



# Continual Contrastive Finetuning Improves Low-Resource Relation Extraction

Wenxuan Zhou<sup>1</sup>, Sheng Zhang<sup>2</sup>, Tristan Naumann<sup>2</sup>, Muhao Chen<sup>1</sup>, Hoifung Poon<sup>2</sup> University of Southern California<sup>1</sup>, Microsoft Research<sup>2</sup>



School of Engineering





### **Motivation**



### Manual annotation for relation extraction is expensive



#### Unlabeled data is abundant and easy to acquire



WikipediA The Free Encyclopedia

6.6M articles



35M citations and abstracts

180K research papers

How to use unlabeled data to improve relation extraction?



2 University of Southern California

### Matching the Blanks (MTB)



#### Harris's distributional hypothesis to relations (Soares et al., 19):

Context linking the same/different entities is more likely to express the same/different relation.





### **MTB-based Pretraining (Soares et al., 19)**



Positive instances: instances with the same entities

Negative instances: instances with different entities

Goal: make embedding of positive/negative instance pairs similar/dissimilar.





### **Continual Contrastive Finetuning**



#### **Previous work**



Continual contrastive finetuning





### **Contrastive Pretraining**



*H*': all relation embedding of all instances except for *t*.  $\tau$ : temperature

#### MTB-based contrastive loss:

Make embedding between *t* and its positive/negative instances to be similar/dissimilar.

$$L_{\text{mtb}} = -\frac{1}{|P|} \sum_{h_1 \in P} \log\left(\frac{e^{\cos(h,h_1)/\tau}}{Z}\right)$$
$$Z = e^{\cos(h,h_1)/\tau} + \sum_{h_2 \in N} e^{\cos(h,h_2)/\tau}$$

#### Self-supervised contrastive loss:

Make the similarity between two embeddings of t to be larger than t and other instances.

$$L_{\text{self}} = -\log\left(\frac{e^{\cos(h,\hat{h})/\tau}}{Z}\right)$$
$$Z = \sum_{h_2 \in H'} e^{\cos(h,h_2)/\tau}$$

Masked language modeling loss  $L_{mlm}$ .



### **Distributional Gap**

Make embedding from the same/different classes similar/dissimilar



#### Classic finetuning objectives:

Softmax classifier with cross-entropy (CE), supervised contrastive loss (SupCon) Distributional gap:

- CE and SupCon are minimized when representations form a single cluster for a class. (Graf et al., 2021)
- Representations from MTB-based pretraining may form multiple clusters for a class.





## **Distributional Gap (Cont.)**



Probing analysis:

Given an MTB-based pretrained RE model, fix the model parameters and fit different classifiers on top of it.

Classifiers:

- Single-cluster: softmax classifier, nearest centroid classifier
- Multi-cluster: kNN



KNN greatly outperforms both softmax and nearest centroid  $\Rightarrow$  MTB-based pretraining generates multi-cluster representations.





t-SNE visualization of relation embedding



#### 10 University of Southern California

- 1. MCCL consistently outperforms CE and SupCon in low-resource settings.
- 2. MTB+CE outperforms PLM+CE, showing that MTB pretraining is effective.
- 70 80 70 60 60 50 50 40 40 30 30 20 20 1% 5% 10% 100% 1% 5% 10% 100% PLM+CE MTB+CE MTB+SupCon MTB+MCCL ■ PLM+CE ■ MTB+CE MTB+MCCL

#### Models: PubmedBERT for BioRED, BERT for Re-DocRED. Evaluation metric: F1



School of Engineering



Re-DocRED

# **Experiments: Main Results**

Datasets: BioRED, Re-DocRED (multi-label)

BioRED

### **Experiments: Additional Analysis**



### Ablation Study BioRED (MCCL)



- 1. All the pretraining objectives are effective.
- Removing MTB leads to the largest drop ⇒ MTB is critical for lowresource RE

#### $F_1$ w.r.t. different % of data



- MCCL consistently outperforms CE when < 20% of data (~80 abstracts) are used.
- 2. MCCL performs similarly to CE with abundant training data.



### Conclusion



- We propose to pretrain the PLMs based on our improved MTB objective and show that it greatly improves PLM performance in low-resource document-level RE.
- 2. We bridge the gap of learning objectives between RE pretraining and finetuning with continual contrastive finetuning and kNN-based inference, helping the RE model leverage pretraining knowledge.
- 3. We design a multi-cluster contrastive learning objective, allowing one relation to form multiple different clusters, thus further reducing the distributional gap between pretraining and finetuning.

