

Large-Scale Label Interpretation Learning for Few-Shot Named Entity Recognition

Jonas Golde

Felix Hamborg

Alan Akbik

Humboldt Universität zu Berlin

goldejon@informatik.hu-berlin.de

{felix.hamborg, alan.akbik}@hu-berlin.de

Abstract

Few-shot named entity recognition (NER) detects named entities within text using only a few annotated examples. One promising line of research is to leverage natural language descriptions of each entity type: the common label PER might, for example, be verbalized as “person entity.” In an initial *label interpretation learning* phase, the model learns to interpret such verbalized descriptions of entity types. In a subsequent *few-shot tagset extension* phase, this model is then given a description of a previously unseen entity type (such as “music album”) and optionally a few training examples to perform few-shot NER for this type. In this paper, we systematically explore the impact of a strong semantic prior to interpret verbalizations of new entity types by massively scaling up the number and granularity of entity types used for label interpretation learning. To this end, we leverage an entity linking benchmark to create a dataset with orders of magnitude of more distinct entity types and descriptions as currently used datasets. We find that this increased signal yields strong results in zero- and few-shot NER in in-domain, cross-domain, and even cross-lingual settings. Our findings indicate significant potential for improving few-shot NER through heuristical data-based optimization.

1 Introduction

Few-shot named entity recognition (NER) refers to identifying and classifying named entities within text by learning from a few annotated examples. A widely adopted strategy in few-shot NER employs transfer learning with pre-trained language models (PLMs) to interpret labels based on their semantic meaning (Yang and Katiyar, 2020; de Lichy et al., 2021; Das et al., 2022; Ma et al., 2022a,b,c; Chen et al., 2023). The main idea is that such models learn to interpret a natural language description of an entity type for use in a word-level decoder. They learn in two phases:

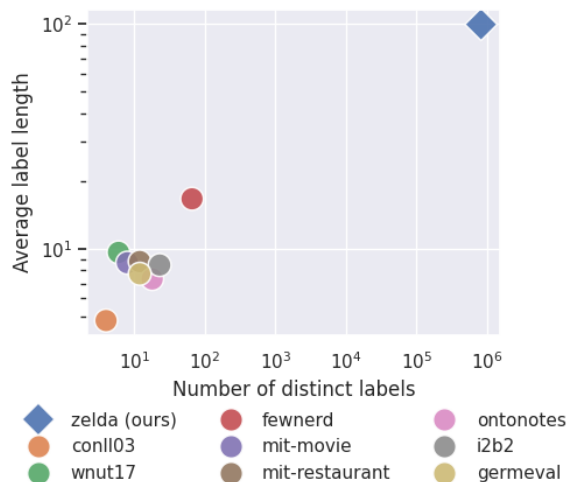


Figure 1: Given existing datasets, few-shot NER methods requiring an initial label interpretation learning are limited regarding entity types and label verbalizations. We propose learning from orders of magnitude more distinct types and more expressive label semantics than current datasets by utilizing ZELDA annotated with WikiData information.

1. a *label interpretation learning* phase on a NER-annotated dataset with a set of entity types and their verbalizations. For instance, the common label PER might be verbalized as "person entity." In this phase, the model learns to associate entity type verbalizations with matching NER annotations.
2. a *few-shot tagset extension phase* in which the model is expanded to previously unseen domains or entity types using only a new verbalization and optionally a few example annotations. For instance, to extend the model to recognize the names of music albums, one would only need to provide a verbalization ("music album") and a few examples.

Limitations. However, as Figure 1 indicates, prior studies used only very limited numbers of distinct entity types for label interpretation learning.

This is an artifact of relying on common NER datasets such as CoNLL-03 (Tjong Kim Sang and De Meulder, 2003), OntoNotes (Pradhan et al., 2012), WNUT-17 (Derczynski et al., 2017), or FewNERD (Ding et al., 2021), which only contain a small number of distinct entity types (between 4 and 66 types). Furthermore, the majority of their entity types have a simple semantic definition, such as “person,” “location,” or “organization,” and occur across several datasets. We hypothesize that these limitations overly constrain the semantic signal that is observed during label interpretation learning, thus constituting a main limiting factor to few-shot NER.

Contributions. With this paper, we introduce a novel approach named LITSET (label interpretation learning by scaling entity types) and systematically investigate the intuition that increasing the number of distinct entity types and their semantic exactness in label interpretation learning introduces a strong semantic prior to understand unseen entities in few-shot settings. To this end, we heuristically create a dataset with orders of magnitude more distinct entity types than commonly employed (cf. Figure 1) and use it for extensive experimentation. In more detail, our contributions are:

- We present experiments to validate our hypothesis on the largest existing NER dataset (FewNERD). We find that few-shot performance increases with label interpretation learning on more distinct entity types and more expressive descriptions (cf. Section 2).
- We derive a dataset with orders of magnitude more granular entity type annotations to massively scale up label interpretation learning. Our approach leverages the recently released entity linking benchmark ZELDA (Milich and Akbik, 2023) and enriches it with type descriptions from WikiData (Vrandečić and Krötzsch, 2014) (cf. Section 3).
- We comprehensively evaluate label interpretation learning on our derived corpus against classical setups for zero- and few-shot NER in in-domain, cross-domain, and cross-lingual settings and transfer it to different model architectures (cf. Section 4).

We find that label interpretation learning on our heuristically derived corpus matches and, in

many cases, significantly outperforms strong baselines. Our findings indicate significant potential for improving few-shot NER through heuristical data-based optimization. We release the generated dataset and source code under the Apache 2 license on Github¹.

2 Validation Experiment for Impact of Entity Types and Label Descriptions

We first conduct an experiment to validate the intuition that a richer training signal for label interpretation learning positively impacts few-shot NER. To this end, we create a set of training datasets for label interpretation learning that each contain the same number of entities but vary in the number of distinct entity types and their label verbalization. We then compare the few-shot NER ability of models trained on each of these datasets.

2.1 Experimental Setup

Definitions. To evaluate few-shot NER, an existing dataset \mathcal{D} is split based on its labels \mathcal{L} : the label interpretation training split \mathcal{D}^{LIT} and a few-shot fine-tuning split \mathcal{D}^{FS} . The corresponding labels of each split \mathcal{L}^{LIT} and \mathcal{L}^{FS} are set such that $\mathcal{L}^{LIT} \cup \mathcal{L}^{FS} = \mathcal{L}$ and $\mathcal{L}^{LIT} \cap \mathcal{L}^{FS} = \emptyset$.

For few-shot tagset extension, we sample a support set \mathcal{S} by k -shot down-sampling \mathcal{D}^{FS} . The support set \mathcal{S} contains each label from \mathcal{L}^{FS} exactly k times. We sample three different support sets using different seeds and report the averaged micro-F1 scores over these iterations.

Dataset. We use FewNERD in our experiment since it is the largest existing dataset w.r.t. the number of distinct entity types (66 types). We set the labels of \mathcal{D}^{LIT} to be the 50 most occurring entity types and the labels of \mathcal{D}^{FS} to be the 16 least occurring. We perform an analysis along two dimensions:

- To measure the impact of more distinct entity types in label interpretation learning, we create 5 versions of the training data containing 3, 5, 10, 30, and all 50 labels, respectively. Importantly, all versions contain the same number of annotations (10k) to ensure an equal entity detection ability.
- To measure the impact of richer verbalizations, we define 3 different labels semantics: (1) a

¹<https://github.com/flairNLP/label-interpretation-learning>

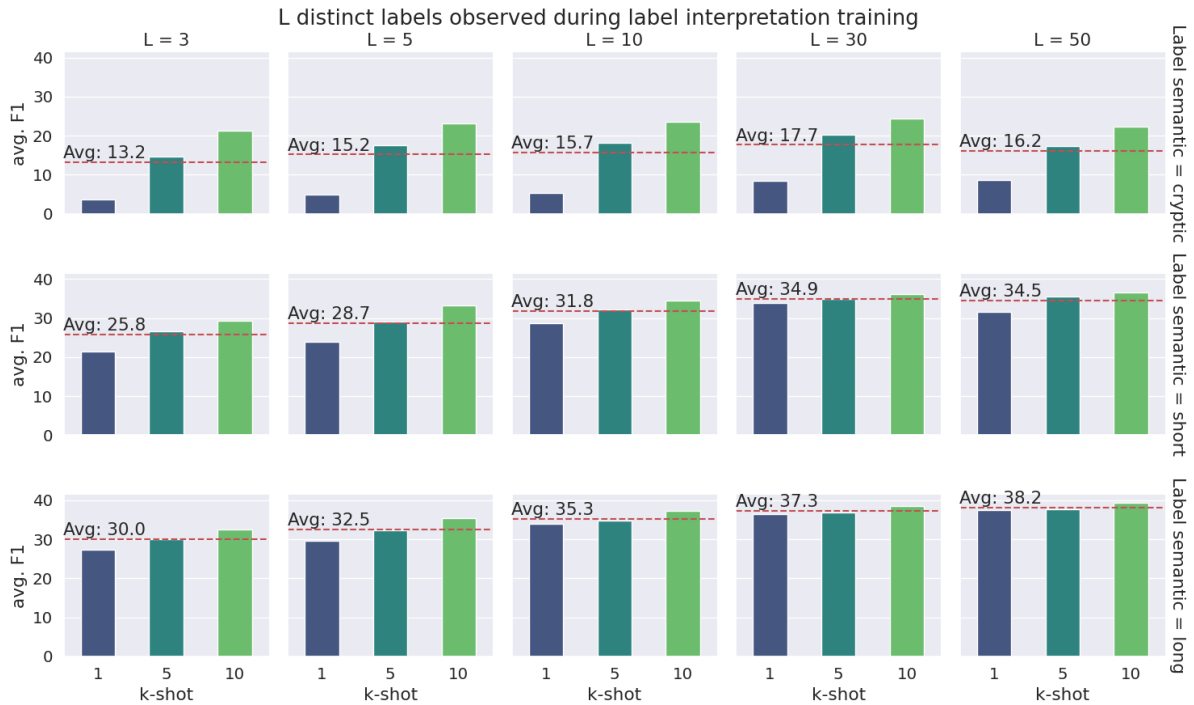


Figure 2: F1 scores for few-shot NER tagset extension on FewNERD depending on how many distinct entity types were seen in label interpretation learning (columns) and how label types were verbalized (rows). We report F1 scores averaged over five seeds. We observe that (1) more distinct labels during label interpretation training and (2) more semantically expressive labels improve the few-shot ability on unseen labels.

"cryptic" unique, random 2-character label, (2) a "short" description as regularly used according to research and (3) a "long" description with examples (cf. Appendix A).

To exclude the respective labels from each split, we follow prior work and mask labels \mathcal{L}^{LIT} in \mathcal{D}^{FS} and \mathcal{L}^{FS} in \mathcal{D}^{LIT} with the 0-token (meaning no named entity).

Few-shot model. We employ the frequently used bi-encoder architecture (Blevins and Zettlemoyer, 2020; Ma et al., 2022a) with two bert-base-uncased transformers (Vaswani et al., 2017) as our backbone architecture.

We argue that this architecture has an essential advantage over approaches using cross-attention such as Li et al. (2020); Halder et al. (2020); Chen et al. (2023). Previously mentioned methods are limited by the input size of the model (e.g., 512 for BERT) because they prepend label verbalizations to the processed sentence. One could overcome this limitation with one forward pass per label-sentence pair. However, both options become computationally expensive with extensive type descriptions or many distinct entity types. The bi-encoder can be easily adapted to handle an arbitrary number of

distinct labels (see Section 3.2).

2.2 Results

Figure 2 shows the results of tagset extension when performing label interpretation learning on FewNERD subsets with different numbers of labels (columns) and different verbalization methods (rows). For each label interpretation learning, we report the average F1-score for tagset extension for 1-shot, 5-shot, and 10-shot learning, respectively.

Improved generalization with more types. We observe that the number of distinct labels seen during label interpretation training increases the generalization in few-shot settings independent of the label semantics used. We find improvements from +3.0 F1 (cf. $L = 3$ vs. $L = 50$, label semantic: cryptic) up to +8.7 F1 (cf. $L = 3$ vs. $L = 50$, label semantic: short) on average in pp.

More expressive descriptions helpful. We also find that increasing the expressiveness of label verbalizations strongly improves the few-shot performance. This observation is independent of the distinct number of labels seen in label interpretation learning, such that we find improvements ranging from +16.8 F1 (cf. label semantics: simple vs. long, with $L = 3$) up to +22.0 F1 (cf. label semantics:

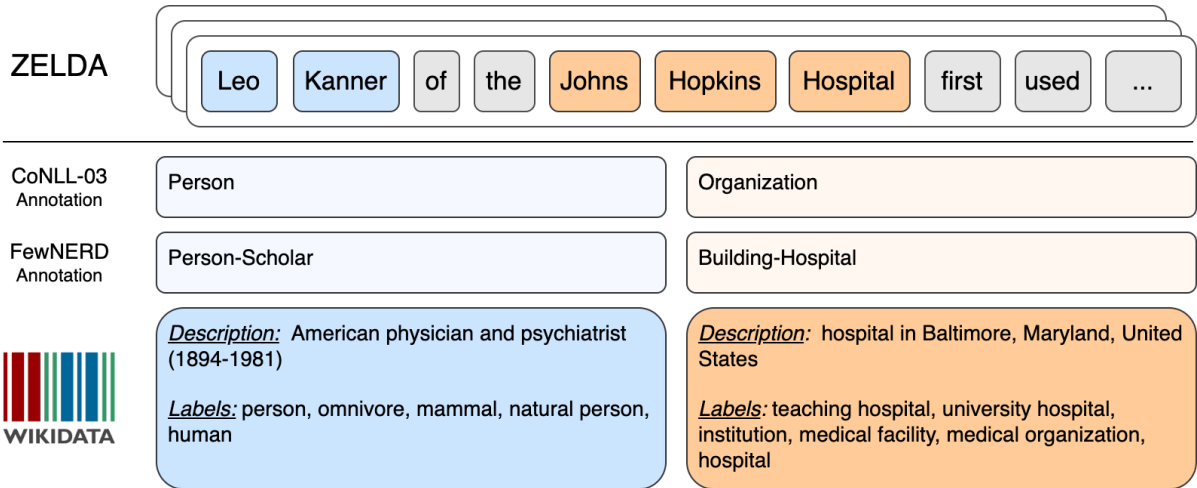


Figure 3: An example annotation of a sentence in ZELDA. WikiData provides precise descriptions and labels about an entity. Annotation types in existing datasets (CoNLL-03, FewNERD) are less informative if not misleading.

simple vs. long, with $L = 50$) on average in pp.

These observations on FewNERD confirm our intuition that a richer training signal in label interpretation learning improves few-shot NER performance. To verify this observation for other models, we repeat this experiment with a pre-trained transformer on sparse latent typing, an objective to sparsely extract sentence-level keywords with diverse latent types, where we make the same observation. These experiments are illustrated in detail in Appendix B.

3 Large-Scale Label Interpretation Learning

As our validation experiment shows a positive impact of increasing the number and expressivity of entity types, we now aim to scale the signal for label interpretation learning to orders of magnitude more entity types. To this end, we heuristically derive a NER-annotated dataset using the recently released entity linking benchmark ZELDA and annotate it with WikiData information (Section 3.1). We also introduce a modified training procedure for the bi-encoder to handle a very large space of entity types that applies to all architectures of its kind (Section 3.2). We call this approach LITSET (label interpretation learning by scaling entity types).

3.1 LITSET Dataset

The task of entity disambiguation is closely related to NER. Here, an already detected entity is disambiguated by linking it to an existing knowledge base such as Wikipedia or WikiData. Existing training and evaluation datasets for entity disambiguation

Dataset	Label length	# Distinct types
CoNLL-03	9.8 ± 2.9	4
WNUT17	8.3 ± 2.8	6
OntoNotes	9.8 ± 8.5	18
FewNERD	17.3 ± 7.6	66
LITSET	99.8 ± 45.4	~817k

Table 1: Average label description length (in characters) and distinct entity types of NER datasets. Label length and distinct entity types for LITSET refers to all annotations as indicated in Figure 3.

thus contain named entities marked with links to entries in the WikiData knowledge base.

One advantage of WikiData is that it contains fine-grained labels and free-form text descriptions of entities in the knowledge base. For instance, the entity "John Hopkins Hospital" (cf. Figure 3) has the free-form description "hospital in Baltimore, Maryland" and belongs to the classes "teaching hospital", "university hospital", and many others. As the Figure shows, these labels are significantly more fine-grained than CoNLL-03 and even FewNERD entity types which simply classify it as an "organization" or a "hospital" respectively.

Deriving the dataset. We leverage the classes and descriptions from WikiData as type annotations in our approach. For each linked entity in the dataset, we retrieve the types and descriptions from WikiData and use them as NER annotations. We refer to Appendix C for a detailed explanation of the fields used.

To best prepare our model for arbitrary labels in a few-shot setting, we sample the annotations to learn to interpret annotations on different hier-

archies. We assume labels to represent high-level types, whereas descriptions are very specific to that entity. Specifically, for each entity x_i , we uniformly sample whether we annotate it with either the description attribute or the labels attribute (cf. Figure 3). When utilizing the labels attribute, we randomly select the number of tags following a geometric distribution with $p = .5$. Subsequently, we uniformly sample tags from the label attribute until the number of tags is reached. Lastly, we concatenate the selected tags for final annotation.

3.2 Backbone Architecture

Due to its simplicity, we conduct our experiments using the widely adopted bi-encoder model. It utilizes two separate transformers to encode tokens and labels, respectively. The first transformer generates embeddings $e_t \in \mathbb{R}^{N \times H}$ for all tokens, where N represents the number of tokens and H denotes the hidden size of the model. The second obtains the [CLS]-token embeddings e_l for the labels converted into natural language. We employ cross-entropy loss and derive final predictions with

$$\hat{y} = \arg \max \text{softmax}(e_t \cdot e_l)$$

However, training a model, including the bi-encoder, with a wide array of distinct classes is non-trivial. With \mathcal{L} denoting the set of labels, the shape of label representations is $e_l \in \mathbb{R}^{|\mathcal{L}| \times H}$. Given that $|\mathcal{L}| \approx 10^6$ (cf. Figure 1), we aim to circumvent the resulting matrix multiplication for two reasons: (1) computational limitations and (2) optimization difficulty. To alleviate these issues, we restrict our consideration to labels present in the current batch \mathcal{L}_b with $|\mathcal{L}_b| \ll |\mathcal{L}|$ for loss calculation.

4 Experiments

We evaluate the impact of label interpretation training in various tagset extension settings. Throughout all experiments, we compare label interpretation learning on LITSET with training on different baseline datasets. We present all hyperparameters used for our experiments in Appendix D. Specifically, we conduct the following experiments:

1. *In-domain transfer*: Identical domain in label interpretation learning and few-shot fine-tuning (cf. Section 4.1).
2. *Cross-domain transfer*: Different domain in label interpretation learning and few-shot fine-tuning (cf. Section 4.2).

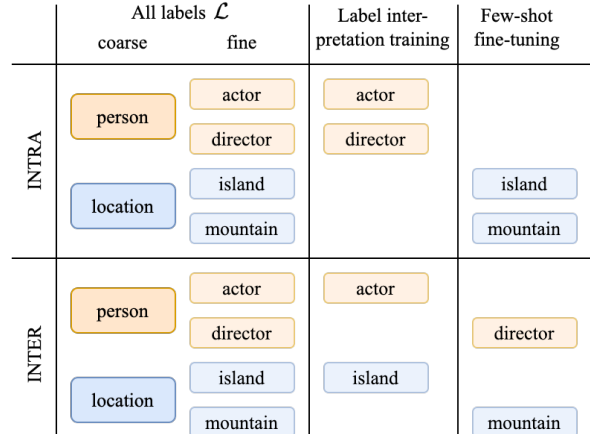


Figure 4: Exemplary illustration on the INTRA and INTER settings of FewNERD experiments.

3. *Transfer to advanced bi-encoders*: Identical to in-domain setting, but we transfer our approach to advanced bi-encoder architectures (cf. Section 4.3).
4. *Cross-lingual transfer*: Identical domain in label interpretation learning and few-shot fine-tuning, but languages differ between both phases (cf. Section 4.4).

Further, we support our experiments by analyzing the impact of different label semantics used between label interpretation learning and few-shot fine-tuning (cf. Section 4.1). At last, we refer to our ablation experiments using (1) different transformers as label encoders and (2) negative sampling (cf. Appendices E and F).

4.1 Experiment 1: In-Domain Transfer

This experiment replicates the most common evaluation setup for few-shot tagset extension, where both \mathcal{D}^{LIT} and \mathcal{D}^{FS} are sourced from the same NER dataset. Our baseline is the default approach of label interpretation learning on \mathcal{D}^{LIT} , which is "in-domain" since it shares the same textual domain and entity annotations are aligned on identical semantic levels as the evaluation data, whereas label interpretation learning on LITSET does not have these advantages.

4.1.1 Experimental Setup

We use OntoNotes and FewNERD in our experiments as they have important properties: OntoNotes covers multiple domains and languages such that we can measure the transferability of our approach. FewNERD comes with two annotation

Evaluation data \mathcal{D}^{FS} for tagset extension from:	Label interpretation learning data \mathcal{D}^{LIT} from:	0-shot	1-shot	5-shot	10-shot	Avg.
OntoNotes	LITSET	8.7 \pm 1.7	21.9 \pm 8.4	40.1 \pm 7.2	<u>48.4</u> \pm 6.2	29.5
	w/ all labels	3.5 \pm 1.3	<u>20.0</u> \pm 9.5	<u>38.4</u> \pm 8.3	46.5 \pm 6.3	<u>27.1</u>
	w/ labels only	0.1 \pm 0.1	14.3 \pm 8.3	29.6 \pm 6.9	37.5 \pm 6.1	20.4
	w/ description only	<u>4.2</u> \pm 1.3	19.8 \pm 8.8	37.5 \pm 7.9	46.2 \pm 5.9	26.9
	OntoNotes (<i>Baseline</i>)	0.2 \pm 0.1	11.2 \pm 9.3	38.3 \pm 12.0	54.9 \pm 7.6	26.2
FewNERD _{INTRA}	LITSET	3.2 \pm 1.0	30.7 \pm 5.3	51.9 \pm 5.2	57.9 \pm 6.2	35.9
	w/ all labels	0.9 \pm 0.4	<u>20.1</u> \pm 5.0	<u>47.7</u> \pm 6.0	<u>54.1</u> \pm 5.9	<u>30.7</u>
	w/ labels only	<u>3.7</u> \pm 0.5	14.3 \pm 8.3	29.6 \pm 7.0	37.5 \pm 6.1	21.3
	w/ description only	1.0 \pm 0.3	19.8 \pm 8.8	37.5 \pm 7.9	46.2 \pm 5.9	26.1
	FewNERD _{INTRA} (<i>Baseline</i>)	5.8 \pm 0.4	8.9 \pm 4.3	31.4 \pm 9.2	38.4 \pm 7.5	21.1
FewNERD _{INTER}	LITSET	24.3 \pm 0.6	39.8 \pm 2.9	<u>49.1</u> \pm 1.9	<u>52.1</u> \pm 1.9	41.3
	w/ all labels	<u>17.6</u> \pm 2.5	36.1 \pm 4.7	47.2 \pm 3.0	50.4 \pm 2.4	37.8
	w/ labels only	2.9 \pm 0.6	24.7 \pm 1.8	37.9 \pm 1.7	42.4 \pm 2.0	27.2
	w/ description only	16.2 \pm 2.0	37.4 \pm 2.9	47.8 \pm 2.2	50.9 \pm 1.9	38.1
	FewNERD _{INTER} (<i>Baseline</i>)	10.6 \pm 0.8	<u>38.4</u> \pm 3.1	50.4 \pm 3.1	53.3 \pm 2.6	<u>38.2</u>

Table 2: Evaluation of zero- and few-shot tagset extension for in-domain settings. We compare the baseline approach of using in-domain data for label interpretation learning against using LITSET. Despite lacking the in-domain advantage of the baselines, training on LITSET matches or significantly outperforms the in-domain baseline in nearly all settings. Best scores are in bold, and 2nd best is underlined.

layers: coarse labels \mathcal{L}^c (8 classes) and fine labels \mathcal{L}^f (66 classes). \mathcal{L}^f are subclasses of the \mathcal{L}^c such that the entity mentions of both annotations are identical, only their surface form differs. Thus, we can evaluate our dataset against FewNERD in two ways: (1) in the INTRA setting in which we split the labels based on coarse annotations, and (2) in the INTER setting in which we split based on the fine annotations (cf. Figure 4).

We split each dataset into two equally sized label sets for both settings. The random split of labels is repeated three times to reduce the impact of randomness. We then perform few-shot fine-tuning runs with three different seeds for each random split.

Comparison with LITSET. To focus solely on understanding the impact of scaling entity types without the influence of increased entity detection, we downsample LITSET to match the number of entity mentions in each baseline dataset. Further, to make a fair comparison, we remove labels from our approach that match those in the baseline labels \mathcal{L}^{FS} and mask them with the 0-token. However, due to our sampling method, LITSET annotations may not always be consistent. Thus, we can only ensure excluding exact overlaps with the few-shot domain.

4.1.2 Results

The experimental results are shown in Table 2, and we find that LITSET substantially improves the few-shot performance in in-domain settings.

Detecting coarse entity types. When performing label interpretation learning on OntoNotes and FewNERD_{INTRA}, we evaluate the model’s ability to identify entirely new concepts (see INTRA in Figure 4). The results in Table 2 show that our approach can effectively leverage its general label interpretation ability to outperform baselines by large margins. We report +14.8 F1 on average in .pp on FewNERD_{INTRA} and +3.3 F1 on OntoNotes. While LITSET consistently outperforms in-domain label interpretation learning on FewNERD (INTRA), this advantage levels off when $k = 10$ on OntoNotes.

Differentiating fine entity types. In this setting, the model is exposed to sub-classes of a coarse category during label interpretation learning (e.g., “actor” is a subclass of “person”, cf. INTER in Figure 4). We observe that all approaches yield improved few-shot generalization in this setting. This finding suggests that transfer to unseen labels is particularly effective when the training includes annotations of high-level categories. With LITSET, we outperform FewNERD_{INTER} in 0- and 1-shot settings (+13.7 F1 and +1.4 F1 on average in .pp) and remain competitive at higher k-shots.

Evaluation data \mathcal{D}^{FS} for tagset extension from:	Label interpretation learning data \mathcal{D}^{LIT} from:	0-shot	1-shot	5-shot	10-shot	Avg.
JNLPBA	LITSET	<u>41.3</u> \pm 2.0	<u>25.4</u> \pm 5.3	51.3 \pm 3.4	57.7 \pm 3.0	43.9
	w/ all labels	42.2 \pm 1.8	22.5 \pm 8.1	<u>49.9</u> \pm 3.8	<u>55.8</u> \pm 2.7	<u>42.6</u>
	FewNERD _{INTER}	8.2 \pm 1.5	29.5 \pm 15.0	46.0 \pm 7.6	49.7 \pm 6.6	33.4
CLUB	LITSET	<u>6.1</u> \pm 0.9	<u>19.4</u> \pm 3.3	<u>25.9</u> \pm 3.7	<u>33.0</u> \pm 2.1	<u>21.1</u>
	w/ all labels	7.3 \pm 0.1	19.9 \pm 2.0	27.6 \pm 4.6	35.1 \pm 3.1	22.5
	FewNERD _{INTER}	1.7 \pm 0.2	16.9 \pm 1.8	25.5 \pm 4.9	32.2 \pm 3.7	19.1

Table 3: LITSET outperforms FewNERD in out-of-domain settings on JNLPBA (bio-medical domain) and CLUB (chemical domain).

Impact of LITSET sampling. We measure the impact of different heuristics for creating LITSET types. To test this, we conduct various experiments using LITSET with (1) only labels, (2) only descriptions, and (3) all label information available (cf. Figure 3). We first find that using only label annotations decreases performance compared to the baselines (cf. FewNERD_{INTER} and OntoNotes), underlining the need for precise label semantics during label interpretation training to obtain a strong few-shot generalization.

When using only the descriptions or all available annotations, we notice that LITSET yields similar performance to their respective baselines, whereas in the FewNERD_{INTRA} setting, substantial improvements are observed compared to the baselines. Again, this emphasizes that learning from detailed label semantics before the few-shot transfer improves the final performance.

At last, we observe that LITSET substantially outperforms all baselines using our sampling technique, which indicates that alternating shorter labels and expressive short descriptions achieves the best generalization.

4.2 Experiment 2: Cross-Domain Transfer

This experiment assesses the performance of LITSET and its corresponding baselines when not only tagsets but also domains of label interpretation learning and few-shot fine-tuning differ. We reuse LITSET and FewNERD_{INTER} models after label interpretation learning from previous experiment and evaluate on out-of-domain datasets JNLPBA (Collier et al., 2004) (bio-medical domain) and the Chemical Language Understanding Benchmark (CLUB) (Kim et al., 2023) (chemical domain) which labels do represent entirely new, domain-specific concepts.

4.2.1 Results

Table 3 shows the results for cross-domain settings. While this setting is identical for LITSET, the baseline now has no advantage of exposure to "in-domain" data during label interpretation training. Further, no additional masking is required since label spaces between JNLPBA and the baseline model are disjoint. Consequently, we do not mask any labels in LITSET to maintain a fair comparison. However, we emphasize that our model may have been exposed to close domain-specific labels during label interpretation training.

LITSET better transfers to new domains. We find that LITSET significantly outperforms FewNERD with average improvements of +10.5 F1 on JNLPBA and +3.4 F1 on CLUB. Further, on JNLPBA, we observe that our sampling approach performs slightly better than using all label information, whereas we observe the opposite when evaluating CLUB. Our approach consistently outperforms FewNERD on CLUB and JNLPBA with higher shots ($k \geq 5$) and achieves an average increase of +34.0 F1 pp. in zero-shot settings on JNLPBA.

Impact of inconsistent annotations. Furthermore, we observe that LITSET underperforms by -4.1 F1 pp. compared to the baseline in 1-shot settings on JNLPBA. Additionally, its performance is inferior even compared to the 0-shot scenario. This indicates the instability of few-shot fine-tuning with LITSET at very low k . Upon further qualitative analysis of the generated dataset, we discovered that annotations from entity linking benchmarks like ZELDA may not be consistently annotated (cf. Appendix G). This inconsistency could be one possible reason for the observed performance drops. However, as k increases, our approach demonstrates the ability to adapt to the target domain.

Model	Tagset extension on \mathcal{D}^{FS}	Label interpretation learning on \mathcal{D}^{LIT}	1-shot	5-shot	10-shot	Avg.
LEAR	FewNERD _{INTRA}	LITSET	16.6 ± 4.2	33.2 ± 9.2	43.4 ± 10.8	31.1
		FewNERD _{INTRA}	13.5 ± 9.2	23.7 ± 11.7	37.0 ± 14.6	24.7
	FewNERD _{INTER}	LITSET	14.1 ± 2.2	38.3 ± 3.3	44.1 ± 2.6	32.2
		FewNERD _{INTER}	27.6 ± 4.6	50.8 ± 3.5	54.8 ± 2.6	44.4
BINDER	FewNERD _{INTRA}	LITSET	18.8 ± 6.2	31.0 ± 4.2	33.8 ± 3.7	27.9
		FewNERD _{INTRA}	2.6 ± 1.3	11.5 ± 5.6	20.7 ± 7.0	11.6
	FewNERD _{INTER}	LITSET	18.6 ± 1.5	27.3 ± 1.8	30.4 ± 2.0	25.4
		FewNERD _{INTER}	6.1 ± 0.9	20.2 ± 3.2	26.6 ± 3.4	17.6

Table 4: Transfer of LITSET to advanced bi-encoder architectures. We outperform baselines when coarse entity types are not learned during label interpretation training. On BINDER, we also improve over in-domain label interpretation learning.

4.3 Experiment 3: Transfer to Advanced Bi-Encoders

This experiment extends our approach to advanced bi-encoder architectures LEAR (Yang et al., 2021) and BINDER (Zhang et al., 2023). Instead of matrix multiplication, LEAR implements a self-attention layer between the token and label encoder, whereas BINDER uses a contrastive loss. The experimental setup is equal to the one from Section 4.1.

4.3.1 Results

The results are shown in Table 4. We find that LITSET with LEAR improves over the corresponding baseline in INTRA settings up to +9.5 F1 on average in pp. at $k = 5$. Notably, both the baseline and our approach exhibit relatively diminished performance compared to results in Section 4.1. However, our approach falls short in INTER settings, confirming our earlier experimental findings. A noteworthy enhancement is discerned at $k=10$ for the baseline in the INTER-setting, suggesting that existing architectures excel in in-domain transfer, particularly when labels closely align. However, in more practical settings (cross-domain and entirely new type concepts), LITSET works well with LEAR.

Further, we surpass baselines in INTRA and INTER settings across all k -shots for BINDER, indicating LITSET also applies to metric-based methods using contrastive objectives. However, to the best of our knowledge, we are the first to evaluate BINDER in such transfer settings. Our evaluation reveals that the overall performance lags behind simpler architectures. We note that BINDER’s contrastive loss is tailored for learning from extensively annotated corpora. Thus, BINDER may

require modifications or extensions for good generalization performance in these transfer scenarios.

4.4 Experiment 4: Cross-Lingual Transfer

In this experiment, we utilize the multilingual xlm-roberta-base model (Conneau et al., 2020) to assess the transferability of LITSET across languages. We use the English version of OntoNotes as the baseline for label interpretation training. ZELDA is also an English corpus. The transfer is done on the Arabic and Chinese versions of OntoNotes. The results are shown in Table 5.

4.4.1 Results

We find strong improvements across all k -shots on the Arabic and Chinese segments of OntoNotes, namely +3.9 F1 and +9.0 F1 on average in pp., respectively. Despite the overlapping domains between label interpretation learning and few-shot fine-tuning on OntoNotes, our model can discern subtle annotation differences across languages. This emphasizes our model’s robust understanding of labels in multilingual scenarios.

Furthermore, we observe that utilizing xlm-roberta-base also improves LITSET’s performance in monolingual settings (cf. Section 4.1). We reduce the previous performance gap at $k = 10$ from -6.5 F1 to -0.5 F1 on average in pp., thereby increasing the overall performance from +3.3 F1 to +6.5 F1.

5 Related Work

Despite advancements achieved through pre-trained word embeddings (Peters et al., 2018; Akbik et al., 2018; Devlin et al., 2019; Liu et al., 2019; Yamada et al., 2020; Raffel et al., 2020), few-shot

Evaluation data \mathcal{D}^{FS} for tagset extension from:	Label interpretation learning data \mathcal{D}^{LIT} from:	0-shot	1-shot	5-shot	10-shot	Avg.
OntoNotes (EN)	LITSET (EN)	9.9 \pm 3.2	27.4 \pm 8.5	46.4 \pm 6.7	55.5 \pm 6.4	34.8
	OntoNotes (EN)	0.3 \pm 0.1	15.9 \pm 8.4	41.1 \pm 15.0	56.0 \pm 12.7	28.3
Ontonotes (AR)	LITSET (EN)	0.0 \pm 0.0	7.2 \pm 6.1	14.8 \pm 6.3	22.0 \pm 5.8	14.7
	Ontonotes (EN)	0.0 \pm 0.0	4.7 \pm 4.7	12.8 \pm 4.8	14.9 \pm 7.9	10.8
Ontonotes (ZH)	LITSET (EN)	3.0 \pm 0.9	22.7 \pm 8.6	37.6 \pm 5.0	42.8 \pm 5.0	26.5
	Ontonotes (EN)	1.6 \pm 0.3	10.8 \pm 5.9	26.2 \pm 6.9	31.2 \pm 7.9	17.5

Table 5: Tag set extension with baseline pre-finetuning and few-shot fine-tuning in the same domain. LITSET outperforms models that are pre-finetuning on in-domain data when pre-finetuning is done on a small number of labels.

NER focuses explicitly on generalizing to previously unseen label categories by leveraging a small number of labeled examples.

Metric learning (Vinyals et al., 2016; Snell et al., 2017) is a common approach for few-shot NER (Fritzler et al., 2019; Wiseman and Stratos, 2019; Ziyadi et al., 2020) and employs a distance metric to learn a shared representation space and assign labels based on class prototypes (Yang and Katiyar, 2020; Hou et al., 2020; Ma et al., 2022a; Han et al., 2023). Additional components like contrastive loss (Das et al., 2022; Layegh et al., 2023) or meta-learning (de Lichy et al., 2021; Ma et al., 2022c; Wang et al., 2022a) often further improve the performance. Our approach aligns with this research by employing the bi-encoder architecture proposed in Ma et al. (2022a) with an adapted loss calculation. However, prior work did not investigate the impact of the dataset used for label interpretation learning. We do so by increasing the training signal with expressive label verbalizations. Thus, our approach may be applied to all prior work that relies on label verbalizations but may require architectural adaptations to accommodate arbitrary labels.

Template-filling and prompting methods with (large) language models (Lewis et al., 2020; Brown et al., 2020; Raffel et al., 2020; Scao et al., 2023; Touvron et al., 2023) have been widely used for few-shot NER (Cui et al., 2021; Ma et al., 2022b; Lee et al., 2022; Kondragunta et al., 2023; Ma et al., 2023). However, these approaches, relying on masked language model (MLM) objectives, may not be directly comparable to our method due to the scale of our labels. In its basic form, the template-based approach requires one forward pass per label or is limited by the model’s maximum sequence length. Additionally, our approach does

not depend on large language models, which are often unavailable or impractical for few-shot NER.

While specific efforts have been made to adapt to tags in few-shot domains (Hu et al., 2022; Ji et al., 2022), these studies evaluated only a limited number of labels. Our approach shares similarities with (Ren et al., 2022) and Chen et al. (2022), where models were pre-trained using event mentions and entity links, respectively. However, our approach differs significantly. In Ren et al. (2022), the pre-training objective targets the latent typing of entities, whereas our approach focuses on explicitly scaling up entity typing of few-shot NER models. Our distinction from Chen et al. (2022) lies in exploring the effectiveness of distantly supervised training in a genuine few-shot context, wherein classes are not observed during label interpretation training.

6 Conclusion

This paper introduces LITSET, a novel approach for label interpretation training with a large-scale set of entity types. We utilize an entity linking dataset annotated with WikiData information, resulting in a dataset with significantly more distinct labels. We conducted a thorough heuristical, data-based optimization of few-shot NER models using LITSET. Our experiments demonstrate that LITSET consistently outperforms various in-domain, cross-domain, and cross-lingual baselines and is transferable to other architectures and transformer models. For example, we surpass FewNERD by +14.7 F1 on average in pp. and Chinese OntoNotes by +9.0 F1 on average in pp. in low-resource settings. Our method and experiments provide valuable insights into the factors influencing the performance of few-shot NER models utilizing label semantics.

Limitations

Our heuristic data-based optimization is an initial exploration of the impact of scaling the number of distinct entity types during label interpretation learning on few-shot capability. Given our focus on this optimization, we select a commonly used backbone architecture and one entity linking dataset. While we achieved substantial improvements in many settings, it is noteworthy that we did not explore all entity linking benchmarks. Thus, applying our approach with different model architectures and entity disambiguation datasets may yield significantly varied results. Further investigation is necessary to understand how these factors interact comprehensively and to develop more generalized few-shot NER models and comparable evaluation settings.

Additionally, achieving 0-shot capability on completely unseen tags remains challenging, especially in languages different from the one used for label interpretation training. This limitation highlights the need for future research and exploring innovative techniques to enhance the adaptability of few-shot NER models in 0-shot scenarios, enabling them to handle diverse domains and situations effectively.

Lastly, concerning LITSET, our best results were obtained by learning solely from in-batch instances. Although this strategy is commonly employed in machine learning, there is substantial related work on learning from negatives, such as contrastive learning. We believe exploring other architectures and loss functions in more detail, including those from contrastive learning, could further improve our method.

Ethics Statement

In our opinion, this work does not raise many ethical problems. One primary concern is that the texts of entity linking datasets serving our approach show signs of bias. If not checked correctly in advance, the model may learn these biases as exemplarily shown in [Haller et al. \(2023\)](#).

Acknowledgements

We thank all reviewers for their valuable comments. Jonas Golde is supported by the German Federal Ministry of Economic Affairs and Climate Action (BMWK) as part of the project ENA (KK5148001LB0). Felix Hamborg is supported

by the WIN program of the Heidelberg Academy of Sciences and Humanities, financed by the Ministry of Science, Research and Arts of the State of Baden-Württemberg, Germany. Alan Akbik is supported by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Emmy Noether grant “Eidetic Representations of Natural Language” (project number 448414230) and under Germany’s Excellence Strategy “Science of Intelligence” (EXC 2002/1, project number 390523135).

References

- Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. [Contextual string embeddings for sequence labeling](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 1638–1649, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Terra Blevins and Luke Zettlemoyer. 2020. [Moving down the long tail of word sense disambiguation with gloss informed bi-encoders](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1006–1017, Online. Association for Computational Linguistics.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared 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 M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher 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](#).
- Jiawei Chen, Qing Liu, Hongyu Lin, Xianpei Han, and Le Sun. 2022. [Few-shot named entity recognition with self-describing networks](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5711–5722, Dublin, Ireland. Association for Computational Linguistics.
- Yanru Chen, Yanan Zheng, and Zhilin Yang. 2023. [Prompt-based metric learning for few-shot NER](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 7199–7212, Toronto, Canada. Association for Computational Linguistics.
- Nigel Collier, Tomoko Ohta, Yoshimasa Tsuruoka, Yuka Tateisi, and Jin-Dong Kim. 2004. [Introduction to the bio-entity recognition task at JNLPBA](#). In *Proceedings of the International Joint Workshop on Natural Language Processing in Biomedicine and its Applications (NLPBA/BioNLP)*, pages 73–78, Geneva, Switzerland. COLING.

- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Leyang Cui, Yu Wu, Jian Liu, Sen Yang, and Yue Zhang. 2021. [Template-based named entity recognition using BART](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 1835–1845, Online. Association for Computational Linguistics.
- 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.
- 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, Online. Association for Computational Linguistics.
- Leon Derczynski, Eric Nichols, Marieke van Erp, and Nut Limsopatham. 2017. [Results of the WNUT2017 shared task on novel and emerging entity recognition](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 140–147, Copenhagen, Denmark. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- 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.
- Elena V. Epure and Romain Hennequin. 2022. [Probing pre-trained auto-regressive language models for named entity typing and recognition](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 1408–1417, Marseille, France. European Language Resources Association.
- Alexander Fritzer, 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, SAC '19*, page 993–1000, New York, NY, USA. Association for Computing Machinery.
- Kishaloy Halder, Alan Akbik, Josip Krapac, and Roland Vollgraf. 2020. [Task-aware representation of sentences for generic text classification](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3202–3213, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Patrick Haller, Ansar Aynedinov, and Alan Akbik. 2023. [Opiniongpt: Modelling explicit biases in instruction-tuned llms](#).
- Chengcheng Han, Renyu Zhu, Jun Kuang, FengJiao Chen, Xiang Li, Ming Gao, Xuezhai Cao, and Wei Wu. 2023. [Meta-learning triplet network with adaptive margins for few-shot named entity recognition](#).
- Yutai Hou, Wanxiang Che, Yongkui Lai, Zhihan Zhou, Yijia Liu, Han Liu, and Ting Liu. 2020. [Few-shot slot tagging with collapsed dependency transfer and label-enhanced task-adaptive projection network](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1381–1393, Online. Association for Computational Linguistics.
- Jinpeng Hu, He Zhao, Dan Guo, Xiang Wan, and Tsung-Hui Chang. 2022. [A label-aware autoregressive framework for cross-domain NER](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 2222–2232, Seattle, United States. Association for Computational Linguistics.
- Bin Ji, Shasha Li, Shaoduo Gan, Jie Yu, Jun Ma, Huijun Liu, and Jing Yang. 2022. [Few-shot named entity recognition with entity-level prototypical network enhanced by dispersedly distributed prototypes](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1842–1854, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Yunsoo Kim, Hyuk Ko, Jane Lee, Hyun Young Heo, Jinyoung Yang, Sungsoo Lee, and Kyu-hwang Lee. 2023. [Chemical language understanding benchmark](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 5: Industry Track)*, pages 404–411, Toronto, Canada. Association for Computational Linguistics.
- Murali Kondragunta, Olatz Perez-de Viñaspre, and Maite Oronoz. 2023. [Improving and simplifying template-based named entity recognition](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop*, pages 79–86, Dubrovnik, Croatia. Association for Computational Linguistics.

- Amirhossein Layegh, Amir H. Payberah, Ahmet Soylu, Dumitru Roman, and Mihhail Matskin. 2023. [Contrastner: Contrastive-based prompt tuning for few-shot ner](#).
- Dong-Ho Lee, Akshen Kadakia, Kangmin Tan, Mahak Agarwal, Xinyu Feng, Takashi Shibuya, Ryosuke Mitani, Toshiyuki Sekiya, Jay Pujara, and Xiang Ren. 2022. [Good examples make a faster learner: Simple demonstration-based learning for low-resource NER](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2687–2700, Dublin, Ireland. Association for Computational Linguistics.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020. [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.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Jie Ma, Miguel Ballesteros, Srikanth Doss, Rishita Anubhai, Sunil Mallya, Yaser Al-Onaizan, and Dan Roth. 2022a. [Label semantics for few shot named entity recognition](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1956–1971, Dublin, Ireland. Association for Computational Linguistics.
- Ruotian Ma, Xin Zhou, Tao Gui, Yiding Tan, Linyang Li, Qi Zhang, and Xuanjing Huang. 2022b. [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.
- Tingting Ma, Huiqiang Jiang, Qianhui Wu, Tiejun Zhao, and Chin-Yew Lin. 2022c. [Decomposed meta-learning for few-shot named entity recognition](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 1584–1596, Dublin, Ireland. Association for Computational Linguistics.
- Yubo Ma, Yixin Cao, YongChing Hong, and Aixin Sun. 2023. [Large language model is not a good few-shot information extractor, but a good reranker for hard samples!](#)
- Marcel Milich and Alan Akbik. 2023. [ZELDA: A comprehensive benchmark for supervised entity disambiguation](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 2061–2072, Dubrovnik, Croatia. Association for Computational Linguistics.
- Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. [Deep contextualized word representations](#). 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 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Olga Uryupina, and Yuchen Zhang. 2012. [CoNLL-2012 shared task: Modeling multilingual unrestricted coreference in OntoNotes](#). In *Joint Conference on EMNLP and CoNLL - Shared Task*, pages 1–40, Jeju Island, Korea. Association for Computational Linguistics.
- 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.
- Nils Reimers and Iryna Gurevych. 2019. [Sentence-BERT: Sentence embeddings using Siamese BERT-networks](#). 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 3982–3992, Hong Kong, China. Association for Computational Linguistics.
- Liliang Ren, Zixuan Zhang, Han Wang, Clare Voss, ChengXiang Zhai, and Heng Ji. 2022. [Language model pre-training with sparse latent typing](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 1480–1494, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Teven Le Scao, Angela Fan, Christopher Akiki, Elie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, and Matthias Gallé et al. 2023. [Bloom: A 176b-parameter open-access multilingual language model](#).
- Jake Snell, Kevin Swersky, and Richard Zemel. 2017. [Prototypical networks for few-shot learning](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. [Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition](#). In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*, pages 142–147.

- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#).
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Oriol Vinyals, Charles Blundell, Timothy Lillicrap, koray kavukcuoglu, and Daan Wierstra. 2016. [Matching networks for one shot learning](#). In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Denny Vrandečić and Markus Krötzsch. 2014. Wiki-data: a free collaborative knowledgebase. *Communications of the ACM*, 57(10):78–85.
- Peiyi Wang, Runxin Xu, Tianyu Liu, Qingyu Zhou, Yunbo Cao, Baobao Chang, and Zhifang Sui. 2022a. [An enhanced span-based decomposition method for few-shot sequence labeling](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5012–5024, Seattle, United States. Association for Computational Linguistics.
- Zihan Wang, Kewen Zhao, Zilong Wang, and Jingbo Shang. 2022b. [Formulating few-shot fine-tuning towards language model pre-training: A pilot study on named entity recognition](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3186–3199, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Sam Wiseman and Karl Stratos. 2019. [Label-agnostic sequence labeling by copying nearest neighbors](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5363–5369, Florence, Italy. Association for Computational Linguistics.
- Ikuya Yamada, Akari Asai, Hiroyuki Shindo, Hideaki Takeda, and Yuji Matsumoto. 2020. [LUKE: Deep contextualized entity representations with entity-aware self-attention](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6442–6454, Online. Association for Computational Linguistics.
- Pan Yang, Xin Cong, Zhenyu Sun, and Xingwu Liu. 2021. [Enhanced language representation with label knowledge for span extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4623–4635, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yi Yang and Arzoo Katiyar. 2020. [Simple and effective few-shot named entity recognition with structured nearest neighbor learning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6365–6375, Online. Association for Computational Linguistics.
- Sheng Zhang, Hao Cheng, Jianfeng Gao, and Hoifung Poon. 2023. [Optimizing bi-encoder for named entity recognition via contrastive learning](#). In *The Eleventh International Conference on Learning Representations*.
- Morteza Ziyadi, Yuting Sun, Abhishek Goswami, Jade Huang, and Weizhu Chen. 2020. [Example-based named entity recognition](#).

Appendix

A FewNERD Label Semantics in Validation Experiment

Tables 6 to 8 show an overview of the label semantics used in our validation experiment.

Original Label	Adapted Label
O	XO
location-GPE	PH
person-politician	EX
organization-education	CE

Table 6: Extract of random two letter labels for FewNERD.

Original Label	Adapted Label
O	XO
location-GPE	geographical political entity
person-politician	politician
organization-education	education

Table 7: Extract of short labels for FewNERD.

Original Label	Adapted Label
O	XO
location-GPE	geographical entity such as cities, states, countries, and political entities
person-politician	politicians such as presidents, senators, and other government officials
organization-education	education institutions such as schools, colleges, and universities

Table 8: Extract of long labels for FewNERD.

B Validation Experiment with Sparse Latent Typing

We perform our validation experiment on the recently released transformer using the sparse latent typing pre-training objective (Ren et al., 2022). The experimental setup, including few-shot splits, is identical to the one in Section 2. The results are depicted in Figure 5.

Similar to the results in Section 2, we observe a better few-shot generalization with more distinct types and increased expressiveness of label verbalizations. However, the overall performance is

higher using the encoder with sparse latent typing pre-training, a dedicated pre-training objective for keyword extraction from sentences. Further, we observe a slight decrease in performance as soon as $L > 30$. This finding indicates that LitSet is transferable to entity-specific pre-trained models.

C WikiData labels

Given all entity mentions from the entity linking dataset, we source various information from WikiData in natural language and annotate those entities with it. In the following, we present the selected attributes along with their respective definitions, which will serve as our labels:

1. `x instance-of y`: Entity x is a particular example and instance of class y . For example, entity `K2` is an instance of a mountain.
2. `y subclass-of z`: Instance y is a subclass (subset) of class z . For example, instance class `volcano` is a subclass of a mountain.
3. `description`: A short phrase designed to disambiguate items with the same or similar labels.

We note that the `instance-of` and `subclass-of` categories commonly encompass multiple tags rather than being limited to a single tag, as demonstrated in the example in Figure 3. We filter out WikiData-related entities such as information or distribution pages because they do not contain any entity-related information.

D Hyperparameters

This section gives a detailed overview of the hyperparameters used throughout all experiments. For our baselines in experiments Sections 2, 4.1, 4.2 and 4.4 and Appendix B we take the same hyperparameters as in (Ma et al., 2022a) for label interpretation learning. An overview is listed in Table 9.

For LITSET in the respective sections, we use a lower learning rate of $1e^{-6}$, which achieved the lowest validation loss on a 5% hold-out split of LITSET.

For few-shot fine-tuning, we use a slightly higher learning rate of $5e^{-6}$ for LITSET while the learning rate for the baselines remains at $1e^{-5}$. We use a maximum of 100 training epochs with early stopping after 5 iterations with no improvements on the training loss. We do not use any validation splits in few-shot fine-tuning for model selection.



Figure 5: K -shot tagset extension on the 16 least occurring labels of FewNERD using the sparse-latent-typing encoder. We sweep over different numbers of distinct entity types and different semantic descriptions observed during label interpretation learning. We find that increasing both dimensions (more distinct types, extensive label verbalizations) contributes to an improved few-shot generalization.

Argument	Value
Learning rate	$1e^{-5}$
Optimizer	AdamW
Scheduler	Linear warm-up (10%)
Training epochs	3
Training batch size	16
Evaluation batch size	16

Table 9: We use S-BERT (all-mpnet-base-v2) and SLT (sparse latent typing) as the label encoder. LITSET transfers to other transformers and outperforms baselines in INTRA settings while remaining competitive in INTER settings with in-domain trained models.

All previous hyperparameters are identical for LEAR and BINDER (cf. Section 4.3), except that we use the recommended learning rate of $3e^{-5}$ for BINDER and early stopping for label interpretation learning (after one epoch with no improvements on the training loss).

E Using Different Transformers as Label Encoder

In this experiment, we investigate whether the all-mpnet-base-v2 sentence trans-

former (Reimers and Gurevych, 2019) and the sparse-latent-typing transformer (Ren et al., 2022) can effectively help to understand label semantics better. Sentence transformers have been trained on a similarity objective, making them intriguing for our model to act as an enhanced label encoder. Sparse latent typing is a pre-training objective designed for extracting keywords from sentences. We present results in Table 10.

We observe that using all-mpnet-base-v2 performs generally worse than plain bert-base-uncased. However, we also observe that using LITSET yields better few-shot generalization in both INTRA and INTER settings and thus confirms that our main findings are transferable to other label encoders. When using SLT encoder, we outperform the baseline by large margins in the INTRA settings but fall slightly short in INTER settings.

F The Impact of Negative Examples

In this experiment, we investigate the impact of integrating negative labels \mathcal{L}^- in each batch. To do so, we additionally sample negative labels from $\mathcal{L} \setminus \mathcal{L}_b$ until the desired number of labels is reached

Transformer	Tagset extension on \mathcal{D}^{FS}	Label interpretation learning on \mathcal{D}^{LIT}	1-shot	5-shot	10-shot	Average
S-BERT	FewNERD _{INTRA}	LITSET	27.6 ± 4.1	49.2 ± 3.4	54.7 ± 4.8	43.8
		FewNERD _{INTRA}	10.7 ± 7.4	37.8 ± 9.8	49.1 ± 8.4	32.5
	FewNERD _{INTER}	LITSET	36.6 ± 2.0	44.3 ± 2.0	47.7 ± 2.1	42.9
		FewNERD _{INTER}	23.4 ± 2.4	42.3 ± 3.8	48.5 ± 3.1	38.1
SLT	FewNERD _{INTRA}	LITSET	27.2 ± 5.8	51.8 ± 4.9	57.2 ± 5.4	45.4
		FewNERD _{INTRA}	6.2 ± 4.9	15.6 ± 4.7	21.9 ± 4.9	14.6
	FewNERD _{INTER}	LITSET	38.6 ± 3.6	49.4 ± 2.5	52.4 ± 2.3	46.8
		FewNERD _{INTER}	40.3 ± 4.1	52.0 ± 3.0	54.9 ± 2.24	49.1

Table 10: We use S-BERT (all-mpnet-base-v2) and SLT (sparse latent typing) as the label encoder. LITSET transfers to other transformers and outperforms baselines in INTRA settings while remaining competitive in INTER settings with in-domain trained models.

Evaluation data \mathcal{D}^{FS} for tagset extension from:	Label interpretation learning data \mathcal{D}^{LIT} from: (/w # max. negative labels per batch)	1-shot	5-shot	10-shot	Average
FewNERD _{INTRA}	LITSET (0)	20.1 ± 5.0	47.7 ± 6.0	54.1 ± 5.9	40.6
	LITSET (64)	20.1 ± 4.8	47.5 ± 5.0	53.2 ± 6.6	40.3
	LITSET (128)	18.9 ± 4.9	46.4 ± 3.9	52.7 ± 5.9	39.3
FewNERD _{INTER}	LITSET (0)	36.1 ± 4.7	47.2 ± 3.0	50.4 ± 2.4	44.6
	LITSET (64)	35.2 ± 4.1	47.4 ± 2.6	50.5 ± 2.4	44.4
	LITSET (128)	34.7 ± 3.3	47.3 ± 2.7	50.4 ± 2.3	44.1

Table 11: The few-shot generalization of LITSET does not improve with a fixed number of labels per batch (we sample additional labels for loss calculation until, e.g., 64 labels are present). We find that the best training setup only uses the labels in the current batch.

and include them for loss calculation. Including negative types could potentially lead to a better generalization in few-shot settings due to the increased signal during loss calculation. We show results in Table 11. We observe that including more labels in each batch harms the performance. While prior work (Epure and Hennequin, 2022; Wang et al., 2022b) has shown that this idea is beneficial in few-shot settings, we find that LITSET works best when only using the labels present in the batch for loss calculation. Since we randomly sample additional labels, it is possible, if not likely, to sample similar labels that are not true negatives and thus not advantageous when using cross-entropy loss.

G Annotation Noise in ZELDA

In some cases, ZELDA is not consistently annotated, which may affect the few-shot fine-tuning performance for settings with very low k . Table 12 shows such an example. We find unique entities, such as proteins, that are not consistently annotated to verify this assumption qualitatively. These in-

consistencies may cause a worse entity detection ability with LITSET than training on consistently annotated datasets. While we show that entity linking benchmarks can be used to obtain a strong label understanding prior, improving the annotation quality or generating a designated label interpretation training dataset remains for future work.

Annotation noise in ZELDA	
annotated	[...] which in turn creates the compound oxyhemoglobin protein .
missing annotation	[...] whereas in oxyhemoglobin O it is a high spin complex.
annotated	GSTK1 promotes adiponectin protein multimerization
missing annotation	[...] ER stress induced adiponectin O downregulation [...]

Table 12: Annotations in the entity linking benchmark may be inconsistent, causing the 1-shot drops on JNLPBA. Since JNLPBA is annotated by humans, it is expected that all sentences are annotated consistently.