

Highlights

Active Relation Discovery: Towards General and Label-aware Open Relation Extraction

Yangning Li, Yinghui Li, Xi Chen, Hai-Tao Zheng, Ying Shen

- We reveal two major shortcomings of previous OpenRE: Insufficient capacity to discriminate between known and novel relations, and requiring additional secondary labeling. We also propose a practical test setup (General OpenRE) to appeal to the information extraction community to focus more on how the model performs in real-world scenarios.
- We propose ARD, a framework that not only adapts to the General OpenRE utilizing relational outlier detection, but also exploits active learning to assign more meaningful and human-readable labels to novel relations. ARD offers a new, feasible and practical perspective for solving OpenRE.
- Extensive experiments on both conventional and General OpenRE settings show that ARD achieve significant improvements in three real-world datasets.

Active Relation Discovery: Towards General and Label-aware Open Relation Extraction

Yangning Li^a, Yinghui Li^a, Xi Chen^b, Hai-Tao Zheng^{a,d,*} and Ying Shen^{c,*}

^aShenzhen International Graduate School, Tsinghua University, Shenzhen, 518055, Guangdong, China

^bPlatform and Content Group, Tencent, Shenzhen, 518055, Guangdong, China

^cSchool of Intelligent Systems Engineering, Sun Yat-Sen University, Guangzhou, 510275, Guangdong, China

^dPeng Cheng Laboratory, Shenzhen, 518055, Guangdong, China

ARTICLE INFO

Keywords:

Open Relation Extraction
Information Extraction
Outlier Detection
Natural Language Processing

ABSTRACT

Open Relation Extraction (OpenRE) aims to discover novel relations from open domains. Previous OpenRE methods mainly suffer from two problems: (1) Insufficient capacity to discriminate between known and novel relations. When extending conventional test settings to a more general setting where test data might also come from seen classes, existing approaches have a significant performance decline. (2) Secondary labeling must be performed before practical application. Existing methods cannot label human-readable and meaningful types for novel relations, which is urgently required by the downstream tasks. To address these issues, we propose the Active Relation Discovery (ARD) framework, which utilizes relational outlier detection for discriminating known and novel relations and involves active learning for labeling novel relations. Extensive experiments on three real-world datasets show that ARD significantly outperforms previous state-of-the-art methods on both conventional and our proposed general OpenRE settings. The source code and datasets will be available for reproducibility.

1. Introduction

Novel relations are cropping up at a rate of tens of thousands per year [44], while most of the rapidly emerging relations are still unlabeled and under-explored, mixed with pre-defined relations. These relations cannot be well handled by supervised RE methods due to the fixed pre-defined relation schema. Therefore, Open Relation Extraction (OpenRE) aims at discovering and extracting potential novel relations from open-domain corpora.

Some recent preliminary studies [49, 50] have noticed the challenge of learning emerging relations and explored methods for OpenRE. Previous works can be divided into two main paradigms: pattern-based and clustering-based methods. Specifically, pattern-based methods [1, 7] utilize statistical or neural approaches to heuristically extract relation patterns from sentences, then clustering-based methods [10, 50] are proposed to aggregate instances representing the same novel relation. However, previous works mainly have two shortcomings in real scenarios:

(1) **The widely used traditional setting can't comprehensively reflect what OpenRE in the real world entails.**

The traditional OpenRE setup is that models are evaluated based on their ability to discriminate among unseen classes, assuming the absence of known relation during the test phase. This test setup is a good measure of the model's ability to learn novel relations, but ignores the model's ability to distinguish between the known and unseen relations. As we all know, the relation distribution in the real world is intricate, mixed with known and unseen relations. Therefore,

it's unrealistic to assume that we will never encounter known relations during the test stage.

In the light of above facts, we loosen the existing setting to a *General OpenRE* setting: test data might also come from known relations. Empirical experiments in Table 4 show that the previous state-of-the-art OpenRE models [19, 50, 55] perform poorly under this setting.

(2) **The results produced by previous OpenRE models require secondary labeling before they can be practically applied.** In other words, for a certainly discovered novel relation, the model cannot assign it a surface name with a specific semantic meaning. As the foundation of a series of downstream tasks, labels with actual meaning are inevitable. However, due to the absence of human knowledge, both pattern-based and clustering-based OpenRE methods lack the ability to name novel relation types as human-readable and meaningful. Pattern-based methods rely heavily on the surface phrase, yet relations between entities are often not directly represented by the span in the sentence. Clustering-based methods merely cluster instances that express the same relations, but do not provide concrete representations of the novel relations. Both methods require manual re-labeling of the novel relations found. This gap between model and practice hinders model application in real-world scenarios.

To address above mentioned issues, we propose the **Active Relation Discovery (ARD)** framework shown in Figure 1. Targeted improvements are made in two aspects: (1) To avoid the model being confused by the set of mixed known and novel relations, we developed a relational outlier detection algorithm that separates known and novel relations by treating novel relations as outliers, performing stably under the General OpenRE setting. (2) To assign meaningful

*Corresponding author

✉ liyn20@mails.tsinghua.edu.cn (Y. Li);

zheng.haitao@sz.tsinghua.edu.cn (H. Zheng)

ORCID(s): 0000-0002-1991-6698 (Y. Li)

labels to novel relations, the incorporation of human knowledge is inevitable. To minimize the labor cost, we propose an active learning algorithm. Specifically, we introduce the *representative instance*, which denotes an instance can offer rich information of unknown relations. Only a handful of representative instances requires manual labeling, and then the model can automatically label the novel relations in a supervised manner.

In summary, our contributions are in three folds:

(1) We reveal two major shortcomings of previous OpenRE, and introduce a new setting called *General OpenRE*, which can realistically measure the capabilities of the OpenRE model.

(2) We propose ARD, a practical framework that not only adapts to the *General OpenRE* utilizing relational outlier detection, but also exploits active learning to assign more meaningful and human-readable labels to novel relations.

(3) We conduct extensive experiments on both conventional and General OpenRE settings to show that our framework can achieve significant improvements in three real-world datasets. Detailed analyses demonstrate the effectiveness of each component of our proposed ARD.

2. Related Work

Open Relation Extraction. Whereas supervised RE [24, 28, 53, 57] relies heavily on manual annotation and the inherent inadequacy of predefined relation schema, OpenRE gains increasing attention. The method of OpenRE can be broadly divided into two categories: pattern-based and clustering-based. Pattern-based approaches extract relation patterns from textual corpora [3, 11, 47, 48]. These methods apply heuristic algorithms to describe relations between marked entities with relation patterns consisting of several key phrases in texts. Due to the ambiguity of relations obtained by the pattern-based methods, the focus of research in recent years has been primarily on clustering-based methods.

Clustering-based method [10, 19, 45, 50, 55] cluster instances in the feature space into novel relation types. [50] enhances unsupervised clustering-based methods by introducing Siamese Network to measure instance similarity. Considering the inherent connection between OpenRE and relation hierarchies, [55] proposes a framework to effectively integrate hierarchy information into relation representations for better novel relation extraction.

As described in Section 1, there are two main problems with the current OpenRE: (1) They focus only on the discrimination of novel relations, supposing that test sets only have novel relations. (2) The model output is not directly usable by downstream tasks. In response, we propose a General OpenRE setup and incorporate outlier detection and active learning into OpenRE.

Active Learning in Relation Extraction. The key idea behind active learning [43] is that the learning algorithm is allowed to ask for true meaningful labels of some selected unlabelled instances. Various criteria [12, 35, 54] have been proposed to choose these instances on traditional supervised

RE tasks. To our best knowledge, we firstly integrate active learning into OpenRE, enabling meaningful tags of the novel relation type with the addition of human knowledge.

Generalized Zero-Shot Learning (GZSL). The motivation for the General OpenRE setting is similar to that of the GZSL. Traditionally, ZSL approaches [38, 56] assume that only the unseen classes are present in the test set. [5] first challenged this implausible setting and proposed the GZSL setting: test data might also come from seen classes. GZSL approaches [20, 37] focus on mitigating the strong bias caused by known classes and preventing novel classes from being categorized as one of the seen classes. While in our General OpenRE setting, we concentrate more on the distinction between known and novel classes.

3. Task Formulation

General OpenRE formulates the task slightly differently from traditional OpenRE setting. The original train set is a large-scale manually annotated corpus $\mathcal{X} = \{x_j^i | r_i \in \mathcal{R}_K\}$, where relations in \mathcal{R}_K are pre-defined as “known relations”. Obviously, we assume that there exists a relation set \mathcal{R}_N that contains “novel relations” in another corpus without annotations. In the real-world scenario, we need to process the dataset whose instances express relations both in \mathcal{R}_K and \mathcal{R}_N , distinguish known and novel relations, then label each instance.

Under this fact, we first consider the *novel relation discovery*, in which we solely focus on the mining of unseen relations. At this stage, we pre-train the model on \mathcal{X} and obtain a trained encoder E . Then for a concrete dataset (test set) $\mathcal{X}' = \{x_j^i, x_j^{i'} | r_i \in \mathcal{R}_K, r_i' \in \mathcal{R}_N\}$. The model will unsupervisedly divide \mathcal{X}' into a “known relation set” \mathcal{X}_K and a “novel relation set” \mathcal{X}_N .

\mathcal{X}_K can be easily labeled for sufficient information obtained from \mathcal{X} . Secondly, we focus on the *annotation of novel relations* \mathcal{X}_N . In this phase, we integrate the intuition of active learning by utilizing limited labor to facilitate the novel relation annotation performance. Our model queries a small set of informative samples in \mathcal{X}_N for manual labeling and then trains a classifier to annotate novel relations.

4. Methodology

4.1. Overview

The overview of the method is illustrated in Figure 1. We will detailedly introduce our work into three components: (1) **Relation representation**, in which we extend to transform semantic relations into low-dimension dense representations. (2) **Relational Outlier Detection**, where the model automatically detects a novel relation set from real-world datasets and feeds them into the active learning stage. (3) **Relational Active Learning**, where the model selects the most informative instances to train a powerful classifier for novel relation.

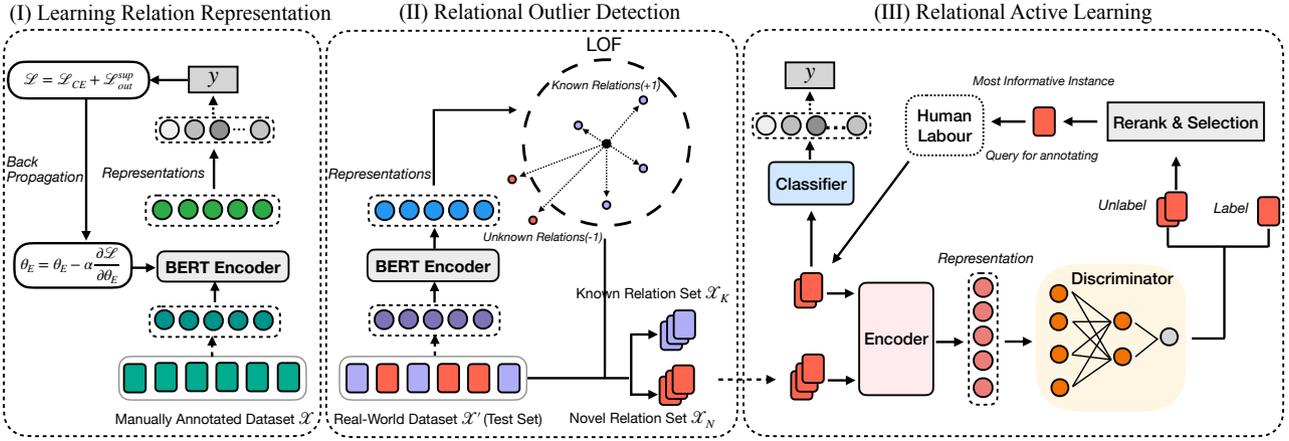


Figure 1: An illustration of our proposed Active Relation Discovery (ARD) framework.

4.2. Relation Representation

Given a dataset $X = \{x_1, \dots, x_n\}$, an instance x is a word (token) sequence $\{w_1, w_2, \dots, w_n\}$ with two marked entities e_h and e_t . We use triplets of relation facts (e_h, r, e_t) to denote that there is a relation r between the marked entity pair. And x^r indicates an instance that expresses the relation r . Specifically, we define four special markers $\langle e_h \rangle$, $\langle /e_h \rangle$, $\langle e_t \rangle$, and $\langle /e_t \rangle$ to locate the head entity and the tail entity. We denote the indices of $\langle e_h \rangle$ and $\langle e_t \rangle$ as $\text{START}(h)$ and $\text{START}(t)$. An instance is represented as:

$$x = \dots, \langle e_h \rangle, w_{\text{START}(h)+1}, \dots, w_{\text{END}(h)}, \langle /e_h \rangle, \dots, \langle e_t \rangle, w_{\text{START}(t)+1}, \dots, w_{\text{END}(t)}, \langle /e_t \rangle, \dots \quad (1)$$

We use pre-trained language model (i.e. BERT [8]) to encode each token w_i to the corresponding representation $\mathbf{h}_i \in \mathbb{R}^d$, where d is denotes the dimension of representation vectors.

For an instance $x_i \in \mathcal{S}$, we use the concatenation of representations of two start positions ($w_{\text{START}(h)}$ and $w_{\text{START}(t)}$) as the representation of the relation:

$$\mathbf{h}_r(x_i) = [\mathbf{h}_{\text{START}(h)}, \mathbf{h}_{\text{START}(t)}], \quad (2)$$

These extra tokens play a similar role like position embeddings in conventional RE tasks [51]. The relation representation $\mathbf{h}_r(x_i)$ will be utilized to predict the relation type r .

As mentioned previously, \mathcal{X} are used to fine-tune the pre-trained language model. Notably, along with the traditional cross-entropy loss, we integrate a supervised contrastive loss $\mathcal{L}_{\text{out}}^{\text{sup}}$ ¹:

$$\mathcal{L}_{\text{out}}^{\text{sup}} = \sum_{i \in I} \frac{-1}{|P(i)|} \sum_{p \in P(i)} \log \frac{\exp(\mathbf{z}_i \cdot \mathbf{z}_p / \tau)}{\sum_{a \in A(i)} \exp(\mathbf{z}_i \cdot \mathbf{z}_a / \tau)}, \quad (3)$$

Here, $P(i) \equiv \{p \in B \setminus \{i\} : \tilde{\mathbf{y}}_p = \tilde{\mathbf{y}}_i\}$ is the set of indices of all positives in the mini-batch B distinct from i . $\mathbf{z}_i = \text{Proj}(\mathbf{h}_r(x_i)) \in \mathbb{R}^{D_p}$, where Proj is a single

¹Scalar temperature parameter τ is 0.1 as in [23]. We refer to [23] for more details.

linear layer outputs vector of size $D_p = 128$. Contrastive loss [23, 25, 26] allows for tighter clustering of intra-class instances and a more dispersed distribution of inter-class instances. The essence behind the employment of contrastive loss is to gain relation representations that are more friendly to outlier detection and active learning. The performance of our relation representation on supervised RE can also be found in Section 5.9.2.

4.3. Relational Outlier Detection

After pre-training, E_θ could encode an instance x into a dense vector $\mathbf{h}_r(x)$ as the relation representation. In the feature space, due to the similarity of the semantics, representations that express the same relation tend to densely gather (forming n separate clusters) and ones that express different relations tend to disperse. Figure 3 illustrated the distribution of different representations. Since the instances express unseen relations have not been pre-trained, in other words, the model has not seen the semantics, the instances are not projected near any clusters. We utilize this property to design local outlier factor (LOF) to reflect the local density of instances in the feature space.

Formally, given any two representations $\mathbf{h}_r(x_i), \mathbf{h}_r(x_j)$ of instances x_i, x_j , we denote $d(\mathbf{h}_r(x_i), \mathbf{h}_r(x_j))$ as the Euclidean distance between them. Then, we define k -th distance, denoted as $d_k(\mathbf{h}_r(x_i))$, to represent the distance from $\mathbf{h}_r(x_i)$ to the k -th nearest neighbour. The reachability distance between $\mathbf{h}_r(x_i)$ and $\mathbf{h}_r(x_j)$ is:

$$\text{rd}_k(\mathbf{h}_r(x_i), \mathbf{h}_r(x_j)) = \max\{d_k(\mathbf{h}_r(x_j)), d(\mathbf{h}_r(x_i), \mathbf{h}_r(x_j))\}, \quad (4)$$

We then compute the density to measure the average distance of reachability distance:

$$\text{den}_k(\mathbf{h}_r(x_i)) = 1 / \frac{\sum_{\mathbf{h}_r(x_j) \in N_k(\mathbf{h}_r(x_i))} \text{rd}_k(\mathbf{h}_r(x_i), \mathbf{h}_r(x_j))}{|N_k(\mathbf{h}_r(x_i))|}, \quad (5)$$

where $N_k(\mathbf{h}_r(x_i))$ denotes all the points within in k -th distance of $\mathbf{h}_r(x_i)$.

The computation of local outlier factor is:

$$\text{LOF}_k(\mathbf{h}_r(x_i)) = \frac{\sum_{\mathbf{h}_r(x_j) \in N_k(\mathbf{h}_r(x_i))} \frac{\text{den}_k(\mathbf{h}_r(x_j))}{\text{den}_k(\mathbf{h}_r(x_i))}}{|N_k(\mathbf{h}_r(x_i))|}, \quad (6)$$

where the larger LOF is, the more likely $\mathbf{h}_r(x_i)$ is an outlier point, i.e., an instance that expresses a novel relation. Our model could unsupervisedly detect the instances with novel relations.

4.4. Relational Active Learning

To this end, the model could divide the real-world dataset into a “known relation set” \mathcal{X}_K and an “novel relation set” \mathcal{X}_N . In view of the fact that \mathcal{X}_K can be conveniently and precisely annotated, we focus on labeling meaningful types for discovered instances in \mathcal{X}_N in this subsection.

To retrieve human-readable labels and avoid subsequent secondary labeling, we need to incorporate human knowledge into the relation learning phase through active learning. Our primary goal is to find a small part of instances with the most information and artificially label them. Then we use the labeled data to train a classifier in a supervised manner. The problem of how to find instances with most information essentially is the problem of how to find the instances that are most likely to express “novel relations”. Inspired by this, we propose the following Relation Active Learning module:

In the beginning, we randomly label a small part of data in \mathcal{X}_N . The labeled dataset is denoted as \mathcal{X}_L and the rest of the unlabeled data is denoted as \mathcal{X}_U . We assume that all the instances x are i.i.d according to a latent distribution $P(x)$. Correspondingly, their labels are distributed by the conditional distribution $P(y|x)$.

Neural Encoder We adopt a neural encoder to learn the distribution of \mathcal{X}_L and \mathcal{X}_U in the latent feature space. Our framework is independent of the choice of neural encoders, in this case, we adopt BERT [8] as the encoder. The goal of the neural encoder is to encode \mathcal{X}_L and \mathcal{X}_U into the same feature space and try to fool a discriminator to correctly predict if the instance is “representative”. The loss function of the encoder is:

$$\mathcal{L}_e = -\mathbb{E}_{x \sim P_{\mathcal{X}_L}} [\log(D_\psi(E_\theta(x)))] - \mathbb{E}_{x \sim P_{\mathcal{X}_U}} [\log(1 - D_\psi(E_\theta(x)))] \quad (7)$$

Discriminator A binary classifier (or a discriminator): $\mathcal{X} \rightarrow \{-1, 1\}$ is adopted to select the most informative samples. We utilize adversarial training to leverage the information of both \mathcal{X}_L and \mathcal{X}_U . The discriminator is adversarially trained to accurately distinguish if the instance expresses a novel relation.² The loss function is a flipped version of the encoder:

$$\mathcal{L}_d = -\mathbb{E}_{x \sim P_{\mathcal{X}_L}} [\log(1 - D_\psi(E_\theta(x)))] - \mathbb{E}_{x \sim P_{\mathcal{X}_U}} [\log(D_\psi(E_\theta(x)))] \quad (8)$$

Naturally, we could jointly optimize the two objective functions by allocate two parameters: $\mathcal{L} = \lambda \mathcal{L}_e + \lambda' \mathcal{L}_d$.

Active learning At each training step, we select k instances with the highest confidence of the discriminator as the most informative instances. Then the instances will be manually annotated and then used to train the classifier. In our experiments, as in most active learning efforts, we use the golden label of the instance as the annotation result. At this point, the discussion of annotations needs to be further developed. Considering the explosive growth of the number of relations, an annotating process that supports online and continual learning of novel relations needs to be designed. Thus, we propose a practical and easy-to-implement annotation procedure. At the start, for each selected instance x_i , the annotator only needs to judge if x_i has the same relation class as any instances of \mathcal{X}_L . x_i will be indexed as a novel relation if it doesn't share the same relation with instances in \mathcal{X}_L , or labeled as one known relation. After the procedure, the labels of relations would be easy to design than before the active learning begins. This manner effectively ensures the ability to continual learning and online learning of our framework, expediently fitting the real situation. Subsequently, \mathcal{X}_L will be fed into a classifier, which is a one-layer MLP [30] with an output layer, optimized by cross-entropy objective function, denoted as \mathcal{L}_c and parameterized by γ :

$$\mathcal{L}_c = \sum_{i \in |\mathcal{X}_L|} -\log p(y_L^{(i)} | x_L^i, \gamma). \quad (9)$$

Algorithm 1 Training for Active Learner, λ, λ', k are hyper-parameters.

Input: Labeled data (\mathcal{X}_L, Y_L) , unlabeled data \mathcal{X}_U , initialized encoder model with θ , discriminator model with ψ , classifier with γ

while not converge **do**

Sample mini-batches (x_L, y_L) from (\mathcal{X}_L, Y_L) Sample mini-batches (x_U) from (\mathcal{X}_U)

Compute \mathcal{L}_e by Eq. 7

Update θ w.r.t \mathcal{L}_e

Compute \mathcal{L}_d by Eq. 8

Update ψ w.r.t \mathcal{L}_d

Select k most informative instances $\{x_1, \dots, x_k\}$ by the output of d

for $i \leftarrow 1$ to k **do do**

if x_i has the same relation as $x_j^r \in \mathcal{X}_L$ **then**

Label x_i with r and append x_i to \mathcal{X}_L

else

Label x_i with a new index and append x_i to \mathcal{X}_L

end if

end for

Update γ w.r.t \mathcal{L}_c

end while

²A single novel relation where it won't be picked to be labeled will eventually be labeled as a novel relation that has already been labeled, but this almost never happens.

5. Experiments

In this section, we verify the performance of the model on three large-scale OpenRE datasets and their variants, and at the same time, a series of auxiliary experiments are carried out to prove the effectiveness of the model. Finally, we give a detailed analysis of the efficacy of our ARD framework.

5.1. Baseline

To demonstrate the effectiveness of our ARD models, we compare our models with three state-of-the-art models: (1) **RSN-CV** [50] employs conditional entropy and virtual adversarial learning to train Siamese Network to measure instance similarity. (2) **SelfORE** [19] utilizes self-training to iteratively learn relation representations and clusters with the weak signals provided by large pretrained language model. (3) **OHRE** [55] integrate hierarchy information into relation representations for better novel relation extraction. For a fair comparison, we substitute all the encoding models in the baseline models with BERT_{LARGE}.

5.2. Datasets and Setting

Datasets Three datasets and their variants are used to evaluate our model: FewRel [15], New York Times Freebase(NYT+FB) [34] and FewRel2.0 [14], the first two of which have been widely used in previous RE works [19, 46, 55]. We follow the division of the datasets from previous works.

FewRel is one of the largest RE dataset. As in the previous work, we use the original train set of FewRel. The dataset contains 80 relation categories and 700 instances of each relation category. Among them, 64 relations are divided into the training set and the remaining 16 relations are chosen as the test set.

NYT+FB dataset aligns entities from the New York Times corpus with Freebase triplets. Following the setting in [46], we filter out sentence with non-binary relations and obtain 41,000 labeled sentences containing 262 relations. The training and test sets comprise 212 and 50 relations respectively.

To verify the cross-domain capability of the model comprehensively, we also use FewRel2.0 dataset whose training and test sets are from completely different domains. As an advanced version of FewRel, FewRel2.0 incorporates knowledge transferring. The test set of FewRel2.0 contains data of 10 relations (100 samples for each relation) in the biomedicine field, and the training set is exactly the same as FewRel. The statistics of the data set are shown in Table 1.

Datasets Processing As described above, in the original OpenRE setting, there are no overlapping relations in the training and test sets. The relations in the test set are all novel relations. To measure the performance of the model in our proposed *General OpenRE* setting, we resample the original dataset and gain two variants: *noisy* and *imbalanced*. In the test sets of the two variants, there exist known relations with different distributions. In other words, the original dataset corresponds to the conventional setting and the noisy and imbalanced variants to the general setting.

Table 1

Statistical results for the dataset. #CLS represents the number of relation types and #SUM stands for the number of samples. In the addition equation $x + y$ in the table, x and y are the statistics for the known and novel relations separately.

Dataset	Setting	Train		Test	
		#CLS	#SUM	#CLS	#SUM
FR	Ori	64	44,800	16	11,200
	Noi	64	40,320	64+16	4,480+11,200
	Imb	64	40,320	64+16	4,480+4,560
NYF	Ori	212	33,990	50	7,010
	Noi	212	30,591	212+50	3,399+7,010
FR2.0	Ori	64	44,800	10	1,000
	Noi	64	40,320	64+10	480+1,000
	Imb	64	40,320	64+10	480+720

Table 2

The discarding probabilities for different relations.

Dataset	Relation ID	P
FewRel	66-73	0.4
	74-77	0.7
	78-81	0.85
FewRel2.0	66-68	0.15
	69	0.2
	70	0.3
	71-72	0.35
	73	0.4
	74-75	0.45

To obtain the noisy variant, we randomly select 40% samples from original training sets. Given that the number of samples for each novel relation is identical in FewRel and FewRel2.0, we further construct the imbalanced variant to explore the performance of the model in the presence of class imbalance. Specifically, we build on the noisy variant by randomly discarding a portion of the samples with different probabilities for each relation class in the test set, yielding class imbalance in test set. The discarding probabilities for different relations are shown in the Table 2.

5.3. Evaluation Settings

Following previous works, we apply instance-level evaluation metrics to evaluate the model, covering B³ [2], V-measure [39] and Adjusted Rand Index(ARI) [22].

For quantitative validation, we divide \mathcal{X}_N into \mathcal{X}_N^{train} and \mathcal{X}_N^{test} , which account for 40% and 60% respectively. The active learning module selects the instance with the most information in \mathcal{X}_N^{train} and trains the relation classifier. In the test phase, we merge \mathcal{X}_K and \mathcal{X}_N^{test} , report metric scores

on it. As the baselines are semi-supervised, \mathcal{X}_N^{train} is also applied to the training of the baseline models to ensure a fair comparison.

For FewRel and NYT+FB, the seminal set size for Active Learning module is 32. The sample size k is 32 and we sample a total of 8 epochs. In other words, a whole of 288 samples is manually labeled. As for FewRel2.0, we choose a smaller sample size: $k = 8$ and keep seminal set size as 32. Finally, 96 informative samples are annotated.

5.4. Implementation Details and Hyper-parameter Choices

To improve the experimental effect, we use BERT_{LARGE} with 300M parameters in the relation representation module. We pre-train the BERT model on 3 epochs, and each epoch costs about 1 GTX 3090 GPU hour. For the discriminator, we constructed a 3-layer fully connected neural network. For active learning, λ and λ' are both 1. For optimization, different models use different optimizers. Specifically, BERT use AdamW [32] with a learning rate of 0.00002, for discriminator, we use Adam with a learning rate of 0.0005, and for task learner of active learning, SGD is utilized. For baseline models, we follow their original setting without modifying any parameters except the division of the dataset.

5.5. Main Experiment

Table 4 shows the quantitative evaluation results on three datasets and their variants, from which we observe that: (1) Our ARD model outperforms state-of-the-art models by a large margin. Specifically, B^3 , V-measure and ARI increased by 10.5, 13.7, and 11.4 respectively compared to OHRE on FewRel. Compared with other semi-supervised methods, the gap is even larger, rises of over 20 are achieved by ARD. This proves that ARD can efficiently discover and learn representations of novel relations at a fraction of the labor cost. (2) A universal and consistent decline in performance of baseline models from the original datasets to noisy variants and then to unbalanced variants. This demonstrates that the *General OpenRE* setting is more challenging and more practical for the real scenario. The $F1$ score for RSN-CV drops dramatically from the original data to the noisy variant by 16.6. In contrast, the ARD model performs better on both the noisy and imbalanced variants than on the original dataset, even with a $F1$ score boosting by 7.2 on FewRel. This indicates the relation discovery procedure and relational active learning is robust in different scenarios. (3) The state-of-the-art models perform poorly on FewRel2.0. This is entirely to be expected, as the instances in the test set are from non-generic and low-resource domains such as biomedicine. ARD, on the other hand, still shows strong stability, confirming the cross-domain capability of the model. Further, to substantiate the applicability of our framework, we deploy ARD to a real medical dataset, as detailed in Section 5.10.

5.6. Analysis on Active Learning

The Efficiency of Active Learning Table 3 shows the results of our active learning approach compared to various active learning baseline models including DBAL [13],

Table 3

$F1$ -measure for various active learning methods on noisy datasets.

Dataset	Model	Epoch					
		#1	#2	#3	#4	#5	#6
FR (Noi)	DBAL	58.8	64.8	71.6	70.4	74.1	76.9
	CoreSet	60.1	61.8	66.1	68.4	70.9	75.4
	SRAAL	61.9	64.7	65.7	69.8	73.7	73.9
	Ours	66.0	69.0	70.5	72.7	75.5	78.5
NYF (Noi)	DBAL	47.4	48.6	51.4	53.3	54.9	55.5
	CoreSet	45.4	49.5	52.0	55.3	56.8	59.2
	SRAAL	50.2	51.9	54.0	55.6	56.2	56.9
	Ours	56.8	62.5	66.6	68.3	69.3	69.9
FR2.0 (Noi)	DBAL	46.9	50.5	51.4	51.7	52.2	53.7
	CoreSet	44.0	45.5	50.3	51.9	53.0	54.2
	SRAAL	45.0	49.7	51.8	52.0	52.8	53.9
	Ours	48.8	51.2	52.4	53.2	53.5	54.5

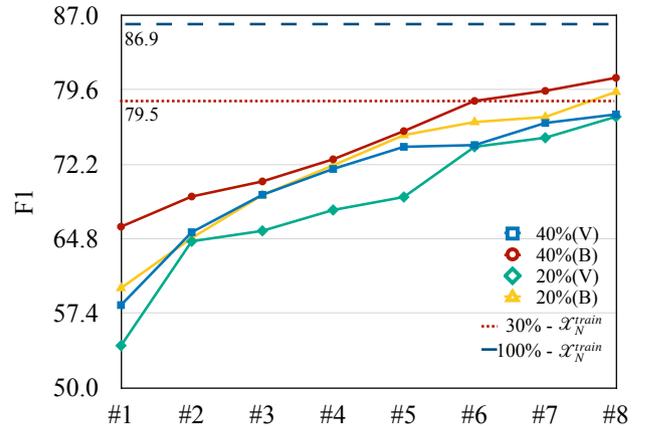


Figure 2: $F1$ -measure on noisy FewRel, (V) denotes the β -VAE and (B) denotes the BERT encoder.

CoreSet [42], SRAAL [52]. It can be observed that in each iteration, our model outperforms the other models, indicating that our method can consistently sample informative samples. In particular, our method performs significantly better on the NYT-FB dataset where the category count is much larger. Compared to the baseline models, our method ensures the information content and category diversity of the selected samples by enabling samples of the present batch to interact implicitly with previously selected samples through the discriminator. Besides, we report first 8 cases selected by discriminator which are regarded as the most representative instances (instances with novel relations) in Section 5.8.

The Impact of Different Encoder and Scope of Query

Figure 2 shows the experimental results on noisy FewRel with different encoders and query ranges. The “query ranges” represents the ratio of \mathcal{X}_N^{train} to \mathcal{X}_N . We also explore the impact of β -VAE [17] and BERT as encoders. From the results we observe that: (1) Generally, the model performance is proportional to the size of \mathcal{X}_N^{train} . However, the results are improved marginally as the number of samples increase.

Table 4

Main results on three original datasets and their variants. Ori, Noi, and Imb stand for original, noisy and imbalanced respectively. Ori corresponds to the conventional setting. Noi, and Imb refer to the general setting. Results are the average of 3 experiments with different random seeds.

Dataset	Model	B ³			V-measure			ARI
		F1	Prec.	Rec.	V	Hom.	Comp.	
FR (Ori)	RSN-CV	59.2 _{±0.5}	55.5	63.4	72.9 _{±0.3}	69.1	77.2	47.2 _{±0.9}
	SelfORE	60.6 _{±0.6}	60.2	61.1	68.4 _{±1.2}	67.5	69.3	56.0 _{±0.7}
	OHRE	63.1 _{±1.0}	54.9	74.1	71.4 _{±0.8}	64.9	79.4	52.7 _{±1.0}
	Ours	73.6_{±0.8}	70.7	76.8	85.1_{±0.1}	84.9	85.3	64.1_{±0.8}
NYF (Ori)	RSN-CV	49.7 _{±0.8}	39.3	67.6	64.2 _{±0.8}	56.7	73.9	35.1 _{±0.2}
	SelfORE	54.9_{±0.7}	52.8	57.2	72.5_{±0.3}	71.6	73.4	56.6_{±0.6}
	OHRE	41.2 _{±0.2}	28.5	74.3	54.1 _{±0.5}	42.7	73.7	26.5 _{±0.6}
	Ours	51.4 _{±0.5}	45.0	60.0	72.3 _{±0.4}	75.0	69.8	45.1 _{±0.7}
FR2.0 (Ori)	RSN-CV	27.7 _{±1.2}	18.2	58.2	48.8 _{±0.2}	39.5	63.7	13.4 _{±0.7}
	SelfORE	36.7 _{±0.8}	26.2	61.3	60.2 _{±0.8}	52.2	71.2	27.3 _{±0.3}
	OHRE	25.1 _{±0.6}	18.5	38.9	15.8 _{±0.3}	14.4	17.6	9.4 _{±0.7}
	Ours	48.8_{±0.7}	43.2	56.1	65.1_{±0.8}	60.8	70.1	34.4_{±0.8}
FR (Noi)	RSN-CV	42.6 _{±0.7}	30.1	72.6	66.6 _{±0.3}	56.4	81.2	28.3 _{±0.3}
	SelfORE	51.3 _{±0.8}	49.3	53.5	56.4 _{±0.5}	55.2	57.7	45.8 _{±0.4}
	OHRE	32.5 _{±0.2}	20.4	79.8	57.5 _{±1.2}	46.1	76.3	26.3 _{±1.2}
	Ours	80.8_{±0.9}	75.7	86.7	90.2_{±0.5}	89.0	91.4	71.3_{±0.5}
NYF (Noi)	RSN-CV	43.0 _{±0.5}	30.6	72.3	59.6 _{±0.1}	50.7	72.3	30.6 _{±0.7}
	SelfORE	48.6 _{±1.2}	43.8	54.6	65.7 _{±0.1}	65.1	66.3	46.2 _{±0.3}
	OHRE	36.4 _{±0.7}	25.1	66.2	48.1 _{±0.3}	38.2	64.6	30.2 _{±0.9}
	Ours	71.3_{±0.2}	60.8	86.2	72.9_{±0.3}	70.5	75.5	51.0_{±0.3}
FR2.0 (Noi)	RSN-CV	27.7 _{±0.2}	32.4	24.2	31.4 _{±0.4}	30.1	32.8	10.2 _{±0.5}
	SelfORE	32.8 _{±1.1}	24.7	48.9	54.8 _{±1.0}	47.3	65.1	27.1 _{±1.0}
	OHRE	25.2 _{±0.5}	15.9	60.5	50.0 _{±1.2}	40.1	66.4	16.2 _{±0.2}
	Ours	55.0_{±0.7}	52.8	57.4	69.3_{±0.7}	65.1	74.0	38.5_{±0.9}
FR (Imb)	RSN-CV	37.2 _{±1.0}	24.6	76.2	65.5 _{±0.6}	55.1	80.8	25.3 _{±0.5}
	SelfORE	48.3 _{±0.2}	44.2	53.5	53.7 _{±0.1}	56.4	51.3	44.2 _{±1.3}
	OHRE	31.0 _{±1.0}	19.8	71.1	56.1 _{±0.6}	44.1	77.6	22.6 _{±1.0}
	Ours	76.5_{±0.6}	74.7	78.4	86.5_{±0.7}	86.8	86.2	67.8_{±0.6}
FR2.0 (Imb)	RSN-CV	26.4 _{±0.4}	20.6	36.8	31.5 _{±0.3}	25.2	41.9	22.7 _{±0.7}
	SelfORE	31.3 _{±1.2}	22.4	52.2	52.9 _{±0.4}	45.8	62.6	25.7 _{±0.9}
	OHRE	22.6 _{±0.6}	13.9	60.7	45.5 _{±1.0}	35.4	63.6	13.4 _{±0.6}
	Ours	52.4_{±0.6}	50.1	54.9	67.4_{±0.2}	63.2	72.2	36.4_{±0.3}

But the model still yields better performance when the query range is 40%. (2) The comparisons between the VAE and the BERT encoder are in line with intuition. Although VAE is intuitive and can be more easily trained, BERT still shows superiority in empirical results.

Compare with Manual Random Selection. In Figure 2, the gain from 288 informative instances (approximately 8% of \mathcal{X}_N^{train}) selected by the active learning is similar to the gain from 30% of instances randomly selected. When trained with the full amount of \mathcal{X}_N^{train} , the F1 is 6.1% higher than ARD while costing 12 times as much in human effort.

Table 5

Comparisons of F1-measure between different sampling strategies on noisy FewRel dataset.

Epoch	Lowest	Random	Highest
#1	57.1	58.7	66.0
#2	57.1	60.4	69.0
#3	57.8	65.6	70.5
#4	57.6	67.2	72.7
#5	57.6	67.7	75.5
#6	58.1	67.7	78.5

Table 6

Average time token to sample once on the corresponding dataset.

Dataset	Time(ms)			
	DBAL	CoreSet	SRAAL	Ours
FewRel	157.0	1145.2	409.3	465.1
NYT+FB	157.9	74.3	132.7	101.9
FewRel2.0	181.4	209.4	418.4	9.2
Average	165.4	476.3	320.1	<u>192.0</u>

The Impact of Different Sampling Strategies In order to prove the effectiveness of the active learning method, we conduct a further ablation experiment. As mentioned above, our sampling strategy is to select the k instances with the highest confidence for manual labeling. In the ablation experiment, we test two other sampling strategies: selecting the k instances with the lowest confidence; randomly selecting k instances by human. The comparison results are shown in Table 5.

It can be seen that after being trained by instances with the highest confidence, the model achieves the most improvement. In contrast, instances with the lowest confidence contribute very little to improving the performance of the model. Even with the continuous increase of training data, the improvement is extremely little. The results prove that the active learning model does select the most informative instances.

Time Efficiency of Relational Active Learning In practice, it is often the time spent on manual annotation that is the time-consuming bottleneck. Nevertheless, the sampling strategy for active learner should also select samples in a time-efficient manner as much as possible. We analysis the time efficiency of different active learning methods. Table 6 shows the average time for different methods to sample once on the corresponding dataset. DBAL is the most competitive baselines in terms of their achieved mean time efficiency. Our method fell marginally behind DBAL, however, our method is outperformed in accuracy by all other methods.

5.7. Analysis on Relational Outlier Detection

The Effects of Relational Outlier Detection ARD employs novel relation discovery module to distinguish between known and novel relations, preserving the active

Table 7

Ablation experiments over novel relation discovery module on noisy datasets.

Dataset	Model	Epoch				
		#1	#2	#3	#4	#5
FR (Noi)	ARD	66.0	69.0	70.5	72.7	75.5
	w/o LOF	62.8	64.4	68.5	70.5	73.4
NYF (Noi)	ARD	56.8	62.5	66.6	68.3	69.3
	w/o LOF	47.7	53.8	57.2	60.3	64.4
FR2.0 (Noi)	ARD	48.8	51.2	52.4	53.2	53.5
	w/o LOF	42.9	44.2	45.6	46.4	48.2

learning module to more efficiently select informative novel relations without being distracted by known relations. To demonstrate the effectiveness and significance of the novel relation discovery, we perform ablation experiments over LOF algorithm on three noisy variants. Table 7 shows the experimental results, and we note that: (1) Despite the robust learning ability of active learning on novel relations, the model performances show different degrees of degradation after the removal of the LOF algorithm. (2) Average of decline of $F1$ scores in each epoch on the FewRel, NYT+FB, and FewRel2.0 datasets is 2.82, 8.02, 6.36 respectively, with the most severe drop on NYT+FB. The phenomenon is intuitive, as the NYT+FB dataset contains the most known relations; the more noise (known relations) there is, the more confused the active learning module becomes about the novel relations. The results demonstrate the novel relation discovery module plays a key role as “noise reduction”.

The Impact of Different Outlier Detection Algorithms

We compare LOF with two different algorithms for the relational outlier detection, including IsolationForest [29], and OneClassSVM [41]. We evaluate the F1-measure of these three algorithms solely on the discovery of novel relations, the results are reported in Table 8. Our LOF algorithm outperforms by large margins, achieving 83.9% F1-measure on FewRel dataset. The principle of the IsolationForest algorithm is to cut data points and isolate data points one by one. Thus the data needs more cuts to be isolated. The main reason for the poor performance of this algorithm is a large amount of the test data. For the same type of new relations, their distribution is relatively dense, and the number of cuts will also increase. Moreover, the dimensions of relation representation are 2048, while IsolationForest has poor processing capabilities for high-dimensional features. Hence, it yields relevant poor results. OneClassSVM aims to learn a tight decision boundary from normal data and treats points outside the decision boundary as abnormal points. In the relational feature space, the distribution of known relations and novel relations are complicated. Thus the OneClassSVM is likely to learn an over-fitting decision boundary, resulting in poor performance.

Table 8

F1-measure on noisy FewRel and FewRel2.0 with different outlier detection algorithms.

Dataset	Method		
	IF	OneClassSVM	LOF
FewRel	64.0	47.3	83.9
FewRel2.0	63.1	54.1	80.3

5.8. Case Study of Active Learning

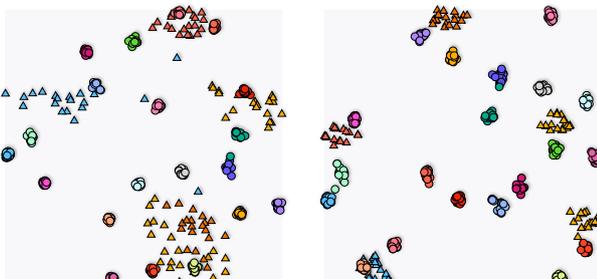
As shown in Table 10, we report 8 cases selected by discriminator in the first iteration on noisy FewRel2.0 dataset, where 64 relations are pre-trained and seen. With the highest confidence, the discriminator successfully selects sentences with unseen relations and guarantees the diversity of relation categories.

5.9. Additional Exploration

5.9.1. Visualization of Relation Representations

In order to intuitively demonstrate the distribution of novel relations relative to known relations and, on the other hand, to illustrate the benefits of introducing contrastive loss, we visualize the relation representation $h_r(x)$ after dimension reduction using t -SNE [33].

As illustrated in Figure 3, instances of the same known relation type are densely clustered with a high local density, while instances of novel relations distribute dispersedly. This fact strongly supports the premise of the LOF algorithm. Also, comparing subfigures 3(a) and 3(b), we observe that contrastive loss firmly constrains the distribution of intra-class instances. In pre-experiments on FewRel, the introduction of contrastive loss boosts the accuracy in distinguishing known and novel relations from 79.3% to 83.9%.



(a) Train using only traditional cross-entropy loss.

(b) Plus contrastive loss.

Figure 3: t -SNE visualization of relation representation. The known and novel relations are distinguished by circular and triangular symbols respectively.

5.9.2. Performance of our Relation Representation on Supervised RE

To demonstrate the effectiveness of the relation representation described in the Methodology section, we conduct a series of experiments on supervised RE task. First,

Table 9

Results on DDI'13 dataset. The first seven rows are the results of the previous SOTA methods, and the bottom results are for ours method in supervised relation learning.

Methods	Pre.	Rec.	F1
SCNN	69.1	65.1	67.0
CNN-bioWE	75.7	64.7	69.8
MCCNN	75.9	65.2	70.2
Joint AB-LSTM	73.4	69.6	71.5
RvNN	74.4	69.3	71.7
Position-aware LSTM	75.8	70.4	73.0
BERE	76.8	71.3	73.9
Ours	92.3	84.4	86.8

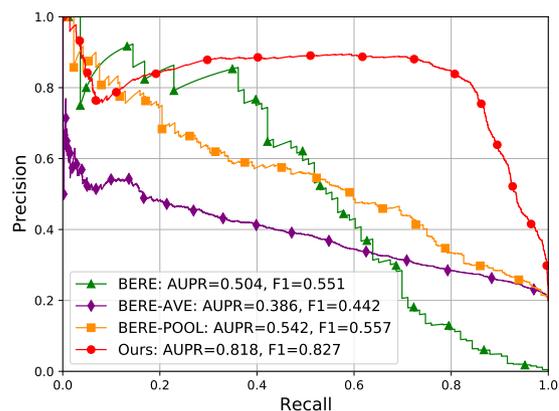


Figure 4: Precision-recall curve of BERE and our model.

we conduct extensive experiments on the biomedical relation extraction benchmark DDI'13 [16] and show in Table 9. We make comparisons with various previous state-of-the-art methods, which fall into two groups according to the neural network architecture: convolutional neural network (CNN) based methods and recurrent neural network (RNN) based methods. For the first group, we report the results of SCNN [58], CNN-bioWE [31] and MCCNN [36], which uses syntax word embeddings, biomedical-related embeddings and multi-channel word embeddings for feature extraction, respectively. For recurrent based networks, we report the results of Joint AB-LSTM [40], Position-aware LSTM [59], RvNN [27] and BERE [18]. Joint AB-LSTM jointly trains two bidirectional LSTM (Bi-LSTM) with different pooling mechanisms: max-pooling for one Bi-LSTM and attentive pooling for the other. Position-aware LSTM adopt position information as attention mechanism for the training of LSTM. RvNN and BERE incorporates parse-tree information to enhance the performance of prediction. Each model is trained on the training dataset to predict a relation class of five pre-defined relation types for the input sequence.

To further evaluate the performance of our representation method on large-scale distantly annotated dataset, we conduct another set of experiments on the DTI dataset. As on the DTI dataset, previous literature has shown the

Table 10

Cases selected by the confidence score of the discriminator and the novel relations, where *red* and *blue* represent the head and tail entities

Selected sentence	Novel relation
Ectopic overexpression of mir-497 promotes chemotherapy resistance in glioma cells by targeting <i>pdcd4</i> , a tumor suppressor that is involved in <i>apoptosis</i> .	<i>Biological process involves gene product</i>
As full-length bid is a weaker apoptogen than <i>tbid</i> , we propose that the phosphorylation of bid by jnks, followed by the accumulation of the full-length protein, delays attainment of <i>apoptosis</i> , and allows the cell to evaluate the stress and make a decision regarding the response strategy.	<i>Biological process involves gene product</i>
Pretreatment with dexamethasone 1 hour before <i>cyclophosphamide injection</i> significantly down-regulated <i>cyclophosphamide</i> induced bladder nuclear factor-u03bab dependent luminescence, ameliorated the grossly evident pathological features of acute inflammation and decreased cellular immunostaining for nuclear factor-u03bab in the bladder.	<i>Ingredient of</i>
Trastuzumab emtansine (<i>t-dm1</i>), an antibody-drug conjugate comprising the cytotoxic agent dm1, a stable linker, and <i>trastuzumab</i> , has demonstrated substantial activity in human epidermal growth factor receptor 2 (her2), -positive metastatic breast cancer, raising interest in evaluating the feasibility and cardiac safety of t-dm1 in early-stage breast cancer (ebc).	<i>Ingredient of</i>
Here we looked for evidence of adult hippocampal <i>neurogenesis</i> using immunohistochemical techniques for the endogenous marker doublecortin (<i>dcx</i>) in 10 species of microchiropterans euthanized and perfusion fixed at specific time points following capture.	<i>Gene plays role in process</i>
Here, we explored the effects of the novel class ii-specific "histone deacetylase inhibitors (hdacis) mc1568 and mc1575 on interleukin-8 (il-8) expression and <i>cell proliferation</i> in cutaneous melanoma cell line <i>gr</i> -m and uveal melanoma cell line ocm-3 upon stimulation with phorbol 12-myristate 13-acetate (pma).	<i>Gene plays role in process</i>
Data indicate that the structurally disordered and abnormally formed ecm of <i>uterine fibroids</i> contributes to <i>fibroid</i> formation and growth.	<i>Classified as</i>
however, individuals heterozygous for both beta "e", "and", beta thalassaemia (hbe/ <i>beta thalassaemia</i>) have a severe clinical disorder which in some cases may approach that seen in <i>homozygous beta thalassaemia</i> and which is by far the commonest form of symptomatic thalassaemia in the indian subcontinent and south-east asia.	<i>Classified as</i>

superiority of BERE compared with CNN-based and RNN-based baselines, we mainly take BERE as the baseline of our experiments. For fairness, we follow the settings of BERE by using precision-recall curve, the area under the precision-recall curve and the F_1 score as the evaluation metrics. We re-run the open-source code of BERE and its two variants: BERE-AVE, BERE-POOL. BERE-AVE adopt the average pooling mechanism to aggregate the semantic information over instances in a bag. BERE-POOL uses the max-pooling strategy. The implementation details of our model on the DTI dataset are identical to the DDI'13 dataset. The precision-recall curve is shown in Figure 4, which indicates the significant performance of our representation method.

5.9.3. Impact of The Size of BERT Model

We change the size of BERT in ARD. The results of this ablation experiment are shown in the Table 11. We can find that the size of BERT is not the key factor to bring gain, and

Table 11

Ablation experiments over the size of BERT model on FewRel datasets.

Data-set	Model	B ³			V-measure			ARI
		F1	Prec.	Rec.	V	Hom.	Comp.	
FR (Ori)	BASE	68.8	55.7	90.0	79.2	73.3	86.2	55.5
	LARGE	73.6	70.7	76.8	85.1	84.9	85.3	64.1
FR (Noi)	BASE	75.7	66.1	88.8	82.2	78.5	86.2	62.2
	LARGE	80.8	75.7	86.7	90.2	89.0	91.4	71.3
FR (Imb)	BASE	70.6	63.4	57.5	81.2	74.2	68.4	59.3
	LARGE	76.5	74.7	78.4	86.5	86.8	86.2	67.8

even if we use BERT_{BASE} as the backbone, the performance of ARD is still considerably higher than that of the baseline model using BERT_{LARGE} in Table 4.

Table 12

Statistical results of dataset for COVID-19.

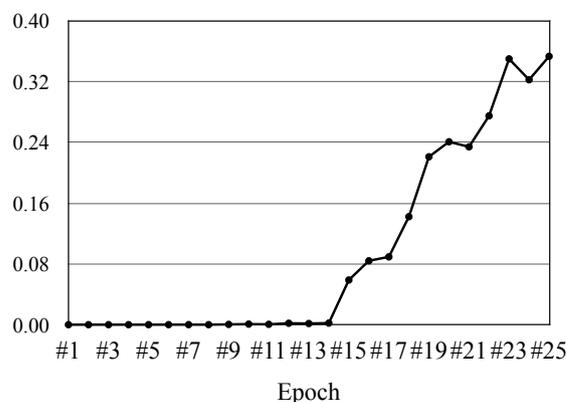
Relation	Count
CHEBI-CHEBI	4680
CHEBI-HP	20455
GO-HP	17254
DO-DO	14430
CHEBI-DO	2415
HP-HP	48236
HP-DO	19770
GO-CHEBI	1285
GO-DO	3615
GO-GO	7303

5.10. Practical Application on Real-world Dataset

We apply the ARD framework in real-world scenarios to verify its practicability. With the increasing number of publications about COVID-19, it is a challenge to extract personalized knowledge suitable for each researcher [9, 21]. [4] aims to build a new semantic-based pipeline for recommending biomedical entities to scientific researchers. In this work, the researchers utilize MER [6] as NER annotation server. As a result, 9,000 articles are automatically annotated with relevant items/concepts for COVID-19. And for further relation extraction task, due to the expensive manual annotation costs, the researchers merely take initial steps towards the results, providing a small sample dataset of ten documents, with all possible relationships between the four types of entities identified by NER pipeline. Thus, we were able to establish ten different types of relations, encompassing the four ontologies (CHEBI, DO, HPO, and GO). We follow the relation types, and apply ARD framework in the results. We take sample size k of 200 and sample 25 epochs. Finally, a total of 139,479 relations between entity pairs are automatically obtained by ARD. The statistical results of the data are shown in Table 12. We also report the confidence of the discriminator in each epoch for \mathcal{X}_U . As can be observed from the Figure 5, the confidence is progressively increasing as the training epoch increases, which indicates that the model is becoming more confident in the classification results. In an ideal case, the confidence should converge toward 0.5.

6. Conclusion and Future Work

The paper proposes Active Relation Discovery (ARD), which aims at accurately discovering and meaningfully annotating new semantic relations under the *General OpenRE* setting. By introducing outlier detection and active learning, ARD solves two problems: (1) *Sufficient capabilities to distinguish between known and novel relations*, with robust performance under General OpenRE settings. (2) *Avoiding Secondary labeling of downstream tasks*. Extensive experiments are conducted to demonstrate the effectiveness of ARD.

**Figure 5:** Discriminator's confidence for \mathcal{X}_U in each epoch.

As a pioneering work in OpenRE, several directions can be further explored: (1) Better methods to discriminate and annotate novel relations in *General OpenRE* setting. (2) Better methods to capture the core relational features for relation representation. (3) Combination with bootstrapping methods to partially replace active learning. (4) Combination with lifelong learning to continuously incorporate novel relations. (5) A universal schema for the standard of active relation learning.

Acknowledgement

This research is supported by National Natural Science Foundation of China (Grant No.62276154 and 62011540405), Beijing Academy of Artificial Intelligence (BAAI), the Natural Science Foundation of Guangdong Province (Grant No. 2021A1515012640), Basic Research Fund of Shenzhen City (Grant No. JCYJ20210324120012033), and Overseas Cooperation Research Fund of Tsinghua Shenzhen International Graduate School (Grant No. HW2021008).

CRedit authorship contribution statement

Yangning Li: Conceptualization of this study, Methodology, Experiments. **Yinghui Li:** Conceptualization of this study, Methodology, Experiments. **Xi Chen:** Revision of the paper, Funding Support. **Hai-Tao Zheng:** Revision of the paper, Funding Support. **Ying Shen:** Investigation process, Experimental verification.

References

- [1] Angeli, G., Premkumar, M.J.J., Manning, C.D., 2015. Leveraging linguistic structure for open domain information extraction, in: Proceedings of ACL-IJCNLP, pp. 344–354. URL: <https://www.aclweb.org/anthology/P15-1034>.
- [2] Bagga, A., Baldwin, B., 1998. Algorithms for scoring coreference chains, in: The first international conference on language resources and evaluation workshop on linguistics conference, Citeseer. pp. 563–566.
- [3] Banko, M., Cafarella, M.J., Soderland, S., Broadhead, M., Etzioni, O., 2007. Open information extraction from the web., in: Proceedings of IJCAI, pp. 2670–2676. URL: <https://patents.google.com/patent/US7877343B2>.

- [4] Barros, M.A., Lamúrias, A., Sousa, D., Ruas, P., Couto, F.M., 2020. Covid-19: A semantic-based pipeline for recommending biomedical entities, in: Proceedings of the 1st Workshop on NLP for COVID-19 (Part 2) at EMNLP 2020.
- [5] Chao, W.L., Changpinyo, S., Gong, B., Sha, F., 2016. An empirical study and analysis of generalized zero-shot learning for object recognition in the wild, in: European conference on computer vision, Springer. pp. 52–68.
- [6] Couto, F.M., Lamurias, A., 2018. Mer: a shell script and annotation server for minimal named entity recognition and linking. *Journal of cheminformatics* 10, 1–10.
- [7] Cui, L., Wei, F., Zhou, M., 2018. Neural open information extraction, in: Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pp. 407–413.
- [8] Devlin, J., Chang, M.W., Lee, K., Toutanova, K., 2019. BERT: Pre-training of deep bidirectional transformers for language understanding, in: Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers), Association for Computational Linguistics, Minneapolis, Minnesota. pp. 4171–4186. doi:10.18653/v1/N19-1423.
- [9] Dong, C., Li, Y., Gong, H., Chen, M., Li, J., Shen, Y., Yang, M., 2021. A survey of natural language generation. *arXiv preprint arXiv:2112.11739*.
- [10] Elsahar, H., Demidova, E., Gottschalk, S., Gravier, C., Laforest, F., 2017. Unsupervised open relation extraction, in: Proceedings of ESWC, pp. 12–16. URL: <https://link.springer.com/chapter/10.1007/978-3-319-70407-4-3>.
- [11] Fader, A., Soderland, S., Etzioni, O., 2011. Identifying relations for open information extraction, in: Proceedings of the conference on empirical methods in natural language processing, Association for Computational Linguistics. pp. 1535–1545.
- [12] Fu, L., Grishman, R., 2013. An efficient active learning framework for new relation types, in: Proceedings of the Sixth International Joint Conference on Natural Language Processing, pp. 692–698.
- [13] Gal, Y., Islam, R., Ghahramani, Z., 2017. Deep bayesian active learning with image data, in: International Conference on Machine Learning, PMLR. pp. 1183–1192.
- [14] Gao, T., Han, X., Zhu, H., Liu, Z., Li, P., Sun, M., Zhou, J., 2019. FewRel 2.0: Towards more challenging few-shot relation classification, 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), Association for Computational Linguistics, Hong Kong, China. pp. 6250–6255. doi:10.18653/v1/D19-1649.
- [15] Han, X., Zhu, H., Yu, P., Wang, Z., Yao, Y., Liu, Z., Sun, M., 2018. Fewrel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation. *arXiv preprint arXiv:1810.10147*.
- [16] Herrero-Zazo, M., Segura-Bedmar, I., Martínez, P., Declerck, T., 2013. The ddi corpus: An annotated corpus with pharmacological substances and drug–drug interactions. *Journal of biomedical informatics* 46, 914–920.
- [17] Higgins, I., Matthey, L., Pal, A., Burgess, C.P., Glorot, X., Botvinick, M.M., Mohamed, S., Lerchner, A., 2017. beta-vae: Learning basic visual concepts with a constrained variational framework, in: International Conference on Learning Representations.
- [18] Hong, L., Lin, J., Li, S., Wan, F., Yang, H., Jiang, T., Zhao, D., Zeng, J., 2020. A novel machine learning framework for automated biomedical relation extraction from large-scale literature repositories. *Nature Machine Intelligence*, 1–9.
- [19] Hu, X., Wen, L., Xu, Y., Zhang, C., Philip, S.Y., 2020. Selfore: Self-supervised relational feature learning for open relation extraction, in: Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 3673–3682.
- [20] Huang, H., Wang, C., Yu, P.S., Wang, C.D., 2019. Generative dual adversarial network for generalized zero-shot learning, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 801–810.
- [21] Huang, S., Ma, S., Li, Y., Li, Y., Lin, S., Zheng, H.T., Shen, Y., 2022. Towards attribute-entangled controllable text generation: A pilot study of blessing generation. *arXiv preprint arXiv:2210.16557*.
- [22] Hubert, L., Arabie, P., 1985. Comparing partitions. *Journal of classification* 2, 193–218.
- [23] Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., Krishnan, D., 2020. Supervised contrastive learning. *Advances in Neural Information Processing Systems* 33.
- [24] Kishimoto, Y., Murawaki, Y., Kurohashi, S., 2020. Adapting bert to implicit discourse relation classification with a focus on discourse connectives, in: Proceedings of The 12th Language Resources and Evaluation Conference, pp. 1152–1158.
- [25] Li, Y., Li, Y., He, Y., Yu, T., Shen, Y., Zheng, H.T., 2022a. Contrastive learning with hard negative entities for entity set expansion, in: Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Association for Computing Machinery, New York, NY, USA. p. 1077–1086. URL: <https://doi.org/10.1145/3477495.3531954>, doi:10.1145/3477495.3531954.
- [26] Li, Y., Zhou, Q., Li, Y., Li, Z., Liu, R., Sun, R., Wang, Z., Li, C., Cao, Y., Zheng, H.T., 2022b. The past mistake is the future wisdom: Error-driven contrastive probability optimization for chinese spell checking, in: Findings of the Association for Computational Linguistics: ACL 2022, pp. 3202–3213.
- [27] Lim, S., Lee, K., Kang, J., 2018. Drug drug interaction extraction from the literature using a recursive neural network. *PLoS one* 13, e0190926.
- [28] Liu, C., Sun, W., Chao, W., Che, W., 2013. Convolution neural network for relation extraction, in: International Conference on Advanced Data Mining and Applications, Springer. pp. 231–242.
- [29] Liu, F.T., Ting, K.M., Zhou, Z.H., 2008. Isolation forest, in: 2008 Eighth IEEE International Conference on Data Mining, IEEE. pp. 413–422.
- [30] Liu, R., Li, Y., Tao, L., Liang, D., Zheng, H.T., 2022. Are we ready for a new paradigm shift? a survey on visual deep mlp. *Patterns* 3, 100520.
- [31] Liu, S., Tang, B., Chen, Q., Wang, X., 2016. Drug-drug interaction extraction via convolutional neural networks. *Computational and mathematical methods in medicine* 2016.
- [32] Loshchilov, I., Hutter, F., 2017. Fixing weight decay regularization in adam. *ArXiv abs/1711.05101*.
- [33] Maaten, L.v.d., Hinton, G., 2008. Visualizing data using t-sne. *Journal of machine learning research* 9, 2579–2605.
- [34] Marcheggiani, D., Titov, I., 2016. Discrete-state variational autoencoders for joint discovery and factorization of relations. *Transactions of the Association for Computational Linguistics* 4, 231–244.
- [35] Qian, L., Hui, H., Hu, Y., Zhou, G., Zhu, Q., 2014. Bilingual active learning for relation classification via pseudo parallel corpora, in: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pp. 582–592.
- [36] Quan, C., Hua, L., Sun, X., Bai, W., 2016. Multichannel convolutional neural network for biological relation extraction. *BioMed research international* 2016.
- [37] Rahman, S., Khan, S., Porikli, F., 2018. A unified approach for conventional zero-shot, generalized zero-shot, and few-shot learning. *IEEE Transactions on Image Processing* 27, 5652–5667.
- [38] Romera-Paredes, B., Torr, P., 2015. An embarrassingly simple approach to zero-shot learning, in: International conference on machine learning, PMLR. pp. 2152–2161.
- [39] Rosenberg, A., Hirschberg, J., 2007. V-measure: A conditional entropy-based external cluster evaluation measure, in: Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL), pp. 410–420.
- [40] Sahu, S.K., Anand, A., 2018. Drug-drug interaction extraction from biomedical texts using long short-term memory network. *Journal of biomedical informatics* 86, 15–24.

- [41] Schölkopf, B., Platt, J.C., Shawe-Taylor, J., Smola, A.J., Williamson, R.C., 2001. Estimating the support of a high-dimensional distribution. *Neural computation* 13, 1443–1471.
- [42] Sener, O., Savarese, S., 2018. Active learning for convolutional neural networks: A core-set approach, in: *International Conference on Learning Representations*.
- [43] Settles, B., 2009. *Active Learning Literature Survey*. Computer Sciences Technical Report 1648. University of Wisconsin–Madison.
- [44] Shi, B., Weninger, T., 2018. Open-world knowledge graph completion, in: *Thirty-Second AAAI Conference on Artificial Intelligence*.
- [45] Shinyama, Y., Sekine, S., 2006. Preemptive information extraction using unrestricted relation discovery, in: *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*, pp. 304–311.
- [46] Simon, E., Guigue, V., Piwowarski, B., 2019. Unsupervised information extraction: Regularizing discriminative approaches with relation distribution losses, in: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 1378–1387.
- [47] Stanovsky, G., Dagan, I., 2016. Creating a large benchmark for open information extraction, in: *Proceedings of EMNLP*, pp. 2300–2305. URL: <https://www.aclweb.org/anthology/D16-1252>.
- [48] Stanovsky, G., Michael, J., Zettlemoyer, L., Dagan, I., 2018. Supervised open information extraction, in: *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 885–895.
- [49] Wu, F., Weld, D.S., 2010. Open information extraction using wikipedia, in: *Proceedings of the 48th annual meeting of the association for computational linguistics, Association for Computational Linguistics*. pp. 118–127.
- [50] Wu, R., Yao, Y., Han, X., Xie, R., Liu, Z., Lin, F., Lin, L., Sun, M., 2019. Open relation extraction: Relational knowledge transfer from supervised data to unsupervised data, 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)*, pp. 219–228.
- [51] Zeng, D., Liu, K., Chen, Y., Zhao, J., 2015. Distant supervision for relation extraction via piecewise convolutional neural networks, in: *Proceedings of EMNLP*, pp. 1753–1762. URL: <https://www.aclweb.org/anthology/D15-1203>.
- [52] Zhang, B., Li, L., Yang, S., Wang, S., Zha, Z.J., Huang, Q., 2020. State-relabeling adversarial active learning, in: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8756–8765.
- [53] Zhang, D., Wang, D., 2015. Relation classification via recurrent neural network. *arXiv preprint arXiv:1508.01006*.
- [54] Zhang, H.T., Huang, M.L., Zhu, X.Y., 2012. A unified active learning framework for biomedical relation extraction. *Journal of Computer Science and Technology* 27, 1302–1313.
- [55] Zhang, K., Yao, Y., Xie, R., Han, X., Liu, Z., Lin, F., Lin, L., Sun, M., 2021. Open hierarchical relation extraction, in: *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 5682–5693.
- [56] Zhang, Z., Saligrama, V., 2015. Zero-shot learning via semantic similarity embedding, in: *Proceedings of the IEEE international conference on computer vision*, pp. 4166–4174.
- [57] Zhao, Y., Wan, H., Gao, J., Lin, Y., 2019. Improving relation classification by entity pair graph, in: *Asian Conference on Machine Learning, PMLR*. pp. 1156–1171.
- [58] Zhao, Z., Yang, Z., Luo, L., Lin, H., Wang, J., 2016. Drug drug interaction extraction from biomedical literature using syntax convolutional neural network. *Bioinformatics* 32, 3444–3453.
- [59] Zhou, D., Miao, L., He, Y., 2018. Position-aware deep multi-task learning for drug–drug interaction extraction. *Artificial intelligence in medicine* 87, 1–8.



Yangning Li received the BEng degree from the Department of Computer Science and Technology, Huazhong University of Science and Technology, in 2020. He is currently working toward a Master's degree with the Tsinghua Shenzhen International Graduate School, Tsinghua University. His research interests include natural language processing and data mining.



Yinghui Li received the BEng degree from the Department of Computer Science and Technology, Tsinghua University, in 2020. He is currently working toward the PhD degree with the Tsinghua Shenzhen International Graduate School, Tsinghua University. His research interests include natural language processing and deep learning.



Xi Chen received his PhD degree in computer science from the Zhejiang University. He is currently the head of the cross-modal algorithm center of Tencent Platform and Content Group and mainly focuses on various applications of NLP.



Hai-Tao Zheng received the bachelor's and master's degrees in computer science from the Sun Yat-Sen University, China, and the PhD degree in medical informatics from Seoul National University, South Korea. He is currently an associate professor with the Shenzhen International Graduate School, Tsinghua University, and also with Peng Cheng Laboratory. His research interests include web science, semantic web, information retrieval, and machine learning.



Ying Shen received the PhD degree in computer science from the University of Paris Ouest Nanterre La Défense, France and the Erasmus Mundus master's degree in natural language processing from the University of Franche-Comté, France and the University of Wolverhampton, U.K. She is currently an associate professor with the School of Intelligent Systems Engineering, Sun Yat-Sen University. Her research interests include natural language processing and deep learning.