

# Self-distilled Transitive Instance Weighting for Denoised Distantly Supervised Relation Extraction

Anonymous ACL submission

## Abstract

The widespread existence of wrongly labeled instances is a challenge to distantly supervised relation extraction. Most of the previous works are trained in a bag-level setting to alleviate such noise. However, sentence-level training better utilizes the information than bag-level training, as long as combined with effective noise alleviation. In this work, we propose a novel Transitive Instance Weighting mechanism integrated with the self-distilled BERT backbone, utilizing information in the intermediate outputs to generate dynamic instance weights for denoised sentence-level training. By down-weighting wrongly labeled instances and discounting the weights of easy-to-fit ones, our method can effectively tackle wrongly labeled instances and prevent overfitting. Experiments on both held-out and manual datasets indicate that our method achieves state-of-the-art performance and consistent improvements over the baselines.

## 1 Introduction

Distantly Supervised Relation Extraction (DSRE) (Mintz et al., 2009) is designed to automatically annotate the sentences mentioning the entity pairs, which enables a significant way for constructing large-scale datasets. However, distant supervision (DS) works under an unrealistic assumption that all sentences mentioning the same entity pair express the same relation. This introduces many noisy (wrongly labeled) instances into the dataset. To tackle this challenge, previous works mostly adopt the bag-level setting as shown at the top of Figure 1, where the vector representations of sentences are aggregated as the bag-level representation using multi-instance learning (MIL) (Riedel et al., 2010), and the prediction is thus produced from the bag representation. The optimization is conducted at the bag level to minimize the loss of bag prediction. Only a small subset of previous works leverage the sentence-level setting (Zhang et al., 2019b; Liu

et al., 2020a) as in the bottom of Figure 1, where the sentence-level predictions are produced and then aggregated into the bag prediction. In fact, sentence-level training can directly optimize the loss from each sentence, enabling higher information utilization than bag-level training. However, sentence-level training is vulnerable to the noise brought by DS, which limits its application. Therefore, sentence-level training should be combined with effective noise-alleviation mechanisms to improve its robustness.

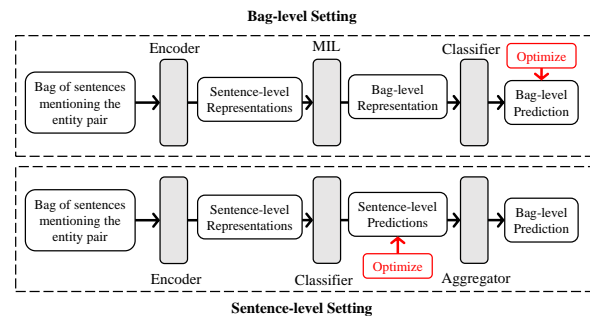


Figure 1: The bag-level and sentence-level pipelines of DSRE.

The mainstream encoders of DSRE models are Piecewise Convolutional Neural Network (PCNN) (Zeng et al., 2015) and Recurrent Neural Network (RNN) (Zhou et al., 2016; Liu et al., 2018) over the years. It is reasonable for most previous works to take the simple encoder as a black box and only utilize its final output during training and inference. However, as large models like BERT (Devlin et al., 2019) becomes popular in recent years, the information within the outputs from their intermediate layers is a non-trivial source of knowledge but is rarely utilized in DSRE. In this work, we apply self-distillation to extract intermediate information as output probabilities and utilize them to denoise from wrong labels. Furthermore, we use soft target selection and set up transitive knowledge passing among the students to alleviate

070 the effects of noisy target probabilities from the  
071 teacher.

072 The instances in DSRE can be roughly divided  
073 into easy, hard and noisy ones. Both easy and hard  
074 instances are correctly labeled but the model learns  
075 from hard instances slower (Huang et al., 2021).  
076 Noisy instances have wrong labels and can be fur-  
077 ther divided into False Positives (FPs) and False  
078 Negatives (FNs). FPs are instances with NA rela-  
079 tion but are wrongly labeled as non-NA relations  
080 by DS, while FNs are non-NA instances wrongly  
081 labeled as NA. We hope to avoid learning from  
082 noisy instances since they contain misleading in-  
083 formation. Moreover, we also need to avoid over-  
084 fitting easy instances to improve the learning of  
085 deeper knowledge. To tackle the above challenges,  
086 we propose a novel **Transitive Instance Weight-**  
087 **ing (TIW)** mechanism for DSRE. Our method  
088 adopts the sentence-level setting in both stages:  
089 fine-tuning and distillation. After fine-tuning the  
090 BERT encoder using a linear classifier (teacher)  
091 in the first stage, we append an auxiliary classi-  
092 fier (student) to each relevant layer and train them  
093 with TIW during distillation. TIW first filters FNs  
094 using binary weights (0 or 1). Then the soft tar-  
095 get probabilities are chosen between the outputs of  
096 the teacher and the previous peer. Finally, the in-  
097 stance weights for the positive (non-NA) instances  
098 are generated by combining two factors: the uncer-  
099 tainty (Liu et al., 2020b) and the soft confidence  
100 score. We apply uncertainty to prevent overfitting  
101 easy instances and use the soft confidence score as  
102 the assessment of learning difficulty, where easy  
103 and hard instances tend to have higher scores than  
104 noisy ones. During filtering and weighting, each  
105 student receives information from both the teacher  
106 and the previous peer, enabling the alleviation of  
107 noise from the teacher and transitive knowledge  
108 passing among the students. The experiments on  
109 both held-out and manual datasets show that our  
110 approach achieves state-of-the-art performance and  
111 consistent improvements over the teacher and the  
112 baselines. We also provide a detailed ablation study  
113 to explore the effects of the modules. Finally, we  
114 analyse the errors and discuss the limitations of our  
115 method.

116 Our contributions are summarized as follows:

- 117 • We are the first to denoise sentence-level  
118 DSRE with dynamic instance weights and har-  
119 ness intermediate knowledge to improve noise  
120 resistance and information utilization.

- We propose a novel Transitive Instance  
Weighting mechanism with multiple func-  
tions, including noise alleviation, overfitting  
prevention, soft target selection and transitive  
knowledge passing.
- Experiment and analysis show that our  
method achieves state-of-the-art performance  
with good generalization and robustness.

## 2 Related Work

Distant supervision (DS) for relation extrac-  
tion (Mintz et al., 2009) enables automatic an-  
notation of large-scale datasets, but its strong as-  
sumption introduces a large number of wrongly  
labeled instances. Following Riedel et al. (2010),  
various multi-instance learning methods are pro-  
posed to denoise from noisy instances, and they  
broadly fall into two categories: instance selec-  
tion (Zeng et al., 2015; Qin et al., 2018; Feng  
et al., 2018) and instance attention (Lin et al., 2016;  
Yuan et al., 2019b,a; Ye and Ling, 2019). Apart  
from multi-instance learning, many of the previ-  
ous works try to improve the effectiveness of train-  
ing. Liu et al. (2017) and Shang et al. (2020)  
try to convert wrongly labeled instances to useful  
information through relabeling. Huang and Du  
(2019) proposes collaborative curriculum learning  
for denoising. Hao et al. (2021) adopts adversarial  
training to filter noisy instances in the dataset. Hao  
et al. (2021) adopts adversarial training to filter  
noisy instances in the dataset. Nayak et al. (2021)  
designs a self-ensemble framework to filter noisy  
instances despite information loss. Li et al. (2022)  
proposes a hierarchical contrastive learning frame-  
work to reduce the effect of noise. Nevertheless, the  
above approaches are trained with bag-level loss,  
leading to lower utilization of information. In our  
work, we adopt sentence-level training to directly  
utilize sentence-level information and effectively  
tackle noise and overfitting using dynamic instance  
weights.

Knowledge distillation (Hinton et al., 2015) is  
an effective way to improve model generalization,  
though it has difficulty in transferring knowledge  
effectively (Stanton et al., 2021). By sharing some  
parameters between teacher and students, self-  
distillation (Zhang et al., 2019a) improves knowl-  
edge transfer from teacher to students. Liu et al.  
(2020b) applies self-distillation on BERT (Devlin  
et al., 2019) to improve inference efficiency. How-  
ever, In our work, we apply self-distillation as the

171 tool to extract intermediate knowledge for denois- 205  
 172 ing and further reduce the noise from the teacher 206  
 173 with transitive information passing among the stu- 207  
 174 dents. 208

175 There are some epoch-level techniques to detect 209  
 176 noisy instances like Swayamdipta et al. (2020) and 210  
 177 Huang et al. (2021). But in sentence-level DSRE 211  
 178 which is highly noisy and contains bias from the 212  
 179 entity mentions (Peng et al., 2020), larger mod- 213  
 180 els like BERT can overfit noisy instances faster, 214  
 181 even before an epoch ends. Therefore, we adopt a 215  
 182 dynamic instance weighting mechanism which is 216  
 183 more suitable for DSRE. 217

### 184 3 Methodology

185 Our model is illustrated in Figure 2. The back- 221  
 186 bone of our model is the BERT encoder on the left, 222  
 187 with a teacher classifier on the top. Each student 223  
 188 contains a subencoder and an auxiliary classifier. 224  
 189 For example, the student 7 has a subencoder end- 225  
 190 ing with the 7th BERT layer and a linear classifier 226  
 191 appended to it. The BERT encoder is fine-tuned 227  
 192 with the teacher classifier on the dataset before dis-  
 193 tillation. As discussed in Jawahar et al. (2019),  
 194 the shallow layers may not be able to encode the  
 195 information needed for the DSRE task. Therefore,  
 196 TIW starts from layer  $L$ , which is empirically set  
 197 and will be called **the head layer** in the rest of the  
 paper.

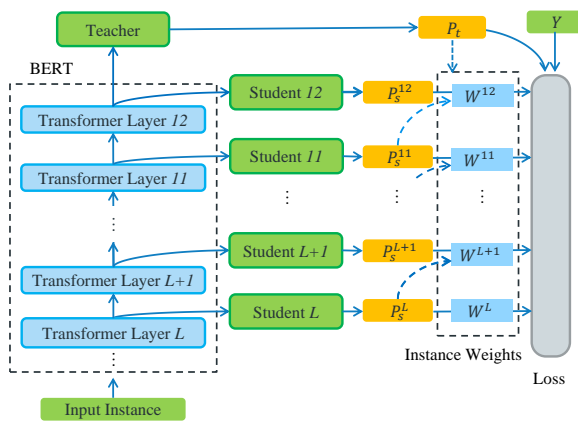


Figure 2: The overall framework of our model. Dotted arrows indicate the generation of instance weight.

#### 198 3.1 Backbone

199 BERT (Devlin et al., 2019) is a powerful 200  
 201 transformer-based pretrained network with broad 202  
 202 applications in natural language processing. Its 203  
 203 intermediate layers encode a rich hierarchy of sen- 204  
 204 tence features, ranging from surface features, and

205 syntactic features, to semantic features (Jawahar 206  
 206 et al., 2019). However, previous BERT applica- 207  
 207 tions in DSRE (Alt et al., 2019; Rao et al., 2022) 208  
 208 only utilize the output from the final layer, neglect- 209  
 209 ing the possibility that hierarchical intermediate 210  
 210 information can be useful in denoising. Therefore, 211  
 211 we apply auxiliary classifiers as in Figure 2 to ex- 212  
 212 tract information from the hierarchical features in 213  
 213 the form of output probabilities and utilize them to 214  
 214 distinguish noisy instances in the distillation stage.

215 Before distillation, we fine-tune the BERT en- 216  
 216 coder on DSRE as in Gao et al. (2021). The 217  
 217 structure of the embedding layer and BERT lay- 218  
 218 ers follow those in the previous works with the 219  
 219 number of transformer layers  $n = 12$  and hidden 220  
 220 size  $d_h = 768$ .

221 Firstly, the input sentence is transformed 222  
 222 to a sequence of vector representations  $s =$  223  
 223  $[w_1, w_2, \dots, w_m]$  by the embedding layer, where 224  
 224  $m$  is the maximum length of the sentence. Then, 225  
 225 BERT conducts layer-wise feature extraction with 226  
 226 the input  $s$ , the output of  $i_{th}$  layer ( $1 \leq i \leq n$ ) is 227  
 227 described as:

$$h_i = BERT_i(s) \quad (1) \quad 228$$

229 where  $BERT_i$  refers to the subencoder containing 229  
 230 transformer layers from the first to the  $i_{th}$ . The 230  
 231 encoder is fine-tuned with a simple feedforward 231  
 232 classifier on the top: 232

$$x_i = [h_i(p_1); h_i(p_2)] \quad (2) \quad 233$$

$$FFN(h_i) = M_2(M_1x_i + b_1) + b_2 \quad (3) \quad 234$$

$$p^t = softmax(FFN_t(h_n)) \quad (4) \quad 235$$

236 where  $M_1 \in R^{d_h \times d_h}$  and  $M_2 \in R^{n_c \times d_h}$  are 237  
 238 weight matrices and  $b_1 \in R^{d_h}$  and  $b_2 \in R^{n_c}$  are 239  
 239 bias terms.  $p_1$  and  $p_2$  are the start positions of the 240  
 240 head entity and tail entity respectively.  $[a : b]$  indi- 241  
 241 cates the concatenation of vectors  $a$  and  $b$ .  $x_i$  is the 242  
 242 entity-aware sentence representation generated by 243  
 243 concatenating the hidden vectors of the entity pair. 244  
 244  $n_c$  is the number of classes and  $p^t$  is the output 245  
 245 probability of the teacher. 246

247 The student  $i$  can be formulated as follows: 247

$$p_i^s = softmax(FFN_i(h_i)) \quad (5) \quad 248$$

249 After fine-tuning, the parameters of the teacher 249  
 250 model including the BERT encoder stay fixed. 250

---

**Algorithm 1** Transitive Instance Weighting

---

**Input:** DS label  $Y$ , teacher’s output probability  $p^t$  and students’  $p^s$  for the instance.

**Output:** The soft target  $p^{tg}$  and the instance weight  $w$  of the instance from the students .

```
1: Initialize  $w_l \leftarrow 1, p_l^{tg} \leftarrow p^t$ 
2: for  $i = l + 1 \rightarrow n$  do
3:   Compute the PoAs of  $i_{th}$  student:  $c_i^t \leftarrow p_i^s \cdot p^t$     $c_i^s \leftarrow p_i^s \cdot p_{i-1}^s$ 
4:   if  $c_i^t > c_i^s$  then  $p_i^{tg} \leftarrow p^t$  else  $p_i^{tg} \leftarrow p_i^s$             $\triangleright$  Soft Target Selection
5:   if  $Y = rel2id(NA)$  then                                            $\triangleright$  False Negative Filtering
6:     if  $Y = argmax_j(p_{i-1}^s(j))$  then  $w_i \leftarrow 1$  else  $w_i \leftarrow 0$ 
7:   else                                                                  $\triangleright$  Positive Weighting
8:     Compute the uncertainty of soft target:  $u_i \leftarrow \sum_{j=1}^{n_c} \frac{p_i^{tg}(j) \log p_i^{tg}(j)}{\log \frac{1}{n_c}}$ 
9:     Compute instance weight:  $w_i \leftarrow max(c_i^t, c_i^s) u_i$ 
10:  end if
11: end for
```

---

### 3.2 Transitive Instance Weighting

The algorithm of TIW is shown in Algorithm 1, where  $rel2id(r)$  is a function that maps the relation class  $r$  to its id for generating the one-hot label. TIW provides dynamic instance weights for each student except the first one (layer  $L$ ), it sets up a transitive way to share knowledge (output probabilities) among the students. Note that we use the last student for the final prediction and the rest of the students aim to provide robust instance weights for the last one.

Most previous works in knowledge distillation directly use the teacher’s output probability as the soft target. However, the teacher can constantly make mistakes if trained with noisy data, as in DSRE. Therefore, as in Line 4 of our algorithm, instead of blindly following the output from the teacher, each student except the first one chooses between the teacher  $p^t$  and the previous peer  $p_{i-1}^s$  to follow. The criterion of choosing is consistency, which can be described as the probability of making the same predictions as each other. We call it the **Probability of Agreement (PoA)** and compute it as the dot product of two probability distributions. The selection of soft targets provides additional referential probability distributions for the learning students and they can switch to a smoother target probability when the output from the teacher is too hard to learn.

In TIW, we adopt different strategies for negative (NA) instances and positive (non-NA) ones because their characteristics are quite different. We conduct **False Negative Filtering (FNF)** as in Lines 5-6 of the algorithm. Since we have sufficient negative in-

stances in the dataset, it is acceptable to avoid more FNs at the cost of slight information loss. Therefore, we assign 0 weight to all the possible FNs and 1 weight to the rest. To correctly identify FNs, we adopt a dynamic approach that if the previous peer agrees with distant supervision and also labels the instance as *NA*, then we classify the instance as a true negative. Otherwise, we assume it to be a false negative that the DS label is unreliable. The student follows the peer’s view in FNF instead of the teacher’s because the teacher already overfits the noisy data and mostly follows the DS label, though the probabilities of label relations may vary.

In order to preserve more information for training, we use soft weights for the positive instances instead of hard filtering. We call it Positive Weighting (PW) and determine the instance weight  $w_i$  of student  $i$  by two factors: uncertainty and the soft confidence score.

The uncertainty term is the normalized entropy as in Liu et al. (2020b) of the chosen soft target. It evaluates how well an instance is fitted so we can leverage it to detect overfitted instances dynamically. Easy instances contain shallow features like *London, UK* indicating a *location/contains* relation, so the model fits them easily and fast. But we do not hope the model becomes overdependent on them and lose focus on deeper features hidden in semantics. Therefore we discount their weights with uncertainty to prevent overfitting.

The maximum between the PoAs from the teacher and the previous peer is the **Soft Confidence (SC)** score which evaluates the learning difficulty of the instance for the student. If the SC score is high, the student successfully follows the

idea of the teacher or the peer, indicating that the instance is easy to learn for the student. If the SC score is low, the student is unable to follow the referential probabilities and the instance may be noisy or very hard to learn.

The instance weight for  $i_{th}$  student ( $l < i \leq n$ ) is computed as the product of the SC score and the uncertainty term, as in Line 9 of the algorithm. Note that during distillation, the student is trained with both soft target distribution and DS labels, as shown in Equation 7. We present the discussions on the SC scores and losses of easy, noisy and hard instances in the following.

Easy instances mostly have high SC scores and are well-fitted by the teacher or the peer, so the optimizations using soft labels and hard labels conform with each other.

Noisy instances are mostly underfitted and very hard to optimize because the soft labels and hard labels are mostly inconsistent. They have low SC scores because the teacher and the students are not likely to provide consistent predictions.

Hard instances are underfitted clean instances with low SC scores at first. However, their soft and hard labels are consistent, leading to smoother optimizations. When clean background knowledge is established by learning from clean instances, learning from hard ones becomes easier so the SC scores of hard instances grow larger.

Based on the above discussions, it is safe to say that both easy and hard instances are faster to fit and tend to have larger SC scores than noisy ones. The uncertainty term only takes effect when easy instances are well-fitted and clean background knowledge is established, so it will not lead to overfitting noisy instances.

To sum up, TIW is robust against noise and overfitting and thus can be combined with sentence-level training to utilize more information for better performance than previous bag-level methods.

### 3.3 Optimization

The teacher model may overfit noisy instances during fine-tuning. Therefore, we apply a dynamic temperature  $\tau$  to the teacher in the following form:

$$\tau_i = 1 + \gamma(1 - u_i) \quad (6)$$

where  $\gamma$  is a hyperparameter empirically set as 3. The idea of  $\tau$  is to further smooth the well-fitted instances to produce softer targets.

The loss function of our model follows the general form of knowledge distillation with the instance weight  $w$  we propose:

$$L = \sum_{i=l}^n w_i (\alpha KL_{\tau_i}(p_i^s, p_i^{tg}) + (1 - \alpha)CE(p_i^s, Y)) \quad (7)$$

where  $\alpha$  is a hyper-parameter empirically set as 0.5.  $KL_{\tau}(p, q)$  computes the KL-divergence between distributions  $p$  and  $q$  with temperature  $\tau$  for the teacher.  $Y$  is the label from distant supervision and  $CE(p, Y)$  is the cross entropy loss with one-hot label obtained from  $Y$ .

## 4 Experiments

In this section, the datasets, settings and hyperparameters are specified first. Then, we present the performance of our model compared with previous baselines and the teacher model. We also conduct an ablation study and error analysis to enable a deeper understanding of the mechanisms.

### 4.1 Datasets and Settings

We use two datasets for evaluation, the widely used **held-out** dataset NYT-10 (Riedel et al., 2010) and recent **manual** dataset NYT-10m (Gao et al., 2021). As a standard dataset for DSRE, NYT-10 is constructed by aligning the relations in Freebase (Bollacker et al., 2008) with the New York Times (NYT) corpus (English). The training set includes sentences from 2005 to 2006, and the test set uses sentences from 2007. NYT-10m is a manual dataset constructed also from NYT corpus, with a human-labeled test set and a new relation ontology. For NYT-10, we divide the dataset into five parts for cross-validation. For NYT-10m, we use the provided validation set. The details of the datasets are shown in Table 1.

Dataset	Train (k)		Test (k)		Rel.
	Sen.	Fac.	Sen.	Fac.	
held-out	522.6	18.4	172.4	2.0	53
manual	417.9	17.1	9.7	3.9	25

Table 1: The statistics of datasets. **Sen.**, **Fac.** and **Rel.** indicate the numbers of sentences, relation facts and relation types (including *NA*) respectively.

In the experiments, we use the *bert-base-uncased* checkpoint with about 110M parameters for initialization as in Han et al. (2019). We apply

the AdamW (Loshchilov and Hutter, 2017) optimizer during distillation and fix the random seed as 42. Apart from the hyperparameters previously mentioned, the batch size is 32 and the learning rate is  $2e - 5$ . The maximum length of sentences  $m$  is 128. The head layer  $L$  is set as layer 7 in our experiments.

We compare the Area Under precision-recall Curve (AUC), the F1 score and the mean of precision at top N predictions (N=100, 200, 300), which is denoted as P@M. Following the *at-least-one* assumption (Riedel et al., 2010), we adopt **ONE** strategy (Zeng et al., 2015) for bag-level evaluation, which takes the maximum score for each relation to generate bag-level predictions. As mentioned in Section 3, We use the output probabilities of the last student as the output of our model. In the appendix, we also display the results from other students and the results using other settings of  $L$ .

## 4.2 Overall Performance

We compare the performance of our model against that of the following baselines:

**PCNN+ATT** (Lin et al., 2016) proposes PCNN with selective attention mechanism.

**RESIDE** (Vashishth et al., 2018) integrates side information into Graph Convolution Networks to improve relation extraction.

**DISTRE** (Alt et al., 2019) extends and fine-tunes GPT on DSRE.

**Intra+inter** (Ye and Ling, 2019) combines intra-bag attention with inter-bag attention to tackle the noisy bags.

**CIL** (Chen et al., 2021) applies contrastive instance learning to reduce noise from DS.

**Teacher** follows the implementation in Gao et al. (2021).

Among the baselines, DISTRE and CIL use pre-trained language models for initialization. CIL adopts the same BERT pretrained encode as ours. The held-out dataset is the mainstream for DSRE evaluation, but it contains wrongly labeled test instances leading to inaccurate evaluation. The manual dataset provides an accurate test set but is limited by its scale in generalization. Therefore, we use both of the datasets for better evaluation.

### 4.2.1 Evaluation on Held-out Dataset

Table 2 shows the experimental results on the held-out dataset. We use the results reported in the papers of previous work. We also plot the precision-recall curves as in Figure 3.

Model	AUC	F1	P@M
PCNN+ATT	33.8	40.7	71.1
RESIDE	41.5	45.7	79.4
DISTRE	42.2	48.6	66.8
Intra+inter	42.3	46.5	84.8
CIL	<u>50.8</u>	<u>52.2</u>	<b>86.0</b>
Teacher	50.6	<u>52.2</u>	83.6
Student	<b>53.9</b>	<b>55.3</b>	<u>84.9</u>

Table 2: The performance (%) of the models on the held-out dataset. The best scores are marked as **bold** and the second best scores are underlined, as in other tables of the experiments.

As shown in the results, our model achieves the best AUC and F1 score among all the compared methods. The P@M of the student is relatively lower than bag-level methods, but still significantly higher than the teacher model. We can see that sentence-level training leads to a slight decline in the P@M due to the existence of noisy sentences but achieves better overall performance on the test set because of its advantage in information utilization. Our method further alleviates noise and overfitting with TIW, thus achieving state-of-the-art performance.

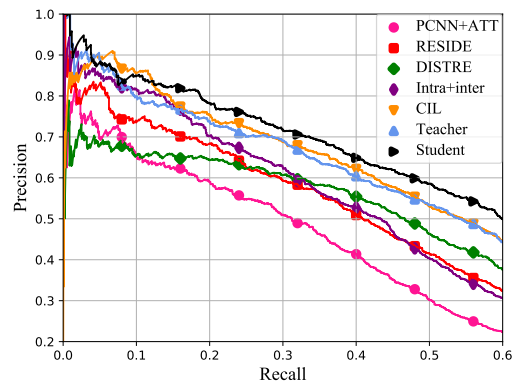


Figure 3: PR curves of the models on the held-out dataset.

### 4.2.2 Evaluation on Manual Dataset

Table 3 shows the experimental results on the manual dataset. We use the original implementations of the baselines. The precision-recall curves are plotted in Figure 4.

In the results, the bag-level methods still perform better at P@M, however, our method outperforms them in AUC and F1 by large margins.

Model	AUC	F1	P@M
PCNN+ATT	57.7	57.0	89.2
Intra+inter	53.6	53.5	<b>91.8</b>
CIL	60.2	58.8	<u>91.7</u>
Teacher	<u>61.3</u>	<u>62.4</u>	84.3
Student	<b>63.9</b>	<b>63.8</b>	90.8

Table 3: The performance (%) of our model and the baselines on the manual dataset.

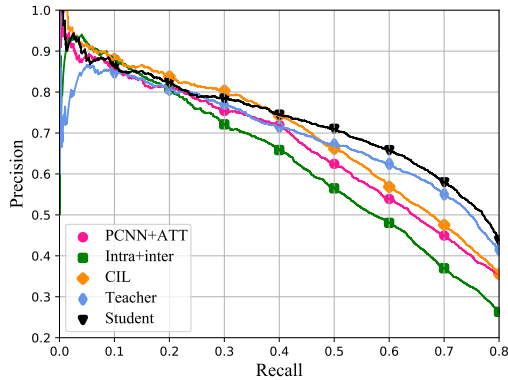


Figure 4: PR curves of the models on the manual dataset.

It shows that previous bag-level methods overfit easy instances, leading to the loss of overall generalization. The student also achieves significant improvements over the teacher, especially in P@M. The results further demonstrate the effectiveness of TIW in improving sentence-level training.

According to Gao et al. (2021), the performance of the model may be inconsistent if evaluated in both the held-out and manual datasets. Good performance on the held-out set may indicate overfitting to the bias from DS. However, our model is robust enough to perform well on both datasets.

### 4.3 Ablation Study

Model	AUC	F1	P@M
Our method	<b>53.9</b>	<b>55.3</b>	<b>84.9</b>
a: - STS	53.2	54.5	83.3
b: - PW	51.9	52.5	<u>84.8</u>
c: - FNF	<u>53.3</u>	<u>54.9</u>	82.5
d: - TIW	52.1	52.6	84.6
e: Probe	50.6	52.5	80.0

Table 4: Ablation study of our method.

As shown in Table 4, all the modules improve the overall performance. Detailed discussions are given below:

*a*: removes Soft Target Selection (STS) and follows the output probabilities from the teacher all the time. The noise from the teacher is not addressed, leading to performance declines.

*b*: removes PW and all the positive instances are treated equally, including the noisy ones. Therefore, the model is heavily affected by noise and the FNF may be inaccurate, leading to further declines in performance. In this case, high P@M indicates that the model overfits easy instances and loses generalization.

*c*: removes FNF. The false negative instances only make up a small part of the dataset, so the effect is relatively small. However, the noise from FNs significantly reduces P@M. We suspect that the fitting of false negatives affects that of true positives. If a false negative  $fn$  has similar syntactic and semantic features to a true positive  $tp$ , fitting  $fn$  is similar to fitting  $tp$  using an incorrect label.

*d*: removes TIW totally and all the instances are weighted as 1. The label smoothness of knowledge distillation is able to alleviate some noise from DS, so there are improvements in performance over *e*. However, the student is still trained with much noise and overfits easy instances, so the overall performance declines significantly.

*e*: is the probing result of 12th layer using the DS label. It shows that without effective denoising mechanisms, simply retraining the classifier does not help in performance.

The above results and discussions further demonstrate the effectiveness of TIW designs in alleviating noise and overfitting.

### 4.4 Error Analysis

For accurate analysis of the errors, we use the test set of the manual dataset for statistical discussions. Each positive label is considered an **item**. The instances with multiple positive labels are considered to have multiple items. We classify the items based on the predictions of the teacher and student, then count the number and percentage of each class as in Table 5. The goal is to explore where the errors of the student come from: a) **from the teacher**, meaning that the knowledge from the teacher is noisy and leads to the student’s errors, or b) **from the student itself**, meaning that the teacher gives correct knowledge but the student fails to follow.

Sentence	Teacher	Student
<u>Carl Friedrich von Weizsäcker</u> was born in <u>Kiel</u> , Germany, on June 28, 1912.	/people/person/place_of_birth	/people/person/place_lived
Presented by <u>Brooklyn College</u> and the office of Borough President <u>Marty Markowitz</u> .	/business/person/company	/people/person/place_lived
Furthermore, the relationship between the central government, dominated by three small <u>Arab</u> tribes living along the Nile, and Darfur's Arabs, who claim a heritage going back to the Prophet <u>Muhammad</u> , is often antagonistic.	/people/person/ethnicity	/people/person/place_of_birth

Figure 5: TCSI examples. The entities are underlined.

Class	Num. of items	Percentage (%)
<i>BC</i>	3,044	78.07
<i>BI</i>	742	19.03
<i>TISC</i>	94	2.41
<i>TCSI</i>	19	0.49

Table 5: Numbers and percentages of different classes of items. *BC* stands for *both correct*, *BI* stands for *both incorrect*, *TISC* stands for *teacher incorrect, student correct* and *TCSI* stands for *teacher correct, student incorrect*.

In the results, the student achieves slightly higher (about 2%) accuracy than the teacher and shows high fidelity with 97.1% of all predictions being the same as the teacher. *BI* represents the student’s errors caused by the errors from the teacher. *TISC* indicates the student’s corrections on the errors from the teacher and *TCSI* represents the errors from the student itself. From the results, we can conclude that almost all (about 97.5%) of the errors come from the teacher, and the corrections made by the student are much more than the errors made by the student itself. This demonstrates the effectiveness of our method in reducing the occurrence of errors and the limitation that it requires a good teacher for good performance.

For further analysis of the student’s errors, we inspect the *TCSI* items and select some representative ones for discussions as in Figure 5. Most of the instances with *place\_of\_birth* relation are correctly classified and the first example should be an easy instance in the form, yet misclassified by the student as *place\_lived*. We observe several similar items and suspect that long and uncommon names like *Carl Friedrich von Weizsäcker* sometimes confuse the student to make conservative predictions, which is the more common relation *place\_lived*. The second example, however, confuses the student with a compound noun *Brooklyn College*. *Brooklyn* appears very often in the dataset in the form of location, making the student believe

that *Brooklyn College* is a location rather than an organization. The third example is mostly related to ambiguity, where the word *Arab* may refer to the Arab people (ethnic group) or the Arab world (location). The latter two examples indicate that the lack of entity-related information may lead to inconsistency between the student and the teacher. The first example shows that the student may be confused to lose focus on key phrases like *was born in*, which may be solved by combining with word-level attention in the future.

## 5 Conclusions and Limitations

In this paper, we propose a novel Transitive Instance Weighting mechanism integrated with self-distillation to denoise from sentence-level training of DSRE. We employ the self-distilled BERT backbone to extract intermediate information for generating reliable instance weights. TIW combines the soft confidence score with uncertainty to tackle noisy instances and alleviate overfitting, it also enables soft target selection and transitive knowledge passing among the students to tackle the noise from the teacher. The experiment results show that our method improves the general resistance to DS noise and prevents overfitting from harming its generalization, thus can achieve state-of-the-art performance and consistent improvements over the baselines on both the held-out and manual datasets.

However, our work still has some limitations. Firstly, since our model is built on the basis of the teacher-student network, the performance of the student is highly affected by the teacher. If the teacher provides too much noisy information, our instance weighting mechanism might not work. Secondly, in some cases, the student fails to follow the correct predictions from the teacher due to ambiguity, lack of information or word-level noise, which indicates that further extension of our method is plausible. Finally, we haven’t explored other instance weighting methods in this paper. There might be better solutions yet to be discovered.



## References

- 609  
610  
611  
612  
613  
614  
615  
616  
617  
618  
619  
620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666
- Christoph Alt, Marc Hübner, and Leonhard Hennig. 2019. [Fine-tuning pre-trained transformer language models to distantly supervised relation extraction](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1388–1398, Florence, Italy. Association for Computational Linguistics.
- Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. 2008. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250.
- Tao Chen, Haizhou Shi, Siliang Tang, Zhigang Chen, Fei Wu, and Yueting Zhuang. 2021. [CIL: Contrastive instance learning framework for distantly supervised relation extraction](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6191–6200, Online. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 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)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jun Feng, Minlie Huang, Li Zhao, Yang Yang, and Xiaoyan Zhu. 2018. Reinforcement learning for relation classification from noisy data. In *Proceedings of the aaai conference on artificial intelligence*, volume 32.
- Tianyu Gao, Xu Han, Yuzhuo Bai, Keyue Qiu, Zhiyu Xie, Yankai Lin, Zhiyuan Liu, Peng Li, Maosong Sun, and Jie Zhou. 2021. [Manual evaluation matters: Reviewing test protocols of distantly supervised relation extraction](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1306–1318, Online. Association for Computational Linguistics.
- Xu Han, Tianyu Gao, Yuan Yao, Deming Ye, Zhiyuan Liu, and Maosong Sun. 2019. [OpenNRE: An open and extensible toolkit for neural relation extraction](#). 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): System Demonstrations*, pages 169–174, Hong Kong, China. Association for Computational Linguistics.
- Kailong Hao, Botao Yu, and Wei Hu. 2021. [Knowing false negatives: An adversarial training method for distantly supervised relation extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9661–9672, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Geoffrey Hinton, Oriol Vinyals, Jeff Dean, et al. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2(7).
- Xiusheng Huang, Yubo Chen, Shun Wu, Jun Zhao, Yuantao Xie, and Weijian Sun. 2021. [Named entity recognition via noise aware training mechanism with data filter](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4791–4803, Online. Association for Computational Linguistics.
- Yuyun Huang and Jinhua Du. 2019. [Self-attention enhanced CNNs and collaborative curriculum learning for distantly supervised relation extraction](#). 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)*, pages 389–398, Hong Kong, China. Association for Computational Linguistics.
- Ganesh Jawahar, Benoît Sagot, and Djamé Seddah. 2019. What does bert learn about the structure of language? In *ACL 2019-57th Annual Meeting of the Association for Computational Linguistics*.
- Dongyang Li, Taolin Zhang, Nan Hu, Chengyu Wang, and Xiaofeng He. 2022. [HiCLRE: A hierarchical contrastive learning framework for distantly supervised relation extraction](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2567–2578, Dublin, Ireland. Association for Computational Linguistics.
- Yankai Lin, Shiqi Shen, Zhiyuan Liu, Huanbo Luan, and Maosong Sun. 2016. Neural relation extraction with selective attention over instances. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2124–2133.
- Tianyi Liu, Xiangyu Lin, Weijia Jia, Mingliang Zhou, and Wei Zhao. 2020a. [Regularized attentive capsule network for overlapped relation extraction](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6388–6398, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Tianyi Liu, Xinsong Zhang, Wanhao Zhou, and Weijia Jia. 2018. Neural relation extraction via inner-sentence noise reduction and transfer learning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2195–2204.
- Tianyu Liu, Kexiang Wang, Baobao Chang, and Zhifang Sui. 2017. A soft-label method for noise-tolerant distantly supervised relation extraction. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1790–1795.



relation classification. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 207–212.

## A Hyperparameter Analysis

There are two key hyperparameters in our experiments, the student selected and the head layer  $L$ . In our best model, we select the last student (12th) for evaluation and set layer 7 as the head layer.

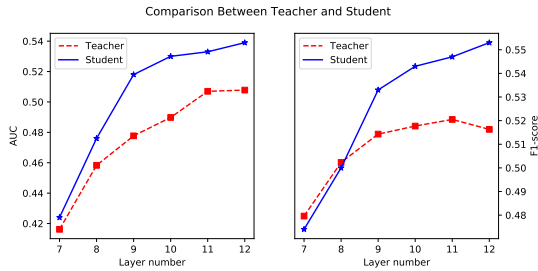


Figure 6: PR curves of the students and auxiliary classifiers of the teacher on the held-out dataset.

As shown in Figure 6, the higher students ( $\geq 9$ ) improve significantly over the teacher. The last student performs the best and the students from 9th to 11th also achieve comparable performances. Lower layers of BERT encode shallower features and the instance weighting in lower students is more affected by noise, so the performances of 7th and 8th students show little advantage over the teacher. With knowledge passed and noise alleviated student by student, the performance gradually improves.

Setting	AUC	F1	P@M
$L = 11$	53.4	55.1	82.8
$L = 10$	53.5	54.9	83.6
$L = 9$	53.6	55.0	84.0
$L = 8$	53.7	55.1	84.7
$L = 7$	<b>53.9</b>	<b>55.3</b>	<b>84.9</b>
$L = 6$	53.8	<b>55.3</b>	84.8
$L = 5$	53.7	55.1	84.6
$L = 3$	53.5	55.0	84.7
$L = 2$	53.5	54.9	84.6
$L = 1$	53.4	54.9	84.5

Table 6: Results of using different head layer  $L$  settings. The best results are marked as **bold**.

To study the effect of head layer  $L$ , we run experiments with  $L$  from 1 to  $n$ . In Table 6, we present the results where  $L = 7$  achieves the best performance. For  $L > 7$ , the head layer is too close to

the top, and TIW filters fewer false negatives. So the P@M declines quickly, which is similar to the effect of removing FNF as in Table 4. For  $L < 7$ , the lower layers of BERT are not able to encode sufficient information for accurate relation extraction, so the lower students are not able to provide reliable instance weights, leading to the transfer of some noise among students. Though other settings are less effective than the best, their performances still dominate the baselines. The above results show that our method is not dependent on the empirical settings of hyperparameters and further demonstrate the effectiveness and robustness of our method.