

DeepTagger: Knowledge Enhanced Named Entity Recognition for Web-Based Ads Queries

Simiao Zuo*, Pengfei Tang, Xinyu Hu, Qiang Lou, Jian Jiao, Denis Charles
{simiaozuo, pengfeitang, xinyuhu, qilou, jian.jiao, cdx}@microsoft.com
Microsoft

Abstract

Named entity recognition (NER) is a crucial task for online advertisement. State-of-the-art solutions leverage pre-trained language models for this task. However, three major challenges remain unresolved: web queries differ from natural language, on which pre-trained models are trained; web queries are short and lack contextual information; and labeled data for NER is scarce. We propose DeepTagger, a knowledge-enhanced NER model for web-based ads queries. The proposed knowledge enhancement framework leverages both model-free and model-based approaches. For model-free enhancement, we collect unlabeled web queries to augment domain knowledge; and we collect web search results to enrich the information of ads queries. We further leverage effective prompting methods to automatically generate labels using large language models such as ChatGPT. Additionally, we adopt a model-based knowledge enhancement method based on adversarial data augmentation. We employ a three-stage training framework to train DeepTagger models. Empirical results in various NER tasks demonstrate the effectiveness of the proposed framework.

1 Introduction

Named Entity Recognition (NER) is the task of classifying each token in an input sequence into predefined categories. For example, in a query such as “hotels in Seattle”, we need to identify “hotels” as a type of product and “Seattle” as a location designator. In the advertisement domain, NER lays the foundation for subsequent applications such as product retrieval (Cheng et al., 2021b), query rewriting (Wen et al., 2019), and attribute value extraction (Zhang et al., 2021). Existing works leverage pre-trained language models (PLMs (Devlin et al., 2019)) for NER. These models are pre-trained on large natural language corpora and contain rich syntactic and semantic information.

There are three major challenges when applying PLMs to NER for web-based ads queries. First, there is a domain shift between web queries and natural language. Most web queries lack

*Corresponding author.

grammatical components such as verbs and subjects. For example, “home furniture bedroom” is an informative web query but not a grammatically correct sentence. Additionally, web queries involve specification attributes and product models that are uncommon in natural language. These properties create a domain shift that hinders the performance of PLMs trained on open-domain natural language.

Second, web queries are short and lack information. For example, in the CoNLL2003 (Tjong Kim Sang, 2002) dataset that contains news articles, the average sequence length is 14.5. However, the average length of ads queries is only 3.9 on a self-collected dataset. This is problematic because the short web queries often lack semantic components for PLMs to make informed predictions. For example, in the natural language sentence “Rabinovich is winding up his term as ambassador”, we can easily infer that “Rabinovich” is a person’s name. However, for a web query such as “credit card square”, PLMs are unlikely to predict that “square” is a brand.

The third problem is label scarcity. NER tasks demand token-level labels, requiring experienced and well-trained human annotators. As a result, domain-specific labeled data of good quality are limited (Jiang et al., 2021; Zhang et al., 2021). Existing works leverage weak supervision to tackle this issue. That is, instead of training human experts to accurately annotate data, automatic tools are used to generate noisy labels. For example, we can match the inputs to external knowledge bases (Liang et al., 2020; Jiang et al., 2021) or semantic rules (Yu et al., 2021; Mukherjee and Awadallah, 2020; Awasthi et al., 2020). However, labels generated by weak supervision can be extremely noisy, rendering the training process unstable (Zuo et al., 2022).

To tackle the above three challenges, we propose DeepTagger, a knowledge-enhanced NER model for web-based ads queries. Specifically, we propose a knowledge enhancement framework that incorporates both *model-free* and *model-based* knowledge enhancement. We further propose a three-stage training framework to train DeepTagger models.

We adopt three model-free knowledge enhancement methods: (1) To address the domain shift issue, we collect large quantities of unlabeled web query data. The syntactic and semantic knowledge in these data can help models adapt to the advertisement domain. (2) We also collect web search results to complement the lack of information in ads queries. Specifically, for each ads query, we retrieve several search results and keep the titles of these results. The web titles provide more context than the query, allowing models to better infer the role of each token (see Table 1 for examples). (3) To alleviate the label scarcity problem, we collect weakly-labeled data from several sources. First, we resort to crowdsourcing platforms to collect inaccurate labels. Second, we leverage the Chain-of-Thoughts prompting (Wei et al., 2022) to automatically generate weak labels using large language models (LLMs) such as ChatGPT and GPT-4 (OpenAI, 2023) (see Figure 1 for an example). We remark that to the best of our knowledge, we are one of the first to augment data using LLMs for NER.

In practice, we find that fine-tuning PLMs on NER tasks still faces severe overfitting. Therefore, we propose a model-based knowledge enhancement method based on adversarial regular-

Table 1: Two examples of web queries and the associated titles of search results (termed *Web4Ads*).

Query	remote it support tools
Web4Ads	Best remote desktop software Best Tools to Easily Perform Remote Tech Support Best Remote Access Software
Query	credit card square
Web4Ads	Square: Solutions & Tools to Grow Your Business Square Payments: Payment Processing Square Review: Fees, Complaints

ization (Miyato et al., 2017, 2018). Specifically, during each training iteration, we find samples on which the model is likely to make erroneous predictions and augment the training data with such samples. The proposed augmentation technique improves model generalization by promoting prediction smoothness.

We train DeepTagger models using a three-stage framework. In **Stage I**, we continue pre-training (Gururangan et al., 2020) a PLM on unlabeled web query data to inject domain knowledge. In **Stage II**, we train the PLM from the previous stage on a large amount of weakly-labeled data. In **Stage III**, we fine-tune the PLM from Stage II using model-based knowledge enhancement on a small amount of strongly-labeled data. We note that similar multi-stage training methods have been shown to be effective in various tasks (Liang et al., 2020; Yu et al., 2021; Jiang et al., 2021; Zuo et al., 2022).

2 Model-Free Knowledge Enhancement

2.1 Search Titles Complement Short Queries

Web queries are inherently short, which presents a problem: the lack of semantic components in these queries often makes it difficult for pre-trained models to perform well. To address this issue, we propose augmenting web queries with search titles. Specifically, we retrieve search results from the search engine for each web query and keep the titles of the results. We refer to these titles as *Web4Ads*.

Intuitively, the *Web4Ads* titles contain richer information than the original web queries, which can help the model better understand the role of each token in the query. Table 1 provides two examples. First, consider the query “remote it support tools”. With only the web query, the model is likely to label the token “it” as *other*. However, the *Web4Ads* titles associate “it” with entities such as “tech” and “software”. Therefore, the model has a better chance of understanding that “it” is a *product*. Second, for the query “credit card square”, even human annotators have trouble

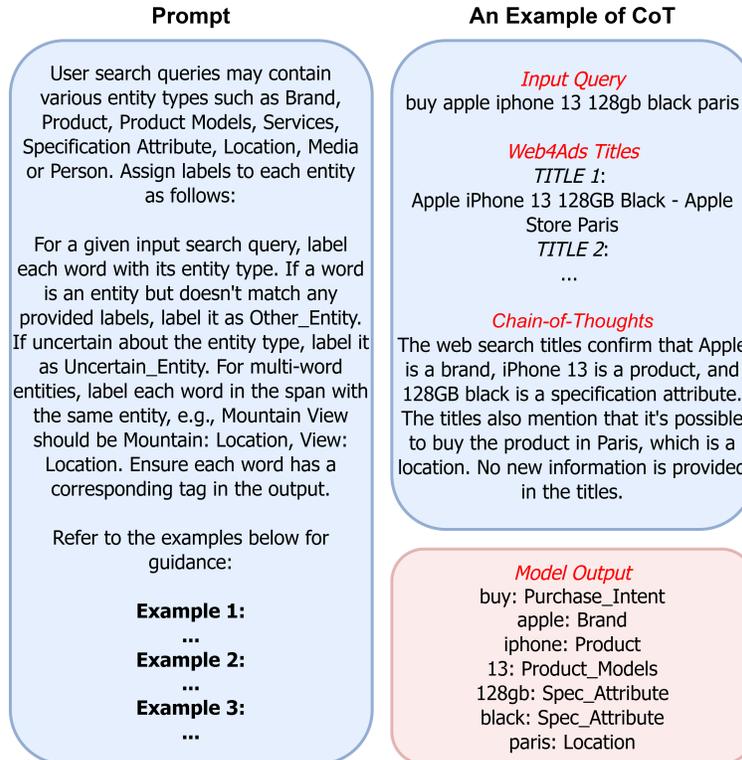


Figure 1: Prompt template for generating weak labels using large language models. The blue chunks are model inputs, and the red chunk contains model outputs.

recognizing that “square” is a company’s name. However, after augmenting the Web4Ads titles, the model can infer that “square” is a *brand* instead of a shape description.

During training, we concatenate a web query with its Web4Ads titles using a “[SEP]” token as the separator, e.g., an example is “Query [SEP] Web4Ads₁ [SEP] Web4Ads₂”. We do not collect labels for the Web4Ads titles. Consequently, we do not compute the supervised loss for tokens that corresponds to the Web4Ads titles. Empirically, our approach avoids introducing additional labeling burden and greatly improves model performance (see Section 6.4 in the experiments for details).

2.2 Augmentation of In-Domain Data

◊ **Unlabeled data.** Web queries differ from natural language in that they often contain uncommon tokens (e.g., specification attributes and product models) and lack grammatical components such as verbs. To address this domain shift issue, we collect large quantities of unlabeled web queries. The domain-specific semantic and syntactic knowledge in such data can be injected into the models.

◊ **Weakly-labeled data.** Experienced and well-trained human annotators are required to generate accurate token-level labels for NER tasks. As a result, strongly (or accurately) labeled data are

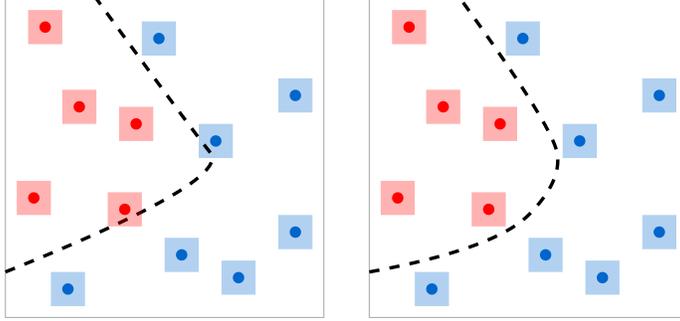


Figure 2: Illustration of decision boundaries without (left) and with (right) adversarial knowledge enhancement. A solid circle indicates a labeled sample, and the square around it indicates its neighborhood. The red and blue colors indicate two difference classes of samples, and the dashed lines are decision boundaries.

often scarce in practice due to cost concerns. To address this issue, we leverage several sources to generate weakly-labeled data.

We collect crowdsourcing data. We remark that even though the data is human-annotated, the quality is often weak due to task difficulty and the absence of properly trained annotators.

We generate weakly-labeled data using large language models (Yoo et al., 2021; Sahu et al., 2022; Liu et al., 2023b). The development of models such as ChatGPT and GPT-4 (OpenAI, 2023) has enabled weak-label generation with nearly no cost. Specifically, we leverage Chain-of-Thoughts prompting (Wei et al., 2022), where we use the Web4Ads titles as intermediate reasoning steps to guide the “thinking” of large language models (see Figure 1 for an example). We thoroughly investigate the effectiveness of several other prompting methods in Section 5.

3 Model-Based Knowledge Enhancement

Fine-tuning pre-trained language models require large quantities of labeled data, which are often unavailable. Even though the proposed model-free knowledge enhancement methods can partially alleviate this issue, in practice, extensive hyper-parameter tuning is still needed to avoid overfitting. We propose a model-based data augmentation method based on adversarial regularization (Miyato et al., 2017, 2018) to reduce over-fitting.

3.1 Virtual Data Augmentation

Given a single datum, if we perturb it by a small noise, the model’s prediction should not change. Such a smoothness assumption promotes the model’s generalization performance (Miyato et al., 2018). Therefore, for each labeled sample, we can augment the training dataset by generating virtual data in its neighborhood (Awasthi et al., 2020).

Figure 2 illustrates the idea. From Figure 2 (left), we see that without augmentation, some data

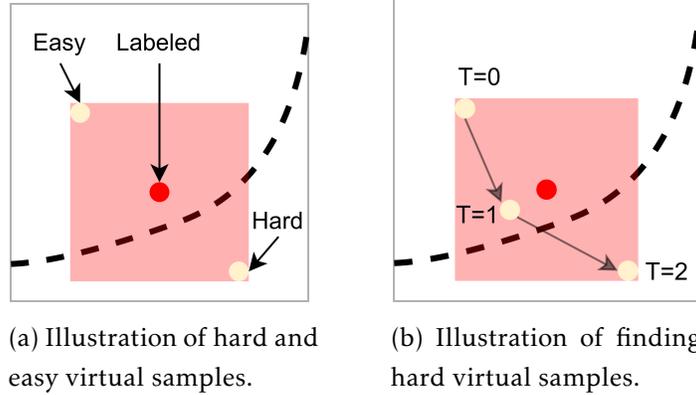


Figure 3: Illustration of virtual samples. The red solid circle indicates a labeled sample, and the square indicates its neighborhood. The light yellow circles are virtual samples.

are very close to the decision boundary. Therefore, a small perturbation to the input data may result in a substantial change to the model’s prediction. Virtual data augmentation encourages the model to make consistent predictions under small perturbations. For example, in Figure 2, if the model predicts a sample as “red”, then it should classify all virtually generated data around the sample (i.e., data in the sample’s neighborhood) as “red”. In this way, the decision boundary becomes smoother, which improves model generalization.

3.2 Model-Based Adversarial Data Augmentation

There are infinitely many virtual data in the neighborhood of a labeled sample. However, not all of the virtual samples are of the same “difficulty”. For example, in Figure 3a, the model correctly classifies the labeled sample (the red circle). In this case, augmenting the data with the *easy* virtual sample (the yellow circle on the top left) will not bring any improvement since the model can already correctly classify it. On the other hand, the model predicts the *hard* virtual sample (the yellow circle on the bottom right) to a different class, which violates our assumption that neighboring data should have the same label. Therefore, augmenting this hard virtual sample will benefit model generalization (see Figure 2).

We find the hard virtual samples via adversarial training. Concretely, denote $f(x, \theta)$ a neural network parameterized by θ , where x is the input. We note that x is continuous and resides in the embedding space. For example, for NER, x is the continuous representation after the embedding layers. Then, We find the hard virtual samples by optimizing the following objective

$$\max_{\|\delta\| \leq \epsilon} \ell_v(x, \delta, \theta), \tag{1}$$

where $\ell_v(x, \delta, \theta) = \text{SymKL}(f(x, \theta), f(x + \delta, \theta))$.

Here, $\text{SymKL}(P, Q) = \frac{1}{2}(\text{KL}(P\|Q) + \text{KL}(Q\|P))$ is the symmetric KL-divergence between two probability distributions; $\|\cdot\|$ is taken as the ℓ_2 norm or the ℓ_∞ norm; and ϵ is a pre-defined radius of

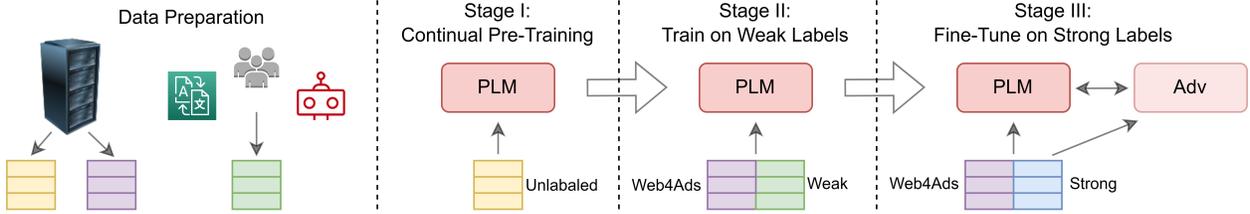


Figure 4: Overall framework of DeepTagger. In data preparation, unlabeled and Web4Ads data are retrieved from search engines; and weak labels are generated from multiple sources. Strongly-labeled data in Stage III are annotated by human experts.

the perturbations. The loss ℓ_v measures discrepancy between the model’s predictions given the clean data x and the perturbed data $x + \delta$. Thus, by maximizing ℓ_v , we can find virtual samples on which the model easily make erroneous predictions (e.g., the hard virtual sample in Figure 3a).

In practice, we solve Equation 1 using projected gradient ascent (Madry et al., 2018). That is, we adopt the following update rule:

$$\delta_{k+1} = \Pi \left(\delta_k + \eta \frac{\nabla_{\delta} \ell_v(x, \delta_k, \theta)}{\|\nabla_{\delta} \ell_v\|_2} \right), \text{ where } \delta_0 \sim \Pi(\mathcal{N}(\mathbf{0}, \mathbf{I})).$$

Here, $\Pi(\cdot)$ is the projection operator onto the $\|\cdot\|$ ball, and η is the learning rate. Figure 3b illustrates finding hard virtual samples. At first ($T = 0$), we randomly initialize a perturbation. Then, we update the perturbation by moving it towards the decision boundary ($T = 1$). Finally ($T = 2$) the model makes a wrong prediction, and we deem the resulting virtual sample *hard*.

To leverage the hard virtual samples during training, we optimize

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \ell(f(x_i, \theta), y_i) + \max_{\|\delta_i\| \leq \epsilon} \ell_v(x_i, \delta_i, \theta),$$

where y_i is the ground-truth label corresponding to x_i , and ℓ is the cross-entropy loss.

We note that virtual data augmentation in Section 3.1 has been shown to improve model generalization (Aghajanyan et al., 2021). However, in practice, adversarial data augmentation works better in terms of both model performance and training stability (Zuo et al., 2021a).

4 Training of DeepTagger

We propose a three-stage training framework to train DeepTagger models. The framework is illustrated in Figure 4.

In **Stage I**, we continue pre-training a pre-trained language model (PLM) on the collected unlabeled web queries. Specifically, we train the PLM using self-supervision objectives such as masked language modeling (Devlin et al., 2019). This can effectively address the domain shift issue by adapting the PLM to the advertisement domain.

Table 2: Dataset statistics. In all the experiments, we report results on the test set sampled from strongly-labeled data.

Language	Type	#Train	#Test	#Classes
En	unlabeled	998.7M	–	11
	crowdsourced	2.6M	–	
	strongly-labeled	167.4k	2.0k	
De	unlabeled	120.8M	–	11
	LLM	100.0k	–	
	strongly-labeled	51.5k	2.2k	
Fr	unlabeled	91.5M	–	11
	LLM	100.0k	–	
	strongly-labeled	50.0k	2.2k	

In **Stage II**, we train the PLM from the previous stage on a large amount of weakly-labeled data, augmented with Web4Ads titles. To prevent the PLM from overfitting to the noise in the weak labels, we adopt an early-stopping strategy (Dodge et al., 2020). We note that other weakly-supervised learning methods such as contrastive learning (Yu et al., 2021) and confidence regularization (Pereyra et al., 2017) can be applied in this stage.

In **Stage III**, we fine-tune the PLM from Stage II on a small amount of strongly-labeled data, augmented with Web4Ads titles. We also apply model-based adversarial data augmentation in this stage. We note that the strongly-labeled data are annotated by human experts, and therefore the amount of them is magnitudes smaller than the weakly-labeled data.

5 Data Preparation

5.1 Data Overview

Table 2 summarizes the collected data. Specifically, we use unlabeled web search queries collected from search engines for continual pre-training. We also collect strongly-labeled data, where the labels are annotated by human experts. However, the quantity of such data is limited because of cost concerns.

Additionally, we collect crowdsourcing data for queries in English. And for other languages (e.g., German and French), we generate weak labels using LLMs (see Section 5.2 for details). We do not collect crowdsourcing data in German and French because of the lack of human annotators.

Table 3: Comparison of different prompting methods. Here, “+” denotes positive examples, “-” denotes negative ones, and “(+,-)” denotes positive examples followed by negative ones.

Method	Few Shots	Embedding	Web4Ads	F1
Prompting	No	No	No	0.62
Demo.	Yes	No	No	0.72
Dyna. Demo.	Yes	EASE (+)	No	0.67
Dyna. Demo.	Yes	EASE (+,-)	No	0.69
Dyna. Demo.	Yes	EASE (-,+)	No	0.70
Dyna. Demo.	Yes	SBERT (+)	No	0.68
CoT	Yes	No	top 3	0.80

5.2 Weak Labels from Large Language Models

Large language models such as ChatGPT excels at straightforward tasks such as text classification (Ding et al., 2022; Chiang and Lee, 2023; Wang et al., 2023; Wu et al., 2023). However, they encounter limitations when dealing with more intricate problems such as NER.

We explore a variety of methods aimed at enhancing prompts for label generation.

- ◊ **Prompting.** In this approach, we only use the prompt (without examples) in Figure 1 (left) to generate labels.
- ◊ **Demonstration.** We utilize a few-shot demonstrations approach (Liu et al., 2023a; Xie et al., 2022) with fixed instances (Figure 1 left). Specifically, we select three fixed examples that ensure full coverage across all categories.
- ◊ **Dynamic demonstration.** In this case, the demonstrations are dynamically retrieved. Specifically, we retrieve three positive and/or three negative examples using two existing embedding models: SBERT (Reimers and Gurevych, 2019) and EASE (Nishikawa et al., 2022). This approach facilitates contrastive in-context learning.
- ◊ **Chain-of-Thoughts.** We enhance the prompt using web information via Chain-of-Thoughts prompting (Wei et al., 2022). Specifically, we augment the Web4Ads titles to intermediate reasoning steps (Figure 1 right).

To evaluate the prompting methods, we generate labels on strongly-labeled test sets (Table 2). Then, we calculate the F1 score of the “Brand” category (there are 11 categories in total) since it plays a crucial role in downstream ads-related tasks.

As shown in Table 3, vanilla *Prompting* yields a brand F1 score of 0.62. By employing few-shot demonstrations (*Demo*), the F1 score increases to 0.72. However, we observe a performance degradation when employing *Dynamic Demonstration* methods with EASE and SBERT. This is because

Table 4: Experimental results on English queries. We report the overall F1 score and the F1 score that correspond to the “Brand” category. The best results are shown in **bold**.

English	Weak Data	Weak Labels	Brand	Overall
Direct training				
BERT	✗	✗	80.77	74.29
Semi-supervised baselines				
Self-Training	✓	✗	81.20	74.76
DRIFT	✓	✗	81.61	75.12
VAT	✓	✗	82.04	75.86
Weakly-supervised baselines				
COSINE	✓	✓	81.34	74.69
NEEDLE	✓	✓	82.27	76.33
Ours				
DeepTagger	✓	✓	83.45	77.94

both embedding models are unsuitable for NER. In particular, SBERT does not incorporate entity-related information, and EASE falls short in capturing the relationships between entities and their contexts. To better leverage the entity information and the augmented web titles, we integrate a Chain-of-Thoughts layer, which effectively guides the LLMs to extract valuable information from Web4Ads titles. This approach results in a substantial increase in the F1 score, raising it to 0.80.

We remark that in reality, usage of LLMs in online deployment is limited due to latency constraints. For example, our DeepTagger model can output token labels within milliseconds, but LLMs require several hundred milliseconds or even seconds.

6 Experiments

6.1 Baselines

We implement all the models using *PyTorch* (Paszke et al., 2019) and the *Huggingface Transformers* (Wolf et al., 2020) code-base. In the experiments, we fine-tune a BERT-base model for English data, and we fine-tune a multilingual version of BERT-base for data in other languages. To facilitate fair comparisons, we conduct unsupervised continual pre-training (i.e., *Stage I* in our framework) for all models.

We compare DeepTagger with several weakly-supervised and semi-supervised learning baselines. We note that for semi-supervised learning baselines, we first fine-tune a BERT model on the strongly-labeled data, and we treat the resulting model as the “teacher”. We do not use the weak labels in the weakly-labeled data, and instead we use the teacher model to generate pseudo-labels

Table 5: Results on German and French queries. We report the overall F1 score and the F1 score that correspond to the “Brand” category. The best results are shown in **bold**.

	German		French		Average	
	Brand	Overall	Brand	Overall	Brand	Overall
BERT	67.11	58.21	69.58	61.45	68.35	59.83
VAT	67.41	58.32	69.62	61.45	68.52	59.89
NEEDLE	67.80	58.41	69.91	61.57	68.86	59.99
DeepTagger	68.59	59.24	70.57	62.63	69.58	60.94

for these data.

◊ *BERT* (Devlin et al., 2019) is where we only train on the strongly-labeled data, i.e., without weakly-labeled data.

◊ *Self-Training* (Rosenberg et al., 2005; Lee et al., 2013) is a classic semi-supervised learning approach. In Self-Training, we simultaneously maintain a teacher and a student model. The teacher generates pseudo-labels, on which the student is trained. The two models are updated alternately.

◊ *DRIFT* (Zuo et al., 2022) adopt a differentiable self-training approach. The method stabilizes conventional self-training by formulating the mean-teacher framework as a Stackelberg game.

◊ *VAT* (Miyato et al., 2017, 2018) is a semi-supervised learning method that employs adversarial training. Specifically, VAT regularizes model training by penalizing the divergence between model predictions on clean and perturbed unlabeled data.

◊ *COSINE* (Yu et al., 2021) is a weakly-supervised learning framework that can efficiently leverage noisily-labeled data. The framework adopts token-level contrastive learning, and also integrates the power of confidence-based sample re-weighting and regularization.

◊ *NEEDLE* (Jiang et al., 2021) is a NER framework that use small strongly-labeled data and large weakly-labeled data. The framework proposes a noise-aware loss function for weakly-supervised learning.

6.2 Main Results

Table 4 summarizes results on English queries; and Table 5 summarizes results on German and French queries. From the results, we see that DeepTagger performs significantly better than all the semi-supervised and weakly-supervised baselines.

We note that in semi-supervised learning, we do not use the weak labels from crowdsourcing platforms and LLMs, which is different from weakly-supervised learning. Because of this, we see that the best-performing semi-supervised learning baseline (VAT) behaves worse than the best-performing weakly-supervised learning baseline (NEEDLE) in terms of both Brand F1 score and

overall F1 score.

Also, we note that in weakly-supervised learning, COSINE adopts a different training framework, where in the last stage the model is trained on both weakly-labeled and strongly-labeled data. In contrast, in DeepTagger and NEEDLE, the model is first trained on weakly-labeled data and then fine-tuned on strongly-labeled data. The results in Table 4 and Table 5 indicate that the weak-then-strong training approach is better than training on both weakly-labeled and strongly-labeled data.

Even though DeepTagger and NEEDLE employ a similar training framework, we see that performance of DeepTagger is significantly better. In NEEDLE, vanilla training (i.e., without data augmentation or modification to the loss function) is used in final fine-tuning on strongly-labeled data. However, because strongly-labeled data are limited, the model can easily overfit to the noise. Our method can reduce overfitting and improve model generalization via adversarial data augmentation.

6.3 Online Deployment

The described system has been deployed to Microsoft Bing Ads for approximately one year. During this period, DeepTagger processes more than one billion queries daily. Compared with the previous NER system, DeepTagger increases revenue by 0.7%, increases Brand detection F1 score by 4.7%, and increases overall entity detection F1 score by 3.2%.

6.4 Analysis

◇ **Effectiveness of continual pre-training.** Recall that in Stage I of DeepTagger, we continue pre-train a BERT model on unlabeled web query data to inject domain knowledge. Table 6 summarizes the results. We see that continual pre-training improves model performance for both vanilla training and DeepTagger. For example, without pretraining, DeepTagger yields a 61.05 overall F1 score on French queries. And after integrating pre-training, the overall F1 score increases to 62.63 (+1.58).

◇ **Effectiveness of adversarial data augmentation.** Table 7 demonstrates model performance with and without adversarial data augmentation. First, we see that even without augmentation, performance of DeepTagger is significantly better than vanilla training. Second, we see that adversarial data augmentation and virtual data augmentation can indeed improve model performance by reducing overfitting. Third, notice that adversarial data augmentation performs better than virtual data augmentation (Section 3.1) because of its ability to find “hard” virtual samples.

◇ **Effectiveness of Web4Ads.** In DeepTagger, we complement the short queries with Web4Ads titles. We verify the effectiveness of such an approach in Figure 5. In Figure 5 (left), we plot the average sequence length with and without Web4Ads titles. We see that the average length of web queries is only 3.9, such that the queries alone might not be informative. After integrating

Table 6: Experimental results with and without continual pre-training. For “BERT”, we directly fine-tune the model on strongly-labeled data without using weak labels.

		BERT		DeepTagger	
		w/o	w/	w/o	w/
English	Brand	79.71	80.77	82.29	83.45
	Overall	73.23	74.29	77.41	77.94
German	Brand	65.87	67.11	67.52	68.59
	Overall	57.08	58.21	57.73	59.24
French	Brand	68.56	69.58	69.16	70.57
	Overall	60.51	61.45	61.05	62.63

Table 7: Experimental results with and without adversarial data augmentation. For “BERT”, we directly fine-tune the model on strongly-labeled data without using weak labels. For “virtual”, we use virtual data augmentation in Section 3.1.

	English		German		French	
	Brand	Overall	Brand	Overall	Brand	Overall
BERT	80.77	74.29	67.11	58.21	69.58	61.45
w/o aug.	84.08	77.10	67.69	58.42	70.47	61.59
virtual aug.	83.47	77.62	67.72	58.72	70.32	62.13
DeepTagger	83.45	77.94	68.59	59.24	70.57	62.63

the Web4Ads titles, the average length increases to 20.8. Therefore, the augmented queries contain richer information. As shown in Figure 5 (right), compared with only using the web queries, model performance significantly increases when we add one Web4Ads title. Also, model performance continues to increase when we further increase the number of Web4Ads titles.

7 Related Works

◊ **Weakly-supervised learning.** In weakly-supervised learning, the labels are noisy and incomplete. For example, we can obtain weak labels by writing semantic rules (Yu et al., 2021; Mukherjee and Awadallah, 2020; Awasthi et al., 2020) or match the data to external knowledge bases (Liang et al., 2020; Jiang et al., 2021). Existing methods aim to denoise the labels by using, for example, soft pseudo-labels (Xie et al., 2016, 2020; Meng et al., 2020; Liang et al., 2020; Yu et al., 2021; Zuo et al., 2022), confidence regularization (Pereyra et al., 2017; Jiang et al., 2021), and labeling-function aggregation (Ratner et al., 2017; Varma and Ré, 2018; Lison et al., 2020).

◊ **Adversarial regularization for natural language processing.** Adversarial training was origi-

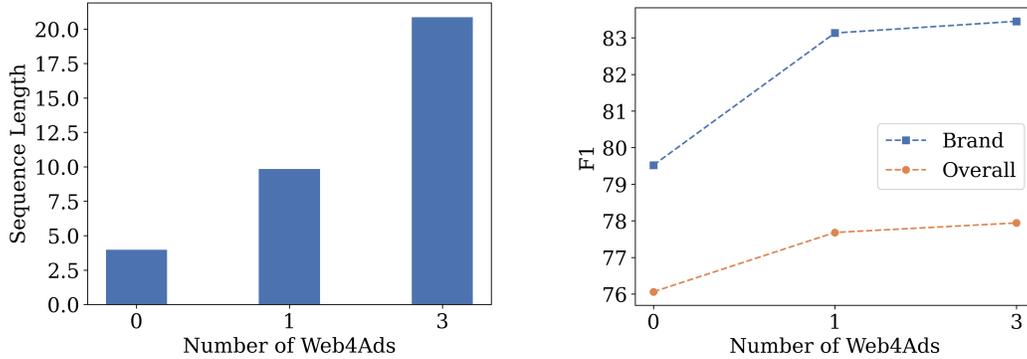


Figure 5: Effectiveness of Web4Ads. Left: average sequence length; Right: model performance on English queries.

nally proposed in computer vision (Szegedy et al., 2014; Goodfellow et al., 2015; Madry et al., 2018), with the goal of training classifiers that are robust to adversarial input images. However, in natural language processing, the goal of adversarial training is to leverage its regularization effect to improve model generalization (Raghunathan et al., 2020). Many works surrounding the efficiency (Zhu et al., 2020; Shafahi et al., 2019; Aghajanyan et al., 2021; Liang et al., 2021; Zuo et al., 2021a) and effectiveness (Cheng et al., 2021a; Liu et al., 2020; Zuo et al., 2021b) of adversarial regularization have been proposed.

8 Conclusion and Discussion

We propose DeepTagger, a knowledge-enhanced NER model for web-based ads queries. DeepTagger leverages both model-free and model-based knowledge enhancement methods. For model-free approaches, we collect unlabeled web queries to inject domain knowledge, and we collect web search titles to complement the short queries. Additionally, we generate weak labels using large language models such as ChatGPT. For model-based approaches, we employ a model-dependent augmentation method based on adversarial training. Extensive experiments in various NER tasks demonstrate the effectiveness of DeepTagger.

In DeepTagger, we adopt a three-stage training framework. In the second stage, we train the model on weakly-labeled data without any modification to the loss function. This is different from previous approaches, e.g., (Yu et al., 2021) uses contrastive learning to suppress noise in weak labels. In practice, we observe only marginal differences when incorporating techniques such as contrastive learning and confidence regularization. We attribute this to the good coverage and quality of weak labels generated by LLMs. We leave further explorations on using LLMs to generate labels as future works.

References

- AGHAJANYAN, A., SHRIVASTAVA, A., GUPTA, A., GOYAL, N., ZETTLEMOYER, L. and GUPTA, S. (2021). Better fine-tuning by reducing representational collapse. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- AWASTHI, A., GHOSH, S., GOYAL, R. and SARAWAGI, S. (2020). Learning from rules generalizing labeled exemplars. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- CHENG, H., LIU, X., PEREIRA, L., YU, Y. and GAO, J. (2021a). Posterior differential regularization with f-divergence for improving model robustness. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online.
- CHENG, X., BOWDEN, M., BHANGE, B. R., GOYAL, P., PACKER, T. and JAVED, F. (2021b). An end-to-end solution for named entity recognition in ecommerce search. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021*. AAAI Press.
- CHIANG, C.-H. and LEE, H.-Y. (2023). Can large language models be an alternative to human evaluations? *ArXiv preprint*, [abs/2305.01937](https://arxiv.org/abs/2305.01937).
- DEVLIN, J., CHANG, M.-W., LEE, K. and TOUTANOVA, K. (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)*. Association for Computational Linguistics, Minneapolis, Minnesota.
- DING, B., QIN, C., LIU, L., BING, L., JOTY, S. and LI, B. (2022). Is gpt-3 a good data annotator? *ArXiv preprint*, [abs/2212.10450](https://arxiv.org/abs/2212.10450).
- DODGE, J., ILHARCO, G., SCHWARTZ, R., FARHADI, A., HAJISHIRZI, H. and SMITH, N. (2020). Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. *ArXiv preprint*, [abs/2002.06305](https://arxiv.org/abs/2002.06305).
- GOODFELLOW, I. J., SHLENS, J. and SZEGEDY, C. (2015). Explaining and harnessing adversarial examples. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.).
- GURURANGAN, S., MARASOVIĆ, A., SWAYAMDIPTA, S., LO, K., BELTAGY, I., DOWNEY, D. and SMITH, N. A. (2020). Don't stop pretraining: Adapt language models to domains and tasks. In *Proceed-*

- ings of the 58th Annual Meeting of the Association for Computational Linguistics. Association for Computational Linguistics, Online.
- JIANG, H., ZHANG, D., CAO, T., YIN, B. and ZHAO, T. (2021). Named entity recognition with small strongly labeled and large weakly labeled data. 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)*. Association for Computational Linguistics, Online.
- LEE, D.-H. ET AL. (2013). Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, vol. 3.
- LIANG, C., YU, Y., JIANG, H., ER, S., WANG, R., ZHAO, T. and ZHANG, C. (2020). BOND: bert-assisted open-domain named entity recognition with distant supervision. In *KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020* (R. Gupta, Y. Liu, J. Tang and B. A. Prakash, eds.). ACM.
- LIANG, X., WU, L., LI, J., WANG, Y., MENG, Q., QIN, T., CHEN, W., ZHANG, M. and LIU, T. (2021). R-drop: Regularized dropout for neural networks. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual* (M. Ranzato, A. Beygelzimer, Y. N. Dauphin, P. Liang and J. W. Vaughan, eds.).
- LISON, P., BARNES, J., HUBIN, A. and TOULEB, S. (2020). Named entity recognition without labelled data: A weak supervision approach. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online.
- LIU, P., YUAN, W., FU, J., JIANG, Z., HAYASHI, H. and NEUBIG, G. (2023a). Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, **55** 1–35.
- LIU, Q., CHEN, N., SAKAI, T. and WU, X.-M. (2023b). A first look at llm-powered generative news recommendation. *ArXiv preprint*, **abs/2305.06566**.
- LIU, X., CHENG, H., HE, P., CHEN, W., WANG, Y., POON, H. and GAO, J. (2020). Adversarial training for large neural language models. *ArXiv preprint*, **abs/2004.08994**.
- MADRY, A., MAKELOV, A., SCHMIDT, L., TSIPRAS, D. and VLADU, A. (2018). Towards deep learning models resistant to adversarial attacks. In *6th International Conference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings*. OpenReview.net.
- MENG, Y., ZHANG, Y., HUANG, J., XIONG, C., JI, H., ZHANG, C. and HAN, J. (2020). Text classification using label names only: A language model self-training approach. In *Proceedings of the*

- 2020 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online.
- MIYATO, T., DAI, A. M. and GOODFELLOW, I. J. (2017). Adversarial training methods for semi-supervised text classification. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- MIYATO, T., MAEDA, S.-I., KOYAMA, M. and ISHII, S. (2018). Virtual adversarial training: a regularization method for supervised and semi-supervised learning. *IEEE transactions on pattern analysis and machine intelligence*, **41** 1979–1993.
- MUKHERJEE, S. and AWADALLAH, A. H. (2020). Uncertainty-aware self-training for text classification with few labels. *ArXiv preprint*, **abs/2006.15315**.
- NISHIKAWA, S., RI, R., YAMADA, I., TSURUOKA, Y. and ECHIZEN, I. (2022). EASE: Entity-aware contrastive learning of sentence embedding. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Seattle, United States.
- OPENAI (2023). Gpt-4 technical report. *arXiv*.
- PASZKE, A., GROSS, S., MASSA, F., LERER, A., BRADBURY, J., CHANAN, G., KILLEEN, T., LIN, Z., GIMELSHEIN, N., ANTIGA, L., DESMAISON, A., KÖPF, A., YANG, E., DEVITO, Z., RAISON, M., TEJANI, A., CHILAMKURTHY, S., STEINER, B., FANG, L., BAI, J. and CHINTALA, S. (2019). Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada* (H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. B. Fox and R. Garnett, eds.).
- PEREYRA, G., TUCKER, G., CHOROWSKI, J., KAISER, Ł. and HINTON, G. (2017). Regularizing neural networks by penalizing confident output distributions. *ArXiv preprint*, **abs/1701.06548**.
- RAGHUNATHAN, A., XIE, S. M., YANG, F., DUCHI, J. C. and LIANG, P. (2020). Understanding and mitigating the tradeoff between robustness and accuracy. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, vol. 119 of *Proceedings of Machine Learning Research*. PMLR.
- RATNER, A., BACH, S. H., EHRENBERG, H., FRIES, J., WU, S. and RÉ, C. (2017). Snorkel: Rapid training data creation with weak supervision. In *Proceedings of the VLDB Endowment. International Conference on Very Large Data Bases*, vol. 11. NIH Public Access.
- REIMERS, N. and GUREVYCH, I. (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)*. Association for Computational Linguistics, Hong Kong, China.
- ROSENBERG, C., HEBERT, M. and SCHNEIDERMAN, H. (2005). Semi-supervised self-training of object detection models.
- SAHU, G., RODRIGUEZ, P., LARADJI, I., ATIGHEHCHIAN, P., VAZQUEZ, D. and BAHDANAU, D. (2022). Data augmentation for intent classification with off-the-shelf large language models. In *Proceedings of the 4th Workshop on NLP for Conversational AI*. Association for Computational Linguistics, Dublin, Ireland.
- SHAFABI, A., NAJIBI, M., GHIASI, A., XU, Z., DICKERSON, J. P., STUDER, C., DAVIS, L. S., TAYLOR, G. and GOLDSTEIN, T. (2019). Adversarial training for free! In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada* (H. M. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. B. Fox and R. Garnett, eds.).
- SZEGEDY, C., ZAREMBA, W., SUTSKEVER, I., BRUNA, J., ERHAN, D., GOODFELLOW, I. J. and FERGUS, R. (2014). Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings* (Y. Bengio and Y. LeCun, eds.).
- TJONG KIM SANG, E. F. (2002). Introduction to the CoNLL-2002 shared task: Language-independent named entity recognition. In *COLING-02: The 6th Conference on Natural Language Learning 2002 (CoNLL-2002)*.
- VARMA, P. and RÉ, C. (2018). Snuba: Automating weak supervision to label training data. In *Proceedings of the VLDB Endowment. International Conference on Very Large Data Bases*, vol. 12. NIH Public Access.
- WANG, J., LIANG, Y., MENG, F., SHI, H., LI, Z., XU, J., QU, J. and ZHOU, J. (2023). Is chatgpt a good nlg evaluator? a preliminary study. *ArXiv preprint*, **abs/2303.04048**.
- WEI, J., WANG, X., SCHUURMANS, D., BOSMA, M., CHI, E., LE, Q. and ZHOU, D. (2022). Chain of thought prompting elicits reasoning in large language models. *ArXiv preprint*, **abs/2201.11903**.
- WEN, M., VASTHIMAL, D. K., LU, A., WANG, T. and GUO, A. (2019). Building large-scale deep learning system for entity recognition in e-commerce search. In *Proceedings of the 6th IEEE/ACM International Conference on Big Data Computing, Applications and Technologies*.
- WOLF, T., DEBUT, L., SANH, V., CHAUMOND, J., DELANGUE, C., MOI, A., CISTAC, P., RAULT, T., LOUF, R., FUNTOWICZ, M., DAVISON, J., SHLEIFER, S., VON PLATEN, P., MA, C., JERNITE, Y., PLU, J., XU, C.,

- LE SCAO, T., GUGGER, S., DRAME, M., LHOEST, Q. and RUSH, A. (2020). Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Online.
- WU, N., GONG, M., SHOU, L., LIANG, S. and JIANG, D. (2023). Large language models are diverse role-players for summarization evaluation. *ArXiv preprint*, **abs/2303.15078**.
- XIE, J., GIRSHICK, R. B. and FARHADI, A. (2016). Unsupervised deep embedding for clustering analysis. In *Proceedings of the 33rd International Conference on Machine Learning, ICML 2016, New York City, NY, USA, June 19-24, 2016* (M. Balcan and K. Q. Weinberger, eds.), vol. 48 of *JMLR Workshop and Conference Proceedings*. JMLR.org.
- XIE, Q., DAI, Z., HOVY, E. H., LUONG, T. and LE, Q. (2020). Unsupervised data augmentation for consistency training. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual* (H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan and H. Lin, eds.).
- XIE, S. M., RAGHUNATHAN, A., LIANG, P. and MA, T. (2022). An explanation of in-context learning as implicit bayesian inference. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- YOO, K. M., PARK, D., KANG, J., LEE, S.-W. and PARK, W. (2021). GPT3Mix: Leveraging large-scale language models for text augmentation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. Association for Computational Linguistics, Punta Cana, Dominican Republic.
- YU, Y., ZUO, S., JIANG, H., REN, W., ZHAO, T. and ZHANG, C. (2021). Fine-tuning pre-trained language model with weak supervision: A contrastive-regularized self-training approach. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. Association for Computational Linguistics, Online.
- ZHANG, D., LI, Z., CAO, T., LUO, C., WU, T., LU, H., SONG, Y., YIN, B., ZHAO, T. and YANG, Q. (2021). Queaco: Borrowing treasures from weakly-labeled behavior data for query attribute value extraction. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*.
- ZHU, C., CHENG, Y., GAN, Z., SUN, S., GOLDSTEIN, T. and LIU, J. (2020). Freelb: Enhanced adversarial training for natural language understanding. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- ZUO, S., LIANG, C., JIANG, H., HE, P., LIU, X., GAO, J., CHEN, W. and ZHAO, T. (2021a). ARCH: Efficient adversarial regularized training with caching. In *Findings of the Association for Com-*

putational Linguistics: EMNLP 2021. Association for Computational Linguistics, Punta Cana, Dominican Republic.

ZUO, S., LIANG, C., JIANG, H., LIU, X., HE, P., GAO, J., CHEN, W. and ZHAO, T. (2021b). Adversarial regularization as stackelberg game: An unrolled optimization approach. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Online and Punta Cana, Dominican Republic.

ZUO, S., YU, Y., LIANG, C., JIANG, H., ER, S., ZHANG, C., ZHAO, T. and ZHA, H. (2022). Self-training with differentiable teacher. In *Findings of the Association for Computational Linguistics: NAACL 2022*. Association for Computational Linguistics, Seattle, United States.