

S²ynRE: Two-stage Self-training with Synthetic data for Low-resource Relation Extraction

Benfeng Xu^{1*}, Quan Wang², Yajuan Lyu³, Dai Dai³
Yongdong Zhang¹⁴ and Zhendong Mao^{14†}

¹University of Science and Technology of China

²Beijing University of Posts and Telecommunications, ³Baidu Inc.

⁴Institute of Artificial Intelligence, Hefei Comprehensive National Science Center

benfeng@mail.ustc.edu.cn, zdmao@ustc.edu.cn

Abstract

Current relation extraction methods suffer from the inadequacy of large-scale annotated data. While distant supervision alleviates the problem of data quantities, there still exists domain disparity in data qualities due to its reliance on domain-restrained knowledge bases. In this work, we propose S²ynRE, a framework of two-stage Self-training with Synthetic data for Relation Extraction. We first leverage the capability of large language models to adapt to the target domain and automatically synthesize large quantities of coherent, realistic training data. We then propose an accompanied two-stage self-training algorithm that iteratively and alternately learns from synthetic and golden data together. We conduct comprehensive experiments and detailed ablations on popular relation extraction datasets to demonstrate the effectiveness of the proposed framework. Code is available at <https://github.com/BenfengXu/S2ynRE>.

1 Introduction

Relation extraction systems aim at discovering relational knowledge between entities by reading from unrestricted texts (Cardie, 1997). Although neural methods, especially pre-trained language models, have greatly advanced the state-of-the-art relation extraction capability (Zeng et al., 2014; Wu and He, 2019), they still require large quantities of training data (Han et al., 2020). However, high-quality annotations are usually very expensive to obtain, making low-resource relation extraction a very practical challenge in many real-world scenarios.

Distant supervision (Mintz et al., 2009), which automatically annotates relational statements by aligning entities with an existing knowledge bases (Bollacker et al., 2008; Vrandečić and Krötzsch, 2014), has been widely explored as an effective way to construct large scale relational

dataset. Many recent works exploit such data in a pretraining stage to learn relational representations (Baldini Soares et al., 2019; Peng et al., 2020; Qin et al., 2021). Although this line of methods have seen certain improvements, they still inevitably raise the concern that the distantly annotated data can vary considerably from downstream tasks both in target schema and in context distributions, thus may not be able to offer optimal transferability. For instance, due to the reliance on existing knowledge bases, current works mostly resort to Wikidata as the source of relational triples and Wikipedia (Vrandečić and Krötzsch, 2014) as the corpus for distant supervision. This circumscribes distant data to only factual knowledge between world entities, while downstream tasks may be of other special interests involving various domains, ranging from semantic relation between nominals (Hendrickx et al., 2009) to chemical-protein interactions (Kringelum et al., 2016).

Meanwhile, recent advances in large-scale pre-trained language models (LLM) (Radford et al., 2019; Brown et al., 2020; Raffel et al., 2020) have demonstrated their great potential in generating realistic texts of various domains (Radford et al., 2019). Accordingly, several very recent works have explored the possibility to exploit LLM as an alternative training data pool (Schick and Schütze, 2021; Vu et al., 2021). However, these attempts are confined to NLI task, while still not effectively explored in the area of relation extraction.

In this paper, we study the construction of synthetic data for relation extraction tasks to simultaneously address both training data scarcity in low resource scenarios and domain disparity in distant supervision. We employ LLM to estimate and adapt to the target domain distribution with only a few training instances, and synthesize a large amount of ones accordingly. The procedure is overall very simple but also carefully designed with two critical choices: 1) we linearize relational state-

*Joint work of USTC and Baidu Inc.

†Corresponding author: Zhendong Mao

ments into natural language sequences where entity pairs are indicated by special marker tokens; 2) we resort to unconditional generation instead of label-conditioned ones, which relaxes the requirements for strict label-semantic correspondence but increases sample availability and diversity.

We experiment with both GPT2 and the recent very large LLMs like GPT-3.5. For standard size generative LMs like GPT2-Large, we first finetune it to adapt to the target domain, while for the capable GPT-3.5 model, we directly apply In-Context Learning (Brown et al., 2020; Xu et al., 2023). We empirically found **finetuned** GPT2-Large produces synthetic data of equivalent quality to **prompted** GPT-3.5. In general, it is observed that with only a few accessible samples, we are able to successfully synthesize a large amount of domain-customized training data with satisfactory quality.

To effectively learn from such synthetic data, we novelly advocate a two-stage self-training algorithm. The approach in general follows the self-training framework (Yarowsky, 1995; Xie et al., 2020), which is widely employed to exploit unlabeled data. Typically, such methods iteratively annotate and learn pseudo labels for unlabeled data to bootstrap the model’s performance. Distinctively, we make a two-stage adaptation where in each of the iterations, the model is firstly trained on synthetic instances, then on golden ones. Such sequential training procedure favors golden data with more importance since they are introduced in the latter stage of the training curriculum.

We refer to our method as **S²ynRE**, a framework of two-stage Self-training with Synthetic data for Relation Extraction. The contributions of this paper is three-fold:

- **Conceptual Contribution** We exploit LLM to generate large quantities of domain adaptive synthetic data for low-resource relation extraction, and challenge the long-prevailing distant supervised methods restricted by KB domain coverages. The proposed solution novelly mitigates the problems of both data scarcity and domain disparity.
- **Technical Contribution** We propose a novel two-stage self-training algorithm to effectively learn from unlabeled synthetic data and golden data together. We demonstrate that this is a non-trivial adaptation that significantly outperforms standard self-training widely employed in semi-supervised learning.

- **Experimental Contribution** We conduct comprehensive experiments on 6 popular relation extraction datasets to investigate, analyze the propose method and make comparisons. We achieve new state-of-the-art for low-resource relation extraction. Compared to standard finetuning baseline, we obtain up to 17.18% absolute improvements, and 11.09% on average across all datasets.

2 Related Works

Relation Extraction Relation extraction is one of the fundamental tasks in natural language processing (Cardie, 1997), where lots of research efforts have been made to advance the state-of-the-art methods (Zeng et al., 2014; Zhou et al., 2016; Zhang et al., 2018; Baldini Soares et al., 2019), as well as the low-resource scenario (Han et al., 2018; Sainz et al., 2021; Dong et al., 2021; Chen et al., 2022). One of the most prominent methods is distant supervision (Mintz et al., 2009), which automatically constructs annotated relational data by aligning corpus with existing knowledge base. Many recent works investigate how to learn effectively with such distant data (Baldini Soares et al., 2019; Peng et al., 2020; Ding et al., 2021; Qin et al., 2021). Generally, they propose various pre-text tasks that pre-train a model to learn relational representation. We will further explain some of these works for comparison in Section 5.3.

Learning from Synthetic Data Built upon massive corpora, pre-trained language models are promising at producing texts of eligible quality, resulting in a surge of research interests in its usage for data augmentation (Feng et al., 2021). One straightforward way is to introduce mask corruptions in the way language models are pre-trained, then collect predictions as augmented data (Kobayashi, 2018; Ng et al., 2020). Later works further developed such technique into conditional augmentation (Wu et al., 2019; Kumar et al., 2020). Nevertheless, these methods are mostly editing existing instances, which limits the diversity and scale of augmented data.

With increasingly powerful LLMs, recent works turn to direct synthesis of new instances (Schick and Schütze, 2021; Wang et al., 2021; Meng et al., 2022; Ye et al., 2022). Different from this work, most of them focus on zero-shot language understanding where no labeled data is available (Schick and Schütze, 2021; Wang et al., 2021; Meng et al.,

2022; Ye et al., 2022). They investigate ways to generate label-conditioned data by prompting LLMs, but these methods can hardly be applied to low-resource or full data scenarios while still preserving effectiveness.

With the existence of labeled data, synthetic data needs to be of higher quality to bring further utility. Several works thus propose to finetune the generator (Anaby-Tavor et al., 2020; Vu et al., 2021; He et al., 2021). There are also explorations for learning from synthetic and golden data together, including threshold-based confidence filtering (Anaby-Tavor et al., 2020), classical semi-supervised learning (He et al., 2021) or restricting the usage of synthetic data within a supplemental intermediate task (Vu et al., 2021).

For structured learning tasks, Ding et al. (2020) similarly formulates NER task data as sequential language. Specifically for relational data synthesis, Papanikolaou and Pierleoni (2020) explore the biomedical domain and Chia et al. (2022) focus on zero-shot setting of triplet extraction. By contrast, Syn²RE distinguishes not only in applied scenario and synthesis strategy, but also in the two-stage learning framework, which is specially designed for improved synthetic data adaptation.

3 Preliminary

This section formulates the task of relation extraction and the baseline models used throughout all experiments.

Task Formulation A typical relation extraction task is defined by a corpus of relational statements and a set of relations, i.e., schema S . Assume the training dataset $\mathcal{D}^{tr} = \{(\mathbf{x}_i, s_i, o_i)\}_{i=1}^N$ and its corresponding labels $\mathcal{Y}^{tr} = \{y_i\}_{i=1}^N$, where \mathbf{x}_i is a sequence of words $\{w_l^i\}_{l=1}^L$, $y_i \in S$, $s_i = [w_{s_{start}} : w_{s_{end}}]$ and $o_i = [w_{o_{start}} : w_{o_{end}}]$ are subject and object entities within the context. The target is to learn a function $f_{\theta}(\mathbf{x}_i, s_i, o_i)$ that predicts the correct relation label y_i .

Baseline Model As S²ynRE is a data-centric framework, we keep the model architecture simple but competitive, which is the vanilla finetuning of pre-trained language models. Instead of autoregressive LMs, we use auto-encoding networks like BERT as they usually perform better on language understanding downstream tasks. Following Baldini Soares et al.’s (2019) comprehensive study of building relation extractors, we inject spe-

cial marker tokens to the input word sequence:

$$\mathbf{x}_{marked} = (\dots, [\text{Sub}], s, [\text{Sub}], \dots, [\text{Obj}], o, [\text{Obj}], \dots) \quad (1)$$

After the encoding process of transformer, the representation \mathbf{h} in corresponding positions will be concatenated for classification:

$$\hat{y} = \text{softmax}(\mathbf{W}^{|\mathcal{S}|}[\mathbf{h}_{[\text{Sub}]}; \mathbf{h}_{[\text{Obj}]}]) \quad (2)$$

where $W^{|\mathcal{S}|}$ is a feedforward network and the predicted categorical distribution \hat{y} will be trained against y using cross-entropy loss.

4 Methodology

4.1 Relational Data Synthesis

Training instances of relation extraction task is of specific structure (\mathbf{x}_i, s_i, o_i) , i.e., the relational statement is expected to be a sentence containing exact two entities as subject and object. Inspired by Paolini et al. (2021), we linearize relational data into marked natural language sequence as in Eq 1. The synthesizer can be built upon any existing LLMs. In this paper, we explore both GPT2 and the even larger GPT-3.5 as two representative LLMs and respectively employ finetuning or in-context learning treatment.

4.1.1 Finetuning for GPT2

The finetuning process is performed in the same autoregressive way as how it is pre-trained:

$$\mathcal{L} = - \sum_{l=1}^{L+4} \log P(w_l | w_0, \dots, w_{l-1}; LLM) \quad (3)$$

where $\{w_l\} = \mathbf{x}_{marked}$, and a <bos> token is prepended as w_0 . Note that we ignore relation labels y in training data and approach it as unconditional generation. This eliminates the noise caused by label-semantic inconsistency, and leaves it to model itself to learn from unlabeled synthetic data.

After the finetuning is completed, we simply prepend the <bos> token to prompt the generation, and repeatedly perform inference using multinomial sampling until we obtain the expected scale of synthetic data \mathcal{D}^{syn} . We show in appendix G that these synthetic data are coherent, realistic, and most importantly, customized to the target domain.

We elaborate on the framework of S²ynRE (see Fig. 2) in this section, including the construction of an LLM-based synthesizer, and the two-stage self-training algorithm.

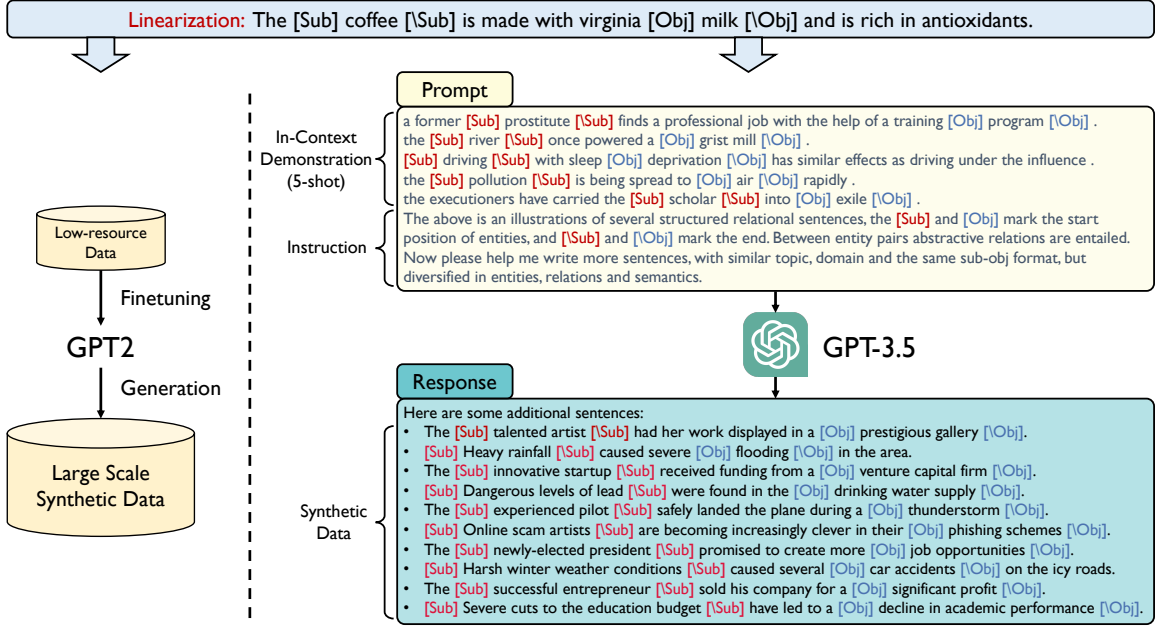


Figure 1: Illustration of relational data synthesis. We either finetune GPT2 or employ in-context learning for even larger LLM like GPT-3.5.

4.1.2 In-Context Learning for GPT-3.5

Even larger LLMs naturally exhibit few-shot learning capabilities, and can be elicited through very few in-context demonstrations (Brown et al., 2020; Xu et al., 2023). We directly prepend 5-shot exemplars and provide a specific instruction asking the LLM to generate more examples. The process is illustrated in Figure 1.

4.2 Two Stage Self-training

Self-training is a widely adopted learning algorithm for semi-supervised learning. Typically, to jointly learn from an unlabeled dataset and a labeled dataset, it iteratively samples from the unlabeled set, assigns them with pseudo labels, merges them with the labeled dataset, and re-trains the model. In this paper, we argue that this design of naive merging is built upon a strong assumption that the unlabeled dataset must be in the exact distribution with the labeled ones, for which the synthetic data does not strictly satisfy.

In S²ynRE, differently, we make a two-stage adaptation: where synthetic data and golden data are trained sequentially (Figure 2). We start from a base model initialized using any auto-encoding language models, e.g., BERT (Devlin et al., 2019), and train it on \mathcal{D}^{tr} to produce a teacher model η , as introduced in Section 3. We first use η to annotate the unlabeled synthetic data \mathcal{D}^{syn} :

$$\hat{y}_i^{syn} = \eta(\mathbf{x}_i^{syn}, s_i, o_i) \quad (4)$$

and we keep $\hat{\mathcal{Y}}^{syn} = \{\hat{y}_i^{syn}\}$ as soft pseudo labels of \mathcal{D}^{syn} , note that here the $\hat{\cdot}$ denotes *soft* as we keep the categorical distribution intact instead of keeping its argmax. Inspired by Li and Qian (2021), to further eliminate fluctuations in pseudo labels, we train multiple teachers using different random seeds, and the pseudo labels annotated by k -th teacher is referred to as \hat{y}_k^{syn} .

We then re-initialize a new student model θ , and apply a two-stage training strategy. In stage-one training, student θ is trained on synthetic data using soft pseudo labels:

$$\theta' \leftarrow \mathcal{L}_{KD}(\theta, \mathcal{D}^{syn}, \{\hat{\mathcal{Y}}_k^{syn}\}_{k=1}^K) \quad (5)$$

This can be seen as a distillation procedure that transfers knowledge from η to θ based on synthetic data \mathcal{D}^{Syn} . And \mathcal{L}_{KD} is calculated as:

$$\mathcal{L}_{KD} = \frac{1}{K} \sum_{k=1}^K D_{KL}(\hat{y}_i^{syn} \parallel \theta(\mathbf{x}_i^{syn}, s_i, o_i)) \quad (6)$$

where D_{KL} is the Kullback-Leibler divergence. Then in stage-two training, we take from θ' , and train it on labeled training dataset:

$$\theta'' \leftarrow \mathcal{L}_{CE}(\theta', \mathcal{D}^{tr}, \mathcal{Y}^{tr}) \quad (7)$$

where \mathcal{L}_{CE} is the standard cross-entropy loss, and θ'' is the resulting model in this iteration. We then use θ'' as the teacher model η for the next iteration to re-annotate \mathcal{D}^{syn} , and this procedure is

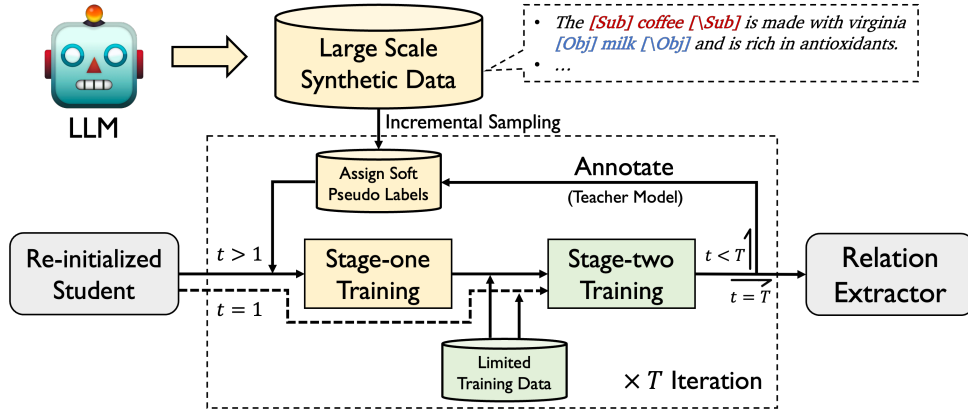


Figure 2: The overall framework of S²ynRE. We iteratively train the student model on both synthetic and golden data via a two-stage self-training strategy. Note that in iteration $t = 1$, stage-two is directly applied. The exemplary instance is sampled from our synthetic data for SemEval.

repeated T times. Following the standard practice of self-training, in each iteration, we incrementally sample $1/T$ more synthetic data from \mathcal{D}^{syn} until in iteration T , where \mathcal{D}^{syn} will be running out of new instances. The entire two-stage self-training process can be formulated as Algorithm 1.

Algorithm 1: Two-stage Self-training.

Input: Golden training dataset \mathcal{D}^{tr} , \mathcal{Y}^{tr} ,
synthetic dataset \mathcal{D}^{syn}

```

/* ===== Iteration 1 ===== */
t = 1;
 $\mathcal{D}_1^{syn} = \emptyset$ ;
Initialize  $\theta$  from auto-encoding LM;
 $\theta_1 \leftarrow Train(\theta, \mathcal{D}^{tr}, \mathcal{Y}^{tr})$ ; // Eq.7
 $\theta_1^{Tea} \leftarrow \theta_1$ ; // assign teacher model
/* ===== Iteration 2~T ===== */
repeat
  t = t + 1;
   $\mathcal{D}_t^{syn} = \mathcal{D}_{t-1}^{syn} \cup \mathcal{D}^{syn}[\frac{t-1}{T} : \frac{t}{T}]$ ;
   $\hat{\mathcal{Y}}_t^{syn} \leftarrow Annotate(\theta_{t-1}^{Tea}, \mathcal{D}_t^{syn})$ ;
  // Eq.4
  Re-initialize  $\theta$  from auto-encoding LM;
  /* stage-one training */
   $\theta'_t \leftarrow Train(\theta, \mathcal{D}_t^{syn}, \hat{\mathcal{Y}}_t^{syn})$ ; // Eq.5
  /* stage-two training */
   $\theta''_t \leftarrow Train(\theta'_t, \mathcal{D}^{tr}, \mathcal{Y}^{tr})$ ; // Eq.7
   $\theta_t^{Tea} \leftarrow \theta''_t$ ; // update teacher model
until performance converges or t reaches
maximum iteration limit T;
Output: Final model  $\theta''_t$ 

```

Dataset	Train	Dev	Test	1% Train	Relation
Semeval	6507	1493	2717	73	19
TACRED	68124	22631	15509	703	42
TACRED-Revisited	68124	22631	15509	703	42
Re-TACRED	58465	19584	13418	570	40
ChemProt	4169	2427	3469	49	13
Wiki80	39200	5600	11200	400	80

Table 1: Numbers of instances in train, dev, test splits and low resource settings.

5 Experiments

5.1 Experimental Settings

We evaluate S²ynRE on popular datasets including **SemEval 2010 Task 8** (Hendrickx et al., 2009), **TACRED** (Zhang et al., 2017), **TACRED-Revisited** (Alt et al., 2020), **Re-TACRED** (Stoica et al., 2021), **ChemProt** (Kringelum et al., 2016) and **Wiki80** (Han et al., 2019). Their statistics are given in Table 1 and we refer to detailed introduction in Appendix A.

For each dataset, we set three different prerequisites of resource availability. Respectively, *FULL* for 100% training data, *LIMITED* for 10% training data and *FEW* for 1% training data. To provide robust and convincing conclusions, we run all experiments (including ablation studies) with 5 different random seeds and report their average. With each random seed, we employ grid search to select the best model as well as the teacher model in each iteration. We use only development set for such selection, and report the corresponding test set score as the final results.

For data synthesis, we use GPT2-Large and GPT-3.5 as the aforementioned LLMs. Specifically, for ChemProt, we use an adapted version of GPT-2 (Pa-

Backbone Init	FT Baseline	S ² ynRE w/ gpt-3.5-turbo ICL	S ² ynRE w/ GPT2 Finetuning
BERT	40.81 \pm 1.62	56.81 \pm 1.51 (+16.00)	57.53 \pm 0.96 (+16.72)
CP	52.77 \pm 1.83	65.39 \pm 0.73 (+12.62)	67.32 \pm 0.54 (+14.55)

Table 2: S²ynRE with Different LLMs. Synthetic data brings significant improvements, and GPT-3.5 with In-Context Learning is just as effective as GPT-2 with Finetuning.

panikolaou and Pierleoni, 2020), which is further trained on 500k PubMed abstracts. When generating, we restrict sequence length to 128, and perform necessary filtering by removing instances that do not conform with the relational structure, i.e., there must exist 4 exact special markers and each start position marker shall appear before its end position marker. The synthesis efficiency is 24.05 instances per second before any filtering. In total, we collect 10,000 samples for *FEW* setting, and 100,000 synthetic samples for *LIMITED* and *FULL* settings. We leave other hyper-parameters to Appendix B.

5.2 Capability of LLM

We first validate the capability of LLMs as data synthesizer and their respective treatment. Considering both affordability and model capability, we use the recently released **gpt-3.5-turbo-0301** API² as even larger LLMs in comparison with GPT2-Large. For in-context learning, we repeatedly sample 5-shot random examples from the golden training set as demonstrations, followed by an instruction that asks LLMs to generate more, with domain, format and diversity constraints. We sent 1,321 queries in total to collect 10,000 synthetic data (the same as our previous experimental settings using GPT-2). Each query produces different synthetic results due to LLM sampling strategy and varied selection and permutation of in-context demonstrations.

The results are shown in Table 2. We empirically found **finetuned** GPT2-Large produces synthetic data of equivalent quality to **prompted** GPT-3.5, while both are effective training data synthesizers that significantly outperforms baselines. In the following experiments, we use GPT2-Large to further verify the proposed S²ynRE framework.

5.3 Main Results

We choose competitive baselines and reproduce them under comparable settings to provide more reliable conclusions. These baseline methods are: **BERT** We finetune BERT model (Devlin et al.,

²<https://openai.com/blog/introducing-chatgpt-and-whisper-apis>

2019) in a straightforward way for relation extraction as explained in Section 3 and implemented in many existing works. This serves as our re-implemented *Finetune Baseline* and will be referred to in the following figures.

MTB (Baldini Soares et al., 2019) pre-trains a relational encoder using matching the blanks task, which is built on the hypothesis that two relational statements containing the same entity pair should express similar relational representations. Note that this is a weaker reliance than distant supervision as it only aligns entities, and does not need relations.

CP (Peng et al., 2020) proposes a contrastive learning pretext task that encourages sentence representations with the same relation to be similar and different ones to be disparate.

ERICA (Qin et al., 2021) further extends distant supervision to document-level corpus, and design similar pretext task that discriminates relational representations across sentences.

We provide an overview of these works regarding various resource usage and requirements in Table 4. The main results are shown in Table 3. On Wiki80, we directly use distant data as they are available in the general wiki domain, we analysis the effects later in Table 5. Under all three settings across five datasets, S²ynRE outperforms the BERT finetune baseline. Specifically for the *FEW* setting, improvements are much more significant, respectively **+17.18**, **+15.47**, **+16.86**, **+8.07**, **+5.59**, and **+3.34**, resulting an average improvements of **+11.09** across all 6 datasets. We further employ CP as a stronger base model to initialize the students, and the performances are even better. This implies that the improvements of S²ynRE are mostly orthogonal with those of the distantly pre-trained methods. In general, S²ynRE_{CP} achieves new state-of-the-art for low resource relation extraction tasks.

²We obtain MTB and CP checkpoints from <https://github.com/thunlp/RE-Context-or-Names> and ERICA checkpoint from <https://github.com/thunlp/ERICA>

Method	SemEval	TACRED	TACRED-Revisited	Re-TACRED	ChemProt	Wiki80
<i>FULL (100% training data)</i>						
BERT	88.86 \pm 0.30	69.27 \pm 0.27	79.24 \pm 0.37	87.75 \pm 0.22	81.66 \pm 0.79	91.54 \pm 0.08
MTB	88.95 \pm 0.31	69.93 \pm 0.40	79.69 \pm 0.32	87.67 \pm 0.37	81.75 \pm 0.86	90.07 \pm 0.97
CP	89.16 \pm 0.17	<u>70.16</u> \pm 0.20	80.08 \pm 0.32	87.95 \pm 0.09	<u>81.77</u> \pm 0.97	90.44 \pm 0.38
ERICA	88.62 \pm 0.24	68.91 \pm 0.75	78.95 \pm 0.86	87.73 \pm 0.31	81.52 \pm 0.43	91.47 \pm 0.13
S²ynRE_{BERT}	89.20 \pm 0.27	70.25 \pm 0.47	<u>79.80</u> \pm 0.29	88.01 \pm 0.24	81.65 \pm 0.60	91.54 \pm 0.14
S²ynRE_{CP}	89.04 \pm 0.32	70.03 \pm 0.27	79.75 \pm 0.49	<u>87.98</u> \pm 0.07	82.15 \pm 0.12	91.33 \pm 0.20
<i>LIMITED (10% training data)</i>						
BERT	82.38 \pm 0.51	59.32 \pm 0.35	66.56 \pm 0.48	80.51 \pm 0.77	68.96 \pm 0.97	85.89 \pm 0.22
MTB	82.56 \pm 0.27	59.45 \pm 0.55	66.48 \pm 0.71	81.15 \pm 0.59	71.44 \pm 1.12	82.42 \pm 2.27
CP	<u>83.80</u> \pm 0.50	<u>62.81</u> \pm 0.39	70.81 \pm 0.58	<u>83.42</u> \pm 0.41	71.89 \pm 1.09	85.86 \pm 0.95
ERICA	82.41 \pm 0.55	58.54 \pm 0.65	66.65 \pm 0.68	80.45 \pm 0.77	69.03 \pm 1.22	86.67 \pm 0.49
S²ynRE_{BERT}	84.01 \pm 0.23	61.26 \pm 0.53	68.62 \pm 0.15	83.28 \pm 0.40	<u>73.62</u> \pm 0.14	85.79 \pm 0.49
S²ynRE_{CP}	84.64 \pm 0.30	62.94 \pm 0.45	<u>70.36</u> \pm 0.75	84.36 \pm 0.32	75.32 \pm 0.92	85.94 \pm 0.95
<i>FEW (1% training data)</i>						
BERT	40.81 \pm 1.62	30.40 \pm 7.74	33.75 \pm 8.68	54.75 \pm 4.52	39.50 \pm 1.47	63.34 \pm 0.76
MTB	45.12 \pm 1.23	36.52 \pm 2.00	40.69 \pm 2.25	58.35 \pm 0.93	41.53 \pm 2.11	62.29 \pm 1.84
CP	53.29 \pm 1.80	<u>49.81</u> \pm 0.59	<u>55.53</u> \pm 0.90	<u>68.03</u> \pm 0.76	43.96 \pm 2.62	80.93 \pm 0.89
ERICA	43.62 \pm 2.33	34.91 \pm 1.40	39.17 \pm 1.69	57.14 \pm 0.83	40.01 \pm 0.86	68.65 \pm 0.95
S²ynRE_{BERT}	<u>57.99</u> \pm 1.08	45.87 \pm 1.07	50.61 \pm 0.99	62.82 \pm 0.52	<u>45.09</u> \pm 0.38	66.68 \pm 0.68
S²ynRE_{CP}	68.03 \pm 0.46	51.91 \pm 0.68	58.48 \pm 0.29	70.21 \pm 0.81	46.23 \pm 0.73	80.93 \pm 0.89

Table 3: Main results. Best performances are **bold**, and the second bests are underlined. We report Accuracy for Chemprot and Wiki80, and Micro-F1 for other datasets. Results for all baseline methods are reproduced with identical hyper-parameter searches for fair comparison¹.

Dataset	Resource Usage	Domain	External Requirements	
			KB Entities	KB Relations
MTB	6,000,000 sent pairs	Wiki	✓	No Requirements
CP	867, 278 sents	Wiki	✓	✓
ERICA	1,000,000 docs	Wiki	✓	✓
S ² ynRE	100,000 sents	Customized		No Requirements

Table 4: Comparison of external resource usage and requirements for different methods.

5.4 Ablation Study

We investigate the advantages of S²ynRE via comprehensive ablations. In accordance with the main claim, all experiments are conducted under the low-resource (*FEW*) setting unless otherwise stated.

Synthetic Data Instead of Distant Data Distant supervision has long been the prevailing solution to automatically construct relational data. We make its comparison against the proposed synthetic data in Table 5. We keep the two-stage self-training algorithm intact, only replace the synthetic data with distant data³. On 5 investigated datasets, distant data can provide appreciable improvements ranging from **+2.06** to **+13.25**, however, synthetic data brings much more significant improvements ranging from **+5.59** to **+17.18**, which clearly demon-

³The distant data is produced and released by Peng et al. (2020), we randomly sample 100,000 instances out of it

Dataset	NA	Distant	Synthetic
SemEval	40.81	49.36 (+ 8.55)	57.99 (+17.18)
ChemProt	39.50	41.56 (+ 2.06)	45.09 (+ 5.59)
TACRED	30.40	42.43 (+12.03)	45.87 (+15.47)
Re-TACRED	54.75	62.34 (+ 7.59)	62.98 (+ 8.23)
TACRED-Revisited	33.75	47.00 (+13.25)	50.61 (+16.86)
Wiki80	63.08	66.68 (+ 3.60)	65.52 (+ 2.44)

Table 5: Comparison between synthetic data and distant data. Inside the parentheses are absolute improvements, **red** means the higher one.

strates the superiority of being domain-customized for target tasks. However, on Wiki80, which very closely follows identical distribution of distant data as both are constructed using distant supervision on wikipedia and wikidata, result shows that synthetic data provides competitive improvements but no longer outperforms distant ones. This verifies the importance and advantage of domain-customized data from an opposite perspective. Nevertheless, real-world scenarios mostly involve distribution beyond the scope of wikipedia, and only the proposed synthetic approach can offer such advantage. We also provide qualitative comparisons for synthetic and distant data in Appendix G to better illustrate the discussed domain disparity.

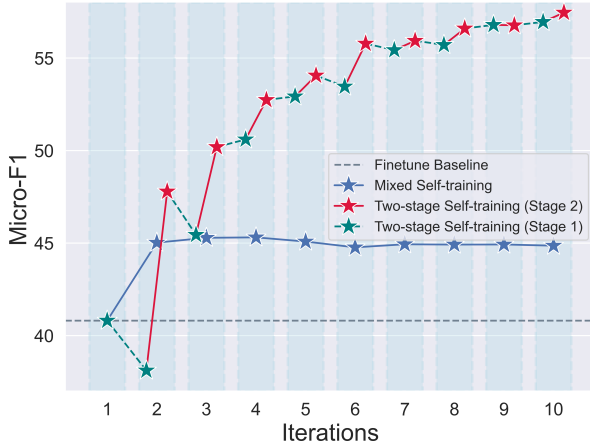


Figure 3: Performance illustration for two-stage self-training compared to classical mixed self-training. Analyzed on SemEval.

Two Stage Self-training Typical self-training algorithms merge the pseudo-labeled data into existing labeled data in each iteration, and minimize the model’s empirical loss on a mixture of both. We refer to such classical implementation as **Mixed Self-training** as opposed to the proposed **Two-stage Self-training**. Fig. 3 compares these two approaches. In each iteration (transparent blue bar), there will be one evaluation for mixed self-training (blue curve), but two evaluation for Two-stage Self-training (teal for stage one, Red for stage two). We observe that in stage-one training, the performance might drop a few compare to its previous iteration, however, it effectively provides a better initialization where the model can further learn from the golden data. Overall, the model can continually bootstrap its performance by learning from synthetic and golden data iteratively and alternately. While in mixed self-training, the golden data are treated equally as synthetic ones, and the model is overwhelmed by large amounts of the latter. Therefore, the improvement quickly saturates to a limited plateau. We also provide illustrations of the bootstrapping performance over iterations on other datasets in Appendix C.

Comparison Under Semi-supervised Setting

Standard semi-supervised setting also investigates low-resource relation extraction by joint learning from both labeled data and unlabeled data. However, they make a strong assumption of identical distribution between unlabeled data and labeled ones, and most existing works actually directly sample from the golden training data and remove the labels to

Method	SemEval	TACRED
MetaSRE	80.09 \pm 0.78	56.95 \pm 0.34
GradLRE	81.69 \pm 0.57	58.20 \pm 0.33
S ² ynRE w/ Golden	84.11 \pm 0.27	59.07 \pm 0.54
S ² ynRE w/ Synthetic	84.01 \pm 0.23	61.26 \pm 0.53

Table 6: Comparison to state-of-the-art methods for semi-supervised setting, including (Hu et al., 2021a) and GradLRE (Hu et al., 2021b). w/ Golden means unlabeled set are sampled from 50% of the golden training data and their original labels are removed accordingly.

Dataset	NA	Conditional Syn	Unconditional Syn
SemEval	40.81	45.26 (+4.45)	57.99
TACRED	30.40	33.34 (+2.94)	45.87
Re-TACRED	54.75	53.03 (-1.72)	62.98
TACRED-Revisited	33.75	37.60 (+3.85)	50.61

Table 7: Comparison between conditional and unconditional synthesis. Inside the parentheses denote the effectiveness comparing to Finetune Baseline.

construct the unlabeled set. We provide comparison with state-of-the-art methods of semi-supervised learning in Table 6 (under the *LIMITED* setting). Results show that 1) the proposed two-stage self-training outperforms other semi-supervised learning algorithms, and 2) synthetic data demonstrates better or comparable performance compared to unlabeled set constructed from golden training data. We attribute the latter to its domain-customized quality and unlimited large-scale quantity.

Unconditional Generation Although a lot of previous works intuitively resort to conditional synthesis, we show that this is not the optimal choice for relation extraction task. We finetune the synthesizer by prepending label-specific prompts: “write a sentence describing relation $V(r)$: ”, where $V(r)$ is the verbalizer for each relation r and we directly use corresponding label strings, e.g., *Component-Whole(e2,e1)*. We synthesize each relation class proportional to its original distribution in golden dataset. As conditional generation provides already labeled data, we can directly finetune the student model instead of self-training. We still train synthetic and golden data sequentially as we empirically found it a better choice. The results show that conditional generation only brings minimum or no benefits. We attribute this to the difficulty of preserving required label semantics for highly abstractive tasks like relation extraction. As a consequence, while these extra amounts of data can still provide certain usability, they also most likely

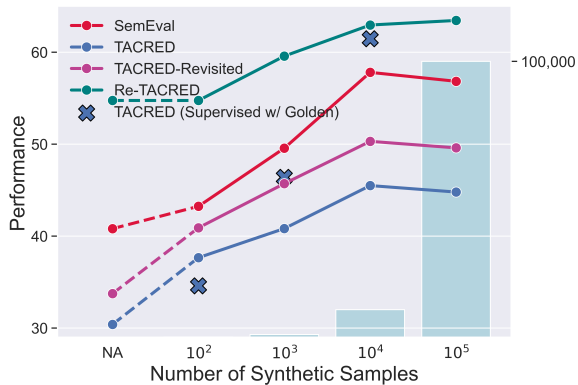


Figure 4: Performances w.r.t. different scales of synthetic data usage.

Scale	Golden	Synthetic
100	98.9	97.8 (- 1.1)
1,000	96.8	88.8 (- 8.0)
10,000	88.6	74.3 (-14.3)

Table 8: Sample diversity (type-token ratio in percentage for 3-grams) of synthetic and golden data w.r.t. different data scales on SemEval.

cause considerable distractions.

Scale of Synthetic Samples Figure 4 investigates the scale of synthetic samples. The improvements are approximately increasing in log scale w.r.t. the number of synthetic samples. The best performance is reached at 10,000, after which if we keep adding more samples, the performance saturates. As the synthesis of data is a repeatedly sampling process, we think exploiting too much data will deteriorate the diversity at the same time. We verify this by evaluating its diversity using type-token ratio (Roemmele et al., 2017; Kumar et al., 2020), which is defined as the ratio of unique n-grams out of all n-grams (see Table 8). We can see that the diversity gap between synthetic and golden data is enlarged when increasing the data scale.

We also report supervised results using additional golden training data to measure the utility of synthetic data. We can achieve two conclusions: 1) the advantage of golden training data are more significant when it is scaled up (10^4). However, this also takes substantially expensive costs. 2) S^2 ynRE approximately achieves the utility of 1,000 additional annotated golden data (10^3), and it only costs several hours of GPU computation to produce according synthetic data (10^4) as needed.

6 Discussion

Distant supervision is the most prevalent solution for low-resource relation extraction, and also the main investigated and compared baseline in this paper. Both distant data and the proposed synthetic data can essentially be recognized as ways of data augmentation to produce sufficient number of additional data. The critical difference which determines the effectiveness lies in their consistency with golden training data, i.e., domain affinity. And in this paper, the superiority of synthetic data is both experimentally proved (Table 5) and qualitatively explained (Appendix G). In conclusion, leveraging LLM to adapt to target domain and generate synthetic data of high utility is in general an performant solution and we hope this novel perspective can further inspire future insights in many related areas that have been greatly impacted by the idea of distant supervision.

7 Conclusion

In this paper, we present S^2 ynRE, a framework of two-stage self-training with synthetic data for relation extraction. We show that synthetic data generated using LLMs can resolve data scarcity in low-resource scenarios and mitigate domain disparity compared to distant supervision. To enable effective learning from such synthetic data, we then propose a novel two-stage self-training algorithm that continually bootstraps model performance by iteratively and alternately training the synthetic and golden data together. The proposed framework brings substantial improvements and achieves new state-of-the-art for low-resource relation extraction. In the future, we expect new possibilities brought by LLMs and will further explore accompanied techniques to exploit their potential.

Ethical Considerations

Synthetic data generated by language models may involve potential ethical risks regarding fairness and bias (Bommasani et al., 2021; Blodgett et al., 2020), which results in further consideration when they are employed in downstream NLP tasks. Although the scope of this paper remains how to produce and leverage such synthetic data to improve relation extraction system, it is worth further investigation to investigate in conjunction with well-established methods that can measure (Nadeem et al., 2021) and mitigate (Nadeem et al., 2021; Gupta et al., 2022) such ethical risks.

Acknowledgements

This work is supported by the National Key Research and Development Program of China No.2020AAA0109400, the National Natural Science Foundation of China under Grant 62222212, U19A0527, 62121002 and 61876223.

References

- Christoph Alt, Aleksandra Gabryszak, and Leonhard Hennig. 2020. [TACRED revisited: A thorough evaluation of the TACRED relation extraction task](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1558–1569, Online. Association for Computational Linguistics.
- Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich, Amir Kantor, George Kour, Segev Shlomov, Naama Tepper, and Naama Zwerdling. 2020. [Do not have enough data? deep learning to the rescue!](#) *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):7383–7390.
- Livio Baldini Soares, Nicholas FitzGerald, Jeffrey Ling, and Tom Kwiatkowski. 2019. [Matching the blanks: Distributional similarity for relation learning](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2895–2905, Florence, Italy. Association for Computational Linguistics.
- Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciBERT: A pretrained language model for scientific text](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3615–3620, Hong Kong, China. Association for Computational Linguistics.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. [Language \(technology\) is power: A critical survey of “bias” in NLP](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476, Online. 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, SIGMOD ’08*, page 1247–1250, New York, NY, USA. Association for Computing Machinery.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, and et al. 2021. [On the opportunities and risks of foundation models](#). *CoRR*, abs/2108.07258.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Claire Cardie. 1997. [Empirical methods in information extraction](#). *AI Magazine*, 18(4):65.
- Xiang Chen, Ningyu Zhang, Xin Xie, Shumin Deng, Yunzhi Yao, Chuanqi Tan, Fei Huang, Luo Si, and Huajun Chen. 2022. [Knowprompt: Knowledge-aware prompt-tuning with synergistic optimization for relation extraction](#). In *Proceedings of the ACM Web Conference 2022, WWW ’22*, page 2778–2788, New York, NY, USA. Association for Computing Machinery.
- Yew Ken Chia, Lidong Bing, Soujanya Poria, and Luo Si. 2022. [Relationprompt: Leveraging prompts to generate synthetic data for zero-shot relation triplet extraction](#). In *Findings of the Association for Computational Linguistics: ACL 2022*.
- 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.
- Bosheng Ding, Linlin Liu, Lidong Bing, Canasai Kruengkrai, Thien Hai Nguyen, Shafiq Joty, Luo Si, and Chunyan Miao. 2020. [DAGA: Data augmentation with a generation approach for low-resource tagging tasks](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*

- (EMNLP), pages 6045–6057, Online. Association for Computational Linguistics.
- Ning Ding, Xiaobin Wang, Yao Fu, Guangwei Xu, Rui Wang, Pengjun Xie, Ying Shen, Fei Huang, Hai-Tao Zheng, and Rui Zhang. 2021. [Prototypical representation learning for relation extraction](#). In *International Conference on Learning Representations*.
- Manqing Dong, Chunguang Pan, and Zhipeng Luo. 2021. [MapRE: An effective semantic mapping approach for low-resource relation extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2694–2704, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Edward Hovy. 2021. [A survey of data augmentation approaches for NLP](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.
- Umang Gupta, Jwala Dhamala, Varun Kumar, Apurv Verma, Yada Pruksachatkun, Satyapriya Krishna, Rahul Gupta, Kai-Wei Chang, Greg Ver Steeg, and Aram Galstyan. 2022. [Mitigating gender bias in distilled language models via counterfactual role reversal](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 658–678, Dublin, Ireland. Association for Computational Linguistics.
- Xu Han, Tianyu Gao, Yankai Lin, Hao Peng, Yaoliang Yang, Chaojun Xiao, Zhiyuan Liu, Peng Li, Jie Zhou, and Maosong Sun. 2020. [More data, more relations, more context and more openness: A review and outlook for relation extraction](#). In *Proceedings of the 1st Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing*, pages 745–758, Suzhou, China. 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.
- Xu Han, Hao Zhu, Pengfei Yu, Ziyun Wang, Yuan Yao, Zhiyuan Liu, and Maosong Sun. 2018. [FewRel: A large-scale supervised few-shot relation classification dataset with state-of-the-art evaluation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4803–4809, Brussels, Belgium. Association for Computational Linguistics.
- Xuanli He, Islam Nassar, Jamie Kiros, Gholamreza Hafari, and Mohammad Norouzi. 2021. [Generate, annotate, and learn: Nlp with synthetic text](#).
- Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2009. [SemEval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals](#). In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions (SEW-2009)*, pages 94–99, Boulder, Colorado. Association for Computational Linguistics.
- Xuming Hu, Chenwei Zhang, Fukun Ma, Chenyao Liu, Lijie Wen, and Philip S. Yu. 2021a. [Semi-supervised relation extraction via incremental meta self-training](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 487–496, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xuming Hu, Chenwei Zhang, Yawen Yang, Xiaohe Li, Li Lin, Lijie Wen, and Philip S. Yu. 2021b. [Gradient imitation reinforcement learning for low resource relation extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 2737–2746, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sosuke Kobayashi. 2018. [Contextual augmentation: Data augmentation by words with paradigmatic relations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 452–457, New Orleans, Louisiana. Association for Computational Linguistics.
- Jens Kringelum, Sonny Kim Kjaerulff, Søren Brunak, Ole Lund, Tudor I Oprea, and Olivier Taboureau. 2016. Chemprot-3.0: a global chemical biology diseases mapping. *Database*, 2016.
- Varun Kumar, Ashutosh Choudhary, and Eunah Cho. 2020. [Data augmentation using pre-trained transformer models](#). In *Proceedings of the 2nd Workshop on Life-long Learning for Spoken Language Systems*, pages 18–26, Suzhou, China. Association for Computational Linguistics.
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2019. [BioBERT: a pre-trained biomedical language representation model for biomedical text mining](#). *Bioinformatics*, 36(4):1234–1240.
- Wanli Li and Tiejun Qian. 2021. From consensus to disagreement: Multi-teacher distillation for semi-supervised relation extraction. *arXiv preprint arXiv:2112.01048*.

- Yu Meng, Jiaxin Huang, Yu Zhang, and Jiawei Han. 2022. [Generating training data with language models: Towards zero-shot language understanding](#). *CoRR*, abs/2202.04538.
- Mike Mintz, Steven Bills, Rion Snow, and Daniel Jurafsky. 2009. [Distant supervision for relation extraction without labeled data](#). In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1003–1011, Suntec, Singapore. Association for Computational Linguistics.
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. [StereoSet: Measuring stereotypical bias in pretrained language models](#). 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 5356–5371, Online. Association for Computational Linguistics.
- Nathan Ng, Kyunghyun Cho, and Marzyeh Ghassemi. 2020. [SSMBA: Self-supervised manifold based data augmentation for improving out-of-domain robustness](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1268–1283, Online. Association for Computational Linguistics.
- Giovanni Paolini, Ben Athiwaratkun, Jason Krone, Jie Ma, Alessandro Achille, RISHITA ANUBHAI, Cicero Nogueira dos Santos, Bing Xiang, and Stefano Soatto. 2021. [Structured prediction as translation between augmented natural languages](#). In *International Conference on Learning Representations*.
- Yannis Papanikolaou and Andrea Pierleoni. 2020. [Dare: Data augmented relation extraction with gpt-2](#). *arXiv preprint arXiv:2004.13845*.
- Hao Peng, Tianyu Gao, Xu Han, Yankai Lin, Peng Li, Zhiyuan Liu, Maosong Sun, and Jie Zhou. 2020. [Learning from Context or Names? An Empirical Study on Neural Relation Extraction](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3661–3672, Online. Association for Computational Linguistics.
- Yujia Qin, Yankai Lin, Ryuichi Takanobu, Zhiyuan Liu, Peng Li, Heng Ji, Minlie Huang, Maosong Sun, and Jie Zhou. 2021. [ERICA: Improving entity and relation understanding for pre-trained language models via contrastive learning](#). 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 3350–3363, Online. Association for Computational Linguistics.
- Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Melissa Roemmele, Andrew S Gordon, and Reid Swanson. 2017. [Evaluating story generation systems using automated linguistic analyses](#). In *Workshop on Machine Learning for Creativity, at the 23rd SIGKDD Conference on Knowledge Discovery and Data Mining*.
- Oscar Sainz, Oier Lopez de Lacalle, Gorka Labaka, Ander Barrena, and Eneko Agirre. 2021. [Label verbalization and entailment for effective zero and few-shot relation extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1199–1212, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Timo Schick and Hinrich Schütze. 2021. [Generating datasets with pretrained language models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6943–6951, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- George Stoica, Emmanouil Antonios Platanios, and Barnabas Poczos. 2021. [Re-tacred: Addressing shortcomings of the tacred dataset](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(15):13843–13850.
- Denny Vrandečić and Markus Krötzsch. 2014. [WikiData: A free collaborative knowledgebase](#). *Commun. ACM*, 57(10):78–85.
- Tu Vu, Minh-Thang Luong, Quoc Le, Grady Simon, and Mohit Iyyer. 2021. [STraTA: Self-training with task augmentation for better few-shot learning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5715–5731, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zirui Wang, Adams Wei Yu, Orhan Firat, and Yuan Cao. 2021. [Towards zero-label language learning](#). *CoRR*, abs/2109.09193.
- Shanchan Wu and Yifan He. 2019. [Enriching pre-trained language model with entity information for relation classification](#). In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, CIKM '19*, page 2361–2364, New York, NY, USA. Association for Computing Machinery.
- Xing Wu, Shangwen Lv, Liangjun Zang, Jizhong Han, and Songlin Hu. 2019. [Conditional bert contextual augmentation](#). In *Computational Science – ICCS 2019*, pages 84–95, Cham. Springer International Publishing.

- Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V. Le. 2020. [Self-training with noisy student improves imagenet classification](#). In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695.
- Benfeng Xu, Quan Wang, Zhendong Mao, Yajuan Lyu, Qiaoqiao She, and Yongdong Zhang. 2023. [\\$k\\$NN prompting: Beyond-context learning with calibration-free nearest neighbor inference](#). In *The Eleventh International Conference on Learning Representations*.
- David Yarowsky. 1995. [Unsupervised word sense disambiguation rivaling supervised methods](#). In *33rd Annual Meeting of the Association for Computational Linguistics*, pages 189–196, Cambridge, Massachusetts, USA. Association for Computational Linguistics.
- Jiacheng Ye, Jiahui Gao, Qintong Li, Hang Xu, Jiangtao Feng, Zhiyong Wu, Tao Yu, and Lingpeng Kong. 2022. [Zerogen: Efficient zero-shot learning via dataset generation](#).
- Daojian Zeng, Kang Liu, Siwei Lai, Guangyou Zhou, and Jun Zhao. 2014. [Relation classification via convolutional deep neural network](#). In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, pages 2335–2344, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- Yuhao Zhang, Peng Qi, and Christopher D. Manning. 2018. [Graph convolution over pruned dependency trees improves relation extraction](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2205–2215, Brussels, Belgium. Association for Computational Linguistics.
- Yuhao Zhang, Victor Zhong, Danqi Chen, Gabor Angeli, and Christopher D. Manning. 2017. [Position-aware attention and supervised data improve slot filling](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing (EMNLP 2017)*, pages 35–45.
- Peng Zhou, Wei Shi, Jun Tian, Zhenyu Qi, Bingchen Li, Hongwei Hao, and Bo Xu. 2016. [Attention-based bidirectional long short-term memory networks for relation classification](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 207–212, Berlin, Germany. Association for Computational Linguistics.

A Datasets

SemEval 2010 Task 8 (Hendrickx et al., 2009) is a widely used testbed for relation extraction, the schema targets at semantic relations between pairs of nominals, which requires certain level of abstractive capabilities. **TACRED** (Zhang et al., 2017) is a large-scale dataset annotated using Amazon Mechanical Turk crowdsourcing. It was initially created for the TAC knowledge base population and mainly covers common relations between people, organizations, and locations based on the TAC KBP scheme. **TACRED-Revisited** (Alt et al., 2020) is a label-corrected version of the TACRED dataset, which motivates from the unresolved challenging cases in original TACRED dataset. **Re-TACRED** (Stoica et al., 2021) further conducted a more comprehensive analysis and re-annotated the entire dataset. Besides, it made alternations to the schema to make it more clear and intuitive, which greatly improved the dataset quality. **ChemProt** (Kringelum et al., 2016) is a bio-domain dataset that extracts 13 kinds of chemical-protein interactions. It is widely used for evaluating domain-specific model capabilities (Lee et al., 2019; Beltagy et al., 2019).

B Experimental Settings

S²ynRE involves three different training processes, respectively the finetuning of LLM, stage-one training, and stage-two training. Except for training steps or epochs, we do not exhaust further search for other hyper-parameters and set them empirically.

For the finetuning of LLM as synthesizer, we set batch size to 64, learning rate to 3e-5. We found that the quality of generated samples is sensitive to the finetuning steps. Considering that the scale of training samples varies from 73 (SemEval 1%) to 68,124 (TACRED 100%) w.r.t. different datasets and different settings, we search steps within different ranges accordingly. The final choices are listed in Table 9.

For stage-one training, we set batch size to 64, learning rate to 3e-5, and fix the training steps as 1500. We save the checkpoint from 500, 1000, and 1500 steps respectively and select the best one. For stage-two training, we set batch size to 16, learning rate to 3e-5, and the epochs are set as Table 10. These epoch settings are empirically chosen in our pilot study to obtain a competitive baseline performance. We set the number of teacher models

K in each iteration to 5 without further searching. We use *bert-base-uncased* to initialize the student model. All experiments are conducted on 40GB A100 machines.

C Performance Over Self-training Iterations

We provide the performance curve w.r.t. iterations in Figure 5. It shows that the iterative training procedure following the classical self-training method is indeed effective. We simply set iteration to 10 as most of the self-training methods did and find it already a robust choice across different datasets.

D Scale of Synthesizer Model

We test S²ynRE with a different scale LLM, i.e., GPT-Small with 117M parameters. The results in Table 11 show that even with such a small size LM, S²ynRE can still bring significant improvements. But in general, larger model unsurprisingly performs better. With the emergence and applicability of increasingly stronger LLMs, we can look forward to further advancement of relation extraction task.

E Ablation on Multi-teacher Distillation

We provide ablation on multi-teacher distillation in Table 12, and demonstrate that the primary improvements come from the utilization of synthetic data and the proposed two-stage self-training. The choices of ablations are:

+Multi-teacher (Naive Ensemble) We train k models on the golden training set and ensemble their predicted logits. This only provides marginal benefits (+3.62, +3.13).

+Synthetic Data We incorporate synthetic data by performing knowledge distillation from multi-teacher to students. This provides rather significant benefits (+6.01, +7.83).

+Two-stage Self-training We further introduce the proposed two-stage self-training method, which brings about the most remarkable improvements (+16.72, +14.55).

F Effects of Domain-Augmented LLM

In the main results (Table 3) we have specifically used a domain-augmented version of GPT2 (GPT2-PubMed) for biomedical task ChemProt. This is our initial choice of design and intuitively should

Setting	SemEval	ChemProt	TACRED	TACRED-Revisited	Re-TACRED
<i>FULL</i>	256	256	512	1024	2048
<i>LIMITED</i>	64	256	256	256	512
<i>FEW</i>	32	32	128	128	128

Table 9: Finetuning steps for LLM under different settings.

Setting	SemEval	ChemProt	Wiki80	TACRED	TACRED-Revisited	Re-TACRED
<i>FULL</i>	{5, 10}	{5, 10}	{5, 10}	2	2	2
<i>LIMITED</i>	{10, 20}	{10, 20}	{10, 20}	5	5	5
<i>FEW</i>	{40, 80}	{40, 80}	{40, 80}	10	10	10

Table 10: Training epochs for stage-two training under different settings.

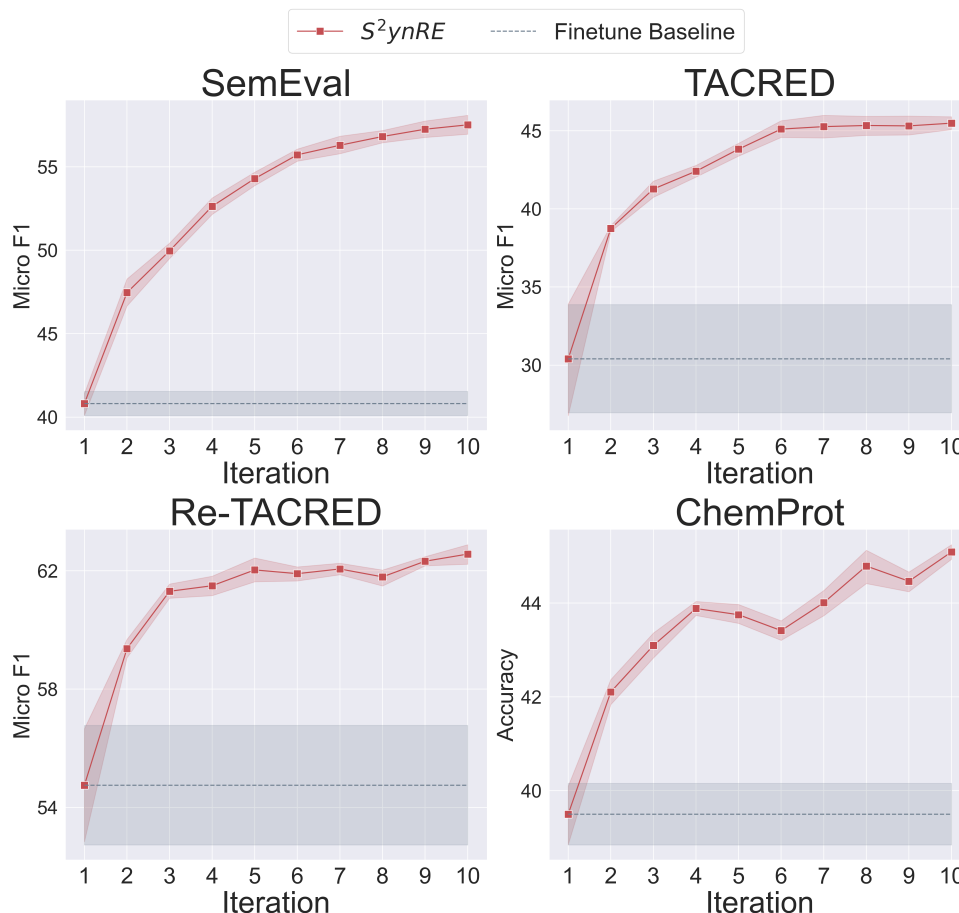


Figure 5: Performance over self-training iterations. Drawn with standard error of mean.

bring better performance. Here we further analysis the effects of such LLM choices in detail. Table 13 provides comparison between GPT2-PubMed and vanilla GPT2, both LLMs can effectively produce synthetic data and bring expected improvements. Nonetheless, we empirically find that vanilla GPT2 would need a bit more finetuning steps to adapt to the target domain (256 steps compared to 32 steps using GPT2-PubMed). In general, the proposed method is rather robust to choices of LLM.

G Case Study

We provide randomly sampled case studies of synthetic data for SemEval, TACRED, and ChemProt in Table 14, 15, and 16 respectively as well as distant data in Table 17. These cases show that LLMs are capable of synthesizing coherent, realistic sentences with relational structure. Most importantly, such synthetic data are customized to target domains with various topics and styles.

Dataset	NA	GPT-2 Small 117M	GPT-2 Large 774M
SemEval	40.81	49.87	57.99
TACRED	30.40	43.95	45.87
TACRED-Revisited	33.75	48.35	50.61
Re-TACRED	54.75	63.51	62.98

Table 11: Performances w.r.t. synthesizer model size.

Method	Perf.	Imp.
BERT Baseline	40.81 \pm 1.62	-
+ Multi-teacher (Naive Ensemble)	44.43	+ 3.62
+ Synthetic Data (Distillation)	46.82 \pm 1.81	+ 6.01
+ Two-stage Self-training (S ² ynRE)	57.53\pm0.96	+16.72
CP Baseline	52.77 \pm 1.83	-
+ Multi-teacher (Naive Ensemble)	55.90	+ 3.13
+ Synthetic Data (Distillation)	60.60 \pm 1.28	+ 7.83
+ Two-stage Self-training (S ² ynRE)	67.32\pm0.54	+14.55

Table 12: Ablation on multi-teacher distillation.

Nevertheless, we also notice several limitations, especially in low-resource scenarios where it’s still challenging to get a good estimation of the target dataset distribution:

- Lack of diversity. For example, instances 2.1, 2.2, 2.3 all start with "*the marmalade*".
- Fragmentary structure. For example, instances 2.4 and 2.8 contain atypically lengthy object.

For pseudo labels, most of the time teacher model confidently assigns one specific label with very high probabilities (> 0.95), but for some other cases, it goes for more than one possible label, such as 1.8, 2.8, 4.1, etc. We attribute this to two possible reasons: 1) the limited capability of the teacher model to accurately recognize all relations, and 2) the imperfections of certain synthetic data, i.e., some synthetic instances do not well align with pre-defined schema and are difficult to be assigned exact relation labels. In these cases, forcing the student to learn from hard labels assigned using argmax might introduce severe noise, while the proposed knowledge distillation process using soft labels in S²ynRE can properly put these imperfect data still into usage.

For distant data, as these instances are produced from wikipedia texts, we can clearly identify that they are quite different from other downstream task data either in content, or in relation schema. This further verifies the superiority of the proposed synthetic in-domain data qualitatively.

BERT Baseline	39.50
S²ynRE w/ GPT2	45.32
S²ynRE w/ GPT2-PubMed	45.09

Table 13: Impacts of domain augmented LLM.

H Potential Limitations

We empirically conclude two limitations for S²ynRE in the hope of inspiring more future research. On one hand, its advantages are less significant when a large amount of annotated data is available. For example, TACRED training set has 68,142 annotated instances. Under this setting, even if we add another 100,000 synthetic samples, the improvement is only +0.98 compared to +22.02 under 1% training set. This means that the quality of synthetic data, although superior to distant ones, is still not as good as golden ones. Thus they can hardly provide identical utility the same as 100,000 golden data. Nevertheless, with the development of LLMs and their powerful generation ability, we look forward to accessing higher-quality synthetic data.

On the other hand, when training data are limited to a few samples (for example, 1% setting for SemEval only includes 73 training instances), even strong LLMs like GPT-2 can not perfectly fit the structure of relational statements within a few steps of finetuning (See Appendix G for illustration of cases). Therefore, many generated sentences may not contain correct subject or object entity markers as requested and have to be discarded. In general, although the formation of marked natural language sequence proposed in this work made such structured synthesis feasible, we look forward to further improving the synthesis efficacy in future works.

⁴https://www.wikidata.org/wiki/Property_talk:P609

Instances		Soft Labels (Top3)	Probs
<i>SemEval FULL</i>			
1.1	the [Sub] mansion [Sub] has been the subject of several [Obj] reports [Obj] on television.	Message-Topic(e2,e1) Component-Whole(e2,e1) Entity-Origin(e2,e1)	0.99956 0.00006 0.00005
1.2	the [Sub] man [Sub] was in the [Obj] building [Obj] at the time.	Other Content-Container(e1,e2) Entity-Origin(e1,e2)	0.99971 0.00005 0.00005
1.3	i had a [Sub] gift [Sub] from the [Obj] hospital [Obj] which was going to cost a lot more than my first.	Entity-Origin(e1,e2) Product-Producer(e1,e2) Product-Producer(e2,e1)	0.99769 0.00110 0.00080
1.4	the video shows the [Sub] person [Sub] getting up from the [Obj] bed [Obj] to put on the robe.	Entity-Origin(e1,e2) Other Cause-Effect(e2,e1)	0.99959 0.00030 0.00002
1.5	the [Sub] cadaver [Sub] was left and kept at the [Obj] museum [Obj] .	Other Entity-Origin(e1,e2) Entity-Destination(e1,e2)	0.94540 0.03551 0.01063
1.6	the [Sub] tumor [Sub] was contained with two instilled [Obj] antibiotics [Obj] ..	Other Cause-Effect(e2,e1) Instrument-Agency(e2,e1)	0.58024 0.40806 0.00442
1.7	it was a [Sub] truck [Sub] that moved the [Obj] furniture [Obj] .	Other Instrument-Agency(e1,e2) Component-Whole(e1,e2)	0.58490 0.37308 0.01200
1.8	he began to set up and operate many of the [Sub] computers [Sub] in the [Obj] store [Obj] .	Component-Whole(e1,e2) Other Content-Container(e1,e2)	0.47224 0.27054 0.24453
<i>SemEval FEW</i>			
2.1	the [Sub] marmalade [Sub] starts with a [Obj] marzipan [Obj] in the centre of a vanilla bean.	Entity-Origin(e2,e1) Entity-Origin(e1,e2) Component-Whole(e2,e1)	0.97080 0.00486 0.00484
2.2	the [Sub] marmalade [Sub] is a blend of [Obj] cherries [Obj] , dulce de leche and cognac that is richly decorated with an intricate series of images of olive branches.	Entity-Origin(e2,e1) Entity-Origin(e1,e2) Component-Whole(e2,e1)	0.98489 0.00257 0.00140
2.3	the [Sub] marmalade [Sub] is a [Obj] blend [Obj] of anise, caster, and grape juice.	Entity-Origin(e2,e1) Entity-Origin(e1,e2) Content-Container(e2,e1)	0.98827 0.00116 0.00086
2.4	the [Sub] cricketers [Sub] have [Obj] struggled to find sponsorship for their \$1.2 million annual home-cooked dinner [Obj] entirely on donated food.	Instrument-Agency(e2,e1) Product-Producer(e2,e1) Other	0.80719 0.07164 0.04683
2.5	there a [Sub] caused by a [Sub] poisoning [Obj] [Obj] .	Cause-Effect(e2,e1) Cause-Effect(e1,e2) Product-Producer(e1,e2)	0.99813 0.00023 0.00020
2.6	the [Sub] troubadour [Sub] starts with a [Obj] snowstorm [Obj] that blankets the streets and then slowly disperses as the temperature drops.	Component-Whole(e2,e1) Entity-Origin(e1,e2) Instrument-Agency(e2,e1)	0.99156 0.00201 0.00085
2.7	the [Sub] water [Sub] is also rich in organic matter [Obj] , mainly cold-water crayfish [Obj] and planktonic foraminifera.	Entity-Origin(e1,e2) Cause-Effect(e1,e2) Instrument-Agency(e1,e2)	0.89010 0.03238 0.01435
2.8	the [Sub] series [Sub] takes its inspiration from a real-life story [Obj] of a young woman who attempted suicide using a water gunslinger [Obj] .	Product-Producer(e2,e1) Product-Producer(e1,e2) Instrument-Agency(e2,e1)	0.53262 0.11620 0.08945

Table 14: Randomly selected cases of synthetic data and the assigned soft labels for SemEval. **Without any cherry picking.**

	Instances	Soft Labels (Top3)	Probs
<i>TACRED FULL</i>			
3.1	The National Union of Students welcomed the move by the [Sub] NUS [Sub] in its opening resolution in a meeting [Obj] today [Obj] .	no_relation org:dissolved org:founded	0.99872 0.00047 0.00038
3.2	It mayn't look it in person, but it's the same thing as playing a show on the radio – it's actually the same thing with radio in general, where if you just hit – [Obj] one [Obj] play through, you 'll get familiar with the basic patterns that make 'em work, and all the subtle nuances – so instead of going out and trying to get “ that ” “ That one ” , ” [Sub] Steve Allen [Sub] told Entertainment Weekly in 1991, “ try doing this, or do that. ”	no_relation per:age org:number_of_employees/members	0.99969 0.00007 0.00004
3.3	It was the third consecutive year the International Skating Union rated [Sub] Skiing Australia [Sub] a gold medal threat, after a bronze medal performance at the 2004 Winter Olympic Games in Nagano and a silver medal performance in the Salt Lake City Games in [Obj] 2006 [Obj] .	no_relation org:founded org:dissolved	0.99901 0.00026 0.00018
3.4	He is survived by [Sub] his [Sub] wife of 63 years, the [Obj] Doris G. Gude [Obj] of Rockville ; a son, Charles Gude Jr. ; five grandchildren ; and three great-grandchildren.	per:spouse no_relation per:other_family	0.91159 0.06497 0.01286
3.5	“ I think these guys have done some amazing work on the set, ” added [Obj] Bryan Fuller [Obj] , whose television show, “ Heroes, ” created another big ensemble cast by including Emmy-nominated actors [Sub] Spencer Pratt [Sub] and Evan Rachel Wood.	no_relation per:other_family per:siblings	0.98786 0.00426 0.00164
3.6	The [Sub] American Family Association [Sub] announced that it is boycotting [Obj] Cathay Pacific [Obj] and is taking a similar stand over the next nine days.	no_relation org:subsidiaries org:member_of	0.95461 0.01223 0.00858
<i>TACRED FEW</i>			
4.1	In addition to his wife, he is survived by four children, William J. Gillette Jr. of Rockville, [Obj] Illinois [Obj] , James P. Gillette of Gilbertsville, Pennsylvania, [Sub] Diana R. [Sub] of Gilbertsville and Michael D. Gillette of Rockville ; 12 grandchildren ; and 12 great-grandchildren.	per:stateorprovinces_of_residence per:siblings org:stateorprovince	0.22273 0.15570 0.12936
4.2	[Sub] Ventura [Sub] 's win brings to eight the number of wins by [Obj] California [Obj] athletes in the 200 meters since 1985.	per:stateorprovinces_of_residence org:stateorprovince_of_headquarters no_relation	0.71593 0.05609 0.03997
4.3	The first episode of [Obj] M*A*S*H [Obj] was broadcast on Saturday, November 2, 1996, on the [Sub] NBC [Sub] network.	no_relation org:alternate_names org:parents	0.99886 0.00007 0.00006
4.4	The [Sub] ICBA [Sub] president, [Obj] Huang Zuocheng [Obj], said in a statement : “ This is a big step forward and will certainly help the whole community of farmers in providing a decent quality food for all. ”	org:top_members/employees org:founded_by org:subsidiaries	0.99060 0.00193 0.00093
4.5	[Sub] Johannesburg [Sub] police chief Inspector-General of Police Lieutenant-general Nathi Nhleko has ordered the arrest of four individuals charged over the grenade attack on a wedding party in [Obj] Johannesburg [Obj] one week ago that left two people - a 27-year-old man and a 41-year-old woman - dead.	per:cities_of_residence org:city_of_headquarters per:city_of_death	0.52491 0.07287 0.05097
4.6	Under the deal, the [Sub] Kuala Lumpur Chamber of Deputies [Sub] has agreed to let foreign [Obj] investors [Obj] buy up to 50 percent of the company, and the government has agreed to give it an additional 10 percent stake once the government approves the deals.	no_relation org:parents org:country	0.99852 0.00034 0.00014

Table 15: Randomly selected cases of synthetic data and the assigned soft labels for TACRED. **Without any cherry picking.**

	Instances	Soft Labels (Top3)	Probs
<i>ChemProt FULL</i>			
5.1	[Sub] Lumiracoxib [Sub] is metabolized to a more potent and selective [Obj] cyclooxygenase-2 [Obj] (COX-2) inhibitor by sequential metabolism.	INHIBITOR SUBSTRATE PRODUCT-OF	0.94689 0.05080 0.00059
5.2	The effect of phenobarbital, a known [Sub] CYP2D6 [Sub] inhibitor, on the pharmacokinetics of [Obj] DEX [Obj] , a substrate of human CYP2D6, in healthy subjects.	INHIBITOR SUBSTRATE ACTIVATOR	0.99792 0.00149 0.00013
5.3	The inhibitory effect of [Sub] pravastatin [Sub] on [Obj] human UGS1 [Obj] mediated by the high affinity UGS2 isoforms EGFR and ErbB2 was also investigated.	INHIBITOR INDIRECT-DOWNREGULATOR DOWNREGULATOR	0.99890 0.00058 0.00017
5.4	Moreover, the [Sub] quinone [Sub] derivative was found to exhibit pronounced [Obj] beta(2)-adrenoceptor [Obj] (beta(2)-AR)/erythrocyte coupling inhibitory effects, in the following order: quinone>diethylglycerol>cis-9,trans-11,12-didehydro-9,trans-11,12- triazol-9-amine (DFTDI)>cis-9,trans-11,12-didehydro-9, cis-9, trans-12, 13-tetrahydro	INHIBITOR ANTAGONIST AGONIST-INHIBITOR	0.99968 0.00010 0.00005
5.5	These data demonstrate that [Sub] troglitazone [Sub] , an inhibitor of [Obj] PTGS2 [Obj] , acts on cells by inhibition of the phosphatidylinositol 3-kinase/Akt/mTOR pathway, which could account for the reduced incidence of osteopetrosis and osteoarthritis that occur in patients receiving this drug.	INHIBITOR INDIRECT-DOWNREGULATOR INDIRECT-UPREGULATOR	0.99984 0.00006 0.00002
5.6	Inhibition of [Sub] PDE11A [Sub] by [Obj] dihydropyridine [Obj] and butyrylcholinesterase inhibitors (BuChE, butyl methylester, and butoxychlor) strongly suggested involvement of cholinergic inhibition at membrane level.	INHIBITOR ACTIVATOR INDIRECT-UPREGULATOR	0.99982 0.00003 0.00003
<i>ChemProt FEW</i>			
6.1	Results show that [Sub] Gossypol [Sub] and buthionine sulfoxane have the most potent inhibitory activities against [Obj] PEPCK [Obj] with IC50 values of 1.46, 1.24 and 0.98 microM, respectively.	INHIBITOR INDIRECT-DOWNREGULATOR AGONIST-INHIBITOR	0.99822 0.00031 0.00023
6.2	Based on the results of this study, it can be concluded that [Sub] sorafenib [Sub] exerted its inhibitory effect on the CSE-induced [Obj] angiogenesis-related phospho-AKT* [Obj] activation through the down-regulation of CSE-induced AKT* phosphorylation.	INHIBITOR ACTIVATOR INDIRECT-DOWNREGULATOR	0.84259 0.11871 0.00990
6.3	Results for [Sub] epinastine [Sub] in the treatment of experimental myasthenia gravis showed IC50 values of 10-11.5 microM against the myasthenia-related enzyme, [Obj] myosin heavy chain [Obj] .	INHIBITOR INDIRECT-DOWNREGULATOR AGONIST	0.98410 0.00765 0.00154
6.4	In a previous study, we have demonstrated that [Sub] sorafenib [Sub] attenuated the growth of C6 glioma cells through [Obj] SRC [Obj] activation.	INHIBITOR INDIRECT-DOWNREGULATOR AGONIST	0.90038 0.04238 0.01377
6.5	Results showed that [Sub] Epinastine [Sub] significantly attenuated the [Obj] l-arginine aminotransferase [Obj] and NADPH oxidase activities in the aorta of MPTP models.	SUBSTRATE INDIRECT-DOWNREGULATOR ACTIVATOR	0.86088 0.04992 0.01992
6.6	Inhibition effect of [Sub] epinastine [Sub] on [Obj] EGFR [Obj] tyrosine kinase activation and its downstream pAKT, ERK, and c-Fos were further investigated.	INHIBITOR INDIRECT-DOWNREGULATOR AGONIST	0.99790 0.00058 0.00029

Table 16: Randomly selected cases of synthetic data and the assigned soft labels for ChemProt. **Without any cherry picking.**

	Instances	Distant Labels
<i>Distant Supervision</i>		
7.1	The tunnel is also part of the UK 's [Sub] National Cycle Route 1 [\Sub] linking Inverness and [Obj] Dover [\Obj] .	P609
7.2	Alfred Faure Alfred - Faure or [Sub] Port Alfred [\Sub] is a permanent French scientific station on [Obj] Île de la Possession [\Obj] (Possession Island) of the subantarctic Crozet Archipelago in the South Indian Ocean .	P709
7.3	He was a respected poet in the [Obj] Latin language [\Obj] , writing under the name of [Sub] Santolius Victorinus [\Sub] .	P1412
7.4	In 1704 , [Sub] Eberhard Ludwig [\Sub] started to build [Obj] Ludwigsburg Palace [\Obj] to the north of Stuttgart , in imitation of Versailles .	P119
7.5	Reports linking full - back [Sub] Fred Speller [\Sub] with Warwick County left the " [Obj] Birmingham [\Obj] Daily Post " " wondering at footballers ingratitude .	P20
7.6	Giovanni di Buiamonte Giovanni di Buiamonte was a [Obj] Florentine [\Obj] nobleman who lived in the late 13th century around the time of Giotto and [Sub] Dante [\Sub] .	P551
7.7	Instead , he left [Sub] Sydney [\Sub] [Obj] Sydney [\Obj] at 1 am on 7 January 1931 , and headed for Blenheim , New Zealand .	P931
7.8	" Nintendo Power " journalist Steve Thomason singled out Sanshiro as the character he would most like to control in the [Obj] Nintendo DSi [\Obj] game " [Sub] Photo Dojo [\Sub] " .	P400
7.9	He provided the vocals for the singing voice of the cub [Obj] Simba [\Obj] in Walt Disney Feature Animation 's 1994 film " [Sub] The Lion King [\Sub] " .	P674
7.10	[Sub] 3rd County of London Yeomanry (Sharpshooters) [\Sub] The 3rd County of London Yeomanry (Sharpshooters) was a Yeomanry regiment of the [Obj] British Army [\Obj] .	P241

Table 17: Randomly selected cases of distant supervision data. The explanation for distant labels can be looked up at the official wikidata website⁴.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Left blank.
- A2. Did you discuss any potential risks of your work?
Left blank.
- A3. Do the abstract and introduction summarize the paper’s main claims?
Left blank.
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Left blank.

- B1. Did you cite the creators of artifacts you used?
Left blank.
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
Left blank.
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Left blank.
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
Left blank.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Left blank.
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Left blank.

C Did you run computational experiments?

Left blank.

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
Left blank.

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Left blank.

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Left blank.

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Left blank.

D Did you use human annotators (e.g., crowdworkers) or research with human participants?

Left blank.

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

Left blank.

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

Left blank.

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

Left blank.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

Left blank.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

Left blank.