

Few-Sample Named Entity Recognition for Security Vulnerability Reports by Fine-Tuning Pre-Trained Language Models

Guanqun Yang¹, Shay Dineen¹, Zhipeng Lin¹, and Xueqing Liu¹

Department of Computer Science,
Stevens Institute of Technology, Hoboken NJ 07030, USA
{gyang16, sdineen1, zlin30, xliu127}@stevens.edu

Abstract. Public security vulnerability reports (e.g., CVE reports) play an important role in the maintenance of computer and network systems. Security companies and administrators rely on information from these reports to prioritize tasks on developing and deploying patches to their customers. Since these reports are unstructured texts, automatic information extraction (IE) can help scale up the processing by converting the unstructured reports to structured forms, e.g., software names and versions [8] and vulnerability types [38]. Existing works on automated IE for security vulnerability reports often rely on a large number of labeled training samples [8,18,48]. However, creating massive labeled training set is both expensive and time consuming. In this work, for the first time, we propose to investigate this problem where only a small number of labeled training samples are available. In particular, we investigate the performance of fine-tuning several state-of-the-art pre-trained language models on our small training dataset. The results show that with pre-trained language models and carefully tuned hyperparameters, we have reached or slightly outperformed the state-of-the-art system [8] on this task. Consistent with previous two-step process of first fine-tuning on main category and then transfer learning to others as in [7], if otherwise following our proposed approach, the number of required labeled samples substantially decrease in both stages: 90% reduction in fine-tuning from 5758 to 576, and 88.8% reduction in transfer learning with 64 labeled samples per category. Our experiments thus demonstrate the effectiveness of few-sample learning on NER for security vulnerability report. This result opens up multiple research opportunities for few-sample learning for security vulnerability reports, which is discussed in the paper. Our implementation for few-sample vulnerability entity tagger in security reports could be found at <https://github.com/guanqun-yang/FewVulnerability>.

Keywords: Software Vulnerability Identification · Few-sample Named Entity Recognition · Public Security Reports

1 Introduction

Security vulnerabilities pose great challenges to software and network systems. For example, there are numerous reported data breaches of Uber, Equifax, Mar-

riott, and Facebook that jeopardize hundreds of millions of customers' information security [15]; vulnerable in-flight automation software used in Boeing 737-MAX are found to be guilty for a series of crashes in 2018 [31]. To allow the tracking and management of security vulnerabilities, people have created public security vulnerability reports and store them in vulnerability databases for reference. For example, the Common Vulnerabilities and Exposures (CVE) database is created and maintained by MITRE ¹; the National Vulnerability Database (NVD) is another database maintained by the U.S. government. These databases have largely facilitated security companies and vendors to manage and prioritize deployment of patches.

To scale up the management of vulnerability entries, existing works have leveraged natural language processing (NLP) techniques to extract information from unstructured vulnerability reports. For example, software name, software version, and vulnerability type. The extraction of such information can help accelerate security administration in the following scenarios: First, after converting the unstructured reports into structured entries such as software name and software version, such information can be directly used for detecting the inconsistencies between different records of vulnerability reports. For example, significant inconsistency was detected between reports from the CVE and NVD databases [8], which may cause system administrators to retrieve outdated or incorrect security alerts, exposing systems under their watch to hazard; inconsistency was also detected between two entries created by the same vendor Samsung in [10]; Second, the extracted entries can be used in other downstream applications. For example, we can leverage different categories of entries (e.g., vendor name, vulnerability types, attacker name) to construct a knowledge ontology for modeling the interplay between security entries [16,2]. The extracted entries have also been used for automatically generating a natural language summarization of the original report by leveraging a template [38].

The scalable processing of vulnerability reports into structured forms relies on the support from an effective information extraction (IE) system. With the massive amount of reports, manual extraction of information is intractable and most of the existing systems rely on machine learning approaches. For example, Mulwad et al. [30] built an SVM-based system which extracted security concepts from web texts. In the recent years, with the rise of deep neural networks, systems based on deep learning were proposed, e.g., by using LSTM [18,48] and bi-directional GRU [8]. Deep learning has been shown to be effective in this task, e.g., with an F1 score of 98% [8].

Nevertheless, a major challenge of applying existing systems (e.g., recurrent neural network like GRU [8] and LSTM [18,48]) is that they often require a large number of labeled training samples to perform well; this number could range from a few thousand [47] to tens of thousand [18]. Labeling a NER dataset specific to computer security domain at this scale is both costly and time-consuming. This problem of annotation is further exacerbated by the fact that new security

¹ <https://www.mitre.org/>

reports are routinely generated. Therefore, it is vital to investigate the possibility to alleviate the burden of labeling huge corpus, i.e., few-sample learning.

In this paper, for the first time, we investigate the problem of information extraction for public vulnerability reports under the setting where only a small number of labeled training samples are available. Our goal is to minimize the number of labeled samples without suffering from a significant performance degradation.

In recent years, significant progress has been made for few-shot learning, e.g., image classification [12], text classification [17], text generation [3], etc. In these tasks, the training set and validation set often only contain ten or fewer examples, and it has been shown that the available small number of examples are all what the model needs to achieve a good performance [5]. To generalize to unseen examples, few-shot learning approaches often leverage meta-learning [12], metric learning [36,6] and model pre-training. For example, the GPT-3 model [3] contains 175 billion parameters and was pre-trained on 45TB of unlabeled text data. It is able to accomplish a variety of generative tasks (e.g., question answering, conversation generation, text to command) when the user only provides a few examples for demonstration. For information extraction, more specifically, named entity recognition (NER), several existing works study its few-shot setting, where only tens of labeled samples are provided during training [44,21]. However, NER systems to date can only achieve a 62% F1 score (evaluated on the CoNLL 2003 dataset [39]). When applied to the security vulnerability reports domain, it is thus unclear whether the performance can be satisfactory enough to security administrators and vendors, as a lower error rate is often required [8]. As a result, in this paper, we investigate the following research questions:

- Is it possible to match the state-of-the-art performance in NER for security vulnerability reports by using only a small number of labeled examples?
- What is the minimum number of labeled examples required for this problem?

To answer these research questions, we conduct an experimental study on *few-sample* named entity recognition for vulnerability reports², which to the best of our knowledge has not been explored by previous works. As an initial study, we focus on the VIEM dataset [8], a dataset of vulnerability reports containing 3 types of entities (i.e., *software name*, *software version*, *outside*) and 13 sub-datasets based on the category of vulnerabilities. For the largest sub-dataset, we find that through fine-tuning pre-trained language models, it is possible to match the state-of-the-art performance (i.e., 0.92 to 0.93 F1 score for different entity types [8]) using only 576 labeled sentences for the training dataset, or 10% of the original training data. For the other 12 categories, the same approach matches the SOTA performance with only 11.2% of the labeled samples. Our experiments

² Here we frame our problem under the term "few-sample" learning instead of "few-shot" learning because our approach generally requires tens to a few hundred labeled training samples, which is more than "few-shot". We adopt the name "few-sample" from [47].

also show that the simple fine-tuning approach works better than state-of-the-art systems for few-shot NER [44].

The main contributions of this paper are:

- We propose, for the first time, to study the problem of named entity recognition (NER) for security vulnerability reports;
- We perform an experimental study by fine-tuning three state-of-the-art language models and evaluating one metric learning-based method for few-sample NER in vulnerability reports, which shows that the fine-tuning approach can achieve the similar performance using only 10% or 11.2% of the original training dataset;
- We discuss multiple future research opportunities of information extraction for security vulnerability reports;

The rest of this paper is organized as follows. Section 2 defines the problem and discusses the domain and data-specific challenges; Section 3 introduces the methods examined; Section 4 describes the experimental steps, datasets, and experimental results; Section 5 analyzes the related work; finally, Section 6 draws the conclusion, analyzes the limitation of this work and discusses several future directions of research for few-sample learning for security vulnerability reports.

2 Problem Definition and Challenges

2.1 Few-Sample Named Entity Recognition

Name entity recognition (NER) is the task of assigning predefined entity names to tokens in the input text. Formally, with predefined tagging set \mathcal{Y} and a sequence of tokens (input sentence) $X_i = [x_{i1}, x_{i2}, \dots, x_{iT}]$, each x_{ik} ($k = 1, 2, \dots, T$) corresponds to a tag $y_{ik} \in \mathcal{Y}$ specifying the its entity type. This gives X_i a tagging sequence $Y_i = [y_{i1}, y_{i2}, \dots, y_{iT}]$ and a typical dataset of NER is of the form $\mathcal{D} = \{(X_i, Y_i)\}_{i=1}^N$.

Notice there exist two distinct definitions for what a few-shot learning algorithm tries to achieve [41]. In the first definition, the goal is to match the performance when only a subset of training data is available; in the second definition, the test set contains unseen classes in the training set and the algorithm needs to generalize to the unseen classes. In this work, we focus on the first setting. We refer readers interested in the details of these two settings to [47].

2.2 Named Entity Recognition for Vulnerability Reports

An example of security vulnerability report (CVE-2015-2384) is shown in Figure 1, which is a snapshot from the NVD website³. In this paper, we focus on

³ <https://nvd.nist.gov>

🚩 CVE-2015-2384 Detail

Current Description

Microsoft **Internet Explorer 11** allows remote attackers to execute arbitrary code or cause a denial of service (memory corruption) via a crafted web site, aka "**Internet Explorer** Memory Corruption Vulnerability," a different vulnerability than CVE-2015-2383 and CVE-2015-2425.

 Software Name(SN)

 Software Version(SV)

 Outside(O)

QUICK INFO

CVE Dictionary Entry:

CVE-2015-2384

NVD Published Date:

07/14/2015

NVD Last Modified:

10/12/2018

Source:

Microsoft Corporation

Fig. 1: An example of public security vulnerability report

the VIEM dataset from [8], which contains 13 categories of vulnerabilities and two types of named entities: software name and software version. That is, the tag set \mathcal{Y} is instantiated as $\{\text{SN}, \text{SV}, \text{O}\}$, where SN, SV, O refer to software name, software version, and non-entities, respectively. For example, as is shown in Figure 1, the report CVE-2015-2384 warrants a security alert where the software *Internet Explorer* of version 11 might expose its vulnerability to attackers. The training, validation, and testing set share the same tag set \mathcal{Y} .

2.3 Data-Specific Challenges

There exists several data-specific challenges in the VIEM dataset, which we summarize as below.

Contextual Dependency When identifying vulnerable software names and versions in vulnerability reports, the tagger has to take the entire context into account. For example, as is shown in Figure 1 (CVE-2015-2384), *Internet Explorer* is tagged as a software name for its first occurrence; in the second occurrence, it appears in the vulnerability name *Internet Explorer Memory Corruption Vulnerabilities* thus is no longer a software name. Software version tag is also context dependent: in the examples below, 1.0 is tagged SV in CVE-2002-2323 and CVE-2010-3487 but O in CVE-2006-4710, this is because when the software name (James M. Snell) is tagged O, the software version should also be tagged O.

- CVE-2002-2323: *Sun PC NetLink 1.0 through 1.2 ... access restriction.*
- CVE-2010-3487: *Directory traversal vulnerability in YelloSoft Pinky 1.0 ... in the URL.*

- CVE-2006-4710: *Multiple cross-site scripting ... as demonstrated by certain test cases of the James M. Snell Atom 1.0 feed reader test suite.*

Data Imbalance The goal of NER system is to identify entities of interest and tag all the other tokens as outsider. It is therefore possible that the entire sentence includes only outsider tokens. As is shown in Table 1, in the general domain NER dataset like CoNLL 2003 and OntoNotes 5, the proportion of such sentences in training set is relatively low. On the other hand, in the vulnerability reports, the proportion of all-0 sentences is high. Such data imbalance problem can cause the classifier to trade off minority group’s performance for that of majority group, leading to a higher error rate for minority samples [23].

Table 1: The proportion of sentences in training set that include only non-entities. The statistics of CoNLL 2003 and OntoNotes 5 are estimated from publicly available dataset releases. In the VIEM dataset, the vulnerability category `memc` is intentionally balanced in [8] to ease transfer learning, making the aggregate (i.e. VIEM (all)) non-entity-only proportion lower.

Dataset	Domain	Non-entity-only sentences
CoNLL 2003 ⁴	General	20.71%
OntoNotes 5 ⁵	General	51.60%
VIEM (all)	Security	60.34%
VIEM (w/o <code>memc</code>)	Security	75.44%

3 Few-Sample NER for Vulnerability Reports

In this section, we introduce the methods investigated in this paper. The first method leverages fine-tuning pre-trained language models (PLM); the second method is a few-shot NER method based on the nearest neighbor approach [44].

3.1 Fine-Tuning Pre-trained Language Models with Hundreds of Training Labels

Fine-Tuning Framework and Selection of Pre-trained Language Models Before we describe our approach, first we introduce our choices for the language model to be fine-tuned. There exists hundreds of candidate models from the HuggingFace model hub⁶. We investigate three models out of all: `bert-large-cased`, `roberta-large`, and `electra-base-discriminator`. BERT

⁴ <https://huggingface.co/datasets/conll2003>

⁵ <https://github.com/Fritz449/ProtoNER/tree/master/ontonotes>

⁶ <https://huggingface.co/models>

was the first pre-trained language model that achieved great performance compared to prior art; RoBERTa and Electra are among the top-ranked language models on the GLUE leaderboard as of Jan 2021 ⁷. For BERT, we choose the cased version `bert-large-cased` rather than the uncased counterpart (i.e. `bert-large-uncased`) as token casings could be informative for the identification of software names (SN). Our implementation is based on the HuggingFace `transformers` library (4.5.1), more specifically, we leverage the `AutoModelForTokenClassification` class for fine-tuning the 13 categories in the VIEM dataset.

Fine-Tuning on the memc Category Following the experimental protocol of VIEM system [8], the PLM in our system is first fine-tuned on the vulnerability category of `memc` (short for "memory corruption") as, from the statistics of the dataset shown in Table 2, this category suffers less from data imbalance described in section 2; fine-tuning on this category could therefore best help PLMs learn representations specific to computer security domain. However, rather than allowing PLMs accessing the entire training dataset of `memc`, we constrain its access to the training data from 1% to 10% of the entire training set.

Transfer Learning to the Other 12 Categories Previous works on fine-tuning language models show that fine-tuning can benefit from multi-task learning or transfer learning, i.e., fine-tuning on a related task before fine-tuning on the objective task [34]. We adopt this approach, i.e., the PLM is first fine-tuned on the `memc` category, then we start from the checkpoint and continue the fine-tuning on the 12 categories. We again limit the number of annotated samples during transfer learning. Throughout the transfer learning experiment, we sample varying number of training samples from each of the 12 categories, and we combine the 12 subsets to form an *aggregate* training set. For validation and testing, however, we evaluate the model with transfer learning on each category *independently*.

Hyperparameter Optimization The hyperparameters for fine-tuning includes the learning rate, batch size, and the number of training epochs. Inappropriately setting one of these hyperparameter might cause performance degradation [28]. As a result, hyperparameter tuning is critical as they could have substantial impact on the model performance.

Setting the Random Seed The fine-tuning performance for PLM are known to be unstable, especially under the typical few-sample setting. Some seemingly benign parameters, like choice of random seed, could introduce sizable variance on the model performance. Indeed, the seed choice could influence weight initialization for task-specific layers when applying PLMs for downstream task [7].

⁷ <https://gluebenchmark.com/>

When these task-specific layers are initialized poorly, coupled with small size of the training set, it is unlikely for the network to converge to a good solution without extensive training. In fact, previous work has shown that some seeds are consistently better than others for network optimization [7]. But it is not possible to know a priori which seed is better than the others.

3.2 Few-shot Named Entity Recognition

After trying the fine-tuning methods above, we find that these methods still require hundreds of training examples, one question is whether it is possible to further reduce the number of training labels, so that only a few training labels are required. For example, the system designed by Yang and Katiyar [44], named **StructShot**, manages to beat competitive systems with a few samples, i.e. only 1 to 5 annotated named entities are required for each named entity class.

The key technique used in **StructShot** is the nearest neighbor search. Specifically, rather than directly training a PLM with token classification objective as we do, their system is first fine-tuned on a high-resource yet distantly related dataset; then they use the model as a feature extractor to obtain token representations of each of tokens in the test set and a small support set, both coming from low-resource task under investigation. With these two sets of representations in hand, test token’s tag is determined through retrieving the most similar token’s tag in the support set. A CRF decoding procedure follows the similarities estimated from previous steps to give the optimal tag sequence.

4 Experiments

4.1 Datasets

We experiment on the VIEM dataset provided by [8]. This dataset is collected from the public vulnerability databases Common Vulnerabilities and Exposure (CVE) and numerous security forums. The goal of this dataset is to identify vulnerable software names and software versions. Therefore, the tag set \mathcal{Y} includes *SN* (software names) and *SV* (software versions). Similar to the other NER datasets, the *O* (outside) tag is used to represent all other tokens except vulnerable software names and versions. The dataset is manually labeled by three security experts in hope of minimizing labeling errors and enforcing consistent annotation standards. The records in dataset range in 20 years and contain all 13 categories listed on CVE website. The statistics of the dataset is shown in Table 2. We only report the average statistics of 12 categories here, and we refer readers interested in detailed statistics of these categories to Table 6. From the statistics of the dataset shown in Table 2, we observe the *dataset imbalance*, a challenge discussed in section 2.3, in different levels.

- **Sample Number Imbalance** The *memc* category has significantly more samples in every split than the other 12 categories.

- **Entity Number Imbalance** The `memc` category has a predominately higher entity proportion in both sentence and token level than the other 12 categories.
- **Entity Type Imbalance** The appearances of `SV` tokens are generally more frequent than `SN`; sometimes this difference could vary up to 4.38% (see valid split of `memc`).

Table 2: Statistic of the VIEM dataset.

Category Split	<code>memc</code>			Average of other 12 categories		
	train	valid	test	train	valid	test
Number of sample	5758	1159	1001	468.00	116.67	556.25
Sentence-level entity proportion	0.5639	0.3287	0.4555	0.2435	0.2677	0.2620
Token-level entity proportion (<code>SN</code>)	0.0613	0.0368	0.0559	0.0214	0.0240	0.0237
Token-level entity proportion (<code>SV</code>)	0.0819	0.0807	0.0787	0.0308	0.0331	0.0331

Preparing for the Few-Sample Dataset Among the 13 vulnerability categories in the dataset official release, only `memc` includes official validation split. In order to enable hyperparameter optimization and model selection on the other 12 categories, we randomly sample 10% of the training samples as the held-out validation set.

As we strive to investigate the performance of PLMs with the reduced number of annotated samples after fine-tuning and transfer learning, we create the training and validation set by sampling a subset of the full dataset. More specifically, for training set:

- **Fine-Tuning** We vary the proportion of sampled data from the `memc` training set from 1% to 10% (equivalent to a sample size from 58 to 576) and investigate the the proportion to measure the adaptation of PLM to the computer security domain.
- **Transfer Learning** We vary the number of samples in the other 12 categories from $\{32, 64, 128, 256\}$ per category. Compared to fine-tuning, we re-train from sampling by proportion considering the fact that there are much smaller training set for the other 12 categories than the `memc` category.

For the validation set, we make sure it has the same size as the training set whenever possible: the infeasible cases are those when the size of the sampled training set exceeds the 10% limit. This comes from the consideration that the validation samples also require annotations and therefore consume labeling budget. Apart from the conversion of text form into data format required by model

training, we do not perform any additional data processing to make sure the comparison is fair with previous system. Throughout the data preparation, we fix the random seed to 42.

4.2 Evaluation Metrics

Following the evaluation metrics used in [8], we use the precision, recall, and F1 score to evaluate our NER system, the definitions of precision, recall and F1 score are standard.

4.3 Experimental Setup

Hyperparameter Optimization Settings We leverage the built-in API for hyperparameter optimization from the HuggingFace library [42]. We leverage the grid search where the search space is as follows: the learning rate is selected from $\{1e-6, 5e-6, 1e-5\}$; the number of training epochs is chosen from $\{3, 5\}$; the batch size, is fixed to 2 per device because we find the resulting additional training iterations favorable as compared to setting it to 4 or 8 in our pilot experiment. All the other hyperparameters values are fixed to the default values in HuggingFace’s `transformers-4.5.1`. Besides the grid search, we also experiment with the Bayesian optimization, but it fails to outperform the grid search with same computational resources. We fix the seed for model training to 42. We randomly restart every experiment 10 times and report the average score. We use the half-precision (i.e. `fp16`) mode to accelerate the training. During the hyperparameter optimization, we select the best checkpoint within each trial. For trials with different epoch numbers, we partition the training iterations so that the same number of model checkpoints are saved during training across different experiments. Due to the class imbalance problem (Section 4.1), measuring checkpoint quality solely based on metric of the SN or SV is suboptimal; therefore, we weigh their F1 score based on the number of their appearances in the groundtruth, i.e. N_{SN} and N_{SV} .

$$\bar{F}_1 = \frac{N_{SN}F_{SN} + N_{SV}F_{SV}}{N_{SN} + N_{SV}}$$

This single evaluation metric provides a tradeoff between metrics on different entities during model selection.

Evaluated Methods Our experiments contain two parts: first, we evaluate the performance of few-sample learning on the `memc` category; second, we evaluate the performance of few-sample learning on the other 12 categories.

For the `memc` category, we fine-tune the models with varying proportions of the `memc` training data from 1% to 10%. We are interested to know whether fine-tuning with a small proportion of samples is feasible for PLMs to reach the performance of competitive system like VIEM.

For the other 12 categories, we compare the following 4 settings:

- **FT** This setting is the direct application of models fine-tuned on `memc` to the other categories. We expect the performance of **FT** to be the lowest across 4 settings.
- **FT+SS** This setting is consistent with the way this system is used in [44]: after the PLM is fine-tuned on the `memc` category, it is coupled with `StructShot` and directly applied to the other 12 categories without any transfer learning.
- **FT+TL** This setting is to apply transfer learning on other categories on the best fine-tuned checkpoint on the `memc` category. With the help of transfer learning, we expect the PLMs have satisfying performance across different vulnerability categories.
- **FT+TL+SS** This setting is built upon the **FT+TL** setting by allowing an additional application of `StructShot` on the model that experiences both fine-tuning and transfer learning. We expect a better performance of this setting over the **FT+TL** setting because of `StructShot`.

4.4 Experimental Results: Fine-Tuning on the `memc` Category

Summary of Main Results The fine-tuning performance is detailed in Table 3. We could observe that the absolute difference of F1 score for both `SN` and `SV` are both within 1% absolute difference of the `VIEM` system. Importantly, the amount of data we use is only 10% of the `VIEM` system.

Table 3: Fine-tuning results for the `memc` category

	SN precision	SN recall	SN f1-score	SV precision	SV recall	SV f1-score
<code>VIEM</code> (full data)	0.9773	0.9916	0.9844	0.9880	0.9923	0.9901
<code>BERT</code> (10% data)	0.9582	0.9808	0.9694	0.9818	0.9848	0.9833
<code>RoBERTa</code> (10% data)	0.9567	0.9808	0.9686	0.9841	0.9844	0.9843
<code>Electra</code> (10% data)	0.9677	0.9754	0.9715	0.9830	0.9890	0.9860

To investigate how fine-tuning helps PLMs adapt to new domain and pick up task-awareness, we show a 2D projection of token-level hidden representation obtained on `Electra` for `memc` training set before and after fine-tuning in Figure 2. We could see that much more compact clusterings of `SN` and `SV` are obtained; fine-tuning informs the PLMs with the underlying patterns of task dataset.

Fine-Tuning with Different Labeled Training Sample Size We vary the number of labeled samples per category from 1% to 10% as used in `VIEM` system to see how different sizes of training set affect the performance of fine-tuning.

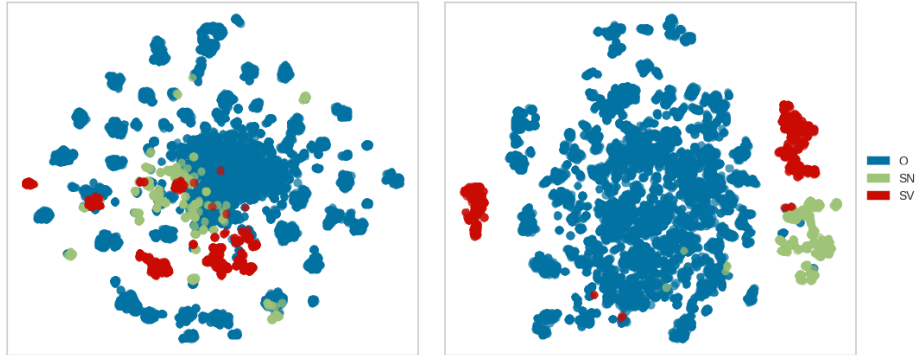


Fig. 2: The t-SNE visualization of `memc` training set in token level before (left) and after (right) fine-tuning. Fine-tuning gives more compact clusters for both SN and SV, illustrating the critical role of fine-tuning.

As is shown in Figure 3, when the proportion of labeled training samples increase from 1% to 10% for fine-tuning, the performance generally improves; but the improvement does not come monotonically. For example, in cases like fine-tuning RoBERTa with 5% of the data, the performance is worse than that of only 3%.

There is also a disparity on the fine-tuning performance of SN and SV. Specifically, the SV could have an F1 of around 95% with as small as 1% of the training samples available across three different models. This does not hold for the SN: under the 1% training data availability, the F1 for SN does not exceed 92% for BERT; for the RoBERTa and Electra, this metric is less than 88%. The disparity on the performance shows that correctly tagging SV is easier than tagging SN. This might arise from the fact that there are naming conventions of software versions for the NER tagger to exploit while similar conventions do not exist for software names.

Another general tendency is the relation between precision and recall: across different models and the tag of interest, the recall is generally higher than precision. This is desirable during the deployment for security applications as incorrectly tagging regular software as vulnerabilities is acceptable, only leading to more manual efforts to double check individual software, while missing vulnerabilities could leave the critical systems unattended.

4.5 Experimental Results: Transfer Learning on the Other 12 Categories

Summary of the Main Results As we could observe from the Table 4, transfer learning with the aggregation of 64 samples per category matches the performance of VIEM system; the absolute difference is less than 1% in F1 score. Despite less than 1% of improvement with respect to our system, the VIEM system is trained with the entire dataset with an aggregation of 7016 samples in

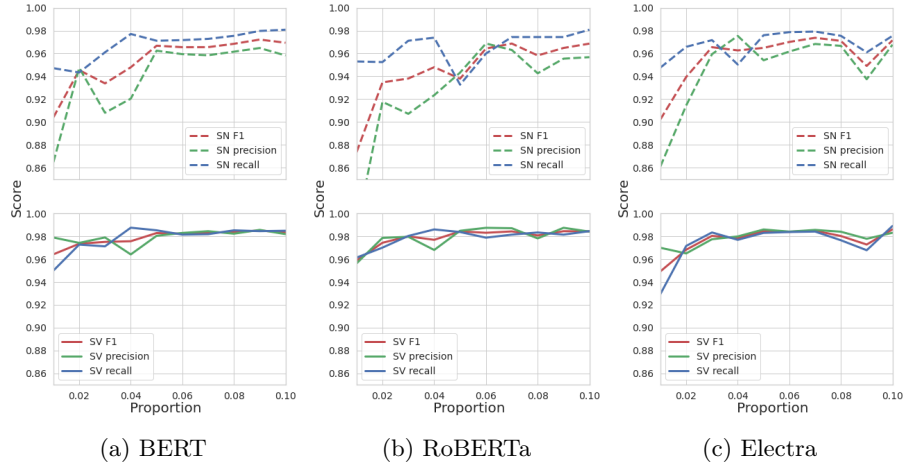


Fig. 3: The influence of training set size on fine-tuning.

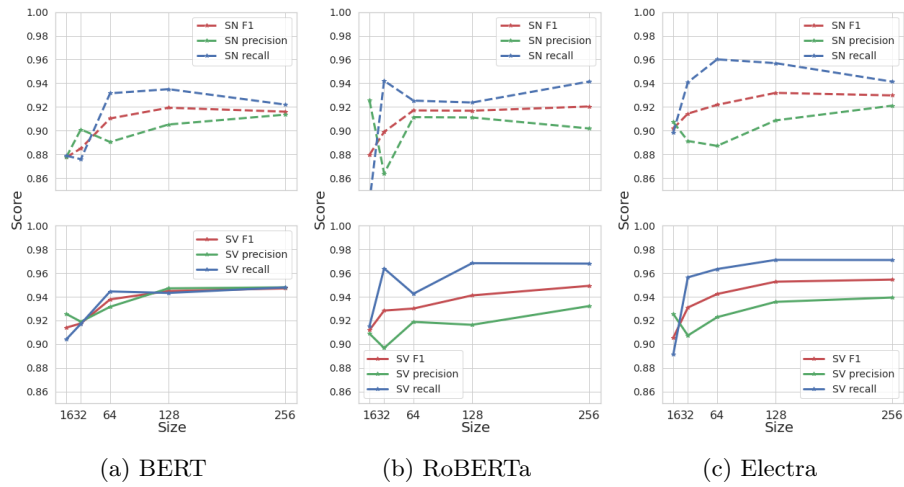


Fig. 4: The influence of training set size on transfer learning.

total (see Figure 6 for details; note the VIEM dataset merges train and valid split as its train split). Comparing the data utilization of our system to VIEM system, we reduce the number of annotations by 88.80%.

StructShot could help with the *precision* of **SN** but struggle with the *recall*. Introducing **StructShot** has caused a degradation in F1 score, which casts the doubt on the applicability of this system to our application. Specifically, the absolute gap in F1 score could be up to 20% (compared to **SN** and **SV**'s F1 scores in **FT+SS** and **FT+TL** for Electra). This gap still exists even with additional transfer learning (**FT+TL+SS**). Specifically,

- Comparing BERT and RoBERTa's **FT+TL+SS** with **FT+TL** in **SN** and **SV**'s precision, adding **SS** on top of either **FT** or **FT+TL** could improve the precision for both **SN** and **SV**. However, this is not always the case: adding **SS** to Electra's **FT+TL** actually degrades both the precision and recall for **SN** and **SV**.
- In cases the precision does improve, as when comparing BERT and RoBERTa's **FT+TL+SS** with **FT+TL**, a minor improvement of precision is often at the cost of a larger degradation of recall, collectively leading to a drop in F1 score; this is especially the case for **SN**.
- Comparing **FT+SS** and **FT+SS+TL** for all three models' precision and recall, additional transfer learning could remedy the degradation caused by the drop of recall discussed above, but this tendency is not reverted.

Transfer Learning with Different Training Sample Size We conduct a similar performance analysis when varying the training sample size as in fine-tuning; we visualize the average performance metrics since we apply transfer learning on fine-tuned model on the other 12 categories. As is shown in Figure 4, when the number of samples range from 32, 64 through 256, more samples do contribute to improvements on transfer learning, but the gains are less clear when we increase the training set size beyond 64 per category: the performance reaches plateau with more than 64 samples (see the F1 score of both **SN** and **SV** across three models).

We could observe a similar performance disparity of **SN** and **SV** as in fine-tuning on the **memc** category (see section 4.4). Specifically, we could use as few as 16 samples per category to achieve an F1 score around 91% while the same performance is not attainable for **SN** until we enlarge the training set to 64 (for BERT and RoBERTa) or even 128 (for Electra) samples per category.

The relation between precision and recall are consistent in transfer learning and fine-tuning; the recalls are both higher than the precision scores in most of the cases.

Training Set Sampling for Transfer Learning with Single Category and 12 Category Aggregate The performance metrics detailed in Table 4 show that we could reach the performance of competitive systems by aggregating the

64 samples from each of the 12 categories, resulting in a training set. In addition, we also evaluate the case where we only use the 64 examples in each category for training (and testing) without aggregation. We could see from the Table 5 that aggregating the training samples is favorable for transfer learning, improving the system performance by 2% to 3%.

Table 4: Transfer learning results averaged through 12 categories.

	SN precision	SN recall	SN f1-score	SV precision	SV recall	SV f1-score
VIEM (w/o transfer)	0.8278	0.8489	0.8382	0.8428	0.9407	0.8891
VIEM (w/ transfer)	0.9184	0.9286	0.9235	0.9382	0.9410	0.9396
BERT						
FT	0.8759	0.7950	0.8318	0.8813	0.9341	0.9060
FT+SS	0.9035	0.6843	0.7623	0.9001	0.7880	0.8183
FT+TL	0.8945	0.9302	0.9116	0.9373	0.9354	0.9355
FT+TL+SS	0.9060	0.7637	0.7766	0.9374	0.9135	0.9235
RoBERTa						
FT	0.8558	0.8510	0.8524	0.8822	0.9308	0.9052
FT+SS	0.8905	0.6810	0.7532	0.9106	0.8489	0.8735
FT+TL	0.8749	0.9409	0.9063	0.9162	0.9561	0.9354
FT+TL+SS	0.8841	0.7853	0.7989	0.9186	0.9364	0.9265
Electra						
FT	0.8656	0.8236	0.8416	0.8852	0.9246	0.9037
FT+SS	0.8692	0.6806	0.7274	0.8355	0.7428	0.7420
FT+TL	0.9003	0.9443	0.9214	0.9264	0.9494	0.9369
FT+TL+SS	0.8933	0.8263	0.8315	0.9035	0.9477	0.9233

Table 5: Comparison of transfer learning results with training set sampled from individual and aggregate of 12 categories. The results are averaged through 12 categories.

	SN precision	SN recall	SN f1-score	SV precision	SV recall	SV f1-score
RoBERTa						
FT+TL (individual)	0.8483	0.8953	0.8688	0.8816	0.9342	0.9059
FT+TL (aggregate)	0.8749	0.9409	0.9063	0.9162	0.9561	0.9354
Electra						
FT+TL (individual)	0.8670	0.9247	0.8940	0.9015	0.9346	0.9169
FT+TL (aggregate)	0.9003	0.9443	0.9214	0.9264	0.9494	0.9369

Analysis on the Two Data-Specific Challenges We discussed in section 4.1 that the VIEM dataset contains the challenges in *context dependency* and

dataset imbalance. Even though we do not explicitly address either of them, our experiments show that they have not largely affected the performance. More specifically,

- **Dataset Imbalance** The average proportion of sentences with entities (SN and SV) in training set is only 24.35%; this might be a concern of performance due to the sparsity of informative tokens. However, as is shown in Table 4, the performance of models reaches or outperforms the VIEM system with an average F1 score of more than 0.9, which shows that the dataset imbalance problem does not hurt the effectiveness of our system.
- **Context Dependency** Our experiments show that transfer learning can help address the context dependency problem. Specifically, for the examples below, the model without fine-tuning misclassifies the tokens *Network Data Loss Prevention*, *Cisco Wireless LAN Controller*, and *PineApp Mail-SeCure* as non-entities (0).
 - CVE-2014-8523 (from `csrf` category): *Cross-site request forgery (CSRF) vulnerability in McAfee Network Data Loss Prevention (NDLP) before 9.3 allows remote ...*
 - CVE-2015-0690 (from `xss` category): *Cross-site scripting (XSS) vulnerability ... on Cisco Wireless LAN Controller (WLC) devices ... aka Bug ID CSCun95178.*
 - CVE-2013-4987 (from `gainpre` category): *PineApp Mail-SeCure before 3.70 allows remote authenticated ...*

Troubleshooting StructShot The unsatisfactory performance of **StructShot** shown in Table 4 alerts caution to which scenario it could be applied. We therefore conduct an error analysis of **StructShot** to pinpoint its weakness. In our controlled experiment, we find that the **StructShot** is sensitive to adversarial examples. More specifically, as is shown in Figure 5, when we manually create the support set by sequentially adding sample one after another, one single sample 0893580d85 could bring down overall performance by over 70%.

One benefit of the **StructShot** reported in [44] is its better capability of assigning non-entity tag (i.e. 0 tag) as these tokens do not carry unified semantic meaning. However, because of the *context dependency* illustrated in Figure 1, in the NER of public security reports, non-entity tokens in one context might be entities in another one, causing confusion for the nearest neighbor-based tagger. We therefore suspect these confusing non-entities are the causes of the observed phenomenon. We monitor the following error cases to validate this hypothesis. Specifically, for CVE-2017-6038, the *Beldn Hirschmann GECKO Lite Managed switch* is originally tagged correctly as SN. However, when an additional

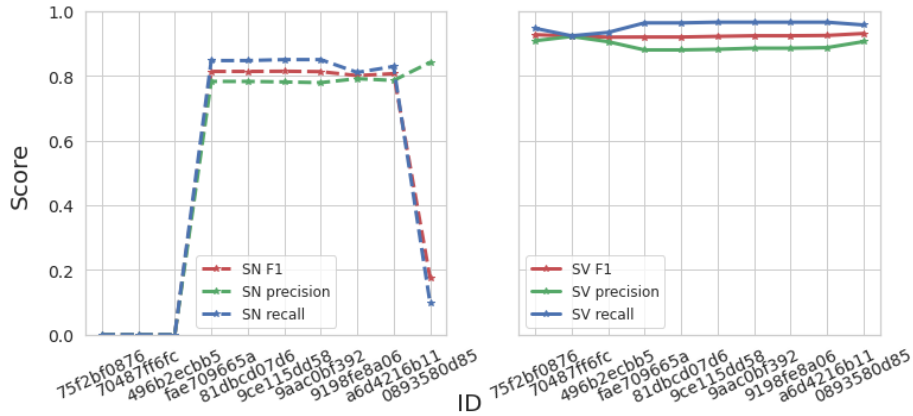


Fig. 5: **StructShot** is sensitive to adversarial samples. When testing on `csrf` category, an additional sample `0893580d85` in the support set could lead to a degradation of 70% of F1 score for SN.

`0893580d85` sample is introduced, all of its predictions are incorrectly reverted to 0. Similarly for CVE-2014-9665, the correct predictions of software names *IB6131* and *Infiniband Switch* are changed to 0 due to the adversarial sample `0893580d85`.

- CVE-2017-6038: A Cross-Site Request Forgery issue was discovered in Belden Hirschmann GECKO Lite Managed switch, Version 2.0.00 and prior versions .
- CVE-2014-9565: Cross-site request forgery (CSRF) vulnerability in IBM Flex System EN6131 40Gb Ethernet and IB6131 40Gb Infiniband Switch 3.4.0000 and earlier.

5 Related Work

5.1 Information Extraction in Public Vulnerability Database

The VIEM system designed by Dong et al. is among the first batch of works that extracts information from public vulnerability database through combination of named entity recognition (NER) and relation extraction (RE) approaches. The extracted information is used to identify vulnerabilities and thereby inconsistencies between two major security databases - CVE and NVD [8].

Prior works of the information extraction in security domain mostly focus on other aspects like dataset curation, data analysis, and visualization. Specifically, Fan et al. create a C/C++ vulnerability dataset under the aid of information provided by the CVE database. Works by Ponta et al. and Wu et al. lead to

similar datasets [35,9]. All of them are facilitated by the CVE database during curation while the main differences are programming language and the dataset scale. Yamamoto et al. try to estimate the severity of vulnerability from the descriptions of CVE reports [43]. They also propose to improve the estimation accuracy by incorporating the relative importance of reports across different years. Pham and Dang implement an interactive tool named `CVEexplorer` that enables the network and temporal view of different vulnerabilities in the CVE database [33].

This work follows the setting of VIEM system but tries to improve the vulnerability identification by minimizing the data labeling cost while matching the performance of existing systems.

5.2 Named Entity Recognition for Computer Security

The role of NER is not limited to serving as the first step of vulnerable software identification as in the VIEM system [8]. Indeed, there has been extensive applications of the NER in miscellaneous computer security tasks.

Feng et al. emphasize the importance of NER applications targeted at software domain for developing security applications like detecting compromised and vulnerable devices [11]. Ye et al. investigate the challenge of NER involving software, and they note that the common word polysemy is one important factor that complicates the task [46]. Pandita et al. identifies software names in the description of mobile applications with NER together with other techniques; then they locate the sentences pertaining to security permissions with extracted information [32]. Sun et al. notice the sizable difference between the source of CVE records named ExploitDB and CVE database in terms of both content completeness and update time. They, therefore, propose to automatically extract the relevant information from ExploitDB descriptions to create CVE records with NER and other technologies [38]. Ma et al. cast the API extraction from software documentation as an NER task, and they leverage an LSTM-based architecture to obtain high-fidelity APIs across different libraries and languages through fine-tuning and transfer learning [27]. Another work by Ye et al. attempts to use the NER to enable entity-centric applications for software engineering [45]. However, when they reduce their training data to 10% of the original dataset, their NER model’s performance drop from an F1 score of 78% to 60%. This drop in performance highlights the difficulty of building NER systems in the few-shot setting.

Our work strives to identify vulnerabilities from public security reports, and therefore treats NER as a vehicle; furthermore, the architecture we leverage also benefits from fine-tuning and transfer learning. However, we face a unique challenge of resource constraints: our system should reach competitive systems’ performance with as few labeled samples as possible.

5.3 Few-Sample Named Entity Recognition

The study of named entity recognition dates back to Message Understanding Conference-6 (MUC-6) held in 1995, where researchers first identified the problem of recognizing named entities such as names, organizations, locations using rule-based or probabilistic approaches [19]. But the adoption of neural approaches and the notion of few-sample NER only happened in recent years.

Lample et al. first propose to utilize neural networks in named entity recognition task and outperforms traditional CRF-based sequence labeling models [24]. Driven by their success, neural network has become the de facto choice for named entity recognition task since then. However, at the same time, the tension between cost of obtaining named entity annotations and data hunger nature of neural model training has gained attention. For example, Buys et al. notice that a graphical model like HMM could outperform its neural counterparts by more than 10% of accuracy in sequence tagging task given the access to same amount of training data [4]. Smaller number of parameters in HMM makes the training easier to converge than a RNN.

In order to tackle this challenge, the notion of few-shot learning, originally proposed by computer vision community [37,40], gets noticed by NLP practitioners. These methods focus on inferring incoming samples' class membership based on prototypes obtained through metric learning. Their success in image classification propels the same methods to be adapted and tested in the NER task. Fritzler et al. are among the first to apply prototypical networks to NER task. However, their model treats each token independently by inferring one entity type at a time, disregarding the context of sentence [14]. Hou et al. extend this work and propose to merge specific named entity types into abstract classes (referred to as "collapsed transitions") to enable domain transfer [20]. Yang and Katiyar borrow the idea of collapsed transitions but simplify the overall architecture [44]. Rather than using dedicated architecture for few-shot learning, they first try to measure distance between tokens based on tokens' contextualized embeddings returned by PLM and then infer token's entity type based solely on nearest neighbor's types. After this step, an optional CRF decoding module is used to account for context of tokens. Their simple model architecture gain state-of-the-art performance compared with previous attempts. More importantly, besides the significant boost in performance metrics, their work shows the potential of pretrained language models (PLM) in reducing annotations by leveraging universal representations of language.

Another line of works consider the alternatives to the prototype-based approaches. Huang et al. show that self-training and active learning could help reduce data annotations [22]. Fries et al. provide empirical evidence of using data programming paradigm to convert rules given by the experts to named entity tags [13]. Li et al. propose a novel training protocol to make model-agnostic meta learning (MAML) framework, previously only available for sentence-level classification [1], applicable to sequence labeling (token-level classification) problem [26,25].

Our work is consistent with previous works in the goal of reducing the number of training labels required. However, different from [44], we discard the tag assignment scheme based on prototypes and propose the direct application of PLMs through fine-tuning and transfer learning. The empirical evidence shows that this design choice attains better and stable performance metrics.

A distinct application scenario of few-sample learning is to test the system by inferring classes not seen in training time. We note that, despite both striving to generalize well to unseen samples during testing, our setting of reducing the required training examples is orthogonal to this goal following the taxonomy provided in [47].

6 Conclusions and Future Work

In this work, we take a first attempt to investigate the problem of few-sample named entity recognition (NER) for public security vulnerability reports. By leveraging fine-tuning of three pre-trained language models (BERT, RoBERTa, Electra), we find that it is possible to match the performance of prior art with much less labeled training data: first, for vulnerability reports in the `memc` category, fine-tuning on PLM can match the score of prior art by using only 10% of their training examples (or 576 labelled sentences); second, for vulnerability reports in the other 12 categories (e.g., `bypass`), we find that through fine-tuning and transfer learning, it is possible to match the scores of [8] by using only 11.2% of their training examples (or an aggregation of 64 labels from each category). As a result, few-sample learning has effectively reduced the training data required for NER for security vulnerability reports.

It is important to notice that our work is just a beginning and there exists great potential for further improvement. We identify future work in the following directions. First, leveraging unlabelled data. This work only considers labelled samples of named entity tags. It is worth exploring how to leverage unlabelled data to further improve the performance. It was shown in [29] that few-shot learning can achieve competitive performance on the GLUE dataset by leveraging models trained on a small amount of labelled data and predictions on a large amount of unlabelled examples as augmented noisy labels. As language models are usually pre-trained in the general domains, it can be expected that learning from unlabelled data in the security domain can help the model more quickly adapt to the in-domain task. Second, leveraging external knowledge base. Our experiment shows that it is more difficult to achieve a good F1 score for software name than software version. This result also meets our expectation, because there is a higher chance for a software name than software version in the testing dataset to be completely unseen. However, the language model may fail to bridge such gap since it is pre-trained on the general domain. To make the language model quickly recognize unseen software names, one approach is to leverage an external knowledge base on top of the few-sample learning model. It is an open question how to leverage the knowledge base to help with the few-sample learning without introducing many mismatched predictions. Third, by

empirically observing the failure cases in transfer learning, we find that there exists some adversarial cases in the 12 sub-datasets (Section 4.5), which results in a dramatic drop in the performance of few-sample learning. One direction of future work is thus to investigate adversarial learning for the few-shot transfer learning to improve its robustness.

References

1. Bao, Y., Wu, M., Chang, S., Barzilay, R.: Few-shot text classification with distributional signatures (Aug 2019)
2. Bridges, R.A., Jones, C.L., Iannacone, M.D., Testa, K.M., Goodall, J.R.: Automatic labeling for entity extraction in cyber security. arXiv preprint arXiv:1308.4941 (2013)
3. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. arXiv preprint arXiv:2005.14165 (2020)
4. Buys, J., Bisk, Y., Choi, Y.: Bridging hmms and rnns through architectural transformations (2018)
5. Chang, J.R., Chen, Y.S.: Pyramid stereo matching network. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5410–5418 (2018)
6. Chen, H., Xia, M., Chen, D.: Non-parametric few-shot learning for word sense disambiguation. arXiv preprint arXiv:2104.12677 (2021)
7. Dodge, J., Ilharco, G., Schwartz, R., Farhadi, A., Hajishirzi, H., Smith, N.A.: Fine-tuning pretrained language models: Weight initializations, data orders, and early stopping. ArXiv **abs/2002.06305** (2020)
8. Dong, Y., Guo, W., Chen, Y., Xing, X., Zhang, Y., Wang, G.: Towards the detection of inconsistencies in public security vulnerability reports. In: USENIX Security Symposium (2019)
9. Fan, J., Li, Y., Wang, S., Nguyen, T.N.: A C/C++ Code Vulnerability Dataset with Code Changes and CVE Summaries. In: Proceedings of the 17th International Conference on Mining Software Repositories. pp. 508–512. ACM, Seoul Republic of Korea (Jun 2020). <https://doi.org/10.1145/3379597.3387501>, <https://dl.acm.org/doi/10.1145/3379597.3387501>
10. Farhang, S., Kirdan, M.B., Laszka, A., Grossklags, J.: An empirical study of android security bulletins in different vendors. In: Proceedings of The Web Conference 2020. pp. 3063–3069 (2020)
11. Feng, X., Li, Q., Wang, H., Sun, L.: Acquisitional rule-based engine for discovering internet-of-thing devices. In: USENIX Security Symposium (2018)
12. Finn, C., Abbeel, P., Levine, S.: Model-agnostic meta-learning for fast adaptation of deep networks. In: International Conference on Machine Learning. pp. 1126–1135. PMLR (2017)
13. Fries, J., Wu, S., Ratner, A., Ré, C.: Swellshark: A generative model for biomedical named entity recognition without labeled data
14. Fritzler, A., Logacheva, V., Kretov, M.: Few-shot classification in named entity recognition task. Proceedings of the 34th ACM/SIGAPP Symposium on Applied Computing (2019)
15. Gaglione Jr, G.S.: The equifax data breach: An opportunity to improve consumer protection and cybersecurity efforts in america. *Buff. L. Rev.* **67**, 1133 (2019)

16. Gao, P., Liu, X., Choi, E., Soman, B., Mishra, C., Farris, K., Song, D.: A system for automated open-source threat intelligence gathering and management. arXiv preprint arXiv:2101.07769 (2021)
17. Gao, T., Fisch, A., Chen, D.: Making pre-trained language models better few-shot learners (Dec 2020)
18. Gasmi, H., Laval, J., Bouras, A.: Information extraction of cybersecurity concepts: An lstm approach. *Applied Sciences* **9**(19), 3945 (2019)
19. Grishman, R., Sundheim, B.: Message understanding conference- 6: A brief history. In: *COLING* (1996)
20. Hou, Y., Che, W., Lai, Y., Zhou, Z., Liu, Y., Liu, H., Liu, T.: Few-shot slot tagging with collapsed dependency transfer and label-enhanced task-adaptive projection network. In: *ACL* (2020)
21. Hou, Y., Che, W., Lai, Y., Zhou, Z., Liu, Y., Liu, H., Liu, T.: Few-shot slot tagging with collapsed dependency transfer and label-enhanced task-adaptive projection network. arXiv preprint arXiv:2006.05702 (2020)
22. Huang, J., Li, C., Subudhi, K., Jose, D., Balakrishnan, S., Chen, W., Peng, B., Gao, J., Han, J.: Few-shot named entity recognition: A comprehensive study. *ArXiv abs/2012.14978* (2020)
23. Johnson, J., Khoshgoftaar, T.: Survey on deep learning with class imbalance. *Journal of Big Data* **6**, 1–54 (2019)
24. Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., Dyer, C.: Neural architectures for named entity recognition
25. Li, J., Chiu, B., Feng, S., Wang, H.: Few-shot named entity recognition via meta-learning. *IEEE Transactions on Knowledge and Data Engineering* pp. 1–1 (2020). <https://doi.org/10.1109/tkde.2020.3038670>
26. Li, J., Shang, S., Shao, L.: MetaNER: Named entity recognition with meta-learning. In: *Proceedings of The Web Conference 2020*. ACM (apr 2020). <https://doi.org/10.1145/3366423.3380127>
27. Ma, S., Xing, Z., Chen, C., Qu, L., Li, G.: Easy-to-deploy api extraction by multi-level feature embedding and transfer learning. *IEEE Transactions on Software Engineering* pp. 1–1 (2019)
28. Mosbach, M., Andriushchenko, M., Klakow, D.: On the stability of fine-tuning bert: Misconceptions, explanations, and strong baselines (Jun 2020)
29. Mukherjee, S.S., Awadallah, A.H.: Uncertainty-aware self-training for few-shot text classification. In: *NeurIPS 2020 (Spotlight)*. ACM (December 2020)
30. Mulwad, V., Li, W., Joshi, A., Finin, T., Viswanathan, K.: Extracting information about security vulnerabilities from web text. In: *2011 IEEE/WIC/ACM International Conferences on Web Intelligence and Intelligent Agent Technology*. vol. 3, pp. 257–260. IEEE (2011)
31. Palmer, C.: The boeing 737 max saga: Automating failure. *Engineering* **6**(1), 2–3 (feb 2020). <https://doi.org/10.1016/j.eng.2019.11.002>
32. Pandita, R., Xiao, X., Yang, W., Enck, W., Xie, T.: Whyper: Towards automating risk assessment of mobile applications. In: *USENIX Security Symposium* (2013)
33. Pham, V., Dang, T.: CVEExplorer: Multidimensional Visualization for Common Vulnerabilities and Exposures. In: *2018 IEEE International Conference on Big Data (Big Data)*. pp. 1296–1301 (Dec 2018). <https://doi.org/10.1109/BigData.2018.8622092>
34. Phang, J., Févry, T., Bowman, S.R.: Sentence encoders on stilts: Supplementary training on intermediate labeled-data tasks. arXiv preprint arXiv:1811.01088 (2018)

35. Ponta, S.E., Plate, H., Sabetta, A., Bezzi, M., Dangremont, C.: A Manually-Curated Dataset of Fixes to Vulnerabilities of Open-Source Software. In: 2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR). pp. 383–387. IEEE, Montreal, QC, Canada (May 2019). <https://doi.org/10.1109/MSR.2019.00064>, <https://ieeexplore.ieee.org/document/8816802/>
36. Snell, J., Swersky, K., Zemel, R.S.: Prototypical networks for few-shot learning. arXiv preprint arXiv:1703.05175 (2017)
37. Snell, J., Swersky, K., Zemel, R.S.: Prototypical networks for few-shot learning (Mar 2017)
38. Sun, J., Xing, Z., Guo, H., Ye, D., Li, X., Xu, X., Zhu, L.: Generating informative cve description from exploitdb posts by extractive summarization. arXiv preprint arXiv:2101.01431 (2021)
39. Tjong Kim Sang, E.F., De Meulder, F.: Introduction to the conll-2003 shared task: Language-independent named entity recognition. In: Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003 - Volume 4. p. 142–147. CONLL '03, Association for Computational Linguistics, USA (2003). <https://doi.org/10.3115/1119176.1119195>, <https://doi.org/10.3115/1119176.1119195>
40. Vinyals, O., Blundell, C., Lillicrap, T., Kavukcuoglu, K., Wierstra, D.: Matching networks for one shot learning (Jun 2016)
41. Wang, Y., Yao, Q., Kwok, J., Ni, L.: Generalizing from a few examples: A survey on few-shot learning. arXiv: Learning (2019)
42. Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T., Louf, R., Funtowicz, M., Brew, J.: Huggingface’s transformers: State-of-the-art natural language processing. ArXiv **abs/1910.03771** (2019)
43. Yamamoto, Y., Miyamoto, D., Nakayama, M.: Text-Mining Approach for Estimating Vulnerability Score. In: 2015 4th International Workshop on Building Analysis Datasets and Gathering Experience Returns for Security (BADGERS). pp. 67–73. IEEE, Kyoto, Japan (Nov 2015). <https://doi.org/10.1109/BADGERS.2015.018>, <http://ieeexplore.ieee.org/document/7809535/>
44. Yang, Y., Katiyar, A.: Simple and effective few-shot named entity recognition with structured nearest neighbor learning. ArXiv **abs/2010.02405** (2020)
45. Ye, D., Xing, Z., Foo, C.Y., Ang, Z.Q., Li, J., Kapre, N.: Software-specific named entity recognition in software engineering social content. 2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER) **1**, 90–101 (2016)
46. Ye, D., Xing, Z., Foo, C.Y., Li, J., Kapre, N.: Learning to extract api mentions from informal natural language discussions. 2016 IEEE International Conference on Software Maintenance and Evolution (ICSME) pp. 389–399 (2016)
47. Zhang, T., Wu, F., Katiyar, A., Weinberger, K.Q., Artzi, Y.: Revisiting few-sample bert fine-tuning (Jun 2020)
48. Zhou, C., Li, B., Sun, X.: Improving software bug-specific named entity recognition with deep neural network. *J. Syst. Softw.* **165**, 110572 (2020)

A Dataset Statistics

Table 6 is the detailed version of Table 2. The valid split for categories other than `memc` is 10% sample of official train set. "Sentence-level entity proportion" refers to the proportion of sentences that have `SN` or `SV`, and "Token-level entity proportion" is proportion of `SN` and `SV` with respect to given dataset split. These proportions reflect the dataset imbalance in different levels.

Table 6: Detailed statistics of VIEM dataset.

Category	Split	Number of sample	Sentence-level entity proportion	Token-level entity proportion (SN)	Token-level entity proportion (SV)
memc	train	5758	0.5639	0.0613	0.0819
	valid	1159	0.3287	0.0368	0.0807
	test	1001	0.4555	0.0559	0.0787
bypass	train	652	0.2239	0.0314	0.0431
	valid	162	0.2469	0.0367	0.0423
	test	610	0.2902	0.0456	0.0531
csrf	train	521	0.2399	0.0207	0.0347
	valid	130	0.2846	0.0251	0.0397
	test	415	0.3181	0.0321	0.0464
dirtra	train	619	0.2359	0.0172	0.0219
	valid	155	0.1871	0.0180	0.0316
	test	646	0.2879	0.0197	0.0220
dos	train	396	0.2273	0.0212	0.0405
	valid	99	0.2020	0.0234	0.0419
	test	484	0.2624	0.0189	0.0331
execution	train	413	0.2639	0.0228	0.0358
	valid	103	0.2718	0.0314	0.0302
	test	639	0.2598	0.0273	0.0357
fileinc	train	546	0.2857	0.0175	0.0185
	valid	137	0.3869	0.0259	0.0222
	test	683	0.3133	0.0206	0.0215
gainpre	train	323	0.2229	0.0243	0.0430
	valid	80	0.3250	0.0357	0.0723
	test	577	0.2114	0.0191	0.0311
https	train	550	0.1891	0.0127	0.0217
	valid	137	0.1241	0.0077	0.0124
	test	411	0.2360	0.0175	0.0304
infor	train	305	0.2459	0.0326	0.0354
	valid	76	0.3158	0.0187	0.0282
	test	509	0.2358	0.0227	0.0348
overflow	train	396	0.2475	0.0217	0.0326
	valid	98	0.2143	0.0185	0.0230
	test	454	0.2819	0.0216	0.0343
sqli	train	538	0.2565	0.0145	0.0141
	valid	134	0.2836	0.0194	0.0151
	test	685	0.2423	0.0171	0.0181
xss	train	357	0.2829	0.0203	0.0289
	valid	89	0.3708	0.0276	0.0386
	test	562	0.2046	0.0219	0.0363