

Mitigating Prototype Shift: Few-Shot Nested Named Entity Recognition with Prototype-Attention

Anonymous ACL submission

Abstract

Nested entities are prone to obtain similar representations in pre-trained language models, posing challenges for Named Entity Recognition (NER), especially in the few-shot setting where prototype shifts often occur due to distribution differences between the support and query sets. In this paper, we regard entity representation as the combination of prototype and non-prototype representations. With a hypothesis that using the prototype representation specifically can help mitigate potential prototype shifts, we propose a Prototype-Attention mechanism in the Contrastive Learning framework (PACL) for the few-shot nested NER. PACL first generates prototype-enhanced span representations to mitigate the prototype shift by applying a prototype attention mechanism. It then adopts a novel prototype-span contrastive loss to reduce prototype differences further and overcome the O-type's non-unique prototype limitation by comparing prototype-enhanced span representations with prototypes and original semantic representations. Our experiments on three English, German, and Russian nested NER datasets show that the PACL outperformed seven baseline models on the 1-shot and 5-shot tasks in terms of F_1 score. Further analyses indicate that our Prototype-Attention mechanism has high generality, enhancing the performance of two baseline models, and can serve as a valuable tool for NLP practitioners facing few-shot nested NER tasks.

1 Introduction

The few-shot Named Entity Recognition (NER) task has gained a lot of attention in recent years as it aims to address the limitations of traditional NER methods that rely on a large number of labeled training instances, which can be both time-consuming and experience-dependent. This task deals with the NER problem using only a few labeled instances. Researchers have made significant progress on this task by applying deep learning models, including pre-trained-model-based (Florez and Mueller,

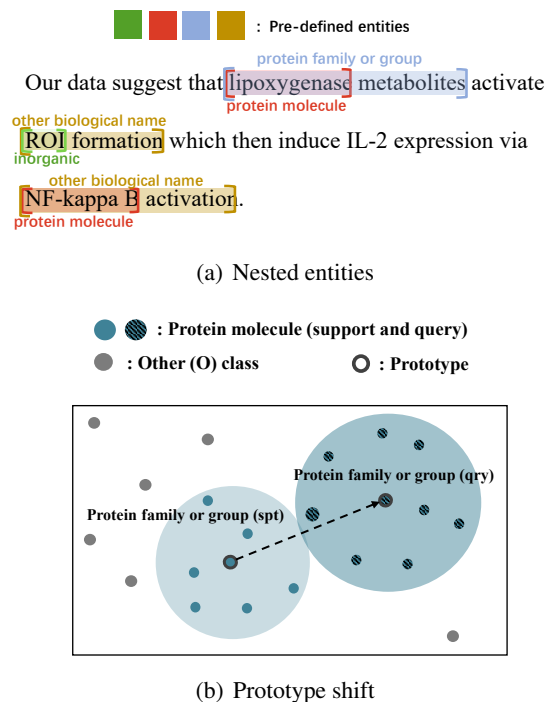


Figure 1: (a) Example of an instance with nested entities from the GENIA dataset. (b) Illustration of prototype shifts, where the prototypes differ due to the distribution difference between the support and query sets.

2019; Hou et al., 2019; Yang et al., 2021; Wang et al., 2022b), metric-learning-based (Snell et al., 2017; Hofer et al., 2018; Yang and Katiyar, 2020), meta-learning-based (Li et al., 2020a; Sung et al., 2018), prompt-tuning-based (Ma et al., 2022; Hou et al., 2022), and contrastive-learning-based (Das et al., 2022) methods.

However, most existing few-shot NER research has focused on flat entities that do not overlap (Ming et al., 2022; Wang et al., 2022b). In reality, many entities share the same words and form nested entities that are part of another entity. This is where the few-shot nested NER task comes in. This task deals with nested entities that share words and are part of another entity. For example, in the

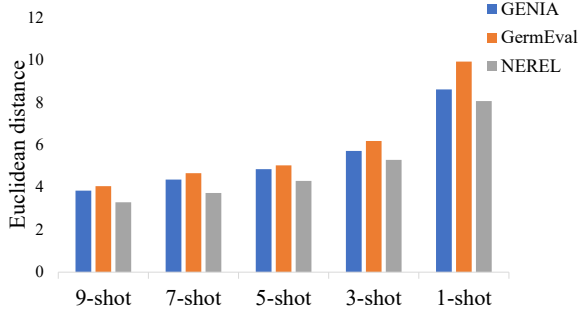


Figure 2: The Euclidean distance of prototype shift between the prototypes in the support set and the query set in the GENIA, GermEval, and NEREL datasets. K-shot denotes the K number of labeled instances in the support set for each type.

GENIA dataset (Kim et al., 2003), about 53.9% of entities are nested. Figure 1 (a) illustrates an instance, that is, a protein molecule entity "lipoxxygenase" is nested within a protein family or group entity "lipoxxygenase metabolites". Due to the overlapped part, nested entities are more likely to obtain similar representations, increasing the difficulty of distinguishing them, especially in the few-shot setting where prototype shifts often occur.

The prototype shift in NER refers to changes in the prototypes between the few-shot labeled data set (support set) and unlabeled data (query set), as exemplified in Figure 1 (b), where a prototype is a representative instance of a specific entity type. The very few labeled data in the support set could hardly represent the whole distribution, resulting in prototype shifts. Figure 2 shows the statistics of prototype shifts in terms of Euclidean distance between the support set and the query set in three nested datasets (GENIA (Kim et al., 2003), GermEval (Benikova et al., 2014), NEREL (Loukachevitch et al., 2021)). We can find that the prototype shift reveals a consistent pattern of increasing Euclidean distance between prototypes as the number of labeled data in the support set decreases. When employing the prototypes derived from the support set for delineating the decision boundaries in the query set, a high frequency of classification errors would be introduced due to prototype shifts. Despite having distinguished nested entities within the support set, they may become interspersed within the query set.

This paper addresses the prototype shift in the few-shot nested NER task. We regard entity representation as the combination of prototype and non-prototype representations. Entities of the same type

should share the same prototype representation. And the non-prototype representation determines the dispersion of the entity distribution. If we could focus more on the prototype representation when learning the entity representation, entities would gather closer around the prototype, and the prototype shift could be reduced. Therefore, we design a prototype-attention mechanism to enhance the prototype representation. Besides, words of the O-type have miscellaneous semantics and cannot be represented by a unique prototype. Therefore, we further design a novel prototype-span contrastive loss. It compares prototype-enhanced span representations with original semantic representations to guarantee the O-type's representations are not enhanced by entity prototypes. It also compares prototype-enhanced span representations with prototypes to reduce prototype differences further.

Our main contributions are as follows:

- We identify the prototype shift challenge in the few-shot learning, particularly in the few-shot nested NER task, and propose a Prototype-Attention Contrastive Learning (PACL) framework to tackle it.
- We devise a unique Prototype-Attention mechanism to generate the prototype-enhanced representation for each span to mitigate the prototype shift between the support and query sets. This mechanism exhibits a high level of generality in enhancing the performance of two baseline models.
- We design a novel prototype-span contrastive loss by comparing prototype-enhanced span representations with prototypes and original semantic representations to reduce prototype differences further and overcome the O-type's non-unique prototype limitation.
- We conduct experiments on three nested NER datasets from three different languages. The results show improvements in PACL over existing nested NER and few-shot NER baselines in terms of F_1 score.

2 Problem Definition

Following the mainstream solutions, we formulate the few-shot nested NER task as a span-based entity classification problem. That is, given an input sentence $x \in \mathcal{X}$ with l tokens, denoted by

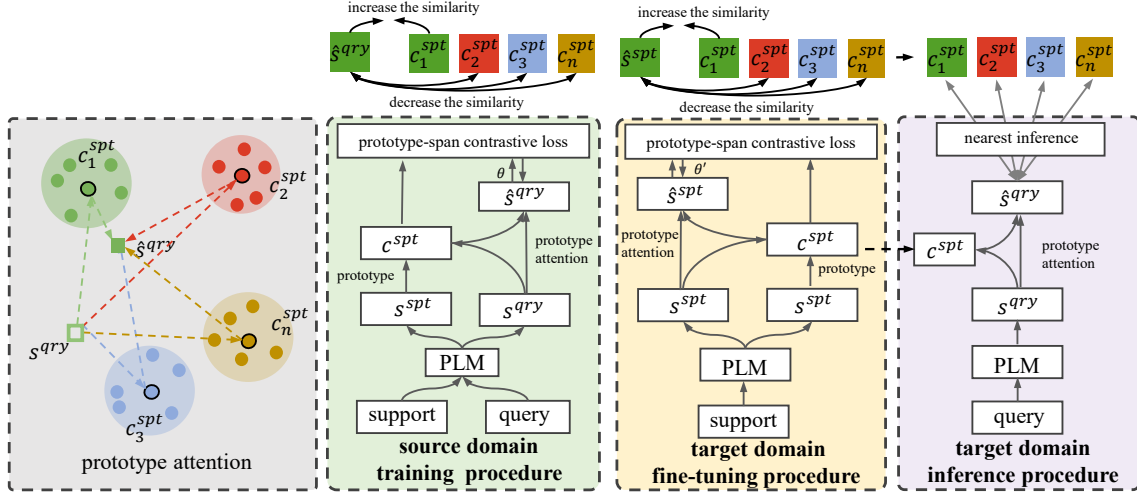


Figure 3: Illustration of our PACL framework and learning procedures. During the training procedure on the source domain, PACL calculates prototypes based on labeled spans of the support set and then utilizes prototype-attention to obtain prototype-enhanced representations for the query set. After that, PACL applies the prototype-span contrastive loss to optimize the representations. During the fine-tuning procedure on the target domain, PACL generates prototype-enhanced representations for the support set to fine-tune the model. Finally, PACL makes inferences on the query set of the target domain based on the nearest neighbor strategy.

$x = \{w_1, \dots, w_l\}$, we generate an entity span set containing all possible spans, and each span s_{pq} is a span of tokens starting from the p^{th} token and ending at the q^{th} token in x , denoted by $s_{pq} = \{w_p, \dots, w_q\}$ ($1 \leq p \leq q \leq l$). Then, we learn a classification model to map each span into an entity label in the label set $E_{\mathcal{X}}$. If we set the task as a K -shot task, then the number of span labels for each entity type used for training is limited to K . Besides, we also apply the meta-learning framework. The formal descriptions are as follows.

Let $\mathcal{D} = \{\mathcal{X}, \mathcal{Y}\}$ denote a dataset with \mathcal{X} and \mathcal{Y} as the sentence set and the corresponding label set, respectively. $\mathcal{D}^{spt} = \{\mathcal{X}^{spt}, \mathcal{Y}^{spt}\}$ and $\mathcal{D}^{qry} = \{\mathcal{X}^{qry}, \mathcal{Y}^{qry}\}$ are disjoint sets sampled from \mathcal{D} for model training and testing, respectively. They are also known as the support set and the query set. Suppose $\mathcal{D}_i = \{\mathcal{X}_i, \mathcal{Y}_i\}$ and $\mathcal{D}_j = \{\mathcal{X}_j, \mathcal{Y}_j\}$ are the source and target domain datasets, respectively. The few-shot nested NER task first samples several subtasks $\{\mathcal{D}_i^{spt}, \mathcal{D}_i^{qry}\}$ from $\mathcal{D}_i = \{\mathcal{X}_i, \mathcal{Y}_i\}$, where $\mathcal{D}_i^{spt} = \{\mathcal{X}_i^{spt}, \mathcal{Y}_i^{spt}\}$, $\mathcal{D}_i^{qry} = \{\mathcal{X}_i^{qry}, \mathcal{Y}_i^{qry}\}$. It then trains a model on these subtasks. After that, it makes adaptations on \mathcal{D}_j , i.e., it fine-tunes the model on $\mathcal{D}_j^{spt} = \{\mathcal{X}_j^{spt}, \mathcal{Y}_j^{spt}\}$ and then predicts the span labels for $\mathcal{D}_j^{qry} = \{\mathcal{X}_j^{qry}\}$. For the K -shot setting, each entity category in \mathcal{X}_i^{spt} and \mathcal{X}_j^{spt} contains K entities.

3 Methodology

This section introduces our PACL framework and then provides details of the prototype-attention mechanism, the prototype-span contrastive loss, and target domain adaption procedures.

3.1 PACL Framework

Figure 3 illustrates our Prototype-Attention Contrastive Learning (PACL) framework and learning procedures.

PACL first applies a Pre-trained Language Model (PLM) to obtain the semantic representation for each span. It then calculates prototypes on the support set and utilizes a novel prototype-attention mechanism to achieve prototype-enhanced representations. After that, PACL optimizes representations by a prototype-span contrastive loss.

During the training procedure on the source domain, PACL utilizes a bunch of subtasks $\{\mathcal{D}_i^{spt}, \mathcal{D}_i^{qry}\}$ to train the model. It generates prototype-enhanced representations for the query set to obtain the adjustment ability for prototype shift. During the fine-tuning procedure on the target domain, PACL utilizes $\mathcal{D}_j^{spt} = \{\mathcal{X}_j^{spt}, \mathcal{Y}_j^{spt}\}$ to fine-tune the model by generating prototype-enhanced representations for the support set. Finally, it predicts the labels for $\mathcal{D}_j^{qry} = \{\mathcal{X}_j^{qry}\}$ by the nearest neighbor strategy.

3.2 Prototype-Attention Mechanism

To mitigate the prototype shift, we propose a prototype-attention mechanism to generate prototype-enhanced representations for the query set based on prototypes obtained from the support set during training on the source domain. This approach improves the span representations in the query set by incorporating more prototype information, which aligns the prototypes of the query set with those of the support set. The detailed procedures are presented below.

We first incorporate a Pre-trained Language Model (PLM) to obtain original span semantic representations. That is, for the sentence x with l tokens, we get all word embeddings, concatenate the start and the end token embeddings of each span, and use a non-linear function to get the span semantic representation s :

$$[h_1, h_2, \dots, h_l] = \text{PLM}([w_1, w_2, \dots, w_l]) \quad (1)$$

$$s = f(h_p \oplus h_q) \quad (2)$$

Where \oplus denotes the concatenation operator, and f is a non-linear function.

For the support set, we calculate the prototype c_k for each entity type k according to span labels except the O-type:

$$c_k = \frac{1}{|s_k|} \sum s_k \quad (3)$$

Where $|s_k|$ denotes the number of spans in type k .

For spans in the query set, we gain the prototype-enhanced representation \hat{s}^{qry} by calculating the attention score between the original span representation s^{qry} and prototypes $\mathcal{C} = [c_1, c_2, \dots]$ in the support set:

$$\hat{s}^{qry} = \text{softmax} \left(\frac{s^{qry} \mathcal{C}^\top}{\sqrt{d_c}} \right) \mathcal{C} + s^{qry} \quad (4)$$

where d_c is the dimension of prototypes. We also include s^{qry} in the attention representation to obtain \hat{s}^{qry} , excluding the O-type spans which cannot be represented by prototypes in \mathcal{C} . This will be further optimized in the next section with prototype-span contrastive loss.

3.3 Prototype-Span Contrastive Loss

The traditional contrastive loss increases span similarities of the same entity type and decreases span similarities between different entity types. This paper aims to address the prototype shift. Therefore,

we want to increase the similarity between spans in the query set and the corresponding prototype in the support set to let the model obtain the ability to mitigate the prototype shift. Besides, the O-type span has miscellaneous semantics and could not be represented by a unique prototype (Fritzler et al., 2019). We also want prototype-enhanced representations of O-type entities close to their original semantic representations. Therefore, we design the following prototype-span contrastive loss based on the circle loss (Sun et al., 2020).

For each span representation \hat{s}^{qry} in the query set, the loss $\mathcal{L}_{\hat{s}^{qry}}$ is calculated by:

$$\mathcal{L}_{\hat{s}^{qry}} = \log(1 + \text{sim}(\hat{s}^{qry}, c^+) * \text{sim}(\hat{s}^{qry}, c^-)) \quad (5)$$

Where c^+ is the corresponding prototype in the support set with the same type as \hat{s}^{qry} , and c^- denotes prototypes in the support set with different types from \hat{s}^{qry} . The similarity function sim is calculated by:

$$\text{sim}(\hat{s}^{qry}, c^+) = e^{-\tau * \phi(\hat{s}^{qry}, c^+)} \quad (6)$$

$$\text{sim}(\hat{s}^{qry}, c^-) = \sum_{c_i^- \in c^-} e^{\tau * \phi(\hat{s}^{qry}, c_i^-)} \quad (7)$$

Where $\phi(\cdot)$ denotes the cosine similarity, τ is the temperature (Wang and Liu, 2021).

When calculating $\text{sim}(\hat{s}^{qry}, c^+)$ for the O-type, we calculate the cosine similarity between the original span representation s^{qry} and the prototype-enhanced representation \hat{s}^{qry} :

$$\phi(\hat{s}_i^{qry}, c_o) = \lambda * \phi(\hat{s}_i^{qry}, s^{qry}) \quad (8)$$

Where λ is a learnable hyperparameter. We calculate the cosine similarity between the prototype-enhanced representation \hat{s}^{qry} and its corresponding prototype in the support set for other entity types.

3.4 Target Domain Adaption

After training the model on the source domain, we make adaption to the target domain, including fine-tuning the model on the support set and making inferences on the query set.

During the fine-tuning procedure, our PACL first generates prototype-enhanced representations \hat{s}^{spt} for spans in the support set by calculating the attention score between the original span representation s^{spt} and the prototypes \mathcal{C} in the support set. After that, PACL fine-tunes the model by utilizing the prototype-span contrastive loss with the input of

\hat{s}^{spt} and \mathcal{C} . Different from using \hat{s}^{qry} as the input in the training procedure, we utilize \hat{s}^{spt} in the fine-tuning procedure since the labels of the query set are unknown. Using \hat{s}^{spt} is to make entities gather around prototypes, and this does not affect our PACL to mitigate the prototype shift. This is because the prototype shift is domain-independent. It is only related to the sampling strategy. Our PACL has already learned prototype shift patterns and acquired the ability to mitigate the prototype shift during training on the source domain.

During the inference procedure, our PACL obtains prototype-enhanced representations \hat{s}^{qry} for spans in the query set according to prototypes \mathcal{C} in the support set. It further applies the nearest neighbor inference for each span according to the maximum similarity with prototypes or its original span representation (O-type).

Note that the support set in the target domain may be too large to load all spans in a single fine-tune procedure. Loading all sentences into a single batch to get prototypes for each type is challenging. Thus, for each sentence in the support set, we leave all entity spans and sample $1/|\mathcal{X}_j^{spt}|$ percentage of O-type spans during fine-tuning, where $|\mathcal{X}_j^{spt}|$ denotes the number of sentences in the support set.

4 Experiments

In this section, we evaluate PACL in few-shot nested NER. After introducing datasets and baseline models, we outline the setup, present results, and analyze them thoroughly.

4.1 Datasets

To evaluate our proposed PACL framework’s performance and generality across languages, we experiment with the Indo-European language family. English is chosen as the source language, and three target languages are selected based on their linguistic proximity to English: English itself, German, and Russian, as obtaining datasets for these languages is feasible.

Dataset	language	Types	Sentences	Entities/Nest entities
GENIA	English	36	18.5k	55.7k / 30.0k
GermEval	German	12	18.4k	41.1k / 6.1k
NEREL	Russian	29	8.9k	56.1k / 18.7k
FewNERD	English	66	188.2k	491.7k / -

Table 1: Datasets used in experiments

As shown in Table 1, the target nested NER

datasets are GENIA¹ in English (Kim et al., 2003), GermEval² in German (Benikova et al., 2014), and NEREL³ in Russian (Loukachevitch et al., 2021). We use a flat NER dataset, FewNERD⁴ in English (Ding et al., 2021), as the source domain dataset to train the model. All these datasets are publicly available under the licenses of CC-BY 3.0 for GENIA, CC-BY 4.0 for GermEval, CC-BY 2.5 for NEREL, and CC-BY-SA 4.0 for FewNERD. We have manually checked to guarantee these datasets are without offensive content and identifiers.

For training in the source domain, We randomly sampled 15,00 5-way 5-shot subtasks from the FewNERD inter-domain subset, among which 500 subtasks are used for validation. We validated the model every 1000 subtasks. When fine-tuning in the target domain, we sampled 32-way, 12-way, and 29-way support sets under 1-shot and 5-shot settings from GENIA, GermEval test subset, and NEREL test subset, respectively. We dropped four entity types in GENIA due to their number of entities being less than 50. After fine-tuning, we made inferences on the left instances in GENIA, GermEval test subset, and NEREL test subset.

4.2 Baselines

We compare our proposed PACL with seven baselines which can be categorized into three groups: 1) Rich-resource nested NER methods including NER-DP (Yu et al., 2020) and TIdentifier (Shen et al., 2021); 2) Metric-based few-shot NER methods including ProtoNet (Snell et al., 2017), NNShot (Yang and Katiyar, 2020), ESD (Wang et al., 2022c), and SpanProto (Wang et al., 2022a); 3) Contrastive-learning-based few-shot NER method CONTaiNER (Das et al., 2021). Appendix A details these baseline models.

4.3 Experimental Settings

We implemented PACL by Huggingface Transformer 4.21.1 and PyTorch 1.12.1. The model is initialized randomly and optimized by AdamW (Loshchilov and Hutter, 2017). We train and fine-tune the model with the learning rate 5e-5. For the text encoder, we use the pre-trained BERT_{base_multilingual} model since the languages of target domain datasets are different. The hidden

¹<http://www.geniaproject.org/genia-corpus>

²<https://sites.google.com/site/germeval2014ner/data>

³<https://github.com/nerel-ds/NEREL>

⁴<https://ningding97.github.io/fewnerd/>

Model/Framework	GENIA (32-way)		GermEval (12-way)		NEREL (29-way)		Average	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
NER-DP	15.26 \pm 2.78	31.89 \pm 4.01	7.12 \pm 2.61	24.89 \pm 3.92	15.86 \pm 5.77	42.25 \pm 2.42	12.75	33.01
TIdentifier	9.73 \pm 5.36	23.90 \pm 4.48	12.26 \pm 8.13	41.11 \pm 4.86	30.06 \pm 7.44	53.29 \pm 5.56	17.35	39.43
CONTAiNER	16.76 \pm 6.00	17.60 \pm 6.61	29.18 \pm 7.05	37.05 \pm 1.01	26.61 \pm 1.75	44.37 \pm 1.27	24.18	33.00
ProtoNet	21.83 \pm 3.39	37.18 \pm 1.81	33.20 \pm 9.00	47.95 \pm 4.06	38.70 \pm 4.62	50.22 \pm 1.28	31.24	45.12
NNShot	25.72 \pm 4.75	33.77 \pm 2.57	28.58 \pm 6.76	41.26 \pm 2.50	38.58 \pm 1.30	46.54 \pm 1.93	30.96	40.52
ESD	19.96 \pm 3.93	25.31 \pm 3.17	34.00 \pm 8.75	34.75 \pm 6.03	28.56 \pm 5.18	47.68 \pm 2.20	27.51	35.91
SpanProto	31.39 \pm 2.86	43.14 \pm 1.37	34.12 \pm 6.64	51.11 \pm 5.89	44.20 \pm 3.55	56.16 \pm 2.15	36.57	50.14
PACL	37.92 \pm 1.97	49.58 \pm 1.82	53.51 \pm 7.70	65.87 \pm 1.80	51.76 \pm 2.85	65.12 \pm 1.97	47.73	60.19

Table 2: F_1 performance on GENIA, GermEval, and NEREL datasets with 1-shot and 5-shot settings (%).

layer of the non-linear function f in equation 2 for getting span semantic representations is set to 512, and the initial value of the learnable hyperparameter λ for the O-type is set to 0.5. We set random seeds ranging from 0 to 10 to get ten results for each setting and report the average and standard deviation values to evaluate all models. We run all the experiments on a single NVIDIA A10 GPU.

4.4 Experimental Results

To evaluate the effectiveness of our PACL, we compare it against state-of-the-art baseline models introduced in 4.2. Table 2 shows their average F_1 results on GENIA, GermEval, and NEREL NER datasets with 1-shot and 5-shot settings over 10 times repeated experiments.

Our PACL achieves superior results compared to other models on 1-shot and 5-shot settings on the GENIA dataset, with F_1 scores of 37.92% and 49.58%, respectively, outperforming the best-performing baseline model, SpanProto, which scored 31.39% and 43.14%, respectively.

Our PACL also demonstrates superior performance on the GermEval dataset, scoring 53.51% on 1-shot and 65.87% on 5-shot, compared to the best-performing baseline model SpanProto on 1-shot with 34.12% and on 5-shot with 51.11%.

Finally, on the NEREL dataset, our PACL again outperforms the other models, scoring 51.76% on 1-shot and 65.12% on 5-shot, compared to the best-performing baseline model SpanProto on 1-shot SpanProto with 44.20% and on 5-shot TIdentifier with 56.16%.

Overall, these results demonstrate the effectiveness of our proposed PACL framework compared to the state-of-the-art baseline models.

4.5 Experimental Analysis

This section presents ablation studies, results on nested and flat entities separately in test datasets,

and the generality of the prototype-attention mechanism.

4.5.1 Ablation Study

		PACL	w/o PA
GENIA	1-shot	37.92 \pm 1.97	34.71 \pm 1.82
	5-shot	49.58 \pm 1.82	48.30 \pm 2.37
GermEval	1-shot	53.51 \pm 7.70	49.74 \pm 7.61
	5-shot	65.87 \pm 1.80	62.29 \pm 2.86
NEREL	1-shot	51.76 \pm 2.85	48.10 \pm 3.69
	5-shot	65.12 \pm 1.97	64.25 \pm 1.21

Table 3: Ablation study of F_1 performance on three datasets (%). “w/o PA” means removing the Prototype-Attention mechanism.

To evaluate the contribution of the designed Prototype-Attention (PA) mechanism to the overall performance of PACL, we conduct the ablation study by removing PA from the PACL. The results in Table 3 suggest that the PA mechanism positively impacts the F_1 score for the GENIA, GermEval, and NEREL datasets. To be specific, the PA mechanism improves the F_1 score by 3.21%, 3.77%, and 3.66% on the GENIA, GermEval, and NEREL datasets with the 1-shot setting, respectively. It also leads to 1.28%, 3.58%, and 0.87% increases on the GENIA, GermEval, and NEREL datasets with the 5-shot setting, respectively.

Overall, the results of this ablation study demonstrate that the PA mechanism enhances performance on various datasets with a particularly pronounced impact in 1-shot settings, as 1-shot settings exhibits larger prototype shifts than 5-shot settings. Appendix B shows how our PACL mitigates the prototype shift.

Note we did not explore the influence of the prototype-span contrastive loss by replacing it with a classical contrastive loss. This is because the classical contrastive loss performs similarity measurement between span and span. In PACL, the span

representations for the query set are enhanced by the PA mechanism while the span representations for the support set are not enhanced. Comparing these two different types of span representations is inappropriate. Therefore, we just study the effectiveness of the PA mechanism.

4.5.2 Nested-Flat Separate Results

In order to more comprehensively demonstrate the efficacy of the outcomes pertaining to nested entities across these datasets, we undertook a process of splitting and filtering exclusively for nested entities. Our proposed PACL framework outperforms all other models in both 1-shot and 5-shot settings on all three datasets. For the nested-only part of three test datasets, PACL achieves an F_1 score of 35.06% in 1-shot and 45.93% in 5-shot on GENIA, 32.82% in 1-shot and 48.52% in 5-shot on GermEval, and 36.70% in 1-shot and 50.89% in 5-shot on NEREL. The other baseline models achieve lower F_1 scores compared to PACL. For the nested part of the query set, our proposed PACL framework could achieve a 6.57% and 9.57% increase in terms of F_1 score in 1-shot and 5-shot settings, respectively. And For the flat part of the query set, our proposed PACL framework could achieve a 9.37% and 8.97% increase in terms of F_1 score in 1-shot and 5-shot settings, respectively. The more specific results are presented in Appendix C.

4.5.3 Generality of Prototype-Attention Mechanism

As the Prototype-Attention (PA) mechanism addresses the fundamental property of the prototype shift phenomenon, we believe it has a high level of generality and can enhance the performance of various models.

		SpanProto	SpanProto w PA
GENIA	1-shot	31.39 \pm 2.86	33.07 \pm 4.47
	5-shot	43.14 \pm 1.37	44.60 \pm 1.99
GermEval	1-shot	34.12 \pm 6.64	44.03 \pm 9.26
	5-shot	51.11 \pm 5.89	54.90 \pm 1.99
NEREL	1-shot	44.20 \pm 3.55	49.04 \pm 3.03
	5-shot	56.16 \pm 2.15	64.17 \pm 1.5

Table 4: F_1 performance before and after integrating the Prototype-Attention (PA) mechanism to SpanProto on three datasets (%).

To assess the generality of the PA mechanism, we conduct experiments by integrating it into the SpanProto and ESD models and comparing the performance before and after integration. As shown in

		ESD	ESD w PA
GENIA	1-shot	19.96 \pm 3.93	25.08 \pm 4.32
	5-shot	25.31 \pm 3.17	35.90 \pm 3.94
GermEval	1-shot	34.00 \pm 8.75	36.08 \pm 6.89
	5-shot	34.75 \pm 6.03	41.95 \pm 7.53
NEREL	1-shot	28.56 \pm 5.18	41.38 \pm 4.93
	5-shot	47.68 \pm 2.20	56.03 \pm 2.47

Table 5: F_1 performance before and after integrating the Prototype-Attention (PA) mechanism to ESD on three datasets (%).

Table 4 and Table 5, the experiment results demonstrate that integrating the PA mechanism into SpanProto and ESD notably improves the F_1 score on GENIA, GermEval, and NEREL datasets in both 1-shot and 5-shot settings.

These findings suggest that the PA mechanism has high generality and can serve as a valuable tool for NLP practitioners looking to improve their models’ performance in few-shot nested NER tasks.

5 Related Work

This section discusses related works on rich-resource nested NER, few-shot NER, and distribution shifts.

5.1 Rich-resource Nested NER

Nested NER aims to recognize entities with nested structures. Most of the current methods for nested NER are established on rich-resource datasets. These methods could be categorized into span-based, hypergraph-based, and layered-based (Wan et al., 2022).

Span-based methods treat sequences of tokens as spans and then label all possible spans by classification models (Shen et al., 2021; Li et al., 2020b; Tan et al., 2021). Hypergraph-based methods analyze the dependence of words in a sentence and then construct a dependency tree (Yu et al., 2020) or other structures (Wang and Lu, 2018; Katiyar and Cardie, 2018) to help identify nested entities. And layered-based methods capture the depth of entity nesting and apply multi-level sequence labeling strategies to recognize nested entities (Wang et al., 2021; Shibuya and Hovy, 2020).

These methods may be stuck in overfitting due to sophisticated models and the limited number of instances for training in the few-shot setting.

5.2 Few-shot NER

Few-shot NER requires recognizing entities with the support of very few labeled instances (Hofer

et al., 2018; Fritzler et al., 2019). Due to limited information contained in the support set, methods for few-shot NER mainly resort to a rich-resource source domain to help train models, resulting in transfer-learning and meta-learning frameworks.

Transfer-learning-based methods train models on a source domain and then transfer models or features to the few-labeled target domain (Yang et al., 2021; Liu et al., 2021). Meta-learning-based methods train models on adequate subtasks to make the model acquire the learning ability on few-shot tasks (de Lichy et al., 2021; Li et al., 2020a). Comparatively speaking, meta-learning-based methods are more widely used in few-shot NER due to their easy adaption to new tasks.

Within the meta-learning framework, various kinds of models are designed. For example, metric-based methods, including ProtoNet (Snell et al., 2017), NNShot (Yang and Katiyar, 2020), and SpanProto (Wang et al., 2022a), measure distances between prototypes in the support set and instances in the query set. Optimization-based methods, such as MAML (Finn et al., 2017) and FEWNER (Li et al., 2020a), train the model by a special optimizer. Model-based methods, such as SNAIL (Mishra et al., 2017) and CNPs (Garnelo et al., 2018), learn the hidden representation of instances on the support set and the query set to make inferences in an end-to-end manner. Contrastive-learning methods, such as CONTaiNER (Das et al., 2022), aims to maximize similarities of the same type and minimize similarities between different types.

These few-shot NER methods mostly focus on flat entities. Few works have discussed the few-shot nested NER setting. Wang converted sequence labeling to span-level matching for the few-shot flat NER and showed their method could handle nested entities (Wang et al., 2022b). However, it is not designed for the few-shot nested NER specifically.

5.3 Distribution Shifts

Distribution shift is a problem of training and testing data following two different distributions. It affects the generalization ability of supervised deep-learning models as the fundamental that these models could work is that training and testing data come from the same distribution. Inspired by real-world challenges, Wiles et al. summarized three distribution shifts: spurious correlation, low-data drift, and unseen data shift (Wiles et al., 2022). There have been some researches aiming to address dis-

tribution shifts in computer vision and general natural language processing tasks (Fang et al., 2020; Tu et al., 2022). To the best of our knowledge, researchers seldom discuss the distribution shift problem in the few-shot NER task. In this paper, we aim to tackle the few-shot nested NER task. Therefore, we rethink the distribution shift problem from the perspective of entity representation distribution and identify the prototype shift since it directly affects entity classification.

6 Conclusion

This paper first identifies the phenomenon of prototype shift that arises when there is a difference in prototypes between the support and query sets. Within the context of few-shot learning tasks, prototype shift is prone to occur since the few labeled instances in the support set could hardly represent the query set. To mitigate this issue in the few-shot nested NER task, we propose the Prototype-Attention Contrastive Learning (PACL) framework combining a prototype-attention mechanism and a prototype-span contrastive loss to enhance prototype representations. The experiments on three English, German, and Russian nested NER datasets demonstrated that PACL outperformed baseline models on the 1-shot and 5-shot settings. Future studies could explore the generality of PACL to other few-shot learning tasks.

7 Limitations

This paper still has several limitations. The first one is about the prototype shift adjustment. It is hard to completely address the prototype shift, while our PACL makes this attempt and achieves inspiring improvement. The second one is about other distribution shifts. Prototype shift is just one kind of distribution shift. Other distribution shifts also need to be identified and addressed to improve the accuracy of the few-shot nested NER task. The third one is about the language used for training. The results validate the performance of PACL across different languages, while the three languages used in this paper belong to the Indo-European family. This may introduce language bias to this language family and cause potential risk.

References

- Darina Benikova, Chris Biemann, and Marc Reznicek. 2014. Nosta-d named entity annotation for german: Guidelines and dataset. In *LREC*, pages 2524–2531.

- Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca Passonneau, and Rui Zhang. 2022. [CONTaiNER: Few-shot named entity recognition via contrastive learning](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 6338–6353, Dublin, Ireland. Association for Computational Linguistics.
- Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca J Passonneau, and Rui Zhang. 2021. Container: Few-shot named entity recognition via contrastive learning. *arXiv preprint arXiv:2109.07589*.
- Cyprien de Lichy, Hadrien Glaude, and William Campbell. 2021. Meta-learning for few-shot named entity recognition. In *Proceedings of the 1st Workshop on Meta Learning and Its Applications to Natural Language Processing*, pages 44–58.
- Ning Ding, Guangwei Xu, Yulin Chen, Xiaobin Wang, Xu Han, Pengjun Xie, Haitao Zheng, and Zhiyuan Liu. 2021. [Few-NERD: A few-shot named entity recognition dataset](#). 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 3198–3213, Online. Association for Computational Linguistics.
- Tongtong Fang, Nan Lu, Gang Niu, and Masashi Sugiyama. 2020. Rethinking importance weighting for deep learning under distribution shift. *Advances in Neural Information Processing Systems*, 33:11996–12007.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR.
- Omar U Florez and Erik Mueller. 2019. Learning to control latent representations for few-shot learning of named entities. *arXiv preprint arXiv:1911.08542*.
- Alexander Fritzier, Varvara Logacheva, and Maksim Kretov. 2019. Few-shot classification in named entity recognition task. In *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing*, pages 993–1000.
- Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. 2018. Conditional neural processes. In *International Conference on Machine Learning*, pages 1704–1713. PMLR.
- Maximilian Hofer, Andrey Kormilitzin, Paul Goldberg, and Alejo Nevado-Holgado. 2018. Few-shot learning for named entity recognition in medical text. *arXiv preprint arXiv:1811.05468*.
- Yutai Hou, Cheng Chen, Xianzhen Luo, Bohan Li, and Wanxiang Che. 2022. [Inverse is better! fast and accurate prompt for few-shot slot tagging](#). In *Findings of the Association for Computational Linguistics: ACL* 2022, pages 637–647, Dublin, Ireland. Association for Computational Linguistics.
- Yutai Hou, Zhihan Zhou, Yijia Liu, Ning Wang, Wanxiang Che, Han Liu, and Ting Liu. 2019. Few-shot sequence labeling with label dependency transfer and pair-wise embedding. *arXiv preprint arXiv:1906.08711*.
- Arzoo Katiyar and Claire Cardie. 2018. [Nested named entity recognition revisited](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 861–871, New Orleans, Louisiana. Association for Computational Linguistics.
- J-D Kim, Tomoko Ohta, Yuka Tateisi, and Jun’ichi Tsujii. 2003. Genia corpus—a semantically annotated corpus for bio-textmining. *Bioinformatics*, 19(suppl_1):i180–i182.
- Jing Li, Billy Chiu, Shanshan Feng, and Hao Wang. 2020a. Few-shot named entity recognition via meta-learning. *IEEE Transactions on Knowledge and Data Engineering*.
- Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020b. [A unified MRC framework for named entity recognition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5849–5859, Online. Association for Computational Linguistics.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2021. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *arXiv preprint arXiv:2107.13586*.
- Ilya Loshchilov and Frank Hutter. 2017. [Fixing weight decay regularization in adam](#). *CoRR*, abs/1711.05101.
- Natalia Loukachevitch, Ekaterina Artemova, Tatiana Batura, Pavel Braslavski, Ilia Denisov, Vladimir Ivanov, Suresh Manandhar, Alexander Pugachev, and Elena Tutubalina. 2021. NEREL: A Russian dataset with nested named entities, relations and events. In *Proceedings of RANLP*, pages 876–885.
- Ruotian Ma, Xin Zhou, Tao Gui, Yiding Tan, Linyang Li, Qi Zhang, and Xuanjing Huang. 2022. [Template-free prompt tuning for few-shot NER](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5721–5732, Seattle, United States. Association for Computational Linguistics.
- Hong Ming, Jiaoyun Yang, Lili Jiang, Yan Pan, and Ning An. 2022. Few-shot nested named entity recognition. *arXiv preprint arXiv:2212.00953*.
- Nikhil Mishra, Mostafa Rohaninejad, Xi Chen, and Pieter Abbeel. 2017. A simple neural attentive meta-learner. *arXiv preprint arXiv:1707.03141*.

726	Yongliang Shen, Xinyin Ma, Zeqi Tan, Shuai Zhang,	Jianing Wang, Chengyu Wang, Chuanqi Tan, Minghui	783
727	Wen Wang, and Weiming Lu. 2021. Locate and label: A two-stage identifier for nested named entity	Qiu, Songfang Huang, Jun Huang, and Ming Gao.	784
728	recognition . In <i>Proceedings of the 59th Annual Meeting</i>	2022a. Spanproto: A two-stage span-based prototyp-	785
729	<i>of the Association for Computational Linguistics</i>	ical network for few-shot named entity recognition .	786
730	<i>and the 11th International Joint Conference on Natural</i>	<i>CoRR</i> , abs/2210.09049.	787
731	<i>Language Processing (Volume 1: Long Papers)</i> ,		
732	pages 2782–2794, Online. Association for Computa-	Peiyi Wang, Runxin Xu, Tianyu Liu, Qingyu Zhou,	788
733	tional Linguistics.	Yunbo Cao, Baobao Chang, and Zhifang Sui. 2022b.	789
734		An enhanced span-based decomposition method for	790
		few-shot sequence labeling . In <i>Proceedings of the</i>	791
735	Takashi Shibuya and Eduard Hovy. 2020. Nested	<i>2022 Conference of the North American Chapter of</i>	792
736	Named Entity Recognition via Second-best Sequence	<i>the Association for Computational Linguistics: Hu-</i>	793
737	Learning and Decoding . <i>Transactions of the Associa-</i>	<i>man Language Technologies</i> , pages 5012–5024, Seat-	794
738	<i>tion for Computational Linguistics</i> , 8:605–620.	tle, United States. Association for Computational	795
		Linguistics.	796
739	Jake Snell, Kevin Swersky, and Richard S Zemel. 2017.	Peiyi Wang, Runxin Xu, Tianyu Liu, Qingyu Zhou,	797
740	Prototypical networks for few-shot learning. <i>arXiv</i>	Yunbo Cao, Baobao Chang, and Zhifang Sui. 2022c.	798
741	<i>preprint arXiv:1703.05175</i> .	An enhanced span-based decomposition method for	799
		few-shot sequence labeling . In <i>Proceedings of the</i>	800
742	Yifan Sun, Changmao Cheng, Yuhang Zhang, Chi Zhang,	<i>2022 Conference of the North American Chapter of</i>	801
743	Liang Zheng, Zhongdao Wang, and Yichen Wei.	<i>the Association for Computational Linguistics: Hu-</i>	802
744	2020. Circle loss: A unified perspective of pair simi-	<i>man Language Technologies</i> , pages 5012–5024, Seat-	803
745	larity optimization. In <i>Proceedings of the IEEE/CVF</i>	tle, United States. Association for Computational	804
746	<i>Conference on Computer Vision and Pattern Recogni-</i>	Linguistics.	805
747	<i>tion</i> , pages 6398–6407.		
		Yiran Wang, Hiroyuki Shindo, Yuji Matsumoto, and	806
748	Flood Sung, Yongxin Yang, Li Zhang, Tao Xiang,	Taro Watanabe. 2021. Nested named entity recog-	807
749	Philip HS Torr, and Timothy M Hospedales. 2018.	nition via explicitly excluding the influence of the	808
750	Learning to compare: Relation network for few-shot	best path . In <i>Proceedings of the 59th Annual Meet-</i>	809
751	learning. In <i>Proceedings of the IEEE conference</i>	<i>ing of the Association for Computational Linguistics</i>	810
752	<i>on computer vision and pattern recognition</i> , pages	<i>and the 11th International Joint Conference on Natu-</i>	811
753	1199–1208.	<i>ral Language Processing (Volume 1: Long Papers)</i> ,	812
		pages 3547–3557, Online. Association for Computa-	813
754	Zeqi Tan, Yongliang Shen, Shuai Zhang, Weiming Lu,	tional Linguistics.	814
755	and Yueting Zhuang. 2021. A sequence-to-set net-		
756	work for nested named entity recognition . In <i>Pro-</i>	Olivia Wiles, Sven Gowal, Florian Stimberg, Sylvestre-	815
757	<i>ceedings of the Thirtieth International Joint Confer-</i>	Alvise Rebuffi, Ira Ktena, Krishnamurthy Dj Dvi-	816
758	<i>ence on Artificial Intelligence, IJCAI-21</i> , pages 3936–	jotham, and Ali Taylan Cemgil. 2022. A fine-grained	817
759	3942. International Joint Conferences on Artificial	analysis on distribution shift . In <i>International Con-</i>	818
760	Intelligence Organization. Main Track.	<i>ference on Learning Representations</i> .	819
761	Jianhong Tu, Ju Fan, Nan Tang, Peng Wang, Chengliang	Guanqun Yang, Shay Dineen, Zhipeng Lin, and Xue-	820
762	Chai, Guoliang Li, Ruixue Fan, and Xiaoyong Du.	qing Liu. 2021. Few-sample named entity recogni-	821
763	2022. Domain adaptation for deep entity resolution .	tion for security vulnerability reports by fine-tuning	822
764	In <i>SIGMOD '22: International Conference on Man-</i>	pre-trained language models. In <i>International Work-</i>	823
765	<i>agement of Data, Philadelphia, PA, USA, June 12 -</i>	<i>shop on Deployable Machine Learning for Security</i>	824
766	<i>17, 2022</i> , pages 443–457. ACM.	<i>Defense</i> , pages 55–78. Springer.	825
767	Juncheng Wan, Dongyu Ru, Weinan Zhang, and Yong	Yi Yang and Arzoo Katiyar. 2020. Simple and effec-	826
768	Yu. 2022. Nested named entity recognition with span-	tive few-shot named entity recognition with struc-	827
769	level graphs . In <i>Proceedings of the 60th Annual Meet-</i>	tured nearest neighbor learning. <i>arXiv preprint</i>	828
770	<i>ing of the Association for Computational Linguistics</i>	<i>arXiv:2010.02405</i> .	829
771	<i>(Volume 1: Long Papers)</i> , pages 892–903, Dublin,		
772	Ireland. Association for Computational Linguistics.	Juntao Yu, Bernd Bohnet, and Massimo Poesio. 2020.	830
		Named entity recognition as dependency parsing . In	831
773	Bailin Wang and Wei Lu. 2018. Neural segmental hy-	<i>Proceedings of the 58th Annual Meeting of the Asso-</i>	832
774	pergraphs for overlapping mention recognition . In	<i>ciation for Computational Linguistics</i> , pages 6470–	833
775	<i>Proceedings of the 2018 Conference on Empirical</i>	6476, Online. Association for Computational Lin-	834
776	<i>Methods in Natural Language Processing</i> , pages 204–	guistics.	835
777	214, Brussels, Belgium. Association for Computa-		
778	tional Linguistics.	A Detail of Baselines	836
779	Feng Wang and Huaping Liu. 2021. Understanding	Detailed information on baseline models is intro-	837
780	the behaviour of contrastive loss. In <i>Proceedings of</i>	duced in this section. We compare our PACL with	838
781	<i>the IEEE/CVF conference on computer vision and</i>	the following seven baseline models:	839
782	<i>pattern recognition</i> , pages 2495–2504.		

- NER-DP (Yu et al., 2020) is a rich-resource-based nested NER method. It applies a bi-affine model to score pairs of start and end tokens for each span to establish dependency parsing for identifying nested entities.
- TIdentifier (Shen et al., 2021) is also a rich-resource-based nested NER method. It applies a Two-stage Identifier (TIdentifier), including a seed span generation module for locating entities and a span proposal module for classifying entities.
- CONTaiNER (Das et al., 2021) is a contrastive-learning-based few-shot NER method. It first obtains entities' Gaussian-distributed embeddings and then optimizes a generalized objective of differentiating between entity types by a contrastive loss function. We adapt it to handle nested entities with the entity span formulation.
- ProtoNet (Snell et al., 2017) is a metric-learning-based few-shot NER method. It applies prototypical networks to learn a metric space for obtaining prototype representations. We also adapt it to handle nested entities with the entity span formulation.
- NNShot (Yang and Katiyar, 2020) is also a metric-learning-based few-shot NER method. It applies structured decoding and nearest-neighbor learning to identify entities. We utilize the entity span formulation to make it handle nested entities.
- ESD (Wang et al., 2022c) is a metric-learning-based few-shot NER method. It formulates the task as a span-level matching problem. To identify entities, it performs span-level procedures, including enhanced span representation, class prototype aggregation, and span conflict resolution.
- SpanProto (Wang et al., 2022a) is a metric-learning-based few-shot NER method. It also applies entity spans to formulate the problem. For identifying entities, it first utilizes a span extractor to recognize candidate entity spans and then applies a mention classifier to determine entity types.

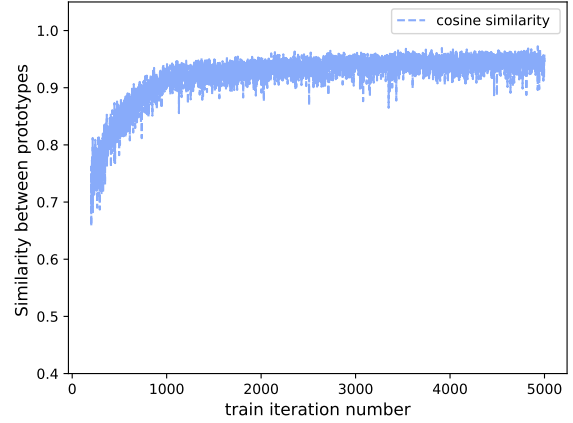


Figure 4: Illustration of the change of the prototype similarity during training.

B Prototype Shift Mitigation by PACL

This paper aims to mitigate prototype shifts, and section 1 has already validated the existence of the prototype shift phenomenon. This section examines how the prototype shift changes by applying our PACL.

We utilize the cosine similarity to denote the prototype differences between the support and query sets to measure the prototype shift. Figure 4 illustrates the change of the prototype similarity with the increase of iteration numbers during training. We could find a consistently increasing trend in prototype similarity, which means the prototype shift is consistently decreasing. This validates the effectiveness of our PACL in mitigating prototype shifts.

C Results on Nested/Flat-Only Entities

We split the query set of each dataset into two subsets: one only contains nested entities and the other one only contains flat entities. We then evaluate the model over 10 times repeated experiments. Table 6 and table 7 show the average F_1 results of nested-only and flat-only entities on GENIA, GermEval, and NEREL NER datasets with 1-shot and 5-shot settings. In the nested part of the query set, our proposed PACL framework achieves a 6.57% and 9.57% increase in terms of F_1 score in the 1-shot and 5-shot settings, respectively. Similarly, in the flat part of the query set, our proposed PACL framework achieves a 9.37% and 8.97% increase in terms of F_1 score in the 1-shot and 5-shot settings, respectively.

Compared to the baseline models, our proposed PACL model achieves the best results not only in

Model/Framework	GENIA (32-way)		GermEval (12-way)		NEREL (29-way)		Average	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
NER-DP	13.94 \pm 3.96	28.16 \pm 2.43	5.84 \pm 3.11	8.12 \pm 4.30	4.83 \pm 2.71	19.87 \pm 3.92	8.20	18.72
TIdentifier	9.19 \pm 5.36	25.15 \pm 5.50	9.22 \pm 6.56	27.79 \pm 3.54	20.13 \pm 7.28	39.93 \pm 5.63	12.85	30.96
CONTAiNER	16.19 \pm 3.65	12.52 \pm 8.11	15.84 \pm 4.70	15.21 \pm 4.92	20.71 \pm 3.94	30.03 \pm 2.04	17.58	19.25
ProtoNet	18.99 \pm 3.29	32.53 \pm 2.24	21.62 \pm 5.85	33.66 \pm 3.23	26.19 \pm 7.67	40.93 \pm 2.36	22.27	35.71
NNShot	24.84 \pm 5.55	30.71 \pm 2.60	27.36 \pm 6.95	28.30 \pm 7.47	28.69 \pm 8.42	42.92 \pm 5.17	26.96	33.98
ESD	18.08 \pm 2.99	22.9 \pm 2.13	19.86 \pm 5.62	22.33 \pm 5.00	24.23 \pm 5.29	30.85 \pm 8.26	20.72	25.36
SpanProto	30.24 \pm 2.77	40.50 \pm 2.04	24.11 \pm 7.57	34.06 \pm 3.43	30.51 \pm 5.86	42.07 \pm 1.29	28.29	38.88
PACL	35.06 \pm 3.52	45.93 \pm 1.93	32.82 \pm 9.39	48.52 \pm 1.90	36.70 \pm 4.88	50.89 \pm 2.99	34.86	48.45

Table 6: **nested** F_1 performance on GENIA, GermEval, and NEREL datasets with 1-shot and 5-shot settings (%).

Model/Framework	GENIA (32-way)		GermEval (12-way)		NEREL (29-way)		Average	
	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot	1-shot	5-shot
NER-DP	12.73 \pm 3.94	26.87 \pm 3.13	6.74 \pm 3.65	15.61 \pm 6.14	9.21 \pm 3.82	34.03 \pm 3.74	9.56	25.50
TIdentifier	14.07 \pm 8.48	26.81 \pm 3.54	12.53 \pm 8.44	42.19 \pm 5.15	32.45 \pm 7.88	55.18 \pm 6.20	19.68	41.39
CONTAiNER	16.10 \pm 2.31	12.12 \pm 7.83	25.44 \pm 8.77	27.70 \pm 9.05	35.71 \pm 2.87	45.08 \pm 2.67	25.75	28.30
ProtoNet	19.42 \pm 3.86	34.21 \pm 2.05	35.79 \pm 8.07	50.45 \pm 3.35	39.81 \pm 6.39	56.62 \pm 2.91	31.67	47.09
NNShot	22.61 \pm 5.06	29.73 \pm 2.94	50.42 \pm 6.92	44.23 \pm 14.07	44.97 \pm 5.68	57.57 \pm 6.01	39.33	43.84
ESD	17.27 \pm 4.41	21.64 \pm 3.38	31.46 \pm 8.46	35.43 \pm 6.42	38.40 \pm 4.08	43.42 \pm 9.94	29.04	33.50
SpanProto	29.16 \pm 3.35	40.73 \pm 1.49	38.97 \pm 9.63	53.66 \pm 3.75	47.51 \pm 3.08	59.15 \pm 1.75	38.55	51.18
PACL	35.96 \pm 2.16	46.62 \pm 2.16	55.77 \pm 7.73	67.63 \pm 2.05	54.38 \pm 2.85	66.19 \pm 2.06	48.70	60.15

Table 7: **flat** F_1 performance on GENIA, GermEval, and NEREL datasets with 1-shot and 5-shot settings (%).

the 1-shot and 5-shot experimental settings but also in both the nested and flat settings. This indicates that our proposed model can effectively classify nested entities compared to other baseline models.