

PeerDA: Data Augmentation via Modeling Peer Relation for Span Identification Tasks*

Weiwen Xu^{1,2,†} Xin Li^{2,‡} Yang Deng¹ Wai Lam¹ Lidong Bing²

¹The Chinese University of Hong Kong

²DAMO Academy, Alibaba Group

{wxwu, wlam}@se.cuhk.edu.hk ydeng@nus.edu.sg

{xinting.lx, l.bing}@alibaba-inc.com

Abstract

Span identification aims at identifying specific text spans from text input and classifying them into pre-defined categories. Different from previous works that merely leverage the Subordinate (SUB) relation (i.e. *if a span is an instance of a certain category*) to train models, this paper for the first time explores the Peer (PR) relation, which indicates that *two spans are instances of the same category and share similar features*. Specifically, a novel **Peer Data Augmentation** (PeerDA) approach is proposed which employs span pairs with the PR relation as the augmentation data for training. PeerDA has two unique advantages: (1) There are a large number of PR span pairs for augmenting the training data. (2) The augmented data can prevent the trained model from over-fitting the superficial span-category mapping by pushing the model to leverage the span semantics. Experimental results on ten datasets over four diverse tasks across seven domains demonstrate the effectiveness of PeerDA. Notably, PeerDA achieves state-of-the-art results on six of them.¹

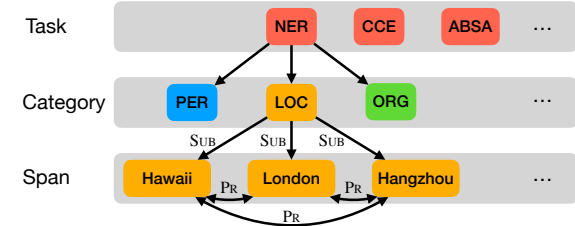
1 Introduction

Span Identification (SpanID) is a family of Natural Language Processing (NLP) tasks with the goal of detecting specific spans from the input text and further classifying them into pre-defined categories (Papay et al., 2020). It serves as the initial step for complex text analysis by narrowing down the search scopes of important spans, which holds a pivotal position in the field of NLP (Ding et al., 2021; Xu et al., 2021). Recently, different

* This work was supported by Alibaba Group through Alibaba Research Intern Program. It was also partially supported by a grant from the Research Grant Council of the Hong Kong Special Administrative Region, China (Project Code: 14200620). † This work was done when Weiwen Xu was an intern at Alibaba DAMO Academy. ‡ Xin Li is the corresponding author.

¹Our code and data are available at <https://github.com/DAMO-NLP-SG/PeerDA>

(a) Relations in SpanID



(b) SpanID in MRC Paradigm

| Context: | Gotta dress up for London fashion week and party in style! | |
|----------------|--|---|
| Original data | SUB Query: | Highlight the parts (if any) related to "LOC". Details: the name of politically or geographically defined locations such as cities, provinces, etc. |
| | Answer: | London |
| Augmented data | PR Query-1: | Highlight the parts (if any) similar to "Hawaii". |
| | Answer: | London |
| | PR Query-2: | Highlight the parts (if any) similar to "Hangzhou". |
| | Answer: | London |

Figure 1: (a) Illustrations of Subordinate (SUB) and Peer (PR) relations in SpanID tasks. (b) The constructions of augmented data with PR relations in MRC paradigm. We use NER here for demonstration purposes.

domain-specific SpanID tasks, such as social media Named Entity Recognition (NER) (Derczynski et al., 2017), Aspect-Based Sentiment Analysis (ABSA) (Liu, 2012), Contract Clause Extraction (CCE) (Chalkidis et al., 2017), Span Based Propaganda Detection (SBPD) (Da San Martino et al., 2019) and Argument Extraction (Cheng et al., 2020), have emerged for various NLP applications.

Precisely, as shown in Figure 1 (a), the process of SpanID can be reinterpreted as extracting **span-category** Subordinate (SUB) relation — *if a span in the input text is an instance of a certain category*. Early works (Chiu and Nichols, 2016) typically tackle SpanID tasks as a sequence tagging problem, where the SUB relation is recognized via predicting the category for each input token under certain context. Recently, to better utilize category semantics, many efforts have been made on reformulating SpanID tasks as a Machine Reading Comprehension (MRC) problem (Liu et al., 2020; Yang et al., 2021). As shown by the example in Figure 1 (b), such formulation first creates a SUB query for each category and then recognizes the SUB relation by

detecting relevant spans in the input text (*i.e.*, context) as answers to the category query.

However, only leveraging the SUB relation in the training data to build SpanID models may suffer from two limitations: 1) **Over-fitting**: With only SUB relation, SpanID models tend to capture the superficial span-category correlations. Such correlations may misguide the models to ignore the semantics of the given span but make predictions based on the memorized span-category patterns, which hurts the generalization capability of the models. 2) **Data Scarcity**: For low-resource scenarios or long-tailed categories, the number of span-category pairs with SUB relation (SUB pairs) could be very limited and insufficient to learn a reliable SpanID model.

In this paper, we explore the **span-span** Peer (PR) relation to alleviate the above limitations. Specifically, the PR relation indicates that *two spans are two different instances of the same category*. The major difference between PR relation and SUB relation is that the former one intends to correlate two spans without giving the categories they belong to. For example, in Figure 1 (a), "Hawaii" and "London" are connected with the PR relation because they are instances of the same category. By jointly recognizing SUB relation and PR relation in the input text, the model is enforced to favor the usage of span semantics instead of span-category patterns for prediction, reducing the risk of over-fitting. In addition, the number of span-span pairs with the PR relation (PR pairs) grows quadratically over the number of SUB pairs. Therefore, we can still construct a reasonable number of training data with PR pairs for categories having insufficient examples.

In this paper, with the aim of leveraging the PR relation to enhance SpanID models, we propose a Peer Data Augmentation (*PeerDA*) approach that treats PR pairs as a kind of augmented training data. To achieve this, as depicted in Figure 1 (b), we extend the usage of the original training data into two views. The first view is the SUB-based training data. It is used to directly solve the SpanID tasks by extracting the SUB relation, which is the typical formulation of MRC-based approaches. The second view is the PR-based training data. It is our augmentation to enrich the semantics of spans by extracting the PR relation in the original training data, where one span is used to identify its peer from the input context. Note that our PR-based

training data can be easily formulated into the MRC paradigm. Therefore, the knowledge learned from such augmentation data can be directly transferred to enhance the model’s capability to capture SUB relation (*i.e.*, the SpanID tasks).

To better accommodate the MRC-style SUB and PR data, we develop a stronger and more memory-efficient MRC model. Compared to the designs in Li et al. (2020b), our model introduces a bilinear component to calculate the span scores and consistently achieves better performance with a 4 times smaller memory consumption. Besides, we propose a margin-based contrastive learning strategy to additionally model the negative spans to the query (*e.g.*, when querying the context in Figure 1 for “ORG” entities, “London” becomes a negative span) so that the spans from different categories are separated more apart in the semantic space.

We evaluate the effectiveness of PeerDA on ten datasets across seven domains, from four different SpanID tasks, namely, NER, ABSA, CCE, and SBPD. Experimental results show that extracting PR relation benefits the learning of semantics and encourages models to identify more possible spans. As a result, PeerDA is a new state-of-the-art (SOTA) method on six SpanID datasets. Our analyses further demonstrate the capability of PeerDA to alleviate scarcity and over-fitting issues.

Our contributions are summarized as follows:

- We propose a novel PeerDA approach to tackle SpanID tasks via augmenting training data with PR relation.
- We conduct extensive experiments on ten datasets, including four different SpanID tasks across seven domains, and achieve SOTA performance on six SpanID datasets.
- PeerDA is more effective in low-resource scenarios or long-tailed categories and thus, it alleviates the scarcity issue. Meanwhile, jointly recognizing the SUB and PR relations makes the MRC model rely less on memorizing the SUB patterns in the training set for inferring the span label, which prevents over-fitting.

2 Related Work

DA for SpanID: DA, which increases the diversity of training data at a low cost, is a widely-adopted solution to address data scarcity (Feng et al., 2021). In the scope of SpanID, existing DA approaches aim to introduce more span-category patterns, including: (1) *Word Replacement* either

replaces or paraphrases some context tokens using simple rules (Wei and Zou, 2019; Dai and Adel, 2020; Xiang et al., 2021) and strong language models (Kobayashi, 2018; Wu et al., 2019; Li et al., 2020a; Yoo et al., 2021), or applies synonym dictionaries or masked language models to replace the labeled tokens with other tokens of the same type (Wei and Zou, 2019; Zhou et al., 2022b). (2) Fine-grained Augmentation Data Generation first trains an auto-regressive language model, and then leverages the model to generate new sentences with entity tags as a special kind of tokens (Ding et al., 2020; Liu et al., 2021b). (3) *Self-training* is to continually train the model on its predicted data (Xie et al., 2019, 2020; Wang et al., 2020; Zhou et al., 2023; Tan et al., 2023), while consistency training also leverages unlabeled data by imposing regularization on the predictions (Zhou et al., 2022a) (4) *Distantly Supervised Training* focuses on leveraging external knowledge to roughly label spans in the target tasks (Bing et al., 2013, 2015; Xu et al., 2023a). Huang et al. (2021) leverage Wikipedia to create distant labels for NER. Chen et al. (2021) transfer data from high-resource to low-resource domains. Jain et al. (2019); Li et al. (2020c) tackle cross-lingual NER by projecting labels from high-resource to low-resource languages, which is particularly common in real applications (Kruengkrai et al., 2020). Differently, the motivation of PeerDA is to leverage the augmented data to enhance models’ capability on semantic understanding by minimizing(maximizing) the distances between semantically similar(distant) spans.

MRC: MRC is to extract an answer span from a relevant context conditioned on a given query. It is initially designed to solve question answering tasks (Hermann et al., 2015), while recent trends have shown great advantages in formulating NLP tasks as MRC problems. In the context of SpanID, Li et al. (2020b); Xu et al. (2022b, 2023b) address the nested NER issues by decomposing nested entities under multiple queries. Mao et al. (2021); Zhang et al. (2021a) tackle ABSA in a unified MRC framework. Hendrycks et al. (2021) tackle CCE with MRC to deal with the extraction of long clauses. Moreover, other tasks such as relation extraction (Li et al., 2019), event detection (Liu et al., 2020, 2021a), and summarization (McCann et al., 2018) are also reported to benefit from the MRC paradigm.

3 PeerDA

Overview of SpanID: Given the input text $\mathbf{X} = \{x_1, \dots, x_n\}$, SpanID is to detect all appropriate spans $\{\mathbf{x}_k\}_{k=1}^K$ and classify them with proper labels $\{y_k\}_{k=1}^K$, where each span $\mathbf{x}_k = \{x_{s_k}, x_{s_k+1}, \dots, x_{e_k-1}, x_{e_k}\}$ is a subsequence of \mathbf{X} satisfying $s_k \leq e_k$ and the label comes from a predefined category set Y (e.g. "Person" in NER).

3.1 Training Data Construction

The training data \mathcal{D} consists of two parts: (1) The SUB-based training data \mathcal{D}^{SUB} , where the query is about a category and the MRC context is the input text. (2) The PR-based training data \mathcal{D}^{PR} is constructed with PR pairs, where one span is used to create the query and the input text containing the second span serves as the MRC context.

3.1.1 SUB-based Training Data

First, we need to transform the original training examples into (query, context, answers) triples following the paradigm of MRC (Li et al., 2020b). To extract the SUB relation between categories and relevant spans, a natural language query Q_y^{SUB} is constructed to reflect the semantics of each category y . Following Hendrycks et al. (2021), we include both category mention $[\text{Men}]_y$ and its definition $[\text{Def}]_y$ from the annotation guideline (or Wikipedia if the guideline is not accessible) in the query to introduce more comprehensive semantics:

$$Q_y^{\text{SUB}} = \text{Highlight the parts (if any)} \\ \text{related to } [\text{Men}]_y. \text{ Details : } [\text{Def}]_y. \quad (1)$$

Given the input text \mathbf{X} as the context, the answers to Q_y^{SUB} are the spans belonging to category y . Then we can obtain one MRC example denoted as $(Q_y^{\text{SUB}}, \mathbf{X}, \{\mathbf{x}_k \mid \mathbf{x}_k \in \mathbf{X}, y_k = y\}_{k=1}^K)$. To guarantee the identification of all possible spans, we create $|Y|$ training examples by querying the input text with each pre-defined category.

3.1.2 PR-based training data

To construct augmented data that derived from the PR relation, we first create a category-wise span set \mathcal{S}_y that includes all training spans with category y :

$$\mathcal{S}_y = \{\mathbf{x}_k \mid (\mathbf{x}_k, y_k) \in \mathcal{D}^{\text{SUB}}, y_k = y\} \quad (2)$$

Obviously, any two different spans in \mathcal{S}_y have the same category and shall hold the PR relation. Therefore, we pair every two different spans in \mathcal{S}_y to create a peer set \mathcal{P}_y :

$$\mathcal{P}_y = \{(\mathbf{x}^q, \mathbf{x}^a) \mid \mathbf{x}^q, \mathbf{x}^a \in \mathcal{S}_y, \mathbf{x}^q \neq \mathbf{x}^a\} \quad (3)$$

For each PR pair (x^q, x^a) in \mathcal{P}_y , we can construct one training example by constructing the query with the first span x^q :

$$Q_y^{\text{PR}} = \text{Highlight the parts (if any)} \\ \text{similar to } x^q. \quad (4)$$

Then we treat the text X^a containing the second span x^a as the MRC context to be queried and x^a as the answer to Q_y^{PR} . Note that there may exist more than one span in X^a satisfying PR relation with x^q , we set all of them as the valid answers to Q_y^{PR} , yielding one training example $(Q_y^{\text{PR}}, X^a, \{x_k^a \mid x_k^a \in X^a, y_k^a = y\}_{k=1}^K)$ of our PeerDA.

Theoretically, given the span set \mathcal{S}_y , there are only $|\mathcal{S}_y|$ SUB pairs in the training data but we can obtain $|\mathcal{S}_y| \times (|\mathcal{S}_y| - 1)$ PR pairs to construct \mathcal{D}^{PR} . Such a large number of augmented data shall hold great potential to enrich spans’ semantics. However, putting all PR-based examples into training would exacerbate the skewed data distribution issue since the long-tailed categories get fewer PR pairs for augmentation and also increase the training cost. Therefore, as the first step for DA with the PR relation, we propose three augmentation strategies to control the size and distribution of augmented data.

PeerDA-Size: This is to increase the size of augmented data while keeping the data distribution unchanged. Specifically, for each category y , we randomly sample $\lambda|\mathcal{S}_y|$ PR pairs from \mathcal{P}_y . Then we collect all sampled PR pairs to construct \mathcal{D}^{PR} , where λ is the DA rate to control the size of \mathcal{D}^{PR} .

PeerDA-Categ: Categories are not evenly distributed in the training data, and in general SpanID models perform poorly on long-tailed categories. To tackle this, we propose PeerDA-Categ to augment more training data for long-tailed categories. Specifically, let y^* denote the category having the largest span set of size $|\mathcal{S}_{y^*}|$. We sample up to $|\mathcal{S}_{y^*}| - |\mathcal{S}_y|$ PR pairs from \mathcal{P}_y for each category y and construct a category-balanced training set \mathcal{D}^{PR} using all sampled pairs. Except for the extreme cases where $|\mathcal{S}_y|$ is smaller than $\sqrt{|\mathcal{S}_{y^*}|}$, we would get the same size of the training data for each category after the augmentation, which significantly increases the exposure for spans from the long-tailed categories.

PeerDA-Both (The final version of PeerDA): To take advantage of the above two strategies, we further propose PeerDA-Both to maintain the data distribution while effectively increasing the size of

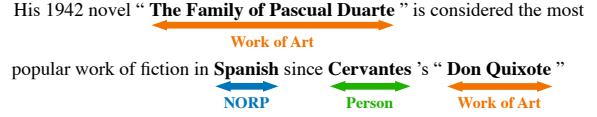


Figure 2: Example of extracting multiple spans in NER.

training data. In PeerDA-Both, we randomly sample $\max(\lambda|\mathcal{S}_{y^*}| + (|\mathcal{S}_{y^*}| - |\mathcal{S}_y|), 0)$ PR pairs from \mathcal{P}_y for each category y to construct \mathcal{D}^{PR} , where $\lambda|\mathcal{S}_{y^*}|$ determines the size of the augmented data, and $|\mathcal{S}_{y^*}| - |\mathcal{S}_y|$ controls the data distribution.

3.1.3 Data Balance

We combine the \mathcal{D}^{SUB} and the \mathcal{D}^{PR} created above as the final training data. Since an input text usually mentions spans from a few categories, when converting the text into the MRC paradigm, many of the $|Y|$ examples are unanswerable. If a SpanID model is trained on this unbalanced data, then the model may favor the majority of the training examples and output an empty span. To balance answerable and unanswerable examples, we follow Hendrycks et al. (2021) to randomly remove some unanswerable examples from the training data.

3.2 Model Architecture

As shown in Figure 2, to achieve the detection of multiple spans for the given query, we follow Li et al. (2020b) to build the MRC model. Compared to the original designs, we further optimize the computation of span scores following a general way of Luong et al. (2015); Xu et al. (2022b).

Specifically, the base model consists of three components: an encoder, a span predictor, and a start-end selector. First, given the concatenation of the query Q and the context X as the MRC input $\bar{X} = \{[\text{CLS}], Q, [\text{SEP}], X, [\text{SEP}]\}$, where $[\text{CLS}]$, $[\text{SEP}]$ are special tokens, the encoder would encode the input text into hidden states H :

$$H = \text{ENCODER}(\bar{X}) \quad (5)$$

Second, the span predictor consists of two binary classifiers, one to predict whether each context token is the start index of the answer, and the other to predict whether the token is the end index:

$$P_{\text{start}} = HW^s \quad P_{\text{end}} = HW^e \quad (6)$$

where $W^s, W^e \in \mathbb{R}^{d \times 2}$ are the weights of two classifiers and d is the dimension of hidden states. The span predictor would output multiple start and end indexes for the given query and context.

Third, the start-end selector matches each start index to each end index and selects the most possible spans from all combinations as the outputs.

| Task | NER | | | | ABSA | | | SBPD | | CCE |
|------------|-------------------|---------------|--------------|-------------------|---------------|---------------|-------------------|-------------------|-----------------|--------------|
| Dataset | OntoNotes5 | WNUT17 | Movie | Restaurant | Weibo | Lap14 | Rest14 | News20 | Social21 | CUAD |
| Domain | <i>mixed</i> | <i>social</i> | <i>movie</i> | <i>restaurant</i> | <i>social</i> | <i>laptop</i> | <i>restaurant</i> | <i>news</i> | <i>social</i> | <i>legal</i> |
| # Train | 60.0k | 3.4k | 7.8k | 7.7k | 1.3k | 2.7k | 2.7k | 0.4k | 0.7k | 0.5k |
| # Test | 8.3k | 1.3k | 2.0k | 1.5k | 0.3k | 0.8k | 0.8k | 75 (<i>dev</i>) | 0.2k | 0.1k |
| # Category | 11 | 6 | 12 | 8 | 4 | 1 / 3 | 1 / 3 | 14 | 20 | 41 |

Table 1: Statistics on the ten SpanID datasets. Note that 1 / 3 denotes that there is 1 category in ATE and 3 categories in UABSA. *dev* denotes that we evaluate News20 on the dev set.

Different from the *concat* way that would create a large $\mathbb{R}^{|\bar{X}| \times |\bar{X}| \times 2d}$ -shape tensor (Li et al., 2020b), we leverage a *general* way following Luong et al. (2015); Xu et al. (2022b) to compute the span score, consuming fewer resources for better training efficiency:

$$P_{s,e} = FFN(\mathbf{H}_s)^T \mathbf{H}_e \quad (7)$$

where *FFN* is the feed-forward network (Vaswani et al., 2017), $P_{s,e}$ denotes the likelihood of $\bar{X}_{s:e}$ to form a possible answer.

3.3 Training Objective

The standard objective is to minimize the cross-entropy loss (CE) between above three predictions and their corresponding ground-truth labels, i.e., $Y_{\text{start}}, Y_{\text{end}}, Y_{s,e}$ (Li et al., 2020b):

$$\mathcal{L}_{mrc} = \text{CE}(\sigma(P_{\text{start}}), Y_{\text{start}}) + \text{CE}(\sigma(P_{\text{end}}), Y_{\text{end}}) + \text{CE}(\sigma(P_{s,e}), Y_{s,e}) \quad (8)$$

where σ is the sigmoid function.

However, these objectives only capture the semantic similarity between the query and positive spans (i.e., the span instances of the query category). In this paper, we propose to explicitly separate the query and its negative spans (i.e., the span instances of other categories) apart with a margin-based contrastive learning strategy, for better distinguishing the spans from different categories.

Specifically, given the MRC input \bar{X} with query of category y , there may be multiple positive spans $\bar{X}^+ = \{\bar{x}_k \in \bar{X}, y_k = y\}$ and negative spans $\bar{X}^- = \{\bar{x}_{k'} \in \bar{X}, y_{k'} \neq y\}$. We leverage the following margin-based contrastive loss to penalize negative spans (Chechik et al., 2010):

$$\mathcal{L}_{ct} = \max_{\substack{\bar{x}_k \in \bar{X}^+ \\ \bar{x}_{k'} \in \bar{X}^-}} \max(0, M - (\sigma(P_{s_k, e_k}) - \sigma(P_{s_{k'}, e_{k'}}))) \quad (9)$$

where M is the margin term, $\max(\cdot, \cdot)$ is to select the larger one from two candidates, and the span score P_{s_k, e_k} can be regarded as the semantic similarity between the query and the target span \bar{x}_k . Note that our contrastive loss maximizes the

similarity difference between the query and the most confusing positive and negative span pairs (*Max-Min*), which we demonstrate to be effective in Sec. 5.3.

Finally, the overall training objective is:

$$\mathcal{L} = \mathcal{L}_{mrc} + \alpha \mathcal{L}_{ct} \quad (10)$$

where α is the balance rate.

4 Experimental Setup

4.1 Tasks

Note that PeerDA is a method for augmenting training data, it is not applied to test sets during evaluation. Therefore, we only use the SUB-based test data to evaluate the models' capability to extract spans according to the category. We conduct experiments on four SpanID tasks from diverse domains, including NER, ABSA, Contract Clause Extraction (CCE), and Span-Based Propaganda Detection (SBPD). The dataset statistics are summarized in Table 1. The detailed task description can be found in Appendix A.1.

NER: It is to detect named entities (i.e. spans) and classify them into entity types (i.e. categories). We evaluate five datasets, including four English datasets: **OntoNotes5**² (Pradhan et al., 2013), **WNUT17** (Derczynski et al., 2017), **Movie** (Liu et al., 2013b), and **Restaurant** (Liu et al., 2013a) and a Chinese dataset **Weibo** (Peng and Dredze, 2015). We use span-level micro-averaged Precision, Recall, and F_1 as evaluation metrics.

ABSA: We explore two ABSA sub-tasks: **Aspect Term Extraction (ATE)** to only extract aspect terms, and **Unified Aspect Based Sentiment Analysis (UABSA)** to jointly identify aspect terms and their sentiment polarities. We evaluate the two sub-tasks on two datasets, including the laptop domain **Lap14** and restaurant domain **Rest14**. We use micro-averaged F_1 as the evaluation metric.

²In order to conduct robustness experiments in Sec. A.3, we use the datasets from Lin et al. (2021) with 11 entity types.

| Methods | OntoNotes5 | | | WNUT17 | | | Movie | | | Restaurant | | | Weibo | | |
|---------|------------|------|----------------|--------|------|----------------|-------|------|----------------|------------|------|----------------|------------|------|----------------|
| | P | R | F ₁ | P | R | F ₁ | P | R | F ₁ | P | R | F ₁ | P | R | F ₁ |
| SOTA | RB-CRF+RM | | | CL-KL | | | T-NER | | | KaNa | | | RoBERTa+BS | | |
| | 92.8 | 92.4 | 92.6 | - | - | 60.5 | - | - | 71.2 | 80.9 | 80.0 | 80.4 | 70.2 | 75.4 | 72.7 |
| | Base | | | | | | | | | | | | | | |
| Tagging | 91.0 | 91.8 | 91.4 | 62.1 | 48.2 | 54.3 | 73.0 | 72.8 | 72.9 | 80.6 | 80.7 | 80.7 | 70.8 | 71.0 | 70.9 |
| MRC | 92.4 | 91.8 | 92.1 | 66.4 | 40.7 | 50.5 | 70.3 | 73.3 | 71.8 | 81.4 | 79.9 | 80.6 | 73.6 | 64.4 | 68.7 |
| PeerDA | 91.9 | 92.6 | 92.4 | 71.1 | 46.9 | 56.5 | 77.9 | 72.3 | 75.0 | 81.3 | 82.8 | 82.1 | 70.0 | 73.3 | 71.6 |
| | Large | | | | | | | | | | | | | | |
| Tagging | 93.0 | 92.3 | 92.6 | 69.4 | 46.2 | 55.4 | 74.2 | 74.0 | 74.1 | 80.9 | 82.0 | 81.4 | 71.4 | 69.2 | 70.3 |
| MRC | 92.8 | 91.8 | 92.3 | 72.4 | 41.7 | 52.9 | 76.7 | 73.2 | 74.9 | 81.6 | 81.7 | 81.7 | 72.2 | 66.8 | 69.4 |
| PeerDA | 92.8 | 93.7 | 93.3 | 70.9 | 48.0 | 57.2 | 78.5 | 73.1 | 75.7 | 81.8 | 82.5 | 82.2 | 73.4 | 71.6 | 72.5 |

Table 2: Performance on NER datasets. The best models are bolded.

SBPD: It aims to detect both the text fragment where a persuasion technique is used (i.e. spans) and its technique type (i.e. categories). We use **News20** and **Social21** from SemEval shared tasks (Da San Martino et al., 2020; Dimitrov et al., 2021). For **News20**, we report the results on its dev set since the test set is not publicly available. We use micro-averaged Precision, Recall, and F₁ as evaluation metrics.

CCE: It is a legal task to detect and classify contract clauses (i.e. spans) into relevant clause types (i.e. categories), such as "Governing Law". We conduct CCE experiments using **CUAD** (Hendrycks et al., 2021). We follow Hendrycks et al. (2021) to use Area Under the Precision-Recall Curve (AUPR) and Precision at 80% Recall (P@0.8R) as the evaluation metrics.

4.2 Implementations

Since legal SpanID tasks have a lower tolerance for missing important spans, we do not include start-end selector (i.e. $CE(P_{s,e}, Y_{s,e})$ and $\alpha\mathcal{L}_{ct}$ in Eq. (10)) in the CCE models but follow Hendrycks et al. (2021) to output top 20 spans from span predictor for each input example in order to extract spans as much as possible. While for NER, ABSA, and SBPD, we use our optimized architecture and objective.

For a fair comparison with existing works, our models utilize BERT (Devlin et al., 2019) as the text encoder for ABSA and RoBERTa (Liu et al., 2019) for NER, CCE, and SBPD. Detailed configurations can be found in Appendix A.2.

4.3 Baselines

Note that our main contribution is to provide a new perspective to treat the PR relation as a kind of training data for augmentation. Therefore, we

compare with models built on the same encoder-only PLMs (Devlin et al., 2019; Liu et al., 2019). We are not focusing on pushing the SOTA results to new heights though some of the baselines already achieved SOTA performance.

NER: We compare with Tagging (Liu et al., 2019) and MRC (Li et al., 2020b) baselines. We also report the previous best approaches for each dataset, including RB-CRF+RM (Lin et al., 2021), CL-KL (Wang et al., 2021), T-NER (Ushio and Camacho-Collados, 2021) KaNa (Nie et al., 2021), and RoBERTa+BS (Zhu and Li, 2022).

ABSA: In addition to the MRC baseline, we also compare with previous approaches on top of BERT. These are SPAN-BERT (Hu et al., 2019), IMNBERT (He et al., 2019), RAEL (Chen and Qian, 2020) and Dual-MRC (Mao et al., 2021).

SBPD: For **News20** we only compare with MRC baseline due to the lack of related work. For **Social21**, we compare with top three approaches on its leaderboard, namely, Volta (Gupta et al., 2021), HOMADOS (Kaczyński and Przybyła, 2021), and TeamFPAI (Hou et al., 2021).

CCE: We compare with (1) MRC baseline, (2) stronger text encoders, including ALBERT (Lan et al., 2019) and DeBERTa (He et al., 2020), (3) the model continually pretrained on contracts: RoBERTa + CP (Hendrycks et al., 2021) and (4) the model leveraged the contract structure: ConReader (Xu et al., 2022a).

5 Results

5.1 Comparison Results

NER: Table 2 shows the performance on five NER datasets. Our PeerDA significantly out-

| Methods | Lap14 | | Rest14 | |
|----------------------|-------------|-------------|-------------|-------------|
| | UABSA | ATE | UABSA | ATE |
| SPAN-BERT | 61.3 | 82.3 | 73.7 | 86.7 |
| IMN-BERT | 61.7 | 77.6 | 70.7 | 84.1 |
| RACL | 63.4 | 81.8 | 75.4 | 86.4 |
| Dual-MRC | 65.9 | 82.5 | 76.0 | 86.6 |
| MRC (<u>Large</u>) | 63.2 | 83.9 | 72.9 | 86.8 |
| PeerDA | 65.9 | 84.6 | 73.9 | 86.8 |

Table 3: Performance on two ABSA subtasks on two datasets. Results are averages F_1 over 5 runs.

| Methods | News20 | | | Social21 | | |
|---------------------|--------|------|-------------|----------|------|-------------|
| | P | R | F_1 | P | R | F_1 |
| Volta | - | - | - | 50.1 | 46.4 | 48.2 |
| HOMADOS | - | - | - | 41.2 | 40.3 | 40.7 |
| TeamFPAI | - | - | - | 65.2 | 28.6 | 39.7 |
| MRC (<u>Base</u>) | 10.5 | 53.5 | 17.6 | 55.8 | 43.5 | 48.9 |
| PeerDA | 21.8 | 31.5 | 25.8 | 49.4 | 70.6 | 58.1 |

Table 4: PeerDA performance on two SBPD datasets.

performs the Tagging and MRC baselines. Precisely, compared to $\text{RoBERTa}_{\text{base}}$ MRC, PeerDA obtains 0.3, 6.0, 3.2, 1.5, and 2.9 F_1 gains on five datasets respectively. When implemented on $\text{RoBERTa}_{\text{large}}$, our PeerDA can further boost the performance and establishes new SOTA on three datasets, namely, **OntoNotes5**, **Movie**, and **Restaurant**. Note that the major improvement of PeerDA over MRC comes from higher Recall. It implies that PeerDA encourages models to give more span predictions.

ABSA: Table 3 depicts the results on ABSA. Compared to previous approaches, PeerDA mostly achieves better results on two subtasks, where it outperforms vanilla MRC by 2.7 and 1.0 F_1 on UABSA for two domains respectively.

SBPD: The results of two SBPD tasks are presented in Table 4. PeerDA outperforms MRC by 8.2 and 9.2 F_1 and achieves SOTA performance on **News20** and **Social21** respectively.

CCE: The results of CCE are shown in Table 5. PeerDA surpasses MRC by 8.7 AUPR and 13.3 P@0.8R and even surpasses the previous best model of larger size ($\text{ConReader}_{\text{large}}$) by 3.2 AUPR, reaching SOTA performance on **CUAD**.

5.2 Analysis on Augmentation Strategies

To explore how the size and category distribution of the augmented data affect the SpanID tasks, we conduct ablation study on the three augmen-

| Methods | #Params | AUPR | P@0.8R |
|--|---------|-------------|-------------|
| $\text{ALBERT}_{\text{xxlarge}}$ | 223M | 38.4 | 31.0 |
| $\text{RoBERTa}_{\text{base}} + \text{CP}$ | 125M | 45.2 | 34.1 |
| $\text{RoBERTa}_{\text{large}}$ | 355M | 48.2 | 38.1 |
| $\text{DeBERTa}_{\text{xlarge}}$ | 900M | 47.8 | 44.0 |
| $\text{ConReader}_{\text{large}}$ | 355M | 49.1 | 44.2 |
| MRC (<u>Base</u>) | 125M | 43.6 | 32.2 |
| PeerDA | 125M | 52.3 | 45.5 |

Table 5: PeerDA performance on CCE.

| Ablation Type | NER | UABSA | SBPD | CCE | Avg. |
|------------------------------|-------------|-------------|-------------|-------------|-------------|
| MRC | 72.7 | 68.1 | 33.3 | 43.6 | 54.4 |
| PeerDA-Size | 74.6 | 69.7 | 38.5 | 48.7 | 57.9 |
| PeerDA-Categ | 74.2 | 69.3 | 40.4 | 51.3 | 58.8 |
| PeerDA-Both (final) | 75.5 | 69.9 | 42.0 | 52.3 | 59.9 |

Table 6: Ablation study on data augmentation strategies. The results (F_1 for NER, UABSA, and SBPD. AUPR for CCE) are averaged of all datasets in each task.

tation strategies mentioned in Sec. 3.1.2, depicted in Table 6. Overall, all of the PeerDA variants are clearly superior to the MRC baseline and the PeerDA-both considering both data size and distribution issues performs the best. Another interesting finding is that PeerDA-Categ significantly outperforms PeerDA-Size on SBPD and CCE. We attribute the phenomenon to the fact that SBPD and CCE have a larger number of categories and consequently, the MRC model is more prone to the issue of skewed data distribution. Under this circumstance, PeerDA-Categ, the variant designed for compensating the long-tailed categories, can bring larger performance gains over MRC model. On the other hand, if the skewed data distribution is not severe (e.g. NER), or the category shows a weak correlation with the spans (i.e. UABSA), PeerDA-Size is more appropriate than PeerDA-Categ.

5.3 Analysis on Model Designs

Calculation of $P_{s,e}$ (Top part of Table 7) Under the same experimental setup ($\text{RoBERTa}_{\text{base}}$, batch size=32, sequence length=192, fp16), using our *general* method (Eq. (7)) to compute span score $P_{s,e}$ greatly reduces the memory footprint by more than 4 times with no performance drop, compared to the original *concat* method. Therefore, our *general* method allows a larger batch size for accelerating the training.

Contrastive Loss (Bottom part of Table 7) After we have settled on the *general* scoring function, we further investigate different methods to compute

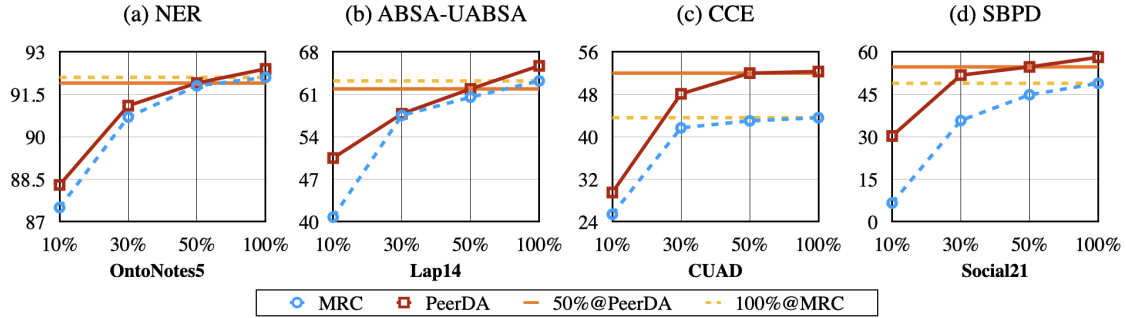


Figure 3: Performance on low-resource scenarios. We select one dataset for each SpanID task and report the test results (AUPR for CCE and F_1 for others) from the models trained on different proportions of the training data.

| Ablation Type | GPU | NER | UABSA | SBPD | Avg. |
|--|--------------|-------------|-------------|-------------|-------------|
| <i>Calculation of $P_{s,e}$</i> | | | | | |
| <i>concat</i> | 1x | 74.5 | 69.2 | 40.3 | 61.3 |
| <i>general (final)</i> | 0.23x | 75.0 | 69.4 | 40.8 | 61.7 |
| <i>Contrastive Loss</i> | | | | | |
| <i>Average</i> | 0.23x | 75.1 | 69.6 | 37.6 | 60.8 |
| <i>Max-Min (final)</i> | 0.23x | 75.5 | 69.9 | 42.0 | 62.4 |

Table 7: Ablation study on model designs. The F_1 scores are averaged of all datasets in each task. The |GPU| column denotes the GPU memory footprint of each variant under the same experimental setup.

| SRC → TGT | RoBERTa _{base} | | RoBERTa _{large} | |
|-----------------------|-------------------------|-------------|--------------------------|-------------|
| | MRC | PeerDA | MRC | PeerDA |
| Onto. → WNUT17 | 43.1 | 46.8 | 44.2 | 46.9 |
| Onto. → Rest. | 1.6 | 5.0 | 2.7 | 11.0 |
| Onto. → Movie | 25.0 | 26.7 | 26.7 | 27.8 |
| Average | 23.3 | 26.2 | 24.5 | 28.6 |

Table 8: F_1 scores on NER cross-domain transfer, where models trained on source-domain training data (SRC) are evaluated on target-domain test sets (TGT).

contrastive loss. We find that the *Average* method, which averages similarity differences between the query and all pairs of positive and negative spans, would affect SpanID performance when the task has more long-tailed categories (i.e. SBPD). While our *Max-Min* (strategy in Eq.(9)) is a relaxed regularization, which empirically is more suitable for SpanID tasks and consistently performs better than the *Average* method.

6 Further Discussions

In this section, we make further discussions to bring valuable insights of our PeerDA approach.

Out-of-domain Evaluation: We conduct out-of-domain evaluation on four English NER datasets, where the model is trained on **OntoNotes5**, the

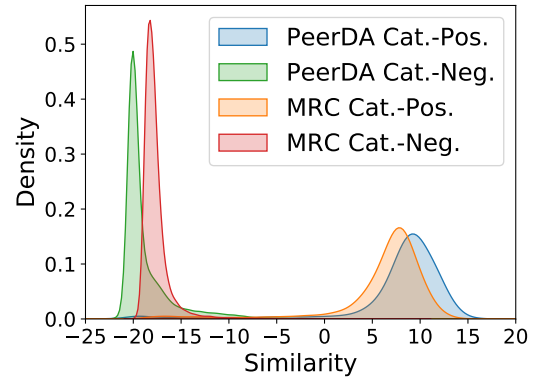


Figure 4: The distribution of similarity score between categories and their corresponding positive/negative spans on **Ontonotes5** test set.

largest dataset among them, and evaluated on the test part of another three datasets. Since these four datasets are from different domains and differ substantially in their categories, this setting largely eliminates the impact of superficial span-category patterns and thus it can faithfully reflect how well the MRC model exploits span semantics for prediction. The results are presented in Table 8. PeerDA can significantly exceed MRC on all three transfer pairs. On average, PeerDA achieves 2.9 and 4.1 F_1 gains over base-size MRC and large-size MRC respectively. These results verify our postulation that modeling the PR relation allows models to weigh more on the semantics for making predictions, and thus mitigates the over-fitting issue.

Semantic Distance: To gain a deeper understanding of the way in which PeerDA enhances model performance, we consider the span score (Eq. 7) as a measure of semantic similarity between a query and a span. In this context, we can create queries for all categories and visualize the similarity distribution between the categories and their corresponding positive and negative spans on **Ontonote5** test set. As shown in Figure 4, we can observe that

| | |
|--------------------|--|
| Multiple Labels | <i>I'm in Atlanta.</i> Gold: ("Atlanta", GPE) PeerDA: ("Atlanta", GPE); ("Atlanta", LOC) (41%) MRC: ("Atlanta", GPE) ("Atlanta", LOC) (3%) |
| Incorrect Label | <i>Why did it take us to get Sixty Minutes to do basic reporting to verify facts?</i> Gold: ("Sixty Minutes", ORG) PeerDA: ("Sixty Minutes", WORK_OF_ART) (37%) MRC: ("Sixty Minutes", WORK_OF_ART) (20%) |
| Missing Prediction | <i>Coming to a retailer near you, PlayStation pandemonium.</i> Gold: ("PlayStation", PRODUCT) PeerDA: \emptyset (19%) MRC: \emptyset (74%) |
| Other Errors | <i>I was guarded uh by the British Royal Marines actually because unfortunately they've had now um uh roadside bombs down there not suicide bombs.</i> Gold: ("the British Royal Marines", ORG) PeerDA: ("Royal Marines", ORG) (3%) MRC: ("Royal Marines", ORG) (3%) |

Table 9: Error analysis of base-sized PeerDA and MRC models on **Ontonotes5** test set. We randomly select 100 examples from the test set and compare the predictions and error percentages of the two models.

the use of PeerDA leads to an increased semantic similarity between spans and their corresponding categories, resulting in higher confidence in the prediction of correct spans. Furthermore, PeerDA has been shown to also create a larger similarity gap between positive and negative spans, facilitating their distinction.

Low-resource Evaluation: We simulate low-resource scenarios by randomly selecting 10%, 30%, 50%, and 100% of the training data for training SpanID models and show the comparison results between PeerDA and MRC on four SpanID tasks in Figure 3. As can be seen, our PeerDA further enhances the MRC model in all sizes of training data and the overall trends are consistent across the above four tasks. When training PeerDA with 50% of the training data, it can reach or even exceed the performance of MRC trained on the full training set. These results demonstrate the effectiveness of our PeerDA in low-resource scenarios.

Error Analysis: In order to know the typical failure of PeerDA, we randomly sample 100 error cases from **Ontonotes5** test set for analysis. As shown in Table 9, there are four major groups:

- *Multiple Labels:* PeerDA would assign multiple labels to the same detected span. And in most cases (35/41), this error occurs among similar categories, such as LOC, GPE, and ORG.
- *Incorrect Label:* Although spans are correctly detected, PeerDA assigns them the wrong categories. Note that MRC even cannot detect many of those spans (23/37). As a result, PeerDA sig-

nificantly improves the model’s capability to detect spans, but still faces challenges in category classification.

- *Missing Prediction:* Compared to MRC, PeerDA tends to predict more spans. Therefore it alleviates the missing prediction issue that MRC mostly suffers.
- *Other Errors:* There are several other errors, such as the incorrect span boundary caused by articles or nested entities.

7 Conclusions

In this paper, we propose a novel PeerDA approach for SpanID tasks to augment training data from the perspective of capturing the PR relation. PeerDA has two unique advantages: (1) It is capable to leverage abundant but previously unused PR relation as additional training data. (2) It alleviates the over-fitting issue of MRC models by pushing the models to weigh more on semantics. We conduct extensive experiments to verify the effectiveness of PeerDA. Further in-depth analyses demonstrate that the improvement of PeerDA comes from a better semantic understanding capability.

Limitations

In this section, we discuss the limitations of this work as follows:

- PeerDA leverages labeled spans in the existing training set to conduct data augmentation. This means that PeerDA improves the semantics learning of existing labeled spans, but is ineffective

to classify other spans outside the training set. Therefore, it would be beneficial to engage outer source knowledge (e.g. Wikipedia), where a variety of important entities and text spans can also be better learned with our PeerDA approach.

- Since PeerDA is designed in the MRC formulation on top of the encoder-only Pre-trained Language Models (PLMs) (Devlin et al., 2019; Liu et al., 2019), it is not comparable with other methods built on encoder-decoder PLMs (Yan et al., 2021b; Chen et al., 2022; Zhang et al., 2021b; Yan et al., 2021a). It would be of great value to try PeerDA on encoder-decoder PLMs such as BART (Lewis et al., 2020) and T5 (Raffel et al., 2020), to see whether PeerDA is a general approach regardless of model architecture.
- As shown in Table 9, although PeerDA can significantly alleviate the *Missing Predictions*, the most prevailing error in the MRC model, PeerDA also introduces some new errors, i.e. *Multiple labels* and *Incorrect Label*. It should be noted that those problematic spans are usually observed in different span sets, where they would learn different category semantics from their peers. Therefore, we speculate that those spans tend to leverage the learned category semantics more than their context information to determine their categories. We hope such finding can shed light on future research to further improve PeerDA.

References

- Lidong Bing, Sneha Chaudhari, Richard Wang, and William Cohen. 2015. [Improving distant supervision for information extraction using label propagation through lists](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*.
- Lidong Bing, Wai Lam, and Tak-Lam Wong. 2013. [Wikipedia entity expansion and attribute extraction from the web using semi-supervised learning](#). In *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*.
- Ilias Chalkidis, Ion Androutsopoulos, and Achilleas Michos. 2017. Extracting contract elements. In *Proceedings of the 16th edition of the International Conference on Artificial Intelligence and Law*.
- Gal Chechik, Varun Sharma, Uri Shalit, and Samy Bengio. 2010. Large scale online learning of image similarity through ranking. *Journal of Machine Learning Research*.
- Shuguang Chen, Gustavo Aguilar, Leonardo Neves, and Thamar Solorio. 2021. [Data augmentation for cross-domain named entity recognition](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Xiang Chen, Lei Li, Shumin Deng, Chuanqi Tan, Changliang Xu, Fei Huang, Luo Si, Huajun Chen, and Ningyu Zhang. 2022. [LightNER: A lightweight tuning paradigm for low-resource NER via pluggable prompting](#). In *Proceedings of the 29th International Conference on Computational Linguistics*.
- Zhuang Chen and Tiejun Qian. 2020. [Relation-aware collaborative learning for unified aspect-based sentiment analysis](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Liying Cheng, Lidong Bing, Qian Yu, Wei Lu, and Luo Si. 2020. [APE: Argument pair extraction from peer review and rebuttal via multi-task learning](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jason P.C. Chiu and Eric Nichols. 2016. [Named entity recognition with bidirectional LSTM-CNNs](#). *Transactions of the Association for Computational Linguistics*.
- Giovanni Da San Martino, Alberto Barrón-Cedeño, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020. [SemEval-2020 task 11: Detection of propaganda techniques in news articles](#). In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*.
- Giovanni Da San Martino, Seunghak Yu, Alberto Barrón-Cedeño, Rostislav Petrov, and Preslav Nakov. 2019. [Fine-grained analysis of propaganda in news article](#). 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)*.
- Xiang Dai and Heike Adel. 2020. [An analysis of simple data augmentation for named entity recognition](#). In *Proceedings of the 28th International Conference on Computational Linguistics*.
- Leon Derczynski, Eric Nichols, Marieke van Erp, and Nut Limsopatham. 2017. [Results of the WNUT2017 shared task on novel and emerging entity recognition](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*.

- Dimitar Dimitrov, Bishr Bin Ali, Shaden Shaar, Firoj Alam, Fabrizio Silvestri, Hamed Firooz, Preslav Nakov, and Giovanni Da San Martino. 2021. [SemEval-2021 task 6: Detection of persuasion techniques in texts and images](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*.
- Bosheng Ding, Linlin Liu, Lidong Bing, Canasai Kruengkrai, Thien Hai Nguyen, Shafiq Joty, Luo Si, and Chunyan Miao. 2020. [DAGA: Data augmentation with a generation approach for low-resource tagging tasks](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Ning Ding, Guangwei Xu, Yulin Chen, Xiaobin Wang, Xu Han, Pengjun Xie, Haitao Zheng, and Zhiyuan Liu. 2021. [Few-NERD: A few-shot named entity recognition dataset](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Edward Hovy. 2021. [A survey of data augmentation approaches for NLP](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*.
- Kshitij Gupta, Devansh Gautam, and Radhika Mamidi. 2021. [Volta at SemEval-2021 task 6: Towards detecting persuasive texts and images using textual and multimodal ensemble](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2020. [Deberta: Decoding-enhanced bert with disentangled attention](#). *arXiv preprint arXiv:2006.03654*.
- Ruidan He, Wee Sun Lee, Hwee Tou Ng, and Daniel Dahlmeier. 2019. [An interactive multi-task learning network for end-to-end aspect-based sentiment analysis](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. 2021. [Cuad: An expert-annotated nlp dataset for legal contract review](#). *arXiv preprint arXiv:2103.06268*.
- Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. [Teaching machines to read and comprehend](#). In *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 1*.
- Xiaolong Hou, Junsong Ren, Gang Rao, Lianxin Lian, Zhihao Ruan, Yang Mo, and Jianping Shen. 2021. [FPAI at SemEval-2021 task 6: BERT-MRC for propaganda techniques detection](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*.
- Minghao Hu, Yuxing Peng, Zhen Huang, Dongsheng Li, and Yiwei Lv. 2019. [Open-domain targeted sentiment analysis via span-based extraction and classification](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Jiabin Huang, Chunyuan Li, Krishan Subudhi, Damien Jose, Shobana Balakrishnan, Weizhu Chen, Baolin Peng, Jianfeng Gao, and Jiawei Han. 2021. [Few-shot named entity recognition: An empirical baseline study](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Alankar Jain, Bhargavi Paranjape, and Zachary C. Lipton. 2019. [Entity projection via machine translation for cross-lingual NER](#). 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)*.
- Konrad Kaczyński and Piotr Przybyła. 2021. [HOMADOS at SemEval-2021 task 6: Multi-task learning for propaganda detection](#). In *Proceedings of the 15th International Workshop on Semantic Evaluation (SemEval-2021)*.
- Sosuke Kobayashi. 2018. [Contextual augmentation: Data augmentation by words with paradigmatic relations](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*.
- Canasai Kruengkrai, Thien Hai Nguyen, Sharifah Mahani Aljunied, and Lidong Bing. 2020. [Improving low-resource named entity recognition using joint sentence and token labeling](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. [Albert: A lite bert for self-supervised learning of language representations](#). *arXiv preprint arXiv:1909.11942*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Kun Li, Chengbo Chen, Xiaojun Quan, Qing Ling, and Yan Song. 2020a. [Conditional augmentation for aspect term extraction via masked sequence-to-sequence generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.

- Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020b. [A unified MRC framework for named entity recognition](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Xiaoya Li, Fan Yin, Zijun Sun, Xiayu Li, Arianna Yuan, Duo Chai, Mingxin Zhou, and Jiwei Li. 2019. [Entity-relation extraction as multi-turn question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*.
- Xin Li, Lidong Bing, Wenxuan Zhang, Zheng Li, and Wai Lam. 2020c. [Unsupervised cross-lingual adaptation for sequence tagging and beyond](#). *arXiv preprint arXiv:2010.12405*.
- Bill Yuchen Lin, Wenyang Gao, Jun Yan, Ryan Moreno, and Xiang Ren. 2021. [RockNER: A simple method to create adversarial examples for evaluating the robustness of named entity recognition models](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Bing Liu. 2012. [Sentiment analysis and opinion mining](#). *Synthesis lectures on human language technologies*.
- Jian Liu, Yubo Chen, Kang Liu, Wei Bi, and Xiaojiang Liu. 2020. [Event extraction as machine reading comprehension](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jian Liu, Yufeng Chen, and Jinan Xu. 2021a. [Machine reading comprehension as data augmentation: A case study on implicit event argument extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Jingjing Liu, Panupong Pasupat, Scott Cyphers, and Jim Glass. 2013a. [Asgard: A portable architecture for multilingual dialogue systems](#). In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*.
- Jingjing Liu, Panupong Pasupat, Yining Wang, Scott Cyphers, and Jim Glass. 2013b. [Query understanding enhanced by hierarchical parsing structures](#). In *2013 IEEE Workshop on Automatic Speech Recognition and Understanding*.
- Linlin Liu, Bosheng Ding, Lidong Bing, Shafiq Joty, Luo Si, and Chunyan Miao. 2021b. [MulDA: A multilingual data augmentation framework for low-resource cross-lingual NER](#). 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)*.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*.
- Ilya Loshchilov and Frank Hutter. 2019. [Decoupled weight decay regularization](#). In *International Conference on Learning Representations*.
- Thang Luong, Hieu Pham, and Christopher D. Manning. 2015. [Effective approaches to attention-based neural machine translation](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*.
- Yue Mao, Yi Shen, Chao Yu, and Longjun Cai. 2021. [A joint training dual-mrc framework for aspect based sentiment analysis](#). *arXiv preprint arXiv:2101.00816*.
- Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. [The natural language decathlon: Multitask learning as question answering](#). *arXiv preprint arXiv:1806.08730*.
- Binling Nie, Ruixue Ding, Pengjun Xie, Fei Huang, Chen Qian, and Luo Si. 2021. [Knowledge-aware named entity recognition with alleviating heterogeneity](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- Yuyang Nie, Yuanhe Tian, Xiang Wan, Yan Song, and Bo Dai. 2020. [Named entity recognition for social media texts with semantic augmentation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Sean Papay, Roman Klinger, and Sebastian Padó. 2020. [Dissecting span identification tasks with performance prediction](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Nanyun Peng and Mark Dredze. 2015. [Named entity recognition for Chinese social media with jointly trained embeddings](#). In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*.
- Maria Pontiki, Dimitris Galanis, John Pavlopoulos, Harris Papageorgiou, Ion Androutsopoulos, and Suresh Manandhar. 2014. [SemEval-2014 task 4: Aspect based sentiment analysis](#). In *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*.
- Sameer Pradhan, Alessandro Moschitti, Nianwen Xue, Hwee Tou Ng, Anders Björkelund, Olga Uryupina, Yuchen Zhang, and Zhi Zhong. 2013. [Towards robust linguistic analysis using OntoNotes](#). In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*.

- Qingyu Tan, Lu Xu, Lidong Bing, and Hwee Tou Ng. 2023. Class-adaptive self-training for relation extraction with incompletely annotated training data. In *Finding of the 61th Annual Meeting of the Association for Computational Linguistics*.
- Asahi Ushio and Jose Camacho-Collados. 2021. **T-NER: An all-round python library for transformer-based named entity recognition**. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. **Attention is all you need**. In *Advances in Neural Information Processing Systems*.
- Xinyu Wang, Yong Jiang, Nguyen Bach, Tao Wang, Zhongqiang Huang, Fei Huang, and Kewei Tu. 2021. **Improving named entity recognition by external context retrieving and cooperative learning**. 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)*.
- Yaqing Wang, Subhabrata Mukherjee, Haoda Chu, Yuanheng Tu, Ming Wu, Jing Gao, and Ahmed Hassan Awadallah. 2020. Adaptive self-training for few-shot neural sequence labeling. *arXiv preprint arXiv:2010.03680*.
- Jason Wei and Kai Zou. 2019. **EDA: Easy data augmentation techniques for boosting performance on text classification tasks**. 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)*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. **Transformers: State-of-the-art natural language processing**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*.
- Xing Wu, Shangwen Lv, Liangjun Zang, Jizhong Han, and Songlin Hu. 2019. Conditional bert contextual augmentation. In *International Conference on Computational Science*.
- Rong Xiang, Emmanuele Chersoni, Qin Lu, Chu-Ren Huang, Wenjie Li, and Yunfei Long. 2021. Lexical data augmentation for sentiment analysis. *Journal of the Association for Information Science and Technology*.
- Qizhe Xie, Zihang Dai, Eduard Hovy, Minh-Thang Luong, and Quoc V Le. 2019. Unsupervised data augmentation for consistency training. *arXiv preprint arXiv:1904.12848*.
- Qizhe Xie, Minh-Thang Luong, Eduard Hovy, and Quoc V Le. 2020. Self-training with noisy student improves imagenet classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Xiaoyu Xing, Zhijing Jin, Di Jin, Bingning Wang, Qi Zhang, and Xuanjing Huang. 2020. **Tasty burgers, soggy fries: Probing aspect robustness in aspect-based sentiment analysis**. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Lu Xu, Lidong Bing, and Wei Lu. 2023a. Sampling better negatives for distantly supervised named entity recognition. In *Finding of the 61th Annual Meeting of the Association for Computational Linguistics*.
- Lu Xu, Zhanming Jie, Wei Lu, and Lidong Bing. 2021. **Better feature integration for named entity recognition**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- Weiwen Xu, Yang Deng, Wenqiang Lei, Wenlong Zhao, Tat-Seng Chua, and Wai Lam. 2022a. **ConReader: Exploring implicit relations in contracts for contract clause extraction**. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*.
- Weiwen Xu, Xin Li, Wai Lam, and Lidong Bing. 2023b. mpmr: A multilingual pre-trained machine reader at scale. In *The 61th Annual Meeting of the Association for Computational Linguistics*.
- Weiwen Xu, Xin Li, Wenxuan Zhang, Meng Zhou, Lidong Bing, Wai Lam, and Luo Si. 2022b. From clozing to comprehending: Retrofitting pre-trained language model to pre-trained machine reader. *arXiv preprint arXiv:2212.04755*.
- Hang Yan, Junqi Dai, Tuo Ji, Xipeng Qiu, and Zheng Zhang. 2021a. **A unified generative framework for aspect-based sentiment analysis**. 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)*.
- Hang Yan, Tao Gui, Junqi Dai, Qipeng Guo, Zheng Zhang, and Xipeng Qiu. 2021b. **A unified generative framework for various NER subtasks**. 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)*.

- Pan Yang, Xin Cong, Zhenyu Sun, and Xingwu Liu. 2021. [Enhanced language representation with label knowledge for span extraction](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Kang Min Yoo, Dongju Park, Jaewook Kang, Sangwoo Lee, and Woomyoung Park. 2021. [GPT3Mix: Leveraging large-scale language models for text augmentation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*.
- Wenxuan Zhang, Yang Deng, Xin Li, Lidong Bing, and Wai Lam. 2021a. [Aspect-based sentiment analysis in question answering forums](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*.
- Wenxuan Zhang, Xin Li, Yang Deng, Lidong Bing, and Wai Lam. 2021b. [Towards generative aspect-based sentiment analysis](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*.
- Wenxuan Zhang, Xin Li, Yang Deng, Lidong Bing, and Wai Lam. 2022. A survey on aspect-based sentiment analysis: tasks, methods, and challenges. *IEEE Transactions on Knowledge and Data Engineering*.
- Ran Zhou, Xin Li, Lidong Bing, Erik Cambria, and Chunyan Miao. 2023. Improving self-training for cross-lingual named entity recognition with contrastive and prototype learning. In *The 61th Annual Meeting of the Association for Computational Linguistics*.
- Ran Zhou, Xin Li, Lidong Bing, Erik Cambria, Luo Si, and Chunyan Miao. 2022a. [ConNER: Consistency training for cross-lingual named entity recognition](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*.
- Ran Zhou, Xin Li, Ruidan He, Lidong Bing, Erik Cambria, Luo Si, and Chunyan Miao. 2022b. [MELM: Data augmentation with masked entity language modeling for low-resource NER](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.
- Enwei Zhu and Jinpeng Li. 2022. [Boundary smoothing for named entity recognition](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*.

A Appendix

A.1 Task Overview

We conduct experiments on four SpanID tasks with diverse domains, including Named Entity Recognition (NER), Aspect Based Sentiment Analysis (ABSA), Contract Clause Extraction (CCE) and Span Based Propaganda Detection (SBPD), to show the overall effectiveness of our PeerDA. The dataset statistics are summarized in Table 1.

NER: It is a traditional SpanID task, where spans denote the named entities in the input text and category labels denote their associated entity types. We evaluate five datasets from four domains:

- **OntoNotes5** (Pradhan et al., 2013) is a large-scale mixed domain NER dataset covering News, Blog and Dialogue. To make a fair comparison in the robustness experiments in Sec. A.3, we use the datasets from Lin et al. (2021), which only add adversarial attack to the 11 entity types, while leaving out 7 numerical types.
- **WNUT17** (Derczynski et al., 2017) is a benchmark NER dataset in social media domain. For fair comparison, we follow the data pre-processing protocols in Nie et al. (2020).
- **Movie** (Liu et al., 2013b) is a movie domain dataset containing movie queries, where long spans are annotated such as a movie’s origin or plot. We use the defaulted data split strategy into train, test sets.
- **Restaurant** (Liu et al., 2013a) contains queries in restaurant domain. Similar to Movie, we use the defaulted data split strategy.
- **Weibo** (Peng and Dredze, 2015) is a Chinese benchmark NER dataset in social media domain. We exactly follow the official data split strategy into train, dev and test sets.

ABSA: It is a fine-grained sentiment analysis task centering on aspect terms (Zhang et al., 2022). We explore two ABSA sub-tasks:

- **Aspect Term Extraction (ATE)** is to extract aspect terms, where there is only one query asking if there are any aspect terms in the input text.
- **Unified Aspect Based Sentiment Analysis (UABSA)** is to jointly extract aspect terms and predict their sentiment polarities. We formulate it

as a SpanID task by treating the sentiment polarities, namely, positive, negative, and neutral, as three category labels, and aspect terms as spans.

We evaluate the two sub-tasks on two datasets, including the laptop domain dataset **Lap14** and restaurant domain dataset **Rest14** from SemEval Shared tasks (Pontiki et al., 2014). We use the processed data from Zhang et al. (2021b).

CCE: It is a legal NLP task to detect and classify contract clauses into relevant clause types, such as "Governing Law" and "Uncapped Liability". The goal of CCE is to reduce the labor of legal professionals in reviewing contracts of dozens or hundreds of pages long. CCE is also a kind of SpanID task where spans are those contract clauses that warrant review or analysis and labels are predefined clause types. We conduct experiments on CCE using **CUAD** (Hendrycks et al., 2021), where they annotate contracts from Electronic Data Gathering, Analysis and Retrieval (EDGAR) with 41 clause types. We follow Hendrycks et al. (2021) to split the contracts into segments within the length limitation of pretrained language models and treat each individual segment as one example. We also follow their data split strategy.

SBPD: It is a typical SpanID task that aims to detect both the text fragment (i.e. spans) where a persuasion technique is being used as well as its technique type (i.e. category labels). We use the **News20** and **Social21** from two SemEval shared tasks (Da San Martino et al., 2020; Dimitrov et al., 2021) and follow the official data split strategy. Note that **News20** does not provide the golden label for the test set. Therefore, we evaluate **News20** on the dev set.

A.2 Implementations

We use Huggingface’s implementations of BERT and RoBERTa (Wolf et al., 2020)³. The hyperparameters can be found in Table 10. We use Tesla V100 GPU cards for conducting all the experiments. We follow the default learning rate schedule and dropout settings used in BERT. We use AdamW (Loshchilov and Hutter, 2019) as our optimizer. The margin term M is set to 0 for NER and ABSA, and 1 for SBPD. The balance rate α is set to 0.1.

³Chinese RoBERTa is from <https://github.com/ymcui/Chinese-BERT-wwm>.

| Dataset | OntoNotes5 | WNUT17 | Movie | Restaurant | Weibo | Lap14 | Rest14 | CUAD | News20 | Social21 |
|---------------|------------|--------|-------|------------|-------|-------|--------|------|--------|----------|
| Query Length | 32 | 32 | 64 | 64 | 64 | 24 | 24 | 256 | 80 | 80 |
| Input Length | 160 | 160 | 160 | 128 | 192 | 128 | 128 | 512 | 200 | 200 |
| Batch Size | 32 | 32 | 32 | 32 | 8 | 16 | 16 | 16 | 16 | 16 |
| Learning Rate | 2e-5 | 1e-5 | 1e-5 | 1e-5 | 1e-5 | 2e-5 | 2e-5 | 5e-5 | 2e-5 | 3e-5 |
| λ | 1 | 1 | 1 | 1 | 1 | 1 | 1 | -0.5 | 0.5 | 1 |

Table 10: Hyper-parameters settings.

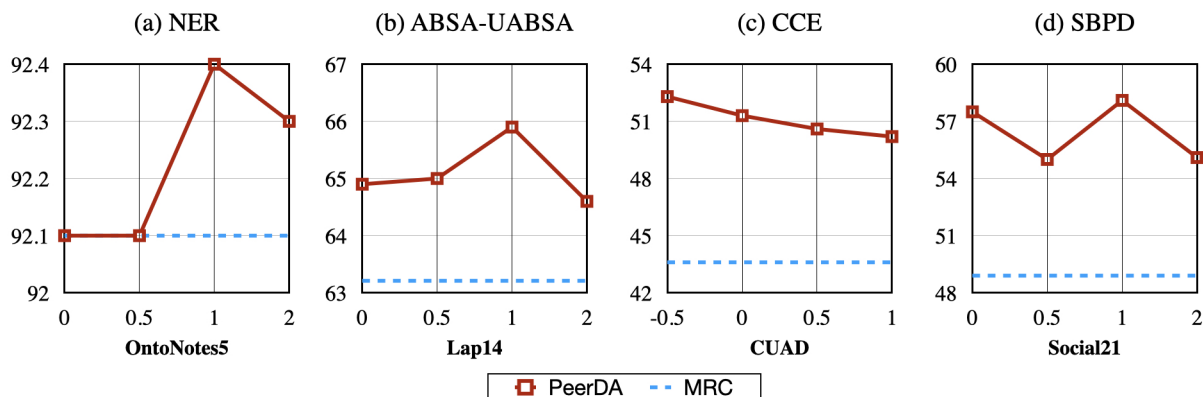


Figure 5: Performance in terms of different DA rates λ . We vary λ to get different volumes of PR-based training data.

| Methods | OntoNotes5 | | | | Lap14 | |
|---------|-------------|-------------|-------------|--------------|-------------|-------------|
| | Ori | Adv. full | Adv. entity | Adv. context | Ori | Adv. |
| Tagging | 89.8 | 56.6 | 61.9 | 83.6 | 62.3 | 44.5 |
| MRC | 90.0 | 55.3 | 61.3 | 83.3 | 63.2 | 46.9 |
| PeerDA | 90.1 | 55.9 | 61.0 | 84.1 | 65.9 | 50.1 |

Table 11: Robustness experiments against adversarial attacks. The results are reported on both original (Ori.) sets and the adversarial (Adv.) sets.

A.3 Robustness

To verify the advantage of PeerDA against the adversarial attack, we conduct robustness experiments using the adversarial dev set of **OntoNotes5** (Lin et al., 2021) on NER and adversarial test set of **Lap14** (Xing et al., 2020) on UABSA. Table 11 shows the performance on the original and the adversarial sets. On **OntoNotes5 full** adversarial set, PeerDA improves the robustness of the model compared to MRC but slightly degrades compared to Tagging. To investigate why this happens, we evaluate each type of adversarial attack independently, including *entity* attack that replaces entities to other entities not presented in the training set and *context* attack that replaces the context of entities. It shows that PeerDA does not work well on *entity* attack because we only use entities in the training set to conduct data augmentation, which

| Methods | OntoNotes5 | Lap14 | CUAD | Social21 |
|----------------|-------------|-------------|-------------|-------------|
| MRC+MenReplace | 91.1 | 63.7 | 45.2 | 50.8 |
| PeerDA | 92.4 | 65.9 | 52.3 | 58.1 |

Table 12: Performance on peer-driven DA approaches.

is intrinsically ineffective to this adversarial attack. This motivates us to engage outer source knowledge (e.g. Wikipedia) into our PeerDA approach in future work. On **Lap14**, PeerDA significantly improves Tagging and MRC by 5.6 and 3.2 F_1 on the adversarial set respectively.

A.4 Peer-driven DA

We compare PeerDA with Mention Replacement (MenReplace) (Dai and Adel, 2020), another Peer-driven DA approach randomly replaces a span mention in the context with another mention of the same category in the training set. The results of four SpanID tasks are presented in Table 12. PeerDA exhibits better performance than MenReplace on all four tasks. In addition, MenReplace would easily break the text coherence as a result of putting span mentions into the incompatible context, while PeerDA can do a more natural augmentation without harming the context.

A.5 Effect of DA Rate

We vary the DA rate λ to investigate how the volume of PR-based training data affect the SpanID models performance.

Figure 5 shows the effect of different λ in four SpanID tasks. PeerDA mostly improves the MRC in all different trials of λ and we suggest that some parameter tuning for λ is beneficial to obtain optimal results.

Another observation is that too large λ would do harm to the performance. Especially on CCE, due to the skewed distribution and a large number of categories, PeerDA can produce a huge size of PR-based training data. We speculate that too much PR-based training data would affect the learning of BL-based training data and thus affect the model's ability to solve a SpanID task, causing the optimal λ to be a negative value. In addition, too much PR-based training data would also increase the training cost. As a result, we should maintain an appropriate ratio of BL-based and PR-based training data to keep a reasonable performance on SpanID tasks.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
Section Limitations
- A2. Did you discuss any potential risks of your work?
Not applicable. Left blank.
- A3. Do the abstract and introduction summarize the paper’s main claims?
Abstract, Section 1
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

Section 4

- B1. Did you cite the creators of artifacts you used?
Section 4
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
Not applicable. Left blank.
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
Section 4
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
Not applicable. Left blank.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
Section 4, Appendix A.1
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
Section 4

C Did you run computational experiments?

Section 5

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
Appendix A.2

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

Section 4, Appendix A.2, Appendix A.5

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

Section 5

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

Section 4, Section A.2

D Did you use human annotators (e.g., crowdworkers) or research with human participants?

Left blank.

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

No response.

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

No response.

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

No response.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

No response.

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

No response.