*, 2018.
- Leblay, J. and Chekol, M. W. Deriving validity time in knowledge graph. In *Companion Proceedings of the The Web Conference 2018*, WWW '18, pp. 1771–1776, Republic and Canton of Geneva, CHE, 2018. International World Wide Web Conferences Steering Committee. ISBN 9781450356404. doi: 10.1145/ 3184558.3191639. URL https://doi.org/10. 1145/3184558.3191639.
- Leetaru, K. and Schrodt, P. A. Gdelt: Global data on events, location, and tone. *ISA Annual Convention*, 2013.
- Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, Ł., Shazeer, N., Ku, A., and Tran, D. Image transformer. arXiv preprint arXiv:1802.05751, 2018.
- Rees-Jones, M., Martin, M., and Menzies, T. Better predictors for issue lifetime. arXiv preprint arXiv:1702.07735, 2017.
- Shaw, P., Uszkoreit, J., and Vaswani, A. Self-attention with relative position representations. *arXiv preprint arXiv:1803.02155*, 2018.
- Soares, D. M., de Lima Júnior, M. L., Murta, L., and Plastino, A. Acceptance factors of pull requests in opensource projects. In *Proceedings of the 30th Annual ACM Symposium on Applied Computing*, pp. 1541–1546, 2015.
- Sun, Z., Deng, Z.-H., Nie, J.-Y., and Tang, J. Rotate: Knowledge graph embedding by relational rotation in complex space. *arXiv preprint arXiv:1902.10197*, 2019.
- Trivedi, R., Dai, H., Wang, Y., and Song, L. Know-evolve: Deep temporal reasoning for dynamic knowledge graphs. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 3462–3471. JMLR. org, 2017.

- Trivedi, R., Farajtabar, M., Biswal, P., and Zha, H. Dyrep: Learning representations over dynamic graphs. In *Inter*national Conference on Learning Representations, 2019.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *Advances in neural information processing systems*, pp. 5998–6008, 2017.
- Yu, Y., Wang, H., Yin, G., and Wang, T. Reviewer recommendation for pull-requests in github: What can we learn from code review and bug assignment? *Information and Software Technology*, 74:204–218, 2016.

# A. Dataset

## A.1. Sampling

Algorithm 1 describes the snowball sampling used to create the GITHUB-SE 1M-NODE dataset. This algorithm aims at preserving the most informative nodes regardless of their types. Moreover, Algorithm 2 describes the temporal sampling used to create the GITHUB 1Y-REPO dataset. This algorithm aims at preserving maximum temporal information regarding particular repositories.

Algorithm 1 Snowball Sampling strategy used for extracting the GITHUB-SE 1M-NODE dataset.

**Require:** set of nodes V, set of edges E, sample size N, growth size S, initial sample size K $L \leftarrow sortDescending(V)$  w.r.t node degree  $Q \leftarrow MaxPriorityQueue()$ for i = 1, ..., K do  $Q.\mathsf{put}(L[i])$ end for  $U \leftarrow Set()$ while size(U) < N do  $V_u \leftarrow Q.top()$  w.r.t node degree  $U.put(V_u)$  $U_s \leftarrow randomSample(E[V_u])$  with size S for i = 1, ..., S do  $Q.\mathsf{put}(U_s[i])$ end for end while

Algorithm 2 Temporal Sampling strategy used for extracting the GITHUB 1Y-REPO dataset.

 $\begin{array}{l} \textbf{Require: set of nodes } V, \text{ size importance factor } W_1, \text{ time} \\ \text{span importance factor } W_2, \text{ sample size } N \\ R \leftarrow extractRelated(V) \left\{ Repository \text{ node type only} \right\} \\ P \leftarrow Array(|R|) \\ \textbf{for } i = 1, ..., |R| \ \textbf{do} \\ P_i \leftarrow calculatePopularity(R_i, W_1, W_2) \\ \textbf{end for} \\ S \leftarrow sorted(P) \\ U \leftarrow Set() \\ \textbf{for } i = 1, ..., |S| \ \textbf{do} \\ \textbf{if } size(U) < N \ \textbf{then} \\ U.union(S_i) \\ \textbf{end if} \\ \textbf{end for} \end{array}$ 

### A.2. Extraction

Table 9 presents the set of extraction rules used to build the KG from raw events each representing a relation type. Although 80 extractions rules are defined in Table 9, the raw events that we used only contained 18 of them. The codes presented in the Relation column of Table 9, when divided on underscore, are interpreted as *a*) the first and the last components respectively represent entity types of the event's subject and object, *b*) *AO*, *CO*, *SE*, *SO*, and *HS* are abbreviations of extracted information from raw payloads <sup>4</sup> serving as distinguishers between different relation types among entities, and *c*) the second to the last component represents the concrete action taken that triggers the event.

# **B.** Model

#### **B.1.** Complexity

Table 4 presents time and space complexity comparison between the existing models and the introduced RT-X model. Notice that, while yielding superior performance, the number of free-parameters introduced in our extension does not increase linearly with the number of entities which is one of the bottlenecks of training large KG embedding models.

### **B.2.** Loss Function

Similar to the self-adversarial negative sampling introduced in Sun et al. (2019), we use weighted negative samples as

$$p(h'_{j}, r, t'_{j} | \{(h_{i}, r, t_{i})\}) = \frac{\exp \eta f_{r}(h'_{j}, t'_{j})}{\sum_{i} \exp \eta f_{r}(h'_{i}, t'_{i})}$$

where  $\eta$  is the sampling temperature and  $(h'_i, r, t'_i)$  is the *i*-th negative sample. Hence, the loss function is defined as

$$\begin{split} L &= -\log \sigma(\omega - d_{r}(h, t)) \\ &- \sum_{i=1}^{n} p(h_{i}^{'}, r, t_{i}^{'}) \log \sigma(d_{r}(h_{i}^{'}, t_{i}^{'}) - \omega) \end{split}$$

where  $\omega$  is a fixed margin and  $\sigma$  is the sigmoid function.

# **C.** Experiments

#### **C.1. Time Prediction**

To evaluate the time prediction queries, we consider the dates in the min-max timestamp range of the set that is being evaluated as the candidate set.

### C.2. Model Selection

The best model is selected using the MRR on validation set and HITS@N with N = 1, 3, 10, Mean Rank (MR), MRR on test set are reported.

<sup>&</sup>lt;sup>4</sup>https://developer.github.com/webhooks/ event-payloads/

DATASET	Туре	#TRAIN	#VALIDATION	#Test
	INTERPOLATED	285,953	3,530	3,531
CITUUD SE 1M NODE	Extrapolated	281,056	2,104	3,276
GITHUB-SE IM-NODE	Standard Extrapolated	275,805	2,104	3,276
	INTERPOLATED	282,597	1,595	1,595
GITHUB-SE 1Y-REPO	Extrapolated	269,789	2,281	1,472
	Standard Extrapolated	252,845	2,281	1,472

Table 3. Details of train, validation, and test splits.

Table 4. Time and space complexity comparison of the models given static embedding dimension  $d_s$ , diachronic embedding dimension  $d_t$ , relative time embedding dimension  $d_r$ , entity set E, and relation set R.

Model	COMPUTATIONAL COMPLEXITY	FREE PARAMETERS COMPLEXITY
ROTATE	$O(d_s)$	$O(d_s( E  +  R ))$
DE-ROTATE	$O(d_s + d_t)$	$O((d_s + d_t)( E  +  R ))$
RA-DE-ROTATE (OURS)	$O(d_s + d_t + d_s d_r + d_r  R )$	$O((d_s + d_t)( E  +  R ) + d_r^2 + d_s d_r)$

Table 5. Performance comparison on standard extrapolated time-conditioned Link Prediction. Results within the 95% confidence interval of the best are bolded.

DATASET	MODEL	HITS@1	HITS@3	HITS@10	$\mathbf{MR}$	MRR
GITHUB-SE 1M-NODE	ROTATE	19.60	38.37	45.54	6437.30	0.2965
	DE-ROTATE	20.97	38.03	45.21	6504.79	0.3005
	RT-DE-ROTATE (OURS)	<b>22.10</b>	38.61	45.54	<b>5782.83</b>	<b>0.3113</b>
GITHUB-SE 1Y-REPO	ROTATE	0.41	1.49	2.45	2259.03	0.0141
	DE-ROTATE	5.16	8.83	16.44	<b>1342.25</b>	0.0911
	RT-DE-ROTATE (OURS)	<b>38.59</b>	<b>40.01</b>	<b>43.27</b>	1613.70	<b>0.4034</b>

### C.3. Negative Sampling

We follow the negative sampling scheme employed by Dasgupta et al. (2018) providing the model with both sets of time-agnostic and time-dependent negative samples.

### C.4. Regularization

We apply L3 regularization parameterized by  $\lambda$  as introduced in Lacroix et al. (2018) on E,  $E_A$ ,  $W_E$ , and  $W_P$ .

### **C.5. Re-ranking Heuristics**

We employed two re-ranking heuristics during the evaluation time for time-conditioned link prediction. First, each entity was only evaluated among entities with the same type. Next, we push down the ranks of entities with prior interactions with the given entity.

#### C.6. Hyperparameters

Initially, we tuned our models using the hyperparameter ranges reported in Table 6 for dropout,  $\eta$ ,  $\omega$ , and  $\alpha$  result-

Table 6. Hyperparameter ranges used for experiments.

Hyperparameter	RANGE
DROPOUT	$\{0.0, 0.2, 0.4\}$
$\eta$	$\{0.5, 1.0\}$
ω	$\{3.0, 6.0, 9.0\}$
$\alpha$	$\{10^{-3}, 10^{-4}, 3 \times 10^{-5}, 10^{-5}\}$
$\lambda$	$\{10^{-3}, 5 \times 10^{-4}, 10^{-4}\}$
$d_s$	$\{128, 96, 64, 32, 0\}$
$d_a$	$\{128, 96, 64, 32, 0\}$
$d_r$	$\{128, 64, 32, 0\}$

Table 7. Average runtime comparison of the models in seconds.

MODEL	SAMPLES	AVG RUNTIME				
Rotate	6400	77s				
DE-ROTATE	6400	80s				
RT-DE-ROTATE	6400	87s				

ing in total of 72 runs. Then, following the best hyperparameters achieved on RotatE and DE-RotatE models, we

Relation (Query)

									Relation	(Context)										
	U_SO_C_P	P_A0_C_R	n_so_c_ic	IC_A0_C_I	U_SE_C_I	I_AO_C_R	R_FO_R	U_SE_O_I	LAO_O_R	U_S0_0_P	P_A0_0_R	U_SE_RO_I	I_AO_RO_R	U_SO_C_PRC	PRC_AO_C_P	U_CO_A_R	U_SO_R_P	P_AO_R_R		
P_AO_R_R	0.23	0.23	0.24	0.19	0.22	0.35	0.13	0.15	0.3	0.24	0.28	0.12	0.15	0.26	0.13	0.24	0.051	0.0067		
U_SO_R_P	0.17	0.031	0.17	0.4	0.11	0.26	0.35	0.22	0.028	0.21	0.074	0.08	0.094	0.22	0.16	0.13	0.048	0.18		
U_CO_A_R	0.36	0.068	0.18	0.018	0.3	0.037	0.14	0.21	0.14	0.15	0.061	0.025	0.0036	0.39	0.02	0.26	0.26	0.017		
PRC_AO_C_P	0.19	0.15	0.13	0.24	0.14	0.55	0.11	0.15	0.21	0.024	0.024	0.18	0.0035	0.055	4.5e-05	0.079	0.28	0.34		<b>—</b> 0.1
U_SO_C_PRC	0.013	0.2	0.041	0.27	0.038	0.26	0.21	0.0032	0.14	0.0084	0.034	0.023	0.13	0.0068	0.16	0.35	0.0055	0.22		
I_AO_RO_R	0.019	0.079	0.053	0.075	0.074	0.15	0.24	0.12	0.13	0.16	0.2	0.074	0.1	0.28	0.078	0.32	0.35	1.1e-05		
U_SE_RO_I	0.046	0.034	0.13	3e-05	7.3e-05	2.2e-05	0.49	0.11	0.11	0.051	0.15	0.091	0.14	0.07	0.19	0.04	0.0097	0.2	-	- 0.2
P_AO_O_R	0.37	0.072	0.14	0.15	0.33	0.28	1.4e-07	0.089	0.023	0.2	0.072	0.07	0.24	0.13	0.069	0.089	0.32	0.017		
U_SO_O_P	8.9e-06	0.048	0.0077	0.36	0.058	0.34	0.071	0.0088	0.25	9.5e-08		0.055	0.032	0.012	0.19	0.3	0.016	0.26		
I_AO_O_R	0.4	0.015	0.077	0.42	0.036	0.27	0.25	0.18	0.27	0.3	0.33	0.054	0.064	0.16	0.16	0.11	0.36	2.1e-07		- 0.3
U_SE_O_I	0.0058	0.25	0.0036	0.058	0.0043	0.12	0.21	0.00079	0.037	0.0015	0.42	0.06	0.066	0.056	0.37	0.041	0.05	0.081		
R_FO_R	0.11	0.2	0.17	0.17	0.37	0.28	6.4e-06	0.19	0.095	0.18	0.2	0.25	0.18	0.3	0.036	0.12	0.37	7.5e-07		
I_AO_C_R	0.32	0.025	0.16	1.7e-06	0.25	0.2	0.082	0.049	0.049	0.23	0.29	0.24	0.25	0.053	0.14	0.021	0.21	8.7e-06		<b>—</b> 0.4
U_SE_C_I	1.7e-06	0.17	0.018	8.7e-07	0.018	0.25	0.17	7.7e-06	5.2e-06	1.3e-06	0.3	2.7e-06	0.23	1.2e-06	0.37	0.31	1e-06	0.25		
IC_AO_C_I	0.16	0.3	0.42	2.6e-06	0.0063	3e-05	0.14	0.21	0.21	0.072	0.26	0.049	0.056	0.32	0.27	0.02	0.35	0.26		
U_SO_C_IC	1.9e-07	0.27	1.4e-06	0.21	1.9e-06	0.054	0.018	1.8e-06	0.26	2.4e-07	0.15	1.6e-06	0.021	8.8e-07	0.19	0.099	1.9e-06	0.26		
P_AO_C_R	0.17	0.042	0.17	0.22	0.12	0.033	7.5e-06	0.23	0.066	0.042	0.042	0.47	0.036	0.16	0.12	0.081	0.097	0.0028		0.5
U_SO_C_P	0.026	0.18	0.026	0.15	0.045	0.0026	0.097	0.011	0.066	0.026	0.026	0.021	0.016	0.013	1.4e-05	0.31	0.0016	0.074		

Figure 1. Example heatmap of the absolute importance scores between relations learned as part of the model.

DATASET	MODEL	HITS@1	HITS@3	HITS@10	$\mathbf{MR}$	MRR
GITHUB-SE 1Y-REPO	ROTATE	1.77	18.27	56.05	10.62	0.1675
	DE-ROTATE	3.46	7.27	60.73	9.35	0.1724
	RT-DE-ROTATE (OURS)	6.18	19.29	55.91	9.40	0.2073
	RANDOM	5.26	15.79	52.63	9.5	0.1867

used dropout = 0.4,  $\eta = 0.5$ ,  $\omega = 6.0$ ,  $\alpha = 3 \times 10^{-5}$ ,  $\lambda = 5 \times 10^{-4}$ , time-agnostic negative ratio = 256, timedependant negative ratio = 32, batch size = 64, warmup steps = 100000, warm-up  $\alpha$  decay rate = 0.1, steps = 200000, and validation steps = 10000 for all experiments.

To make a fair comparison, we chose a base embedding size of 128 for all experiments. Subsequently, we only report on the combinations of static embedding dimension  $d_s$  values and diachronic embedding dimension  $d_t$  values presented in Table 6 where  $d_s + d_t = 128$ . We evenly distribute  $d_t$ among all diachronic embeddings to prevent giving models a distinct advantage in terms of free-parameters. As for the relative time embedding dimension  $d_r$ , we report on all the combinations in Table 6 with  $d_s$  and  $d_t$  respecting the stated restriction resulting in total of 17 experiments per dataset.

### C.7. Runtime

Table 7 presents the average runtime of each model for every 100 steps with batch size set to 64. All experiments were carried on servers with 16 CPU cores, 64GB of RAM, and a NVIDIA V100/P100 GPU.

### **C.8. Standard Error**

We use standard error to calculate confidence intervals and detect statistically indistinguishable results.

# **C.9. Relations Importance Matrix**

Figure 1 presents the importance matrix between relations, i.e.  $W_P$ , learned as part of the RT-X model. From this figure, it is evident that the learned matrix is not symmetric, indicating that the model learns different importance scores conditioned on the query relation.

EVENT TYPE	HEAD	RELATION (CODE)	TAIL				
COMMIT COMMENT	USER	ACTOR (U_AO_CC)	COMMIT COMMENT				
Fork	REPOSITORY	Fork (R_FO_R)	Repository				
	USER	CREATED (U_SO_C_IC) EDITED (U_SO_E_IC)	ISSUE COMMENT				
ISSUE COMMENT		CREATED (U.SOLDIC)					
	ISSUE COMMENT	EDITED (IC_AO_E_I) DELETED (IC_AO_D_I)	REPOSITORY				
		OPENED (U_SE_O_I)					
		EDITED (U_SE_E_I) DELETED (U_SE_D_I)					
		PINNED (U_SE_P_I)					
		UNPINNED (U_SE_UP_I)					
	USER	REOPENED (U_SE_RO_I)	Issue				
		ASSIGNED (U_SE_A_I)					
		LOCKED (U_SE_LO_I)					
		UNLOCKED (U_SE_ULO_I)					
Incura	Larp	ASSIGNED (U_AO_A_I)	LOOVE				
ISSUES	USER	UNASSIGNED (U_AO_UA_I)	ISSUE				
		EDITED (LAO_C_R)					
		DELETED (I_AO_D_R)					
		PINNED (I_AO_IP_R)					
	Issue	CLOSED (I_AO_C_R)	PEROSITORY				
	1550E	REOPENED (I_AO_A_R)	REPOSITORI				
		UNASSIGNED (I_AO_UA_R)					
		LOCKED (I_AO_LO_R)					
		TRANSFERRED (I_AO_ULO_R)					
Member	USER	REMOVED (U_CO_E_R)	REPOSITORY				
		EDITED (U_CO_R_R)					
		CREATED (U_SO_C_PRC)	Deve Devere Devere German				
	USER	DELETED (U_SO_E_PRC)	PULL REQUEST REVIEW COMMENT				
PULL REQUEST REVIEW COMMENT		CREATED (PRC_AO_C_P)	PULL REQUEST				
	PULL REQUEST REVIEW COMMENT	EDITED (PRC_AO_E_P) DELETED (PRC_AO_D_P)					
		SUBMITTED (U_SO_S_PR)					
	USER	EDITED (U_SO_E_PR)	PULL REQUEST REVIEW				
PULL REQUEST REVIEW		SUBMITTED (PR_AO_S_P)					
	PULL REQUEST REVIEW	EDITED (PR_AO_E_P)	PULL REQUEST				
		DISMISSED (PR_AO_D_P)					
		ASSIGNED (U_SO_A_P)					
		REVIEW REQUESTED (U_SO_RR_P)					
		REVIEW REQUEST REMOVED					
	11	OPENED (U_SO_RKK_P)	Dev D				
	USER	EDITED (U_SO_E_P)	PULL REQUEST				
		CLOSED (U_SO_C_P) READY FOR REVIEW (U SO RFR P)					
		LOCKED (U_SO_L_P)					
		UNLOCKED (U_SO_UL_P) REOPENED (U_SO_R_P)					
DULL PROUEST		SYNCHRONIZE (U_SO_S_P)					
PULL REQUEST	USER	ASSIGNED (U_AO_A_P)	PULL REQUEST				
	Lapp	REVIEW REQUESTED (U_RRO_A_P)	Drift Drovper				
	USER	REVIEW REQUEST REMOVED (U_RRO_R_P)	PULL REQUEST				
		UNASSIGNED (P_AO_UA_R)					
		REVIEW REQUESTED (P_AO_RR_R)					
	PULL REQUEST	PULL REOUEST (P_AO_C.RK)					
		EDITED (P_AO_E_R)	3511041				
		READY FOR REVIEW (P-AO_RFR_R)					
		LOCKED (P_AO_LL_R)					
		REOPENED (P_AO_UL_K)					
		SYNCHRONIZE (P_AO_S_R)					
Push	USER	SENDER (U_SO_C)	REPOSITORY				
STAR	USER	CREATED (U_HS_A_R)	REPOSITORY				
0.1111	COLK	DELETED $(U_HS_R_R)$	ALL ODITORI				

# Table 9. Extraction rules used to build the KG from raw events.