

A Comparative Study of Pre-trained Encoders for Low-Resource Named Entity Recognition

Anonymous ACL submission

Abstract

Pre-trained language models (PLM) are effective components of few-shot named entity recognition (NER) approaches when augmented with continued pre-training on task-specific out-of-domain data or fine-tuning on in-domain data. However, their performance in low-resource scenarios, where such data is not available, remains an open question. We introduce an encoder evaluation framework, and use it to systematically compare the performance of state-of-the-art pre-trained representations on the task of low-resource NER. We analyze a wide range of encoders pre-trained with different strategies, model architectures, intermediate-task fine-tuning, and contrastive learning. Our experimental results across ten benchmark NER datasets in English and German show that encoder performance varies significantly, suggesting that the choice of encoder for a specific low-resource scenario needs to be carefully evaluated.

1 Introduction

Pre-trained language models (PLM) have been shown to be very effective few-shot learners for a wide range of natural language processing tasks (Brown et al., 2020; Gao et al., 2021), as they capture semantically and syntactically rich representations of text via self-supervised training on large-scale unlabeled datasets (Peters et al., 2018; Devlin et al., 2019). Recent research in few-shot named entity recognition (NER) has leveraged such representations, e.g. for metric learning on task-specific out-of-domain¹ data (Fritzler et al., 2019; Yang and Katiyar, 2020), optionally augmented by continued pre-training with distantly supervised, in-domain data (Huang et al., 2020). However, there has been no systematic comparison of the NER performance of such representations in low-resource scenarios without task-specific out-of-domain data

¹Out-of-domain and in-domain refer to NER-specific data with disjoint label spaces, i.e. $\mathcal{Y}_{out} \neq \mathcal{Y}_{in}$.

and very limited in-domain data; a prevalent setting in many practical applications.

In this paper we conduct a comparative study to answer the following research questions: How well do representations learnt by different pre-trained models encode information that benefits these low-resource scenarios? What can we observe for different categories of encoders, such as encoders trained with masked language modeling, versus encoders that are additionally fine-tuned on downstream tasks, or optimized with contrastive learning? How do they perform across different datasets and languages? We present an evaluation framework inspired by few-shot learning to evaluate representations obtained via different pre-training strategies, model architectures, pre-training data, and intermediate-task fine-tuning in low-resource NER scenarios of varying difficulty (see Figure 1).

We find that the choice of encoder can have significant effects on low-resource NER performance, with F1 scores differing by up to 25% between encoders, and simply picking an encoder of the BERT family at random will usually not yield the best results for a given scenario. We observe that while BERT in general performs adequately, ALBERT and RoBERTa outperform BERT by a large margin in many cases, with ALBERT being especially strong in very low-resource settings with only one available labeled example per class.

In summary, the main contributions of this study are: (1) a systematic performance evaluation of a wide range of encoders pre-trained with different strategies, such as masked language modeling, task-specific fine-tuning, and contrastive learning on the task of low-resource named entity recognition; (2) an evaluation on ten benchmark NER datasets in two languages, English and German; (3) an encoder-readout evaluation framework that can be easily extended with additional scenarios, encoders, datasets, and readout approaches, and that we make available as open source at [anonymized-url](#).

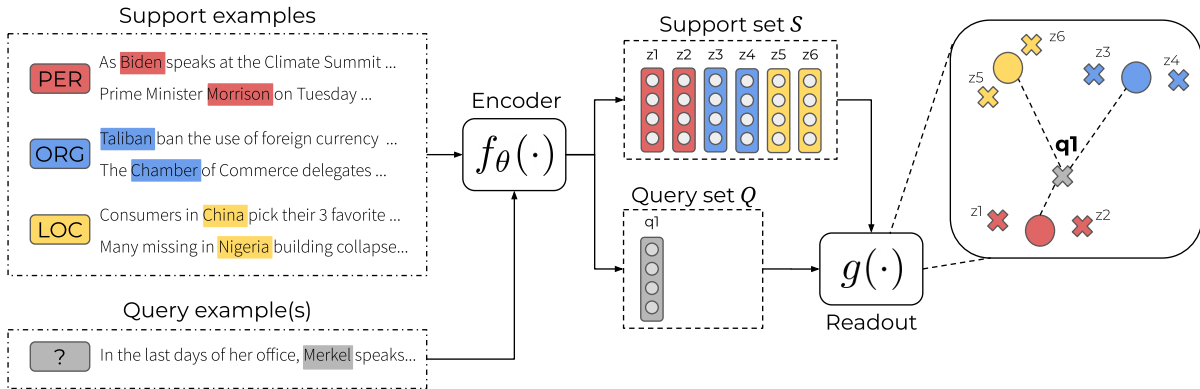


Figure 1: Encoder-readout evaluation framework. For each of the N classes, we randomly sample K support tokens including their sentence context, and an unlabeled query token with sentential context. The encoder $f_{\theta}(\cdot)$ provides an embedding (or representation) for each token, and the readout module $g(\cdot)$ assigns a class to a query token by comparing its representation q_j to the representations $\{z_1, \dots, z_{N \times K}\}$ of the support tokens. Depending on the readout approach, the c -th class in \mathcal{S} is represented either by its prototype embedding (as shown in the example) or by its set of associated token embeddings, e.g. for nearest neighbor classification. In this example q_1 representing *Merkel* would be assigned the class *PER* based on the closest class prototype embedding (red circle).

2 Encoder Evaluation Framework

To simulate low-resource NER scenarios of varying difficulty, we draw inspiration from the evaluation of few-shot learning methods. We first give a formal definition of the few-shot NER task, and then introduce the encoder evaluation framework itself.

2.1 Few-shot NER task definition

NER is typically formulated as a sequence labeling problem, where the input is a sequence of tokens $\mathbf{X} = \{x_1, x_2, \dots, x_T\}$ and the output is the corresponding T -length sequence of entity type labels $\mathbf{Y} = \{y_1, y_2, \dots, y_T\}$. In contrast, few-shot learning is cast as an episodic N -way K -shot problem, where in each episode, N classes are sampled with K examples each to construct a support set $\mathcal{S} = \{\mathbf{X}_i, \mathbf{Y}_i\}_{i=1}^{N \times K}$ for learning, and K' examples per class are sampled to create a query set $\mathcal{Q} = \{\mathbf{X}_j, \mathbf{Y}_j\}_{j=1}^{N \times K'}$ for evaluation ($\mathcal{S} \cap \mathcal{Q} = \emptyset$). In a sequence labeling problem like NER, samples are typically sentences, due to the importance of contextual information for token classification, but care has to be taken to ensure that the sampled sentences contain no other entities. In particular, there should be no entity overlap between the support and the query sets (Ding et al., 2021).

2.2 Encoder-Readout Framework

Our framework consists of two modules, an encoder $f(\cdot)$ and a readout module $g(\cdot)$, as shown in Figure 1. The encoder provides an embedding $z = f_{\theta}(x)$ of a token x , where θ denotes the pa-

rameters of the encoder. The readout module is responsible for assigning a class to each token x' in the query set \mathcal{Q} given the support set \mathcal{S} . Depending on the readout approach, the c -th class in \mathcal{S} is represented either by its prototype embedding or by its associated set of token embeddings, e.g. for nearest neighbor classification. The decision is made by comparing the embedding $q = f_{\theta}(x')$ with each of the N class prototypes built from the support set \mathcal{S} , or with each of the token-level embeddings.

3 Experiments

We illustrate the evaluation framework using a representative set of encoders pre-trained with different strategies. We then give details of the readout approaches, the datasets we used, and all other experimental settings.

3.1 Encoders

We group encoders into four categories, depending on their type of pre-training:

PLM These models are pre-trained on a large general corpus in a self-supervised manner without any task-specific fine-tuning. We consider six representative encoders for English: BERT cased and uncased (Devlin et al., 2019), SpanBERT (Joshi et al., 2020), XLNet (Yang et al., 2019), ALBERT (Lan et al., 2020) and RoBERTa (Liu et al., 2019), and three encoders for German: deepset’s BERT, GottBERT (Scheible et al., 2020) and XLM-RoBERTa (Conneau et al., 2020).²

²HuggingFace model identifiers for these and all other

Language	Dataset	Domain	# Entity types	Entity tag set
English	CoNLL-2003 _{EN}	News	4	LOC, MISC, ORG, PER
	OntoNotes 5.0	News, Dialogue	18	CARDINAL, DATE, EVENT, MONEY, ...
	Few-NERD _{coarse}	General	8	art, building, event, product, ...
	Few-NERD _{fine}	General	66	art-film, product-car, other-law, ...
	WNUT-17	Social Media	6	corporation, creative-work, group, ...
	WikiAnn	General	3	LOC, ORG, PER
	WikiGold	General	4	LOC, MISC, ORG, PER
Zhang et al.	e-Commerce	4	ATTRIBUTE, BRAND, COMPONENT, PRODUCT	
German	CoNLL-2003 _{DE}	News	4	LOC, MISC, ORG, PER
	GermEval 2014	General	12	LOC, LOCderiv, LOCpart, ORG, ...
	Smartdata	News, General	16	DISASTER-TYPE, DISTANCE, LOCATION, ...

Table 1: Statistics of the evaluated datasets

Fine-tuned PLM Recent research has shown that intermediate-task training can result in significant performance gains on the target task even in low-resource settings (Vu et al., 2020; Poth et al., 2021). We evaluate three BERT encoders that are fine-tuned on token-level, sentence-level, and document-level intermediate tasks, respectively: BERT_{POS} for part-of-speech tagging, BERT_{MNLI}, fine-tuned on the MultiNLI dataset (Williams et al., 2018), and BERT_{SQuAD} for extractive question answering (Rajpurkar et al., 2016). Evaluating these encoders may allow us to observe whether the representation granularity induced by the tasks they were fine-tuned on has an effect on NER performance: While token-level part-of-speech tag information is a staple feature of classic NER approaches (Finkel et al., 2005), it is less clear if encoders trained on tasks that require conceptual representations (and possibly understanding) of sentence- and document-length context, learn entity representations useful for NER.

PLM fine-tuned on NER We also experiment with BERT_{CoNLL}, a BERT model fine-tuned on the CoNLL-2003 NER dataset. As this model’s hidden representations have been adapted to NER, we expect it to exhibit better performance than the other representations. The most interesting question of using this model is whether its representations transfer to NER datasets with non-CoNLL tagsets.

PLM with contrastive learning For each of the English PLM encoders, we apply contrastive learning to learn representations with better separability. The idea of contrastive learning is to pull positives closer and push negatives away in the representation space during the pre-training phase (Rethmeier and Augenstein, 2021). We use the loss function

models are listed in Appendix A.

proposed by Chopra et al. (2005):

$$\mathcal{L}_{CL}(x_i, x_j; \theta) := \mathbb{1}_{y_i=y_j} \cdot \|f_{\theta}(x_i) - f_{\theta}(x_j)\| + \mathbb{1}_{y_i \neq y_j} \cdot \max(0, \epsilon - \|f_{\theta}(x_i) - f_{\theta}(x_j)\|).$$

To guarantee that this label-aware contrastive learning conforms to the few-shot setting, we construct positive/negative pairs from the support set: Given an N -way K -shot support set, for each of the N classes we construct 1 positive pair and K negative pairs.³

3.2 Readout approaches

We analyze three variants for the readout approach:⁴ (1) **Logistic Regression (LR)**, a linear classification algorithm that can be extended to multinomial logistic regression to deal with multi-class (N -way) settings, such as the one discussed here. (2) **k-Nearest Neighbor (NN)**, a non-parametric classification method adopted in metric space. As proposed in STRUCTSHOT (Yang and Katiyar, 2020), we set $k = 1$ to find the exact nearest token in the support set. (3) **Nearest Centroid (NC)** works similar to NN, but instead of computing the distance between the query and every instance in the embedding space, we represent each class by the centroid of all token embeddings belonging to this class, and assign the query to the class with the nearest centroid.

3.3 Datasets

In order to provide a comprehensive evaluation, we evaluate all encoders on a range of

³One extra example per class is needed for $K = 1$ to build one positive pair for this class. This extra example is involved only in the contrastive learning phase and not introduced to the encoding and readout steps.

⁴Computational details of the readout approaches can be found in Appendix B.

205 datasets covering different languages and domains, 253
206 including seven English benchmarks: CoNLL- 254
207 2003 (Tjong Kim Sang and De Meulder, 2003), 255
208 Few-NERD (Ding et al., 2021), OntoNotes 256
209 5.0 (Weischedel et al., 2013), WikiAnn (Pan 257
210 et al., 2017), WNUT-17 (Derczynski et al., 2017), 258
211 WikiGold (Balasuriya et al., 2009), and the dataset 259
212 of Zhang et al. (2020). For German, we se-
213 lected the following three datasets: CoNLL-
214 2003 (Tjong Kim Sang and De Meulder, 2003),
215 Smartdata (Schiersch et al., 2018) and GermEval
216 2014 (Benikova et al., 2014). Table 1 lists the do-
217 mains and tagset details of each dataset.

218 3.4 Experimental settings / Hyperparameters

219 **Datasets** We use the BIO tagging schema by de-
220 fault and the IO schema only when BIO is not pro-
221 vided by the original dataset (in case of Few-NERD,
222 OntoNotes 5.0 and WikiGold). WikiGold and
223 the dataset of Zhang et al. (2020) do not provide
224 train/test splits, we therefore use the full dataset
225 to sample support and query sets. For all other
226 datasets, test splits are used for sampling.⁵

227 **General settings** For each dataset, we evaluate
228 our methods under three few-shot scenarios: 5-way
229 1-shot, 5-way 5-shot and 5-way 10-shot. To pro-
230 duce accurate performance estimates, we sample
231 600 episodes for each scenario and report the mean
232 token-level micro-F1 score over all episodes, av-
233 eraged over all positive classes, and excluding the
234 'O' class.

235 **Encoders** Max-length is fixed at 128. We use
236 randomly initialized, static embeddings as the base-
237 line encoder (*Random*). For contrastive learning,
238 we use the Adam optimizer and set the learning
239 rate to be 5×10^{-5} and the number of epochs to be
240 1 across all encoders.

241 **Readout approaches** We L2-normalize the en-
242 coder embeddings before feeding them to the read-
243 out model. For NN and NC classification, Eu-
244 clidean distance serves as the similarity metric be-
245 tween tokens. For LR, an L2-penalty is applied to
246 the coefficients. All reported results use LR as the
247 default readout method, unless specified otherwise,
248 as we found LR to perform best on average (see
249 Section 4.4).

250 **Framework implementation** We implement
251 our low-resource NER encoder evaluation frame-
252 work using the HuggingFace Transformers li-

brary (Wolf et al., 2020), Hydra (Yadan, 2019),
and PyTorch (Paszke et al., 2019). Additional sce-
narios, encoders, and datasets can be easily added
simply by creating new experiment configurations.
Adding new readout methods is also a simple mat-
ter of a few lines of code. We make our code base
available as open source at [anonymized-url](#).

260 4 Results and Discussion

261 4.1 Comparison of PLM encoders

262 We first analyze PLM encoders which have not
263 been fine-tuned on any task.

264 **English results** Table 2 presents the experimen-
265 tal results of English-language encoders for differ-
266 ent scenarios and datasets. For all scenarios and
267 datasets, the PLM encoders outperform the ran-
268 domly initialized baseline by a large margin. As ex-
269 pected, the NER classification performance of the
270 encoders increases with higher K , i.e. with more
271 instances per class in the support set. Overall, the
272 level of performance across various datasets of this
273 encoder-only approach to low-resource NER is sur-
274 prisingly good: We observe that ALBERT achieves
275 a token-level F1 score of $F1 = 72.8$ on CoNLL-
276 2003, XLNet a score of $F1 = 85.7$ on Few-NERD
277 fine-grained, and RoBERTa a score of $F1 = 83.8$
278 on OntoNotes 5.0. While these results are not di-
279 rectly comparable to those of state-of-the-art, fully
280 supervised approaches due to the differences in the
281 evaluation setup, they are achieved essentially fine-
282 tuning-free, and with much fewer labeled instances
283 per class.

284 **Encoder analysis** The best-performing en-
285 coders, on average and across datasets, are AL-
286 BERT, RoBERTa, and BERT. ALBERT is by far
287 the best encoder for $K = 1$, but the other encoders
288 achieve comparable performance or outperform
289 ALBERT for $K \geq 5$. Even though ALBERT is
290 an order of magnitude smaller in terms of its num-
291 ber of parameters than either BERT or RoBERTa,
292 it provides very competitive embeddings in our
293 evaluation setup. As can be expected, BERT_{cased}
294 consistently outperforms BERT_{uncased} for datasets
295 with tag sets where casing provides useful informa-
296 tion for NER (e.g. CoNLL, WikiGold), but does not
297 necessarily perform better if the tag set contains en-
298 tity types whose instances use lower-case spelling.
299 XLNet achieves mixed results, mainly depending
300 on the dataset – on CoNLL-2003, WikiAnn and
301 WNUT-17, its F1 scores are significantly lower for
302 all scenarios than those of the best encoder, while

⁵For Few-NERD, we use the test data from the "super-
vised" split.

Dataset	K	Random	BERT↓	BERT↑	ALBERT↓	RoBERTa↑	SpanBERT↑	XLNet↑
CoNLL-2003 _{EN}	1	9.52	21.96	22.04	33.03 †	21.71	18.39	18.49
	5	12.53	60.94	62.17	68.33 †	64.49	43.22	44.82
	10	13.71	66.11	68.79	72.76	72.09	49.79	52.43
OntoNotes 5.0	1	18.66	42.71	45.09	50.45 †	42.74	34.30	38.40
	5	19.73	74.68	77.70	77.66	78.70	65.64	72.60
	10	18.88	80.92	82.70	82.10	83.80 †	74.14	78.38
Few-NERD _{coarse}	1	12.12	25.99	28.52	35.67 †	28.12	23.34	25.93
	5	15.59	53.85	56.04	59.14	58.66	45.50	52.32
	10	16.04	59.44	63.20	63.30	65.52 †	52.65	61.94
Few-NERD _{fine}	1	21.14	49.74	48.50	54.27 †	51.27	39.13	47.02
	5	21.00	80.12	79.26	78.08	81.70	71.93	82.73
	10	20.62	84.07	83.21	81.17	84.95	78.39	85.73
WNUT-17	1	18.86	25.71	25.67	28.47 †	25.43	23.14	24.36
	5	19.11	51.56	50.58	55.12	54.59	42.29	42.26
	10	18.52	58.77	60.37	60.41	63.93 †	48.84	49.74
WikiAnn	1	12.07	24.53	25.92	32.63 †	24.80	22.67	22.06
	5	15.64	48.33	52.29	53.11 †	51.34	40.60	36.81
	10	16.95	54.84	59.48	59.10	60.83	46.44	44.19
WikiGold	1	3.71	18.40	21.30	32.30 †	20.63	14.90	18.01
	5	10.02	49.19	55.54	55.87	56.08	41.07	45.44
	10	11.62	55.85	63.91	61.23	64.84	48.09	53.85
Zhang et al.	1	13.49	37.39	36.82	41.23 †	38.79	25.83	31.25
	5	17.08	63.19	62.17	62.73	66.44 †	49.08	57.69
	10	16.21	67.45	67.09	66.61	70.16 †	54.80	63.79

Table 2: Token-level micro-F1 scores of PLM encoders and a random baseline for 5-way K -shot scenarios, with logistic regression readout. † denotes scores with significant difference to the next-best encoder’s score ($\alpha = 0.05$). † and ‡ indicate cased and uncased models.

Dataset	K	Random	BERT↑	Gott-BERT↑	XLNet↑
CoNLL-2003 _{DE}	1	12.53	29.42	26.27	30.65
	5	15.38	65.98	58.37	65.22
	10	16.00	71.43	64.77	71.18
GermEval 2014	1	17.52	25.89	24.08	27.24
	5	20.70	61.79 †	54.06	58.51
	10	18.33	71.18 †	60.30	65.37
Smartdata	1	26.12	51.12	49.96	53.17
	5	23.52	82.50 †	79.30	80.89
	10	21.55	86.01	83.10	85.66

Table 3: Token-level micro-F1 scores of German PLM encoders and a random baseline under 5-way K -shot scenarios, with logistic regression readout. † denotes scores with a significant difference to the next-best encoder’s score ($\alpha = 0.05$). † indicates cased models.

on Few-NERD fine-grained, XLNet achieves the best score of all encoders. SpanBERT on average shows the worst performance of all encoders, with F1 scores in most scenarios several percentage points lower than even those of XLNet. This suggests that SpanBERT’s span-level masking and training with a span boundary objective produce token-level embeddings that are less well separable

by the logistic regression classifier.

Dataset analysis On a per-dataset basis, we can observe the following from Table 2: On CoNLL-2003, ALBERT outperforms the next-best encoder BERT_{cased} for $K = 1$ by 11% F1, and achieves a best score of $F1 = 72.8$ for $K = 10$, closely followed by RoBERTa. XLNet’s and SpanBERT’s F1 scores are more than 20% lower than those of ALBERT for $K = 5$ and $K = 10$. On Few-NERD with coarse labels, ALBERT is again the best encoder at $K = 1$. For $K = 10$, RoBERTa achieves $F1 = 65.5$, but the other encoders except for SpanBERT perform almost as well. Using the fine-grained labels of Few-NERD, all encoders achieve around 80% F1 score. The overall picture is similar for OntoNotes 5.0 and the dataset of Zhang et al., with ALBERT being the best encoder at $K = 1$ and RoBERTa outperforming the other encoders at $K = 10$. BERT and XLNet show competitive performance to ALBERT and RoBERTa, yielding slightly lower F1 scores in all scenarios. This trend is also confirmed for the remaining datasets, WikiAnn, WNUT-17 and WikiGold, with ALBERT and RoBERTa being the

Dataset	K	BERT \downarrow	B _{POS} \downarrow	B _{MNLI} \downarrow	B _{SQuAD} \downarrow
CoNLL-2003 _{EN}	1	21.96	43.01 \dagger	22.29	35.05
	5	60.94	65.72	61.34	65.94
	10	66.11	68.46	64.71	68.50
OntoNotes 5.0	1	42.71	50.85 \dagger	42.99	47.83
	5	74.68	66.17	75.29	76.37
	10	80.92	68.02	80.94	79.68
Few-NERD _{coarse}	1	25.99	34.70	26.08	35.07
	5	53.85	49.88	52.52	59.77 \dagger
	10	59.44	52.78	58.17	63.09 \dagger
Few-NERD _{fine}	1	49.74	43.97	46.71	51.17
	5	80.12 \dagger	63.08	77.14	78.58
	10	84.07 \dagger	66.43	81.26	81.58
WNUT-17	1	25.71	32.04 \dagger	25.12	29.04
	5	51.56	44.90	48.50	51.05
	10	58.77 \dagger	49.11	56.30	54.58
WikiAnn	1	24.53	32.92	23.35	33.33
	5	48.33	43.54	46.94	55.93 \dagger
	10	54.84	45.70	53.47	63.37 \dagger
WikiGold	1	18.40	37.46 \dagger	20.33	30.80
	5	49.19	55.54 \dagger	50.86	53.96
	10	55.85	55.62	55.81	57.99 \dagger
Zhang et al.	1	37.39	45.67 \dagger	37.29	40.90
	5	63.19	59.58	62.98	61.01
	10	67.45	60.61	66.23	61.95

(a) Micro-F1 scores of BERT, and fine-tuned BERT_{POS}, BERT_{MNLI} and BERT_{SQuAD}.

Dataset	Overlap	K	BERT \downarrow	B _{CoNLL} \downarrow
CoNLL-2003 _{EN}	1.00	1	21.96	90.46 \dagger
		5	60.94	94.73 \dagger
		10	66.11	94.40 \dagger
WikiGold	1.00	1	18.40	68.83 \dagger
		5	49.19	81.40 \dagger
		10	55.85	84.68 \dagger
WikiAnn	0.75	1	24.53	55.15 \dagger
		5	48.33	67.22 \dagger
		10	54.84	71.34 \dagger
Few-NERD _{coarse}	0.50	1	25.99	53.25 \dagger
		5	53.85	70.04 \dagger
		10	59.44	72.66 \dagger
WNUT-17	0.25	1	25.71	44.96 \dagger
		5	51.56	63.99 \dagger
		10	58.77	69.76 \dagger
OntoNotes 5.0	0.16	1	42.71	58.99 \dagger
		5	74.68	76.21 \dagger
		10	80.92 \dagger	77.75
Few-NERD _{fine}	0	1	49.74	59.36 \dagger
		5	80.12	79.70
		10	84.07 \dagger	82.00
Zhang et al.	0	1	37.39	49.22 \dagger
		5	63.19	65.40 \dagger
		10	67.45	66.13

(b) Micro-F1 scores of BERT and BERT_{CoNLL}. The datasets are listed in descending order of tag set overlap with CoNLL-2003, as measured by Jaccard Index.

Table 4: Token-level micro-F1 scores of fine-tuned encoders under 5-way K -shot scenarios, with LR readout. \dagger denotes scores with significant difference to the next-best encoder’s score ($\alpha = 0.05$). \downarrow indicates uncased models.

strongest contenders, and BERT often catching up in terms of F1 scores with increasing K .

German results Table 3 shows the results of German-language encoders and the random baseline on three evaluation datasets. Similar to the English results, we observe that: (i) BERT, GottBERT and XLM-RoBERTa all benefit from more support instances, and outperform the random baseline by a large margin. (ii) XLM-RoBERTa shows the best performance across datasets in one-shot settings, whereas BERT outperforms the other encoders for $K \geq 5$. (iii) GottBERT’s encodings yield features that are less useful for low-resource NER, resulting in worse performance than the other two encoders in all scenarios.

On CoNLL-2003, BERT achieves a micro-F1 score of 71.4 at $K = 10$, XLM-R a competitive score of 71.2, while GottBERT only achieves $F1 = 64.8$. Similar performance differences between the three encoders can be observed for the other two datasets at $K = 5$ and $K = 10$. At $K = 1$, XLM-R consistently outperforms BERT

and GottBert, with GottBERT showing the worst performance. The results show that BERT, a model trained with less, but likely quality training data (Wikipedia, OpenLegalData, News) produces representations that are more suited for low-resource NER in most of the evaluated settings, compared to GottBERT (145GB of unfiltered web text), and XLM-RoBERTa (≈ 100 GB filtered CommonCrawl data for German).

4.2 Fine-tuned encoders

Fine-tuned PLM The next group of encoders we analyze are encoders fine-tuned on an intermediate task, in our case POS tagging, NLI, and QA. Results are shown in Table 4a. We can see that using a BERT encoder fine-tuned on POS tagging significantly improves F1 scores at $K = 1$ for all datasets except Few-NERD fine-grained, on average by about 9 points. However, for $K \geq 5$, BERT_{POS}’s performance is significantly worse than that of BERT for the majority of datasets, except CoNLL-2003 and WikiGold.

The BERT_{MNLI} model’s performance is compet-

itive with the base BERT model’s, with no statistically significant differences. Fine-tuning on this sentence-level task, which is rather unrelated to NER, hence seems to have neither negative nor positive effects on the resulting token embeddings.

Embeddings obtained from BERT_{SQuAD}, fine-tuned on document-level span extraction, outperform BERT in most settings, often with statistical significance. However, on some datasets (e.g. WNUT-17, Few-NERD_{fine}), BERT_{SQuAD}’s scores are lower than BERT’s for $K \geq 5$. Compared to the other fine-tuned encoders, BERT_{SQuAD} performs better in general for $K \geq 5$. Its good performance may be attributed to the fact that approximately 41.5% of the answers in the SQuAD dataset correspond to common entity types, and another 31.8% to common noun phrases (Rajpurkar et al., 2016).

The observations for these three encoders coincide with the intuition, that the more relevant the knowledge encoded by the intermediate task is w.r.t. the target task, the more likely an improvement on the target task becomes.

PLM fine-tuned on NER Table 4b shows the results obtained for BERT_{CoNLL}, an encoder that was fine-tuned on CoNLL-2003. As can be expected, this encoder performs very well on the CoNLL-2003 test set, with large F1 gains in all scenarios. For most of the other datasets, F1 scores are also significantly improved for all settings of K , especially with a large tagset overlap. These results coincide with the intuition that the higher the tagset overlap, the larger the improvement. However, we note that some of these datasets are constructed from other data sources, e.g. web and social media texts, which indicates some transferability of the CoNLL-2003-tuned representations. Even for datasets where there is little or no overlap (OntoNotes 5.0, Zhang et al.), there are at least some gains at $K = 1$. However, at $K = 10$, the performance of the embeddings obtained from BERT_{CoNLL} is significantly worse than that of the base BERT model.

4.3 PLM with contrastive learning

Table 5 compares the results of English encoders before and after contrastive learning. In general, results are mixed: For ALBERT and SpanBERT, using CL improves F1 scores in most cases, often with significant differences, whereas for BERT, RoBERTa and XLNET, the base encoders mostly exhibit (marginally) better performance.

Encoder analysis We observe that ALBERT benefits the most from contrastive learning, with significant F1 gains in 5 out of 12 comparisons, followed by SpanBERT (3), XLNet (1), BERT (1) and RoBERTa (0). Surprisingly, it achieves slightly higher F1-scores on Few-NERD coarse-grained and significantly higher F1-scores on WikiGold in all three scenarios. For 1-shot scenario on CoNLL-2003, ALBERT also gets a large F1 increase by 3.68%, the best improvement among all encoders.

Dataset analysis Few-NERD coarse-grained and WikiGold show better compatibility with contrastive learning, with 11 and 8 F1 improvements out of 15 comparisons after contrastive learning, respectively, compared with CoNLL-2003 (6) and OntoNotes 5.0 (4). Specifically, all five encoders have F1 gains on Few-NERD dataset in the one-shot scenario.

4.4 Readout approaches

Finally, Table 6 compares the different readout approaches on the CoNLL-2003 and OntoNotes 5.0 datasets, using ALBERT. For $K \geq 5$, Logistic Regression outperforms Nearest Centroid and Nearest Neighbor classification, while for one-shot scenarios Nearest Neighbor performs best. NC is outperformed by LR and NN in all scenarios but 5-shot on OntoNotes 5.0. This suggests that with very few samples, the raw token embedding information, as used by NN, is a better representation of a class than the averaged embeddings as produced by LR and CN, but with more samples, weighted embeddings obtained with LR are more useful.

5 Related Work

Few-shot NER Recent work on few-shot NER has primarily focused on integrating additional knowledge to support the classification process. Fritzler et al. (2019) are the first to use pre-trained word embeddings for this task. Yang and Katiyar (2020) extend a Nearest Neighbor token-level classifier with a Viterbi decoder for structured prediction over entire sentences. Huang et al (2020) propose to continue pre-training of a PLM encoder with distantly supervised, in-domain data, and to integrate self-training to create additional, soft-labeled training data. Recently, Gao et al. (2021) and Ma et al. (2021) investigate methods for making PLMs better few-shot learners via prompt-based fine-tuning. While these approaches extend standard few-shot learning algorithms in promising di-

Dataset	K	BERT↓		ALBERT↓		RoBERTa↑		SpanBERT↑		XLNet↑	
		w/o CL	CL	w/o CL	CL	w/o CL	CL	w/o CL	CL	w/o CL	CL
CoNLL-2003 _{EN}	1	21.96	23.87 †	33.03	36.71 †	21.71	22.57	18.39	17.61	18.49	18.25
	5	60.94	60.55	68.33	66.85	64.49	62.45	43.22	44.23	44.82	45.93
	10	66.11	65.03	72.76	70.66	72.09	70.17	48.79	49.82	52.43	49.25
OntoNotes 5.0	1	42.71	42.89	50.45	51.38	42.74	41.66	34.30	32.95	38.40	38.64
	5	74.68	74.02	77.66	76.65	78.70	75.29	65.64	64.29	72.60	70.66
	10	80.92	80.36	82.10	81.47	83.80	82.51	74.14	74.72	78.38	75.99
Few-NER _{coarse}	1	25.99	27.42	35.67	38.16 †	28.12	29.10	23.34	23.40	25.93	26.35
	5	53.85	52.97	59.14	59.71	58.66	55.75	45.50	46.03	52.32	54.91 †
	10	59.44	59.89	63.30	64.53	65.52	62.86	52.65	55.47 †	61.94	61.45
WikiGold	1	18.40	16.85	32.30	34.05 †	20.63	19.90	14.90	15.39	18.01	19.13
	5	49.19	49.19	55.87	57.67 †	56.08	53.91	41.07	42.92 †	45.44	44.21
	10	55.85	56.87	61.23	62.68 †	64.84	63.05	48.09	50.93 †	53.85	52.26

Table 5: Token-level micro F1-scores of PLM encoders without and with contrastive learning (CL) for 5-way K -shot scenarios, with logistic regression readout. † denotes scores with a significant ($\alpha = 0.05$) improvement after contrastive learning. † and ‡ indicate cased and uncased models.

Dataset	K	LR	NC	NN
CoNLL-2003 _{EN}	1	33.03	35.21	40.76 †
	5	68.33 †	61.53	62.24
	10	72.76 †	62.65	67.79
OntoNotes 5.0	1	50.45	51.52	52.72
	5	77.66 †	72.46	71.04
	10	82.10 †	73.49	76.11

Table 6: Micro-F1 scores of ALBERT for 5-way K -shot scenarios, comparing Logistic Regression (LR), Nearest Centroid (NC) and Nearest Neighbor (NN) readout approaches.

reactions, none of them directly investigate the contribution of different pre-trained representations. As such, our analysis complements these works. Das et al. (2021) present a contrastive pre-training approach for few-shot NER that uses in-domain data to fine-tune token embeddings before few-shot classification. In contrast, we only consider contrastive examples from the sampled few-shot set to conform to the low-resource setting.

Encoder comparisons In parallel to our work, Pearce et al. (2021) compare different Transformer models on extractive question answering and, similar to our results, find RoBERTa to perform best, outperforming BERT. However, they did not reproduce the strong performance we achieved with ALBERT and, unlike our results, found XLNet to be consistently outperforming BERT. Cortiz (2021) compare Transformer models for text-based emotion recognition and also found RoBERTa to perform best with XLNet being (shared) second, again outperforming BERT.

Our work can also be viewed as a kind of prob-

ing task (Conneau et al., 2018; Belinkov and Glass, 2019; Tenney et al., 2019), since we analyze how much information about named entities is preserved in the pre-trained representations, as measured by a linear classifier.

6 Conclusion

We presented a systematic, comparative study of pre-trained encoders on the task of low-resource named entity recognition. We find that encoder performance varies significantly depending on the scenario and the mix of pre-training and fine-tuning strategies. This suggests that the choice of encoders for a particular setting in current state-of-the-art low-resource NER approaches may need to be carefully (re-)evaluated. We also find that PLM encoders achieve reasonably good token classification performance on many English and German NER datasets with as little as 10 examples per class, in a fine-tuning-free setting. In particular, ALBERT turned out to be a very strong contender in one-shot settings, whereas RoBERTa often outperforms other PLMs in settings with more examples. For German, BERT shows the best average performance across scenarios, with XLM-R being more useful in one-shot settings.

One obvious direction for future work is to evaluate additional encoders, in particular models that are pre-trained in an entity-aware manner (Peters et al., 2019; Zhang et al., 2019). While our analysis is limited to NER, the encoder-readout framework can easily be adapted to evaluate other low-resource classification tasks.

References

Dominic Balasuriya, Nicky Ringland, Joel Nothman, Tara Murphy, and James R. Curran. 2009. [Named entity recognition in Wikipedia](#). In *Proceedings of the 2009 Workshop on The People’s Web Meets NLP: Collaboratively Constructed Semantic Resources (People’s Web)*, pages 10–18, Suntec, Singapore. Association for Computational Linguistics.

Yonatan Belinkov and James Glass. 2019. [Analysis methods in neural language processing: A survey](#). *Transactions of the Association for Computational Linguistics*, 7:49–72.

Darina Benikova, Chris Biemann, and Marc Reznicek. 2014. [NoSta-D named entity annotation for German: Guidelines and dataset](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC’14)*, pages 2524–2531, Reykjavik, Iceland. European Language Resources Association (ELRA).

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Sumit Chopra, Raia Hadsell, and Yann LeCun. 2005. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, volume 1, pages 539–546. IEEE.

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.

Alexis Conneau, Germán Kruszewski, Guillaume Lample, Loïc Barrault, and Marco Baroni. 2018. [What you can cram into a single \backslash\\$&!#* vector: Probing sentence embeddings for linguistic properties](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2126–2136, Melbourne, Australia. Association for Computational Linguistics.

Diogo Cortiz. 2021. [Exploring transformers in emotion recognition: a comparison of bert, distillbert, roberta, xlnet and electra](#). *CoRR*, abs/2104.02041. 589
590
591

Sarkar Snigdha Sarathi Das, Arzoo Katiyar, Rebecca J. Passonneau, and Rui Zhang. 2021. [Container: Few-shot named entity recognition via contrastive learning](#). *CoRR*, abs/2109.07589. 592
593
594
595

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. 596
597
598
599
600
601
602

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. 603
604
605
606
607
608
609
610
611

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. 612
613
614
615
616
617
618
619
620

Jenny Rose Finkel, Trond Grenager, and Christopher Manning. 2005. [Incorporating Non-local Information into Information Extraction Systems by Gibbs Sampling](#). In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics, ACL ’05*, pages 363–370, Stroudsburg, PA, USA. Association for Computational Linguistics. 621
622
623
624
625
626
627

Alexander Fritzler, Varvara Logacheva, and Maksim Kretov. 2019. [Few-shot classification in Named Entity Recognition Task](#). *Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing - SAC ’19*, pages 993–1000. ArXiv: 1812.06158. 628
629
630
631
632

Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. [Making pre-trained language models better few-shot learners](#). 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 3816–3830, Online. Association for Computational Linguistics. 633
634
635
636
637
638
639
640

Jiaxin Huang, Chunyuan Li, Krishan Subudhi, Damien Jose, Shobana Balakrishnan, Weizhu Chen, Baolin Peng, Jianfeng Gao, and Jiawei Han. 2020. [Few-Shot Named Entity Recognition: A Comprehensive Study](#). *arXiv:2012.14978 [cs]*. ArXiv: 2012.14978. 641
642
643
644
645

646	Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S. Weld, Luke Zettlemoyer, and Omer Levy. 2020. SpanBERT: Improving pre-training by representing and predicting spans . <i>Transactions of the Association for Computational Linguistics</i> , 8:64–77.	704
647		705
648		706
649		707
650		
651	Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations . In <i>8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020</i> . OpenReview.net.	708
652		709
653		710
654		711
655		712
656		713
657		714
658	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 . <i>CoRR</i> , abs/1907.11692.	715
659		716
660		717
661		718
662		719
663		720
664	Ruotian Ma, Xin Zhou, Tao Gui, Yiding Tan, Qi Zhang, and Xuanjing Huang. 2021. Template-free prompt tuning for few-shot ner . <i>CoRR</i> , abs/2109.13532.	721
665		722
666		723
667		724
668	Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1946–1958, Vancouver, Canada. Association for Computational Linguistics.	725
669		726
670		727
671		728
672		
673	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. Pytorch: An imperative style, high-performance deep learning library . In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, <i>Advances in Neural Information Processing Systems 32</i> , pages 8024–8035. Curran Associates, Inc.	729
674		730
675		731
676		732
677		733
678		734
679		735
680		736
681		737
682		
683		
684		
685		
686	Kate Pearce, Tiffany Zhan, Aneesh Komanduri, and Justin Zhan. 2021. A comparative study of transformer-based language models on extractive question answering . <i>CoRR</i> , abs/2110.03142.	738
687		739
688		740
689		741
690		742
691		743
692		744
693	Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep Contextualized Word Representations . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.	745
694		746
695		747
696		748
697		749
698		750
699		
700	Matthew E. Peters, Mark Neumann, Robert Logan, Roy Schwartz, Vidur Joshi, Sameer Singh, and Noah A. Smith. 2019. Knowledge enhanced contextual word representations . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 43–54, Hong Kong, China. Association for Computational Linguistics.	704
701		705
702		706
703		707
	Clifton Poth, Jonas Pfeiffer, Andreas Rücklé, and Iryna Gurevych. 2021. What to pre-train on? Efficient intermediate task selection . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10585–10605, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	708
		709
		710
		711
		712
		713
		714
	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	715
		716
		717
		718
		719
		720
	Nils Rethmeier and Isabelle Augenstein. 2021. A primer on contrastive pretraining in language processing: Methods, lessons learned and perspectives . <i>CoRR</i> , abs/2102.12982.	721
		722
		723
		724
	Raphael Scheible, Fabian Thomczyk, Patric Tippmann, Victor Jaravine, and Martin Boeker. 2020. Gottbert: a pure german language model . <i>CoRR</i> , abs/2012.02110.	725
		726
		727
		728
	Martin Schiersch, Veselina Mironova, Maximilian Schmitt, Philippe Thomas, Aleksandra Gabryszak, and Leonhard Hennig. 2018. A German Corpus for Fine-Grained Named Entity Recognition and Relation Extraction of Traffic and Industry Events . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).	729
		730
		731
		732
		733
		734
		735
		736
		737
	Ian Tenney, Patrick Xia, Berlin Chen, Alex Wang, Adam Poliak, R. Thomas McCoy, Najoung Kim, Benjamin Van Durme, Sam Bowman, Dipanjan Das, and Ellie Pavlick. 2019. What do you learn from context? Probing for sentence structure in contextualized word representations . In <i>International Conference on Learning Representations</i> .	738
		739
		740
		741
		742
		743
		744
	Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition . In <i>Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003</i> , pages 142–147.	745
		746
		747
		748
		749
		750
	Tu Vu, Tong Wang, Tsendsuren Munkhdalai, Alessandro Sordani, Adam Trischler, Andrew Mattarella-Micke, Subhransu Maji, and Mohit Iyyer. 2020. Exploring and predicting transferability across NLP tasks . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7882–7926, Online. Association for Computational Linguistics.	751
		752
		753
		754
		755
		756
		757
		758

- 759 Ralph Weischedel, Martha Palmer, Mitchell Marcus,
760 Eduard Hovy, Sameer Pradhan, Lance Ramshaw, Ni-
761 anwen Xue, Ann Taylor, Jeff Kaufman, Michelle
762 Franchini, et al. 2013. Ontonotes release 5.0
763 ldc2013t19. *Linguistic Data Consortium, Philadel-*
764 *phia, PA*, 23.
- 765 Adina Williams, Nikita Nangia, and Samuel Bowman.
766 2018. [A broad-coverage challenge corpus for sen-](#)
767 [tence understanding through inference](#). In *Proceed-*
768 *ings of the 2018 Conference of the North American*
769 *Chapter of the Association for Computational Lin-*
770 *guistics: Human Language Technologies, Volume*
771 *1 (Long Papers)*, pages 1112–1122. Association for
772 Computational Linguistics.
- 773 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien
774 Chaumond, Clement Delangue, Anthony Moi, Pier-
775 ric Cistac, Tim Rault, Rémi Louf, Morgan Funtow-
776 icz, Joe Davison, Sam Shleifer, Patrick von Platen,
777 Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu,
778 Teven Le Scao, Sylvain Gugger, Mariama Drame,
779 Quentin Lhoest, and Alexander M. Rush. 2020.
780 [Transformers: State-of-the-art natural language pro-](#)
781 [cessing](#). In *Proceedings of the 2020 Conference on*
782 *Empirical Methods in Natural Language Processing:*
783 *System Demonstrations*, pages 38–45, Online. Asso-
784 ciation for Computational Linguistics.
- 785 Omry Yadan. 2019. [Hydra - a framework for elegantly](#)
786 [configuring complex applications](#). Github.
- 787 Yi Yang and Arzoo Katiyar. 2020. [Simple and effective](#)
788 [few-shot named entity recognition with structured](#)
789 [nearest neighbor learning](#). In *Proceedings of the*
790 *2020 Conference on Empirical Methods in Natural*
791 *Language Processing (EMNLP)*, pages 6365–6375,
792 Online. Association for Computational Linguistics.
- 793 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Car-
794 bonell, Russ R Salakhutdinov, and Quoc V Le. 2019.
795 Xlnet: Generalized autoregressive pretraining for
796 language understanding. *Advances in neural infor-*
797 *mation processing systems*, 32.
- 798 Hanchu Zhang, Leonhard Hennig, Christoph Alt,
799 Changjian Hu, Yao Meng, and Chao Wang.
800 2020. [Bootstrapping named entity recognition in E-](#)
801 [commerce with positive unlabeled learning](#). In *Pro-*
802 *ceedings of The 3rd Workshop on e-Commerce and*
803 *NLP*, pages 1–6, Seattle, WA, USA. Association for
804 Computational Linguistics.
- 805 Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang,
806 Maosong Sun, and Qun Liu. 2019. [ERNIE: En-](#)
807 [hanced language representation with informative en-](#)
808 [tities](#). In *Proceedings of the 57th Annual Meet-*
809 *ing of the Association for Computational Linguis-*
810 *tics*, pages 1441–1451, Florence, Italy. Association
811 for Computational Linguistics.

A Additional Training Details

We used a single RTX A6000-GPU for all experiments. The average runtime per scenario (dataset, encoder) for 600 episodes was approximately 1 minute (1-shot), 3 minutes (5-shot) and 6 minutes (10-shot). Contrastive pre-training was also performed on the same single RTX A6000-GPU, and took approximately 1 hour of GPU-time, including hyperparameter search.

For contrastive pre-training, the following hyperparameters were manually tuned: learning rate in $[2 \times 10^{-5}, 5 \times 10^{-5}]$, the number of epochs in $[1, 2, 5]$. We used the most occurrences of F1-gains across all encoders and scenarios on CoNLL-2003 dataset as criterion for hyperparameter selection.

All pre-trained models evaluated in this study were used as they are available from HuggingFace’s model hub, without any modifications. Table 7 lists the model identifiers. We used HuggingFace’s dataset hub for all datasets except the dataset by Zhang et al. (2020), which is used here with the permission of the authors.

Model	HuggingFace ID
BERT↓	bert-base-uncased
BERT↑	bert-base-cased
ALBERT	albert-base-v2
RoBERTa	roberta-base
SpanBERT	SpanBERT/spanbert-base-cased
XLNET	xlnet-base-cased
BERT DE	bert-base-german-cased
GottBERT	uklfr/gottbert-base
XLM-R	xlm-roberta-base
BERT _{POS}	vblagoje/bert-english-uncased-finetuned-pos
BERT _{MNLI}	textattack/bert-base-uncased-MNLI
BERT _{SQuAD}	csarron/bert-base-uncased-squad-v1
BERT _{CoNLL}	dslim/bert-base-NER-uncased

Table 7: HuggingFace model identifiers of evaluated encoders

B Readout approaches

Logistic Regression (LR) is a linear classification algorithm that can be extended to multinomial logistic regression to deal with multi-class (N -way) settings, such as the one discussed here. The probability that query token x' belongs to the c -th class is given by:

$$\Pr(y = c) = \frac{\text{score}(x', c)}{\sum_{i=1}^N \text{score}(x', i)} \quad (1)$$

$$\text{score}(x', i) := \exp(W_i \cdot f_{\theta}(x')),$$

where W is a matrix of N rows learned from the support set \mathcal{S} , and W_i denotes the i -th row of W . $\text{score}(\cdot)$ serves as the metric to measure the affinity between token x' and the prototype of class c , and the prediction is given by

$$y^* = \arg \max_{c \in \{1, \dots, N\}} \text{score}(x', c).$$

k-Nearest Neighbor (NN) is a non-parametric classification method adopted in metric space. As proposed in STRUCTSHOT (Yang and Katiyar, 2020), we set $k = 1$ to find the exact nearest token in the support set. Given a query token x' ,

$$y^* = \arg \min_{c \in \{1, \dots, N\}} d_c(x') \quad (2)$$

$$d_c(x') := \min_{x \in \mathcal{S}_c} d(f_{\theta}(x'), f_{\theta}(x)),$$

where \mathcal{S}_c is the set of support tokens whose tags are c , and d denotes the distance between two embeddings in the representation space.

Nearest Centroid (NC) works similar to NN. In contrast, for each query token x' , instead of computing the distance between $f_{\theta}(x')$ and every instance in the embedding space, we represent each class by the centroid c_c of all embeddings belonging to this class, and assign token x' to the class with the nearest centroid:

$$y^* = \arg \min_{c \in \{1, \dots, N\}} d(f_{\theta}(x'), c_c) \quad (3)$$

$$c_c = \frac{1}{|\mathcal{S}_c|} \sum_{x \in \mathcal{S}_c} f_{\theta}(x).$$

C Entity tag sets of English datasets

We list the full entity tag sets for all English benchmarks. Overlap entity tags with CoNLL-2003_{EN} are highlighted with underline.

C.1 CoNLL-2003_{EN}

LOC, MISC, ORG, PER.

C.2 OntoNotes 5.0

CARDINAL, DATE, EVENT, FAC, GPE, LANGUAGE, LAW, LOC, MONEY, NORP, ORDINAL, ORG, PERCENT, PERSON, PRODUCT, QUANTITY, TIME, WORK_OF_ART.

C.3 Few-NERD_{coarse}

art, building, event, location, organization, other⁶, person, product.

⁶Few-NERD_{coarse} sets non-entity as 'O' and various entity types as 'other'. Therefore, we treat 'other' as 'MISC' in this case.

873 **C.4 Few-NERD_{fine}**

874 art-broadcastprogram, art-film, art-music, art-
875 other, art-painting, art-writtenart, building-
876 airport, building-hospital, building-hotel,
877 building-library, building-other, building-
878 restaurant, building-sportsfacility, building-
879 theater, event-attack/battle/war/militaryconflict,
880 event-disaster, event-election, event-other,
881 event-protest, event-sportsevent, location-
882 GPE, location-bodiesofwater, location-island,
883 location-mountain, location-other, location-
884 park, location-road/railway/highway/transit,
885 organization-company, organization-education,
886 organization-government/governmentagency,
887 organization-media/newspaper, organization-other,
888 organization-politicalparty, organization-religion,
889 organization-showorganization, organization-
890 sportsleague, organization-sportsteam, other-
891 astronomything, other-award, other-biologything,
892 other-chemicalthing, other-currency, other-
893 disease, other-educationaldegree, other-god,
894 other-language, other-law, other-livingthing,
895 other-medical, person-actor, person-artist/author,
896 person-athlete, person-director, person-other,
897 person-politician, person-scholar, person-soldier,
898 product-airplane, product-car, product-food,
899 product-game, product-other, product-ship,
900 product-software, product-train, product-weapon

901 **C.5 WNUT-17**

902 corporation, creative-work, group, location, person,
903 product.

904 **C.6 WikiAnn**

905 LOC, ORG, PER.

906 **C.7 WikiGold**

907 LOC, MISC, ORG, PER.

908 **C.8 Zhang et al.**

909 ATTRIBUTE, BRAND, COMPONENT, PROD-
910 UCT.