

# Alleviating Sparsity of Open Knowledge Graphs with Ternary Contrastive Learning

Qian Li<sup>1,2</sup>, Shafiq Joty<sup>2,3</sup>, Daling Wang<sup>1\*</sup>, Shi Feng<sup>1</sup> and Yifei Zhang<sup>1</sup>

<sup>1</sup> Northeastern University, China

<sup>2</sup> Nanyang Technological University, Singapore

<sup>3</sup> Salesforce Research

feiwangyuzhou@foxmail.com, srjoty@ntu.edu.sg  
{wangdaling, fengshi, zhangyifei}@cse.neu.edu.cn

## Abstract

Sparsity of formal knowledge and roughness of non-ontological construction make sparsity problem particularly prominent in Open Knowledge Graphs (OpenKGs). Due to sparse links, learning effective representation for few-shot entities becomes difficult. We hypothesize that by introducing negative samples, a contrastive learning (CL) formulation could be beneficial in such scenarios. However, existing CL methods model KG triplets as binary objects of entities ignoring the relation-guided ternary propagation patterns and they are too generic, i.e., they ignore zero-shot, few-shot and synonymy problems that appear in OpenKGs. To address this, we propose TernaryCL, a CL framework based on ternary propagation patterns among head, relation and tail. TernaryCL designs *Contrastive Entity* and *Contrastive Relation* to mine ternary discriminative features with both negative entities and relations, introduces *Contrastive Self* to help zero- and few-shot entities learn discriminative features, *Contrastive Synonym* to model synonymous entities, and *Contrastive Fusion* to aggregate graph features from multiple paths. Extensive experiments on benchmarks demonstrate the superiority of TernaryCL over state-of-the-art models.

## 1 Introduction

Open Knowledge Graphs (OpenKGs) structure textual facts in the form of (*subject entity*, *relation*, *object entity*) without depending on an ontology schema (Fader et al., 2011; Gashteovski et al., 2019). They benefit knowledge-intensive tasks, such as question answering (Sun et al., 2019) and dialogue systems (Dinan et al., 2019). Representation learning of OpenKGs aims to learn implicit embeddings of entities and relations (Gupta et al., 2019; Broscheit et al., 2020), and has become an indispensable step in the applications of OpenKGs.

Due to the rigidity of grammatical patterns and roughness of non-ontological construction, a

common challenge in representation learning of OpenKGs is the sparsity problem, where a large portion of entities have few- or zero-shot links. In the two standard OpenKGs: ReVerb20K and ReVerb45K (Vashishth et al., 2018), the degree of 55% and 89% entities is less than 3. As such, few- and zero-shot entities in these OpenKGs do not get enough training, resulting in poor generalization.

We hypothesize that negative samples could help existing sparse links to learn discriminative features in the form of a negative feedback. Being popular in self-supervised learning, *contrastive learning* (CL) aims to learn representations by contrasting negative samples with the positive ones (He et al., 2020; Gao et al., 2021; Zhu et al., 2021). Although existing CL in Graphs (Velickovic et al., 2019) and CuratedKGs (Ahrabian et al., 2020; Wang et al., 2022) have shown promising results, they could not effectively tackle the sparsity and synonymy problems of OpenKGs. First, they model KG triplets as binary objects of entities, ignoring the relation-guided ternary propagation patterns where entities propagate to multiple neighbor entities through multiple relations. Second, they do not specifically tackle the sparsity issues of zero- and few-shot entities. Finally, they can not address the synonymy problem in OpenKGs where multiple entities with different surface forms have the same meaning.

To alleviate the above problems, we propose TernaryCL, a contrastive learning framework based on ternary relational patterns among head, relation and tail. TernaryCL uses the following key ideas: (1) **Contrastive Entity** learns ternary discriminative features of different entities under the same (head entity, relation)-pair, which alleviates the sparsity issue by negative entities (Fig. 1B); (2) **Contrastive Relation** learns ternary discriminative features of different relations under the same (head entity, tail entity)-pair, which alleviates the sparsity by negative relations (Fig. 1C); (3) **Contrastive Self** designs *self* positive samples to give zero- and

\*Corresponding author

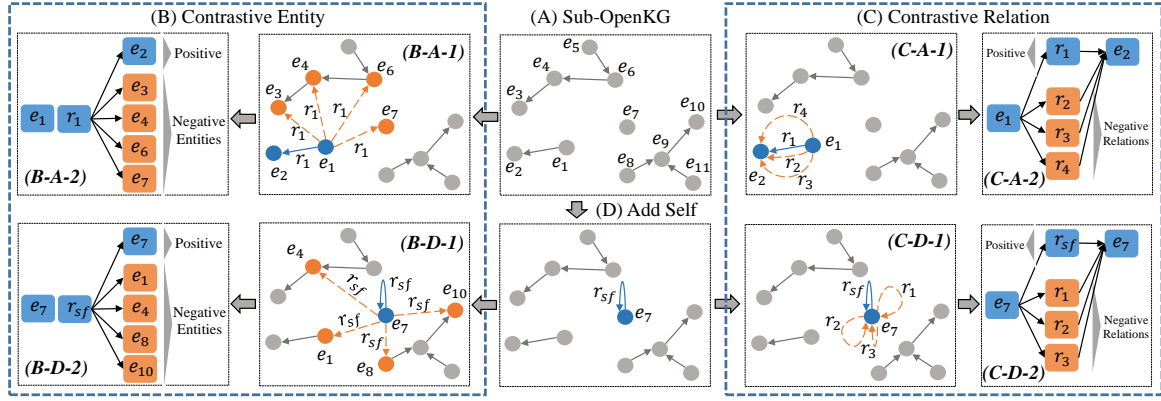


Figure 1: Framework of TernaryCL for alleviating sparsity of OpenKGs. (A) A given sub-OpenKG. (B) Contrastive Entity generates negative entities (yellow) and contrasts them with a positive entity (blue). (C) Contrastive Relation generates negative relations (yellow) and contrasts them with a positive relation (blue). (D) Contrastive Self constructs a positive sample by adding a *self* relation  $r_{sf}$  (blue) to the entity (blue), generates negative entities (yellow) in (B-D) and negative relations (yellow) in (C-D), and contrasts them with the *self* positive sample.

few-shot entities chances to learn discriminative features (Fig. 1D); (4) **Contrastive Synonym** constructs *synonymous* positive samples to automatically aggregate synonymous entities (Fig. 2); (5) **Contrastive Fusion** extends the ternary propagation pattern from the above 1-to-1 to 1-to-N: a head entity propagates to multiple tail entities through multiple relations. To gain insights into the method, we analyze the gradients of the above components.

We show that TernaryCL can learn effective embeddings from the KG itself without relying on pretrained language models or external information. This makes it light-scale and easy to apply to downstream tasks. We perform extensive evaluation to understand its performance on full-shot and at different sparsity levels, and for few- and zero-shot types. Our results show that TernaryCL can significantly outperform state-of-the-art baselines.

In summary, our key contributions are:

- To the best of our knowledge, this is the first work to do contrastive learning with relation and synonymous views over OpenKGs.
- We propose a ternary contrastive learning model, TernaryCL, to learn representations from complex propagation patterns of OpenKGs. It subtly generates negative samples from the perspective of entities and relations, and uses contrastive learning to incorporate structural information. It does not rely on external resources making it easy to scale and to apply to downstream tasks.
- We introduce *Contrastive Self* to solve the problem of learning for zero- and few-shot entities,

and propose *Contrastive Synonym* to aggregate signals from synonymous entities.

- We perform extensive experiments to show the superiority of our method over state-of-the-art baselines. We release our code at <https://github.com/feiwangyuzhou/TernaryCL>.

## 2 Related Work

OpenKGs represent factual knowledge in structured forms, as triples of *head-relation-tail* or  $(h, r, t)$ . They are extracted with OpenIE tools (Fader et al., 2011; Gashtevski et al., 2019), and generally do not rely on specification of ontology. Although OpenKGs have the advantage that they can be easily bootstrapped to new domains, because of the sparsity of formal grammatical patterns and non-ontological construction, many relevant facts are often missing from such OpenKGs. This makes them difficult to be used effectively in downstream tasks (Chandras and Talukdar, 2021).

Representation learning of KGs devotes to learning informative features of entities and relations, which can be used for other KG-related downstream tasks. General representation learning models over KGs focus on inducing structural features with linear (Bordes et al., 2013), bilinear (Wang et al., 2014; Lin et al., 2015), complex (Yang et al., 2015; Trouillon et al., 2016) or convolutional (Dettmers et al., 2018; Nguyen et al., 2018) operations, while OpenKG-specific models enhance the embeddings with side information (Gupta et al., 2019) and pretrained language models (Chandras and Talukdar, 2021). However, these methods

are limited in alleviating the sparsity issue. We propose a contrastive learning method that is more effective and does not rely on external resources.

Contrastive learning aims to learn effective representation by pulling close neighbors and pushing apart non-neighbors (Hadsell et al., 2006; Gao et al., 2021), which has achieved great success in vision (He et al., 2020), text (Gao et al., 2021) and graph (Zhu et al., 2021). Contrastive learning methods in Graphs attempt to leverage a contrastive loss at node (Velickovic et al., 2019), graph (Sun et al., 2020) and multi-view levels (Hassani and Ahmadi, 2020; Zhu et al., 2021). Ahrabian et al. (2020) and Wang et al. (2022) introduce contrastive learning into CuratedKGs by designing negative sampling strategies. However, existing contrastive learning on Graphs and CuratedKGs only model binary objects. In contrast, we propose contrastive learning to study ternary patterns in OpenKGs.

### 3 Preliminaries

- Let  $\mathcal{G} = (\mathcal{E}, \mathcal{R})$  be an **OpenKG** and  $(h, r, t)$  be a triple in  $\mathcal{G}$  with  $h, t \in \mathcal{E}$  being the head and tail entities, and  $r \in \mathcal{R}$  being the relation between them. Entities and relations are non-empty word sequences; let  $w_h = \{w_{h,i}\}_{i=1}^{|w_h|}$  and  $w_r = \{w_{r,i}\}_{i=1}^{|w_r|}$  represent word sequences of entity  $h$  and relation  $r$  respectively. Representations of entities and relations are denoted as  $\mathbf{E} \in \mathbb{R}^{|\mathcal{E}| \times D}$  and  $\mathbf{R} \in \mathbb{R}^{|\mathcal{R}| \times D}$  respectively, with  $D$  being the embedding dimension. We will use  $\mathbf{h} \in \mathbf{E}$  and  $\mathbf{r} \in \mathbf{R}$  to denote the embeddings of entity  $h$  and relation  $r$ , respectively.
- **Link prediction** is used as a downstream task to verify the effectiveness of representation learning. Link prediction in OpenKGs is to predict the answers for the two questions: (1) predicting the tail  $Q_t=(h, r, ?)$  and (2) predicting the head  $Q_h=(?, r, t)$ . For each such question, the possible answer entities can be one or more, because there could be multiple entities with the same meaning but different textual forms in an OpenKG (Broscheit et al., 2020). For example, for question  $Q_t=(\text{"NBC-TV"}, \text{"has office in"}, ?)$ , we expect all answers from the set of entities  $\{\text{"New York"}, \text{"NYC"}, \text{"New York City"}\}$ .
- A **zero-shot** entity (relation) is an entity (relation) without any link in the KG. A **few-shot** entity (relation) is an entity (relation) with few links in the KG, e.g., an one-shot entity has only one link. A zero-shot (few-shot) triple is a triple that contains at least one zero-shot (few-shot) entity or relation.

## 4 Proposed TernaryCL Model

We propose TernaryCL to alleviate the sparsity of OpenKGs in representation learning (Fig. 1). In the following, we first introduce a simple *Ternary Similarity* function to compute a similarity score for each  $(h, r, t)$  triple in an OpenKG by considering both textual and structural information (§4.1). We present *Contrastive Entity* in §4.2 to learn embeddings of different entities with the same  $(h, r)$ -pair, followed by *Contrastive Relation* in §4.3 to learn embeddings of different relations with the same  $(h, t)$ -pair. We describe *Contrastive Self* in §4.4 that constructs a positive sample  $(h, r_{sf}, h^+)$  and contrasts it with negative samples to give zero- and few-shot entities chances to learn informative features. We present *Contrastive Synonym* in §4.5 to aggregate signals from multiple synonymous entities. *Contrastive Fusion* in §4.6 further extends propagation patterns from the above 1-1 to 1-N. Finally, the training procedure is described in §4.7.

### 4.1 Ternary Similarity

For a triple  $(h, r, t) \in \mathcal{G}$ , word sequence of head entity  $h$  is  $w_h = \{w_{h,i}\}_{i=1}^{|w_h|}$ , of relation  $r$  is  $w_r = \{w_{r,i}\}_{i=1}^{|w_r|}$ , and of tail entity  $t$  is  $w_t = \{w_{t,i}\}_{i=1}^{|w_t|}$ . We encode each of these sequences with a text encoder (Enc) such as BiGRU (Cho et al., 2014) and BERT (Devlin et al., 2019).

$$\mathbf{i}^w = \text{Enc}(w_i) \text{ for } w_i \in \{w_h, w_r, w_t\} \quad (1)$$

This yields the *textual* embeddings of  $h$ ,  $r$  and  $t$  as  $\mathbf{h}^w$ ,  $\mathbf{r}^w$  and  $\mathbf{t}^w$ , respectively. For BiGRU, we concatenate the last states in the forward and backward directions to get the sequence representation. For BERT, we take the [CLS] representation.

Then, we focus on exploiting potential connections between entities and relations. We use a two-dimensional convolutional network (Dettmers et al., 2018) to learn potential connections between a head entity  $h$  and a relation  $r$  as follows:

$$\varphi(h, r) = \rho(\text{Linear}(\rho(\text{Conv2d}_\omega([\hat{\mathbf{h}}; \hat{\mathbf{r}}]))) \quad (2)$$

where  $\rho$  represents a ReLU activation, and  $\hat{\mathbf{h}}$  and  $\hat{\mathbf{r}}$  denote a reshaping of  $[\mathbf{h} + \mathbf{h}^w]$  and  $\mathbf{r}^w$  respectively, with  $\mathbf{h} \in \mathbf{E}$ .<sup>1</sup> The reshaping operation converts a vector  $\mathbf{v} \in \mathbb{R}^D$  from one-dimension  $D$  to two-dimensions  $\mathbf{v} \in \mathbb{R}^{D_1 \times D_2}$ , where  $D = D_1 \cdot D_2$ ;

<sup>1</sup>Recall that  $\mathbf{h} \in \mathbf{E}$  and  $\mathbf{r} \in \mathbf{R}$  are the embedding parameters of entity  $h$  and relation  $r$ , which are initialized randomly and updated during training. For relations, adding  $\mathbf{r} \in \mathbf{R}$  with the textual embedding  $\mathbf{r}^w$  worsen the performance.

$[\hat{\mathbf{h}}; \hat{\mathbf{r}}] \in \mathbb{R}^{(2D_1) \times D_2}$  represents the concatenation of the reshaped embeddings of  $\hat{\mathbf{h}}$  and  $\hat{\mathbf{r}}$ . The  $\text{Conv2d}_\omega$  symbol denotes a two-dimensional convolutional layer with filters  $\omega$ . This layer returns a feature map tensor  $\mathcal{C} \in \mathbb{R}^{C_1 \times C_2 \times C_3}$ , where  $C_1$  is the number of feature maps of dimensions  $C_2 \times C_3$ .  $\mathcal{C}$  is then reshaped into a vector  $\mathbb{R}^{C_1 \cdot C_2 \cdot C_3}$ , and projected back to  $\mathbb{R}^D$  with a Linear layer. Through the 2d convolution module, potential embeddings of entity  $h$  and relation  $r$  are jointly encapsulated.

Finally, we compute a similarity score for each triple  $(h, r, t)$  with a dot product similarity function:

$$\beta(h, r, t) = \varphi(h, r) \cdot \mathbf{t} \quad (3)$$

where  $\mathbf{t} \in \mathbf{E}$ . When predicting the head  $h$  based on pair  $(r, t)$ , we reverse the relation by adding a special symbol, and obtain a new triple  $(t, r_{\text{rev}}, h)$ . The similarity score for  $(t, r_{\text{rev}}, h)$  is computed in a similar fashion as above following Eq. (1)-Eq. (3).

## 4.2 Contrastive Entity

Contrastive Entity (Fig. 1B) alleviates sparsity of OpenKGs from the perspective of *negative* entities, and induces discriminative features of different entities with the same  $(h, r)$ -pair. The contrastive score for a triplet  $p_e = (h, r, t^+)$  is:

$$S(h, r, t^+) = -\log \frac{e^{\beta(h, r, t^+)/\tau}}{\sum_{n \in \{p_e, \mathcal{N}_e\}} e^{\beta(n)/\tau}} \quad (4)$$

where  $\tau$  is a temperature hyperparameter,  $\beta(\cdot)$  is the similarity score as in Eq. (3), and  $p_e = (h, r, t^+)$  is a true triple in the OpenKG;  $\mathcal{N}_e = \{(h, r, t_j^-)\}_{j=1}^{|\mathcal{N}_e|}$  is a set of negative samples, where a negative entity  $t_j^-$  is selected from a candidate entity list defined by:  $\mathcal{E} - \mathcal{E}(h, r)$  with  $\mathcal{E}(h, r)$  being the entity list of true answers (tail entities), that is,  $t_i \in \mathcal{E}(h, r)$  if the triple  $(h, r, t_i) \in \mathcal{G}$ . To analyse how this contrastive loss affects the learning, we perform gradient analysis. It can be shown that the gradients with respect to the head entities are:

$$\begin{aligned} -\frac{\partial S(h, r, t^+)}{\partial \mathbf{h}} &= \frac{\varphi'(h, r)}{\tau A} \left( \left[ \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right] \mathbf{t}^+ \right. \\ &\quad \left. - \sum_{(h, r, t_j^-) \in \mathcal{N}_e} \left[ e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right] \mathbf{t}_j^- \right) \end{aligned} \quad (5)$$

where  $A$  is a normalization constant (see Appendix for a derivation). This is consistent with our intuition, where positive entity  $t^+$  gives positive feedback while negative entities  $t_j^-$  give negative feedback. The gradients with respect to the relations  $(-\frac{\partial S(h, r, t^+)}{\partial \mathbf{r}})$  has a similar form as Eq. (5).

Similarly, we can derive the gradients for the positive  $t^+$  and negative tail entities  $t_j^-$  as:

$$-\frac{\partial S(h, r, t^+)}{\partial \mathbf{t}^+} = \frac{\varphi(h, r)}{\tau A} \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \quad (6)$$

$$-\frac{\partial S(h, r, t^+)}{\partial \mathbf{t}_j^-} = -\frac{\varphi(h, r)}{\tau A} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \quad (7)$$

For a few-shot entity  $t$ , when it appears as a positive (Eq. (6)) or negative (Eq. (7)) sample headed by entity  $h$  of high degree, it gets sufficient gradients to learn informative representations. Similarly, through Eq. (7), the parameters of a zero-shot entity get updated when it appears as a negative sample with a  $(h, r)$ -pair. As discussed later in §4.4, with Contrastive Self, zero-shot entities get updated with their own contrastive losses. Overall, through the negative samples, training of few-shot, zero-shot and other entities (with many links) gets more balanced, while with existing approaches, zero-shot entities generally do not get trained as they are disconnected from the rest of the OpenKG.

## 4.3 Contrastive Relation

Entities in a KG propagate information to neighboring entities through one or more relations, so the features of relations are as important as that of entities. Contrastive Relation (Fig. 1C) alleviates the sparsity from the perspective of *negative* relations, and captures potential features of different relations with the same  $(h, t)$ -pair. The contrastive score for a positive relation  $p_r = (h, r^+, t)$  is:

$$S(h, r^+, t) = -\log \frac{e^{\beta(h, r^+, t)/\tau}}{\sum_{n \in \{p_r, \mathcal{N}_r\}} e^{\beta(n)/\tau}} \quad (8)$$

where  $p_r$  is a true triple in the OpenKG, and  $\mathcal{N}_r = \{(h, r_j^-, t)\}_{j=1}^{|\mathcal{N}_r|}$  is a set of negative samples, where a negative relation  $r_j^-$  is sampled from a candidate relation list  $\mathcal{R} - \mathcal{R}(h, t)$  with  $\mathcal{R}(h, t)$  being a relation list that satisfies the condition:  $r_i \in \mathcal{R}(h, t)$  if the triple  $(h, r_i, t) \in \mathcal{G}$ . The gradients have the same form as above (see Appendix A.1). In particular, the gradients that a tail entity gets is:

$$\begin{aligned} -\frac{\partial S(h, r^+, t)}{\partial \mathbf{t}} &= \frac{1}{\tau B} \left( \left[ \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \right] \varphi(h, r^+) \right. \\ &\quad \left. - \sum_{(h, r_j^-, t) \in \mathcal{N}_r} \left[ e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \right] \varphi(h, r_j^-) \right) \end{aligned} \quad (9)$$

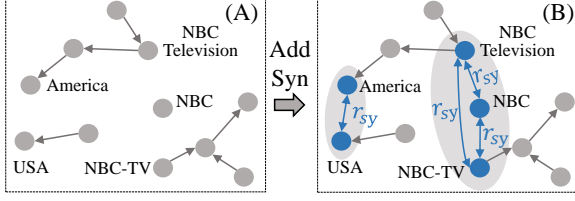


Figure 2: Examples of Contrastive Synonym. (A) A given OpenKG, where two entities “America”, “USA” are synonymous, three entities “NBC Television”, “NBC”, “NBC-TV” are synonymous. (B) Add a  $sy$  relation ( $r_{sy}$ ) between each pair of synonyms.

Different from Eq. (6), tail entities ( $t = t^+$ ) get contrastive gradients, because the signals from the positive ( $r^+$ ) and negative ( $r_j^-$ ) samples are in opposite direction in Eq. (9).

#### 4.4 Contrastive Self

As described in §4.2 and §4.3, parameters of a few-shot entity get updated with contrastive gradients when it acts as a head (Eq. (5)) or tail entity (Eq. (9)). Parameters of a zero-shot entity are also updated when it appears as a negative entity (Eq. (7)), but not from the perspective of contrastive. This means zero-shot entities have no chance to learn discriminative representations, because they have no links with the rest of OpenKG.

In view of this, we propose to construct a positive sample  $p_{sf} = (h, r_{sf}, h^+)$  by adding a *self* relation (Fig. 1D), where  $h$  and  $h^+$  are the same entity but with different embeddings: the embedding of  $h$  is  $[\mathbf{h} + \mathbf{h}^w]$  and the embedding of  $h^+$  is  $\mathbf{h} \in \mathbf{E}$ . For such a positive sample, negative samples  $\mathcal{N}_e = \{(h, r_{sf}, h_j^-)\}_{j=1}^{|\mathcal{N}_e|}$  are generated with negative entities by selecting  $h_j^-$  strategically from an entity list  $\mathcal{E} - h$  (Fig. 1B-D). Similarly, negative relation samples  $\mathcal{N}_r = \{(h, \{r_{sf}\}_j^-, h^+)\}_{j=1}^{|\mathcal{N}_r|}$  are generated with negative relations by selecting  $\{r_{sf}\}_j^-$  strategically from a relation list  $\mathcal{R} - r_{sf}$  (Fig. 1C-D). The contrastive scores for  $p_{sf} = (h, r_{sf}, h^+)$  can then be computed with Eq. (4) and Eq. (8).

Through the *self* positive sample, parameters of zero-shot entities can have the chances to be updated from the contrastive perspective (Eq. (5) and Eq. (9)), where  $h^+$ ,  $r_{sf}^+$  give positive feedback while  $h_j^-$ ,  $\{r_{sf}\}_j^-$  give negative feedback.

#### 4.5 Contrastive Synonym

We propose Contrastive Synonym based on the synonymity characteristic of an OpenKG, where multiple entities with different surface forms have

the same meaning, e.g., “NBC-TV”, “NBC” and “NBC Television” mean the same thing. These synonymous entities are designed as positive samples of each other. We identify synonyms automatically by a fuzzy matching with the IDF-token-overlap tool (Galárraga et al., 2014) and a semantic matching with Eq. (3) based on the embeddings obtained by Eq. (1). We put the extracted synonymous entities for  $h$  into  $\mathbb{S}(h)$  and construct a positive sample  $p_{sy} = (h, r_{sy}, h_i^+)$  by adding a  $sy$  relation for each entity  $h_i^+ \in \mathbb{S}(h)$  (Fig. 2). We follow the same strategy as in §4.2 and §4.3 for sampling negative entity and relation samples, and compute the contrastive score for each  $p_{sy} = (h, r_{sy}, h_i^+)$  with Eq. (4) and Eq. (8).

#### 4.6 Contrastive Fusion

The contrastive learning of the above modules has the form of 1-to-1, where a head entity propagates information to a tail entity through a relation, i.e., it involves only one positive sample. However, a head entity can connect to multiple tail entities through multiple relations. To model this multi-propagation pattern, we extend the 1-to-1 to 1-to-N, where multiple positive samples are considered.

We design two sets of positive 1-to-N patterns:  $p_e = \{(h, r, t_j^+)\}_{j=1}^{|p_e|}$  where a head entity  $h$  connects to multiple tail entities through a relation  $r$ , and  $p_r = \{(h, r_j^+, t)\}_{j=1}^{|p_r|}$  where a head entity  $h$  connects to a tail entity  $t$  through multiple relations. We generate negative samples  $\mathcal{N}_e = \{(h, r, t_j^-)\}_{j=1}^{|\mathcal{N}_e|}$  for  $p_e$  as in §4.2, and negative samples  $\mathcal{N}_r = \{(h, r_j^-, t)\}_{j=1}^{|\mathcal{N}_r|}$  for  $p_r$  as in §4.3. Let  $p_f = \{p_e, p_r\}$  and  $\mathcal{N}_f = \{\mathcal{N}_e, \mathcal{N}_r\}$ . With this, we design two types of contrastive scores to learn 1-to-N patterns.

$$S_a(p_f) = -\log \frac{\sum_{n^+ \in p_f} e^{(\beta(n^+)/\tau)}}{\sum_{n^- \in \{p_f, \mathcal{N}_f\}} e^{(\beta(n^-)/\tau)}} \quad (10)$$

$$S_b(p_f) = -\sum_{n^+ \in p_f} \log \frac{e^{(\beta(n^+)/\tau)}}{\sum_{n^- \in \{n^+, \mathcal{N}_f\}} e^{(\beta(n^-)/\tau)}} \quad (11)$$

In Eq. (10), all positive samples are trained together and normalized to a unified space, while in Eq. (11), different positive samples are trained separately and the independent similarity score of all positives are accumulated. From the perspective of negative samples, negative entity  $\mathcal{N}_e$  and relation  $\mathcal{N}_r$  samples are merged to train with the positive samples, which is different from Eq. (4) and Eq. (8) which only have one type of negatives.

Dataset	Ent	Rel	Clust	Valid	Test	Train (Sparsity Levels)					Few-Shot		Zero-Shot	
						100%	80%	60%	40%	20%	Ent	Rel	Ent	Rel
<b>ReVerb20K</b>	11.1	11.1	10.8	1.6	2.4	15.5	12.4	9.3	6.2	3.1	2.8 / 1.0	2.5 / 0.8	8.1 / 1.1	8.5 / 1.0
<b>ReVerb45K</b>	27.0	21.6	18.6	3.6	5.4	36.0	28.8	21.6	14.4	7.2	7.7 / 1.8	5.1 / 1.5	18.8 / 2.5	16.4 / 1.8

Table 1: Dataset statistics (in 1000s). Ent, Rel, Clust show the no. of entities, relations, entity clusters. Valid, Test, Train show the no. of triples in valid, test, train sets. In Train,  $x\%$  represents sparsity levels. In Few- and Zero-Shot, left of / represents the no. of entities or relations, and right of / represents the no. of related triples in its Test set.

## 4.7 Training Procedure

We train TernaryCL in Pretrain and Finetune stages, where Pretrain stage aims to learn discriminative representations with Eq. (10) or Eq. (11), and Finetune stage aims to optimize parameters for a target task. Note that TernaryCL is a general framework that can be applied to diverse KG types, such as OpenKGs, CuratedKGs and TemporalKGs, and to diverse applications, such as link prediction, relation prediction, node classification, relation extraction and other ternary tasks. In addition, the principle of Contrastive Self can be applied to dynamic tasks, e.g., when a new entity gets added to a KG, it can transfer embeddings of existing entities to the new entity. Since our focus in this work is the sparsity of OpenKGs, we design experiments with link prediction tasks in OpenKGs.

For a test triple  $(h_i, r_i, t_i)$ , the jointly encapsulated representation  $\varphi(h_i, r_i)$  of entity  $h_i$  and relation  $r_i$  is matched with the embeddings of all the entities in  $\mathbf{E}$  to predict  $\hat{Y}_i \in [0, 1]^{|\mathcal{E}|}$  as:

$$\hat{Y}_i = \text{sigmoid}(\varphi(h_i, r_i) \cdot \mathbf{E}^\top) \quad (12)$$

We use a binary cross-entropy loss defined as:

$$-\frac{1}{|\mathcal{E}|} \sum_{j=1}^{|\mathcal{E}|} Y_{i,j} \log \hat{Y}_{i,j} + (1 - Y_{i,j}) \log(1 - \hat{Y}_{i,j}) \quad (13)$$

where  $Y_{i,j} = 1$  if  $t_j \in \mathcal{E}(h_i, r_i)$  otherwise  $Y_{i,j} = 0$  with  $\mathcal{E}(h, r)$  being the set of all true tail entities for a pair  $(h, r)$ . As mentioned earlier, when predicting the head  $h_i$  based on a pair  $(r_i, t_i)$ , we reverse the relation to obtain a reversed triple  $(t, r_{\text{rev}}, h)$ , then train it with Eq. (12) and Eq. (13).

## 5 Experiments

### 5.1 Datasets and Experiment Setup

We use ReVerb20K and ReVerb45K OpenKG benchmarks (Vashishth et al., 2018), which are constructed through ReVerb (Fader et al., 2011). Table 1 presents their statistics. ReVerb45K with 27K entities and 21.6K relations is larger and sparser

than ReVerb20K with 11.1K entities and 11.1K relations. Entity clusters are gold canonicalized clusters, extracted through the Freebase entity linking tools (Gupta et al., 2019). Entities in an entity cluster have the same meaning. Usually, entity clusters are only used for evaluation.

To evaluate on sparsity, we design two sets of controlled experiments: the first is at different sparsity levels and the second is on few- and zero-shot samples. For the first, we construct *train sets* at different sparsity granularity  $\{100\%, 80\%, 60\%, 40\%, 20\%\}$  by respectively removing  $\{0\%, 20\%, 40\%, 60\%, 80\%\}$  of the links from the original train set. We use the same original data for validation and testing. For the second, we evaluate on few-shot or zero-shot samples separately, where few-shot refers to 3-, 2- and 1-shot. Few- or zero-shot entities (relations) are extracted as per the definition in §3, e.g., a 3-shot entity has three links in the test KG. For the test set of a few-shot (zero-shot) entity (relation), its related triples are extracted from the original test set. For training, we use the 20% train set, because it is sparse enough to contain more few-shot (zero-shot) samples.

### 5.2 Evaluation Metrics and Baselines

To evaluate on a single test triple, we use Mention Ranking or MR (Gupta et al., 2019), which is the minimum ranking position of the answer entities. To evaluate over all test triples, we use the three most widely used ways to integrate the individual MR scores: (a)  $H@N$ : proportion of MR scores not higher than  $N$ , (b)  $AR$ : average of all MR scores, and (c)  $ARR$ : compute the reciprocal of each MR score, and average all reciprocals. A model with better performance should have higher  $H@N$  and  $ARR$  scores and a lower  $AR$  score.<sup>2</sup>

To show the effectiveness of our approach, we compare TernaryCL against several strong base-

<sup>2</sup>In previous work,  $ARR$  and  $AR$  are called  $MRR$ , and  $MR$  respectively, where  $M$  stands for *Mean*. However, since *Mention Ranking* is abbreviated as  $MR$ , to prevent confusion, we use *Average* in the names in stead of *Mean*.

Type	Model	ReVerb20K						ReVerb45K					
		AR ↓	ARR	H@1	H@10	H@50	H@100	AR ↓	ARR	H@1	H@10	H@50	H@100
General	TransE	1497	13.3	2.2	29.6	43.0	49.2	2222	15.8	9.3	25.9	37.1	43.2
	DistMult	4569	1.9	1.3	2.7	5.2	7.1	5782	8.5	7.7	9.7	12.0	13.6
	ComlEx	4376	2.0	1.4	3.0	5.6	7.7	5173	8.9	7.5	11.3	16.0	18.9
	ConvE	1085	25.5	19.9	35.8	50.1	57.2	2483	22.1	16.6	32.4	43.3	47.9
	ConvTransE	1080	26.1	20.5	35.9	50.0	57.1	2490	23.4	17.9	33.8	44.4	48.8
OpenKG	CaReTransE	950	30.3	23.2	42.8	58.4	64.6	2414	19.5	7.8	37.5	47.5	51.4
	CaReConvE	801	31.6	25.6	42.9	56.7	63.4	1589	29.7	23.4	41.3	53.6	58.7
OpenKG+PLM	SimKGC	538	29.7	23.3	41.7	58.7	65.7	1080	27.1	20.4	39.8	53.2	59.2
	OKGITBert	524	35.1	27.5	49.5	65.9	72.7	<b>735</b>	<b>33.7</b>	<b>26.7</b>	47.1	59.8	65.2
	OKGITRob	594	35.8	28.4	49.2	65.4	72.1	849	33.4	26.5	46.4	58.8	63.9
<b>Our</b>	<b>TernaryCL</b>	<b>393</b>	<b>38.9</b>	<b>30.5</b>	<b>54.6</b>	<b>69.2</b>	<b>75.5</b>	767	33.3	25.3	<b>48.7</b>	<b>63.0</b>	<b>68.3</b>

Table 2: Results on ReVerb20K and ReVerb45K in the standard full data (100% train) setup. Best scores are made bold. Columns with ↓ denote lower is better, otherwise higher is better.

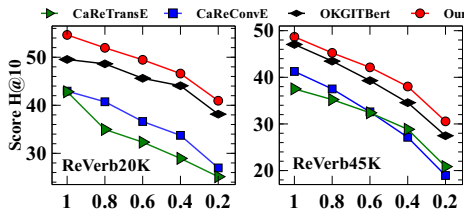


Figure 3: Results at sparsity levels on ReVerb20K (Left) and ReVerb45K (Right).

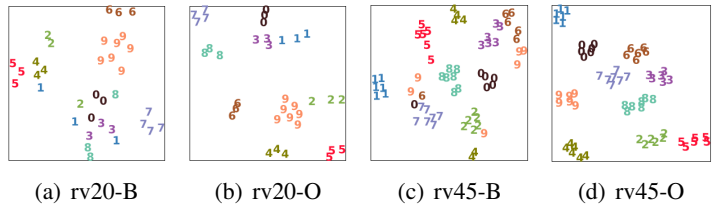


Figure 4: Visualization of Baseline(B) and Our(O) on ReVerb20K (rv20) and ReVerb45K(rv45). Same numbers denote similar entities.

lines which fall in three groups:

- **General:** This applies general KG models to OpenKGs. It includes TransE (Bordes et al., 2013), DistMult (Yang et al., 2015), ComplEx (Trouillon et al., 2016), ConvE (Dettmers et al., 2018) and ConvTransE (Shang et al., 2019).
- **OpenKG:** This comprises the baselines designed specifically for OpenKGs. It includes CaReTransE and CaReConvE (Gupta et al., 2019).
- **OpenKG+PLM:** This group uses pretrained language models (PLM) to improve the representation learning. It includes SimKGC (Wang et al., 2022), OKGITBert and OKGITRob (Chandrasah and Talukdar, 2021). SimKGC is a state-of-the-art contrastive model for CuratedKGs, we apply it to OpenKGs. OKGITBert and OKGITRob are state-of-the-art models for OpenKGs.

**For our TernaryCL** TernaryCL is implemented based on the PyTorch library with a single GeForce RTX 2080 GPU. The number of parameters is 18.06M for ReVerb20K and 33.63M for ReVerb45K. We run the model five times and report the maximum of results. For the training settings, the optimizer is set to Adam, the embedding size is set to 300. The entity and relation embeddings

are initialized randomly, and the word vectors are initialized with the GloVe embeddings. Default fusion strategy (§4.6) is  $S_b(p_f)$  and text encoder (§4.1) is BiGRU.

We tune our model with the grid search to select the optimal hyper-parameters based on the performance on the validation dataset. The list of hyperparameters are from two aspects: (1) Hyperparameters for Finetune stage: batch size  $\in \{32, 64, 128, 256, 512\}$ , learning rate  $\in \{1e-3, 1e-4, 8e-5, 5e-5, 1e-5\}$ . (2) Hyperparameters for Pretrain stage: learning rate  $\in \{1e-3, 1e-4, 5e-5, 1e-5\}$ , temperature regulation value  $\in \{0.1, 0.05, 0.01\}$ . The results of hyperparameters are shown in Fig. 6.

**For Baselines** The results of baselines are reproduced with open source implementations. Concretely, TransE, DistMult and ComplEx are reproduced with public code in <sup>3</sup>. The code for ConvE is in <sup>4</sup> and for SimKGC is in <sup>5</sup>. The code for ConvTransE is implemented by us and public in our code. We use the grid search technique to select the optimal values of hyperparameters for above baselines. For OpenKG and OpenKG+PLM baselines, CaReTransE, CaReConvE are reproduced with the

<sup>3</sup><https://github.com/uma-pil/kgc>

<sup>4</sup><https://github.com/mallabiisc/CaRE>

<sup>5</sup><https://github.com/intfloat/SimKGC>

Model	Few-Shot Entity				Few-Shot Relation				Zero-Shot Entity				Zero-Shot Relation			
	AR↓	ARR	H@1	H@10	AR↓	ARR	H@1	H@10	AR↓	ARR	H@1	H@10	AR↓	ARR	H@1	H@10
CaReTransE	1215	10.6	1.2	25.6	2022	10.7	1.7	26.4	3286	7.3	0.0	20.9	2071	10.4	2.2	24.8
CaReConvE	2196	24.4	19.6	33.1	2663	20.6	16.2	29.3	3304	14.2	11.1	20.8	2557	22.3	18.1	30.4
OKGITBert	<b>285</b>	29.4	22.2	43.1	1696	26.2	19.3	38.8	3465	20.1	15.1	29.8	1775	27.4	20.7	40.8
TernaryCL	364	<b>33.1</b>	<b>25.0</b>	<b>48.6</b>	<b>863</b>	<b>27.7</b>	<b>20.5</b>	<b>41.8</b>	<b>1473</b>	<b>20.7</b>	<b>15.5</b>	<b>31.1</b>	<b>868</b>	<b>30.1</b>	<b>23.1</b>	<b>43.2</b>
CaReTransE	2610	12.6	6.7	21.4	3323	13.6	7.8	23.0	5776	10.2	8.8	12.4	3368	11.8	6.4	20.3
CaReConvE	2739	19.2	15.5	26.2	3379	14.8	10.9	22.8	5662	4.3	3.0	6.6	3429	12.5	8.8	19.2
OKGITBert	<b>644</b>	23.5	<b>18.2</b>	33.7	2211	20.0	14.5	31.0	4952	7.9	5.4	12.5	2349	<b>18.6</b>	<b>13.2</b>	29.2
TernaryCL	954	<b>23.6</b>	16.9	<b>35.8</b>	<b>1460</b>	<b>21.3</b>	<b>14.9</b>	<b>32.8</b>	<b>2626</b>	<b>13.2</b>	<b>9.9</b>	<b>19.5</b>	<b>1539</b>	18.3	12.3	<b>29.9</b>

Table 3: Results of few-/zero-shot entities/relations on ReVerb20K (Top) and ReVerb45K (Below) with 20% train.

public code in <sup>4</sup>, OKGITBert and OKGIT+Rob are reproduced with the public code in <sup>6</sup>. The optimal values of hyperparameters for this four baselines are consistent with that in their paper.

### 5.3 Results

• **Full-data evaluation.** We first present the results on the standard full-data setup in Table 2, for which we use the original train set (100%). We notice that TernaryCL achieves substantial improvements in comparison to the baselines. For ReVerb20K, it outperforms all the baselines by a good margin across all the metrics – ARR increases by 3.1 point and  $H@1, 10, 50, 100$  increase by 2.1, 5.1, 3.3, 2.8 points, respectively. For ReVerb45K, its performance is better than General and OpenKG baselines with sizeable margins in all metrics, notably 3.6 point in ARR and 1.9, 7.4, 9.4, 9.6 points in  $H@1, 10, 50, 100$ . It also outperforms OpenKG+PLM baselines in all metrics for ReVerb20K and in 3 of 6 metrics for ReVerb45K.

Overall, TernaryCL with a simple structure can achieve better performance through innovative training methods at lower costs than OpenKG+PLM baselines that use a relatively complex structure and large PLMs which is memory and compute intensive. This can make TernaryCL a favorable choice against the baselines.

• **Targeted evaluation with sparsity.** We now evaluate whether TernaryCL can alleviate the sparsity problem from the perspective of different sparsity levels. Fig. 3 reports the results on the (original) test set for different train sets with varying sparsity level. We observe that TernaryCL achieves the best scores at all sparsity levels on both datasets. Taking ReVerb20K as an example, it gives improvements of 3.3, 3.8, 2.5 and 2.8 points over OKGITBert at 80%, 60%, 40%, 20%, respectively. Note

<sup>6</sup><https://github.com/Chandrasahd/OKGIT>

that OKGITBert enhances learning with large-scale PLMs, which possess a lot of commonsense and factual knowledge within its huge parameter space by training on large data. In contrast, TernaryCL does not use side information or PLMs, but gives significant performance gains over OKGITBert, even when the sparsity level is high (e.g., 20%).

• **Targeted evaluation for few- and zero-shot.**

We now analyze whether TernaryCL can alleviate sparsity on test sets with few- or zero-shot entities (or relations). As shown in Table 3, CaReTransE performs poorly, especially, its  $H@1$  score is 0.0 for zero-shot entities on ReVerb20K. OKGITBert achieves better results than CaReTransE and CaReConvE on most metrics, which shows that pretrained knowledge introduced by Bert could be helpful to few- and zero-shot entities and relations. The proposed TernaryCL, without using any PLM or extra information, outperforms all baselines by a good margin on most metrics.

• **Visualization evaluation.** To qualitatively show that TernaryCL can make entities with the same meaning closer in the vector space, we show t-SNE visualization (Van der Maaten and Hinton, 2008) in Fig. 4. For this, we train TernaryCL on the original train set (100%) and compare with the state-of-the-art baseline OKGITBert. For visualization, we selected 10 entity clusters from the test set, where each cluster has more than three entities. We observe that entities with the same meaning (same number) are closer in TernaryCL than in OKGITBert. These qualitative results are consistent with the above quantitative results, which further verify the effectiveness of the proposed TernaryCL.

### 5.4 Ablation Study and Analysis

• **Fusion strategy.** The ablations for the fusion strategies (§4.6) are shown in Fig. 5(a), where  $S_a(p_e)$ ,  $S_a(p_r)$ ,  $S_b(p_e)$ ,  $S_b(p_r)$  use only one type



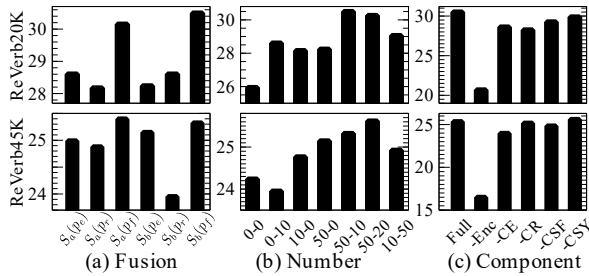


Figure 5: Results ( $H@1$ ) for (a) fusion strategies, (b) no. of negative entities and relations, and (c) components

of 1-to-N pattern (denoted by the subscript of  $p$ ), and  $S_a(p_f)$  and  $S_b(p_f)$  are our two variants.  $S_a(p_f)$ , which fuses both 1-to-N types, outperforms its counterparts  $S_a(p_e)$  and  $S_a(p_r)$ . We see the same phenomenon with  $S_b(p_f)$ , which outperforms  $S_b(p_e)$  and  $S_b(p_r)$ . This proves that both 1-to-N types carry different semantic features, and integration of them yields additive gains. When we compare our two variants,  $S_b(p_f)$  shows better performance than  $S_a(p_f)$  on ReVerb20K while similar performance on ReVerb45K. Beside  $H@1$ , we also compare the  $H@10,50,100$  scores on ReVerb45K; these scores of  $S_a(p_f)$  are 48.3, 62.5, 67.9, which are lower than that of  $S_b(p_f)$  in Table 2.  $S_b(p_f)$  contrasts a positive sample with both negative entities and relations, which is able to capture more discriminative features for the positive sample.

• **No. of negative entities and relations.** To investigate how the number of negative entities and relations could affect the performance of TernaryCL, we design several collocations of negative entities and negative relations (entity-relation): 0-0, 0-10, 10-0, 50-0, 50-10, 50-20, 10-50. Results in Fig. 5(b) show that TernaryCL achieves the best results on ReVerb20K at 50-10 and on ReVerb45K at 50-20 which have the bigger no. of negative entities and the smaller no. of negative relations. It also reveals that negative entities are more important than negative relations in improving performance.

• **Component.** We now probe the role of each component in our model (Fig. 5(c)). -Enc denotes removing the text encoder (§4.1) from the full model. Similarly, -CE, -CR, -CSF, and -CSY represent removing the Contrastive Entity (§4.2), Contrastive Relation (§4.3), Contrastive Self (§4.4), Contrastive Synonym (§4.5), respectively. Compared to the full model, performance of -Enc decreases prominently meaning that the text encoder plays a crucial role in improving performance.

The performance degradation for -CE, -CR, -CSF proves that Contrastive Entity, Contrastive Relation, Contrastive Self are also important components of the model. Removing Contrastive Synonym (-CSY) decreases the performance on ReVerb20K while has similar performance on ReVerb45K. Beside  $H@1$ , we also computed  $H@10,50,100$  scores on ReVerb45K; -CSY gets 48.1, 61.9, 67.7, which are lower than ones of the full model (Table 2).

Note that the roles of BiGRU (-Enc) and CL (-CE, -CR, -CSF, -CSY) are different – BiGRU is a basic component to encode textual features, while CL is a training mechanism to capture discriminative features. The main take home messages from our experiments and ablations are: (1) as an encoder, Bi-GRU is more effective than BERT; (2) CL gives further significant improvements irrespective of the encoder.

## 6 Conclusion

In this work, we have provided empirical insights about the sparsity problem of OpenKGs, and proposed TernaryCL, a contrastive learning framework to alleviate the sparsity by introducing contrastive entity, relation, self, synonym and fusion methods. Through extensive experiments and comprehensive analysis, we have shown that TernaryCL outperforms state-of-the-art baselines by a good margin.

## Limitations

Text encoder is an important component in the proposed TernaryCL model, where the default is BiGRU. Pretrained language models (PLMs) have achieved great results on most NLP tasks, so we also attempted to use the PLM BERT instead of BiGRU. However, performance of the model with BERT (both tuned and not tuned) as an encoder is weaker than the one with BiGRU. The results are shown in Appendix (A.3). This shows that usage of PLMs may not benefit TernaryCL further. In future, we would like to explore other large-scale PLMs like RoBERTa (Liu et al., 2019) and ELECTRA (Clark et al., 2020).

## Acknowledgements

The work was supported by National Natural Science Foundation of China (62172086, 61872074, 62272092) and Chinese Scholarship Council.

## References

- Kian Ahrabian, Aarash Feizi, Yasmin Salehi, William L. Hamilton, and Avishek Joey Bose. 2020. Structure aware negative sampling in knowledge graphs. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020*, pages 6093–6101. Association for Computational Linguistics.
- Antoine Bordes, Nicolas Usunier, Alberto García-Durán, Jason Weston, and Oksana Yakhnenko. 2013. Translating embeddings for modeling multi-relational data. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pages 2787–2795.
- Samuel Broscheit, Kiril Gashteovski, Yanjie Wang, and Rainer Gemulla. 2020. Can we predict new facts with open knowledge graph embeddings? A benchmark for open link prediction. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL, Online, July 5-10, 2020*, pages 2296–2308.
- Chandrasah and Partha P. Talukdar. 2021. OKGIT: open knowledge graph link prediction with implicit types. In *Findings of the Association for Computational Linguistics: ACL/IJCNLP, Online Event, August 1-6*, pages 2546–2559. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1724–1734. ACL.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. 2020. ELECTRA: pre-training text encoders as discriminators rather than generators. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Tim Dettmers, Pasquale Minervini, Pontus Stenetorp, and Sebastian Riedel. 2018. Convolutional 2d knowledge graph embeddings. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 1811–1818.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. Wizard of wikipedia: Knowledge-powered conversational agents. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Anthony Fader, Stephen Soderland, and Oren Etzioni. 2011. Identifying relations for open information extraction. In *Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing, EMNLP, 27-31 July 2011, John McIntyre Conference Centre, Edinburgh, UK, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1535–1545.
- Luis Galárraga, Jeremy Heitz, Kevin Murphy, and Fabian M. Suchanek. 2014. [Canonicalizing open knowledge bases](#). In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, CIKM 2014, Shanghai, China, November 3-7, 2014*, pages 1679–1688. ACM.
- Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. SimCSE: Simple contrastive learning of sentence embeddings. In *Conference on Empirical Methods in Natural Language Processing*, pages 6894–6910, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Kiril Gashteovski, Sebastian Wanner, Sven Hertling, Samuel Broscheit, and Rainer Gemulla. 2019. OPIEC: an open information extraction corpus. In *1st Conference on Automated Knowledge Base Construction, AKBC, Amherst, MA, USA, May 20-22*.
- Swapnil Gupta, Sreyash Kenkre, and Partha P. Talukdar. 2019. Care: Open knowledge graph embeddings. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP, Hong Kong, China, November 3-7, 2019*, pages 378–388.
- Raia Hadsell, Sumit Chopra, and Yann LeCun. 2006. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2006), 17-22 June 2006, New York, NY, USA*, pages 1735–1742. IEEE Computer Society.
- Kaveh Hassani and Amir Hosein Khas Ahmadi. 2020. Contrastive multi-view representation learning on graphs. In *Proceedings of the 37th International Conference on Machine Learning, ICML, 13-18 July, Virtual Event, volume 119 of Proceedings of Machine Learning Research*, pages 4116–4126. PMLR.
- Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross B. Girshick. 2020. Momentum contrast for unsupervised visual representation learning. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA*,

- June 13-19, 2020, pages 9726–9735. Computer Vision Foundation / IEEE.
- Yankai Lin, Zhiyuan Liu, Maosong Sun, Yang Liu, and Xuan Zhu. 2015. Learning entity and relation embeddings for knowledge graph completion. In *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*, pages 2181–2187.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.
- Dai Quoc Nguyen, Tu Dinh Nguyen, Dat Quoc Nguyen, and Dinh Q. Phung. 2018. A novel embedding model for knowledge base completion based on convolutional neural network. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, New Orleans, Louisiana, USA, June 1-6, 2018, Volume 2 (Short Papers)*, pages 327–333.
- Chao Shang, Yun Tang, Jing Huang, Jinbo Bi, Xiaodong He, and Bowen Zhou. 2019. End-to-end structure-aware convolutional networks for knowledge base completion. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 3060–3067.
- Fan-Yun Sun, Jordan Hoffmann, Vikas Verma, and Jian Tang. 2020. Infograph: Unsupervised and semi-supervised graph-level representation learning via mutual information maximization. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Haitian Sun, Tania Bedrax-Weiss, and William W. Cohen. 2019. Pullnet: Open domain question answering with iterative retrieval on knowledge bases and text. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, pages 2380–2390. Association for Computational Linguistics.
- Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. 2016. Complex embeddings for simple link prediction. In *Proceedings of the 33rd International Conference on Machine Learning, ICML, New York City, NY, USA, June 19-24, 2016*, volume 48 of *JMLR Workshop and Conference Proceedings*, pages 2071–2080.
- Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-sne. *Journal of machine learning research*, 9(11).
- Shikhar Vashishth, Prince Jain, and Partha P. Talukdar. 2018. CESI: canonicalizing open knowledge bases using embeddings and side information. In *Proceedings of the 2018 World Wide Web Conference on World Wide Web, WWW, Lyon, France, April 23-27, 2018*, pages 1317–1327.
- Petar Velickovic, William Fedus, William L. Hamilton, Pietro Liò, Yoshua Bengio, and R. Devon Hjelm. 2019. Deep graph infomax. In *7th International Conference on Learning Representations, ICLR, New Orleans, LA, USA, May 6-9*. OpenReview.net.
- Liang Wang, Wei Zhao, Zhuoyu Wei, and Jingming Liu. 2022. Simkgc: Simple contrastive knowledge graph completion with pre-trained language models. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, pages 4281–4294. Association for Computational Linguistics.
- Zhen Wang, Jianwen Zhang, Jianlin Feng, and Zheng Chen. 2014. Knowledge graph embedding by translating on hyperplanes. In *Proceedings of the Twenty-Eighth AAAI Conference on Artificial Intelligence, July 27-31, 2014, Québec City, Québec, Canada*, pages 1112–1119.
- Bishan Yang, Wen-tau Yih, Xiaodong He, Jianfeng Gao, and Li Deng. 2015. Embedding entities and relations for learning and inference in knowledge bases. In *3rd International Conference on Learning Representations, ICLR, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Yanqiao Zhu, Yichen Xu, Feng Yu, Qiang Liu, Shu Wu, and Liang Wang. 2021. Graph contrastive learning with adaptive augmentation. In *WWW '21: The Web Conference 2021, Virtual Event / Ljubljana, Slovenia, April 19-23, 2021*, pages 2069–2080. ACM / IW3C2.

## A Appendix

### A.1 Gradient Update

**Gradients of Contrastive Entity** Taking the Eq. (4) as an example, we give the detailed reasoning steps of gradient about  $h, r, t^+, t^-$ .

$$\begin{aligned}
-\frac{\partial S(h, r, t^+)}{\partial h} &= \frac{\partial}{\partial h} \left( \log \frac{e^{\beta(h, r, t^+)/\tau}}{\sum_{n \in \{p_e, \mathcal{N}_e\}} e^{\beta(n)/\tau}} \right) \\
&= \frac{\partial}{\partial h} \left( \frac{\beta(h, r, t^+)}{\tau} - \log \sum_{n \in \{p_e, \mathcal{N}_e\}} e^{\frac{\beta(n)}{\tau}} \right) \\
&= \frac{\partial}{\partial h} \left( \frac{\beta(h, r, t^+)}{\tau} - \log \left( e^{\frac{\beta(h, r, t^+)}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\beta(h, r, t_j^-)}{\tau}} \right) \right) \\
&= \frac{\partial}{\partial h} \left( \frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau} - \log \left( e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right) \right) \\
&= \frac{\varphi'(h, r)}{\tau} \mathbf{t}^+ - \frac{e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} \frac{\varphi'(h, r)}{\tau} \mathbf{t}^+ + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \frac{\varphi'(h, r)}{\tau} \mathbf{t}_j^-}{e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}}} \\
&= \frac{\varphi'(h, r)}{\tau} \left( \frac{\sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}}}{e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}}} \right) \mathbf{t}^+ \\
&\quad - \frac{\sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \mathbf{t}_j^-}{e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}}} \\
&= \frac{\varphi'(h, r)}{\tau A} \left( \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right) \mathbf{t}^+ \\
&\quad - \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \mathbf{t}_j^- \tag{14}
\end{aligned}$$

$$A = e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \tag{15}$$

Parameters of entity  $h$  can be updated from contrastive with the gradient in Eq. (14):

$$\begin{aligned}
\mathbf{h}^{l+1} &= \mathbf{h}^l - \eta \left( \frac{\partial S(h, r, t^+)}{\partial h} \right) \\
&= \mathbf{h}^l + \frac{\varphi'(h, r) \eta}{\tau A} \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \mathbf{t}^+ \\
&\quad - \frac{\varphi'(h, r) \eta}{\tau A} \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \mathbf{t}_j^- \tag{16}
\end{aligned}$$

$$\frac{\partial S(h, r, t^+)}{\partial r} = \frac{\partial S(h, r, t^+)}{\partial h} \tag{17}$$

$$\begin{aligned}
-\frac{\partial S(h, r, t^+)}{\partial t^+} &= \frac{\partial}{\partial t^+} \left( \frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau} - \log \left( e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right) \right) \\
&= \frac{\varphi(h, r)}{\tau} - \frac{e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} \frac{\varphi(h, r)}{\tau}}{e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}}} \tag{18}
\end{aligned}$$

$$\begin{aligned}
-\frac{\partial S(h, r, t^+)}{\partial t_j^-} &= \frac{\partial}{\partial t_j^-} \left( \frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau} - \log \left( e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right) \right) \\
&= \frac{\varphi(h, r)}{\tau A} \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \\
-\frac{\partial S(h, r, t^+)}{\partial t_j^-} &= \frac{\partial}{\partial t_j^-} \left( \frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau} - \log \left( e^{\frac{\varphi(h, r) \cdot \mathbf{t}^+}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right) \right) \\
&\quad + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} \left( \frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau} - \log \left( e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \right) \right) \\
&= \frac{-e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \frac{\varphi(h, r)}{\tau}}{e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} + \sum_{(h, r, t_j^-) \in \mathcal{N}_e} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}}} \\
&= -\frac{\varphi(h, r)}{\tau A} e^{\frac{\varphi(h, r) \cdot \mathbf{t}_j^-}{\tau}} \tag{19}
\end{aligned}$$

**Gradients of Contrastive Relation** Taking the Eq. (8) as an example, we give the detailed reasoning steps of gradient about  $h, t, r^+, r^-$ .

$$\begin{aligned}
-\frac{\partial S(h, r^+, t)}{\partial h} &= \frac{\partial}{\partial h} \left( \frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau} - \log \left( e^{\frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau}} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \right) \right) \\
&= \frac{\varphi'(h, r^+)}{\tau} \mathbf{t} - \frac{e^{\frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau}} \frac{\varphi'(h, r^+)}{\tau} \mathbf{t} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \frac{\varphi'(h, r_j^-) \mathbf{t}}{\tau}}{e^{\frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau}} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}}} \\
&= \frac{\varphi'(h, r^+)}{\tau B} \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} - \frac{1}{B} \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \frac{\varphi'(h, r_j^-) \mathbf{t}}{\tau} \tag{20}
\end{aligned}$$

$$\begin{aligned}
-\frac{\partial S(h, r^+, t)}{\partial t} &= \frac{\partial}{\partial t} \left( \frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau} - \log \left( e^{\frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau}} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \right) \right) \\
&= \frac{\varphi(h, r^+)}{\tau} - \frac{e^{\frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau}} \frac{\varphi(h, r^+)}{\tau} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \frac{\varphi(h, r_j^-)}{\tau}}{e^{\frac{\varphi(h, r^+) \cdot \mathbf{t}}{\tau}} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}}} \\
&= \frac{\varphi(h, r^+)}{\tau B} \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} - \frac{1}{B} \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\frac{\varphi(h, r_j^-) \cdot \mathbf{t}}{\tau}} \frac{\varphi(h, r_j^-)}{\tau} \tag{21}
\end{aligned}$$

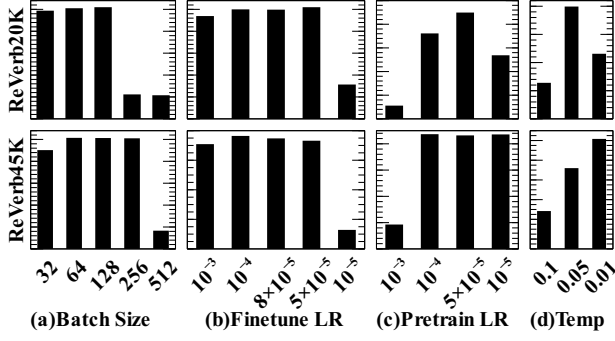


Figure 6: Results of hyperparameters on ReVerb20K (upper) and ReVerb45K (Below).

$$\begin{aligned}
& -\frac{\partial S(h, r^+, t)}{\partial r_j^-} \\
&= \frac{\partial}{\partial r_j^-} \left( \frac{\varphi(h, r^+) \cdot t}{\tau} - \log \left( e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} \right) \right. \\
&\quad \left. + \sum_{(h, r_j^-, t) \in \{p_r, \mathcal{N}_r - (h, r_j^-, t)\}} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} \right) \\
&= \frac{-e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} \frac{\varphi'(h, r_j^-) t}{\tau}}{e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} + \sum_{(h, r_j^-, t) \in \{p_r, \mathcal{N}_r - (h, r_j^-, t)\}} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)}} \\
&= -\frac{\varphi'(h, r_j^-) t}{\tau B} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} \quad (22)
\end{aligned}$$

$$\begin{aligned}
& -\frac{\partial S(h, r^+, t)}{\partial r^+} = \frac{\partial}{\partial r^+} \left( \log \frac{e^{\left( \frac{\beta(h, r^+, t)}{\tau} \right)}}{\sum_{n \in \{p_r, \mathcal{N}_r\}} e^{\left( \frac{\beta(n)}{\tau} \right)}} \right) \\
&= \frac{\partial}{\partial r^+} \left( \frac{\beta(h, r^+, t)}{\tau} - \log \sum_{n \in \{p_r, \mathcal{N}_r\}} e^{\left( \frac{\beta(n)}{\tau} \right)} \right) \\
&= \frac{\partial}{\partial r^+} \left( \frac{\beta(h, r^+, t)}{\tau} - \log \left( e^{\left( \frac{\beta(h, r^+, t)}{\tau} \right)} \right) \right. \\
&\quad \left. + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\beta(h, r_j^-) \cdot t}{\tau} \right)} \right) \\
&= \frac{\partial}{\partial r^+} \left( \frac{\varphi(h, r^+) \cdot t}{\tau} - \log \left( e^{\left( \frac{\varphi(h, r^+) \cdot t}{\tau} \right)} \right) \right. \\
&\quad \left. + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} \right) \\
&= \frac{\varphi'(h, r^+) t}{\tau} - \frac{e^{\left( \frac{\varphi(h, r^+) \cdot t}{\tau} \right)} \frac{\varphi'(h, r^+) t}{\tau}}{e^{\left( \frac{\varphi(h, r^+) \cdot t}{\tau} \right)} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)}} \quad (23) \\
&= \frac{\varphi'(h, r^+) t}{\tau} \left( \frac{\sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)}}{e^{\left( \frac{\varphi(h, r^+) \cdot t}{\tau} \right)} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)}} \right) \\
&= \frac{\varphi'(h, r^+) t}{\tau B} \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)}
\end{aligned}$$

$$B = e^{\left( \frac{\varphi(h, r^+) \cdot t}{\tau} \right)} + \sum_{(h, r_j^-, t) \in \mathcal{N}_r} e^{\left( \frac{\varphi(h, r_j^-) \cdot t}{\tau} \right)} \quad (24)$$

Model	AR ↓	ARR	H@1	H@10	H@50	H@100
<i>BERT-F</i>	754	28.1	22.1	39.1	53.8	61.2
<i>BERT-U</i>	1315	23.7	18.7	33.2	45.9	51.8
<i>BiGRU</i>	<b>393</b>	<b>38.9</b>	<b>30.5</b>	<b>54.6</b>	<b>69.2</b>	<b>75.5</b>
<i>BERT-F</i>	1165	28.8	21.6	42.8	55.3	60.4
<i>BERT-U</i>	2131	18.0	12.5	28.9	41.3	46.6
<i>BiGRU</i>	<b>767</b>	<b>33.3</b>	<b>25.3</b>	<b>48.7</b>	<b>63.0</b>	<b>68.3</b>

Table 4: Results of Textual Technology on ReVerb20K (upper) and ReVerb45K (below).

## A.2 Text Encoder

Results of textual technologies to encoder sequences in §4.1: BiGRU and BERT, are shown in Fig. 5, where *BERT-F* fixes the parameters of Bert, *BERT-U* unfixes them. By observing the results, performance of model with BERT (both fixed and unfixed parameters) as encoder is weaker than that with BiGRU. This shows that BiGRU is more capable of capturing sequential information in KGs, and pretrained language model can not give any help to TernaryCL. State-of-the-art baseline OKGIT also points out that pretrained language models can not predict the correct entities on the top. It could be inadvisable to introduce the pretrained language model into the proposed model due to the difficulty of training and reproduction. The intention of constructing a KG is to convert text into structures for easy calculation, that is, KG itself is rich in world knowledge. So, our TernaryCL pays attention to mining own structure features in an effective way.