

MCL-NER: Cross-Lingual Named Entity Recognition via Multi-view Contrastive Learning

Ying Mo¹, Jian Yang^{1*}, Jiahao Liu², Qifan Wang³, Ruoyu Chen⁴,
Jingang Wang², Zhoujun Li^{1*},

¹Beihang University; ²Meituan; ³Meta AI

⁴Beijing Information Science and Technology University

{moying, jiaya, lizj}@buaa.edu.cn; {liujiahao12, wangjingang02}@meituan.com;
wqfcr@fb.com; chenruoyu@bistu.edu.cn

Abstract

Cross-lingual named entity recognition (CrossNER) faces challenges stemming from uneven performance due to the scarcity of multilingual corpora, especially for non-English data. While prior efforts mainly focus on data-driven transfer methods, a significant aspect that has not been fully explored is aligning both semantic and token-level representations across diverse languages. In this paper, we propose Multi-view Contrastive Learning for Cross-lingual Named Entity Recognition (MCL-NER). Specifically, we reframe the CrossNER task into a problem of recognizing relationships between pairs of tokens. This approach taps into the inherent contextual nuances of token-to-token connections within entities, allowing us to align representations across different languages. A multi-view contrastive learning framework is introduced to encompass semantic contrasts between source, codeswitched, and target sentences, as well as contrasts among token-to-token relations. By enforcing agreement within both semantic and relational spaces, we minimize the gap between source sentences and their counterparts of both codeswitched and target sentences. This alignment extends to the relationships between diverse tokens, enhancing the projection of entities across languages. We further augment CrossNER by combining self-training with labeled source data and unlabeled target data. Our experiments on the XTREME benchmark, spanning 40 languages, demonstrate the superiority of MCL-NER over prior data-driven and model-based approaches. It achieves a substantial increase of nearly +2.0 F_1 scores across a broad spectrum and establishes itself as the new state-of-the-art performer.

Introduction

Cross-lingual named entity recognition (CrossNER) suffered from significant performance degradation in low-resource languages with limited data. In response to this challenge, the advent of multilingual pre-trained models (Devlin et al. 2018; Conneau et al. 2019; Yang et al. 2020; Chi et al. 2020) has prompted the development of model-based approaches (Keung, Lu, and Bhardwaj 2019; Bari, Joty, and Jwalapuram 2020; Wu et al. 2020a,c). These methods aim to facilitate knowledge transfer from languages with ample resources to those with fewer resources. Furthermore,

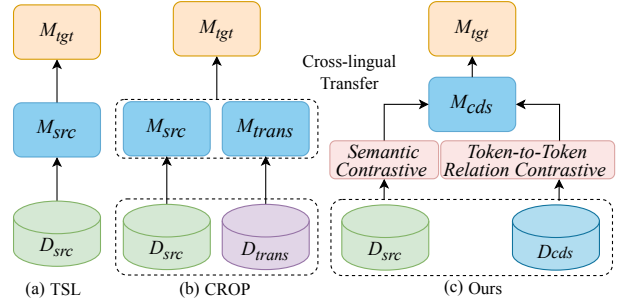


Figure 1: Illustration of MCL-NER vs. existing methods TSL (Wu et al. 2020a) and CROP (Yang et al. 2022a). D_{src} , D_{trans} and D_{cds} are the source, translated-source and code-switched data respectively. M_* represents the trained models from the corresponding data. Our model leverages multi-view contrastive learning to bridge the gap between cross-lingual semantic and token-to-token representations.

recent research (Wu et al. 2020b; Yang et al. 2022a; Ashok and Lipton 2023; Zhou et al. 2022, 2023) unifies model-based and data-based transfer techniques to enhance CrossNER. The efficacy of these methods hinges on both the inherent cross-lingual capabilities of the pre-trained models and the quality of the synthetic data generated through phrase-level and sentence-level translation.

Along the research line of leveraging the cross-lingual pre-trained model, previous works (Mayhew, Tsai, and Roth 2017; Xie et al. 2018; Wu et al. 2020b; Chen et al. 2021; Yang et al. 2022a) perform phrase-level and sentence-level translation and annotate corresponding target entities. These methods can be broadly classified into two groups, as illustrated in Figure 1. Within the first group, methods like TSL (Wu et al. 2020a) generate soft labels in the target language based on the model trained using the source language. These soft labels are then utilized to facilitate the training of a NER model for the target language. In contrast, techniques belonging to the second group, such as CROP (Yang et al. 2022a), leverage both the source model and translation data to fine-tune the NER model for the target domain. Despite their effectiveness, most existing studies focus primarily on aligning semantic spaces, often disregarding the crucial cross-lingual syntactic context encompassing token-to-token relationships within entities. Yet, these

relations play a pivotal role in cross-lingual NER learning, given that the structure of token-to-token relationships within bilingual sentences should exhibit similar patterns. Consequently, formulating a CrossNER framework capable of capturing token-level contextual nuances across diverse languages emerges as a significant challenge.

To address the above limitation, in this work we propose a novel multi-view contrastive learning framework for cross-lingual named entity recognition. We reformulate the problem of entity recognition into token-to-token relation classification. A multi-view contrastive learning with source-codeswitch semantic contrastive and token-to-token relation contrastive is employed to train the cross-lingual NER model. Specifically, the token-to-token relation aspect concentrates on both intra-entity and extra-entity relationships. In the context of intra-entity relations, two tokens’ adjacency is scrutinized, along with determining if they signify the start-end relationship, i.e., whether they respectively represent the initial and concluding words of the same entity. Conversely, extra-entity relations between tokens signify that they do not pertain to the same entity. To improve the representation of entities across different languages, we introduce sentence-level semantics source-codeswitch contrastive learning. Within this framework, code-switching sentences incorporate both source and target tokens concurrently, promoting more effective cross-lingual knowledge transfer. Further, the multilingual NER model can be iteratively trained on the distilled multilingual corpora from the initial multilingual model, which makes full use of unlabeled training data to improve the performance of the CrossNER. We conduct experiments on the XTREME benchmark covering 40 languages and then test on the CoNLL benchmark of 4 languages. Experimental results demonstrate that our method significantly outperforms most cross-lingual sequence labeling and span-based methods. Extensive probing experiments further analyze how our method can benefit the token-level cross-lingual representation by encouraging the consistency of different languages.

Our main contributions are summarized as follows:

- We develop a multi-view contrastive learning approach with both sentence-level and token-level aligning, simultaneously enhancing the semantic and entity representation of different languages.
- We introduce the code-switched data together with the source data to jointly conduct the cross-lingual transfer, which effectively captures the relationships between tokens in entities, improving the CrossNER performance.
- We conduct comprehensive experiments on two benchmarks, demonstrating competitive cross-lingual NER performance, establishing new state-of-the-art results on most of the evaluated cross-lingual transfer pairs (XTREME-40 and CoNLL).

Related work

Cross-lingual NER Cross-language named entity recognition (CrossNER) achieves significant progress in recent years due to the development of pre-trained language models (Ni, Dinu, and Florian 2017; Mayhew, Tsai, and Roth 2017;

Xie et al. 2018; Wu and Dredze 2019; Yu, Bohnet, and Poesio 2020; Hu et al. 2020; Wu et al. 2020b,a; Liu et al. 2021a; Han et al. 2022; Zhou et al. 2022; Yang et al. 2022a). The CrossNER methods is mainly divided into two categories: model-transfer and data-transfer methods. Model-transfer methods (Xie et al. 2018) generally use language features to train a NER model on the labeled source language data and then directly uses it on the target language data. These features include aligned word representations (Ni, Dinu, and Florian 2017; Wu and Dredze 2019; Li et al. 2021), Wikifier features (Mayhew, Tsai, and Roth 2017), and meta-learning (Wu et al. 2020d). Data-transfer methods (Wu et al. 2020b,a; Zhou et al. 2022; Yang et al. 2022c,b,a) construct pseudo-labels for target language data by translating data (Wu et al. 2020b; Yang et al. 2021) from the source language typically. Further, self-training can continue to solve the lack of target data based on the existing trained CrossNER model. But these methods usually use sequence tagging and do not consider the impact of token-to-token relationships in NER.

Contrastive learning Contrastive learning is used in computer vision for image classification (Chen et al. 2020; He et al. 2020; Khosla et al. 2020) and is now widely used in various tasks (Chuang et al. 2020; Giorgi et al. 2020; You et al. 2020; Hou et al. 2020; Gao, Yao, and Chen 2021; Das et al. 2021; Chen et al. 2022). In NLP, (Giorgi et al. 2020; Gao, Yao, and Chen 2021; Chen et al. 2022) proposes to enhance semantic representation and a pre-training model based on contrastive learning. Hou et al. (2020) applied contrastive learning to slot filling and then Das et al. (2021) proposed CONTaiNER for few-shot NER combining contrastive learning with Gaussian distribution. The contrastive learning can effectively pull the distance between positive samples and push the distance between negative samples to achieve better recognition results.

Problem Formulation

CrossNER aims to extract entities from target language sentences and assign them to predefined categories, without labeled data specific to the target language. Specifically, given the target sentence $X = \{x_1, \dots, x_N\}$, each token is associated with a corresponding label $t = t_1, \dots, t_N$ (Tags follow the BOI schema, such as B-LOC, I-LOC, O), turning this task into a sequence labeling task (Wu et al. 2020b; Yang et al. 2022a). Building on the insights from prior research (Tang et al. 2022; Ye et al. 2022; Li et al. 2022b; Shang, Huang, and Mao 2022; Zhu and Li 2022; Li et al. 2022a; Mo et al. 2023), our approach tackles the CrossNER problem by establishing relationships between all tokens within multilingual sentences. This strategy serves to simultaneously bridge the semantic and syntactic gaps across different languages, thus enabling more effective cross-lingual NER.

We reformulate the CrossNER task as a token pair relation classification problem. Given the source sentence $X = \{x_1, \dots, x_N\}$ of N tokens, the relations between token pairs in X are extracted into a relation set/matrix R , which can be categorized into intra-entity R_{in} and extra-entity R_{extra} relations. For an entity $e = (x_s, \dots, x_e)$ spanning from the s -th token to the e -th token, $R(x_s, x_e)$ denotes the start-end

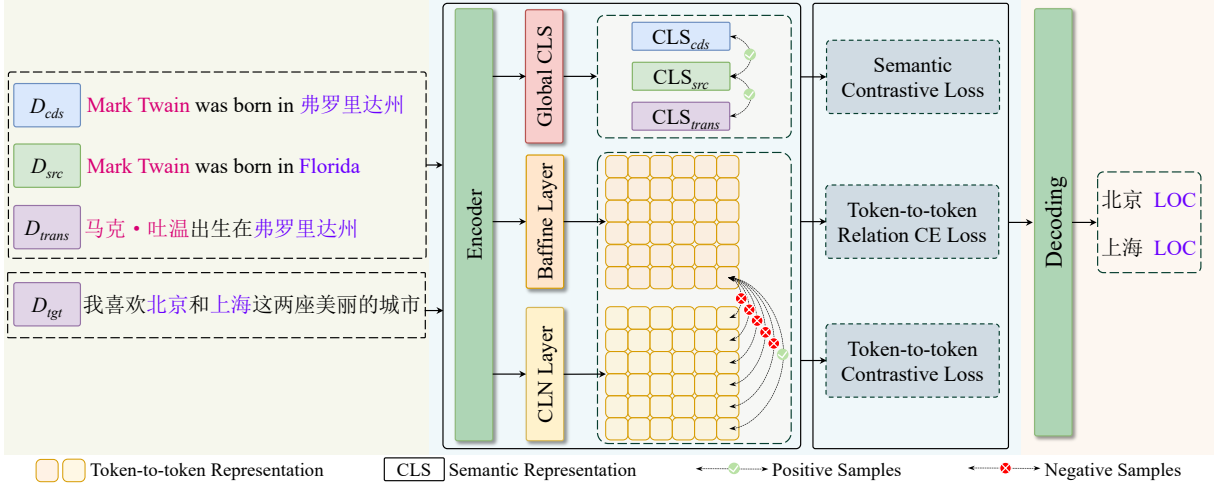


Figure 2: Overview of the proposed mCL-NER. It transforms the sequence labeling problem to the classification of the token pair relation with semantic contrastive learning and token-to-token relation contrastive learning.

relation. Additionally, $R(e_i, e_{i+1})$ ($s \leq i < e$) signifies the neighbor relation between consecutive tokens within the entity. Here, $R(\cdot, \cdot)$ denotes the relation between two tokens. The intra-entity R_{in} encompasses both start-end relations ($R(x_s, x_e)$) and neighbor relations ($R(x_i, x_{i+1})$) within an entity. Moreover, relations such as $R(e_i, e_j)$ ($s \leq i < e \wedge j \notin [s, e]$) represent the extra-entity relations.

$$P(\mathcal{R}|X) = \prod_{1 < i, j < N \wedge i \neq j} P(R(x_i, x_j) | x_i, x_j; \Theta_{ner}) \quad (1)$$

where the relation set \mathcal{R} contains relations between the token x_i and x_j in the sentence. Θ_{ner} is the NER model parameter.

In the cross-lingual setting, we amalgamate both source and target languages within a single sentence. Subsequently, we employ contrastive learning techniques to narrow the disparity between the source relation $R(x_i, x_j)$ and the corresponding aligned relations $R(y_a, y_b)$ in the target language. Simultaneously, we refine the representations of $x_{i,j}$ and $y_{a,b}$ through semantic contrastive learning. This multi-view contrastive learning process can be conceptualized as:

$$\min(|R(x_i, x_j) - R(y_a, y_b)| + |F(x_i, x_j) - F(y_a, y_b)|) \quad (2)$$

where $|R(x_i, x_j) - R(y_a, y_b)|$ denote the relation distance and $|F(x_i, x_j) - F(y_a, y_b)|$ is the representation distance. $R(\cdot)$ extracts relation representation and $F(\cdot)$ obtains the semantic representation.

Cross-Lingual Multi-view Contrastive Learning

As shown in the overall architecture of Figure 2, we resolve the cross-lingual NER task by distinguishing the relation between tokens and tokens for better modeling of intra-entity and extra-entity. We design a semantic contrastive objective to cluster the representations with the same meaning in different languages and then introduce token-to-token relation contrastive learning for gathering similar token-to-token relations within entities across different languages. Following

the previous work (Wu et al. 2020b; Yang et al. 2022a), the initial model is jointly trained on labeled source language data and its code-switch counterparts (where the translated data is a special case of code-switched data). Further, the target raw data is annotated by the previous model to construct the synthetic corpora for the subsequent iteration of optimizing via self-training.

Semantic Contrastive Learning

For zero-resource CrossNER without target annotated data, effectively representing the same entity with the cross-lingual pre-trained model in different languages becomes crucial to enhance the performance of the CrossNER model. To address this challenge, we enrich the original source sentences by randomly substituting source phrases with target translations to create code-switched sentences comprised of both source and target tokens. This approach provides assistance for our proposed model to align the cross-lingual context of different languages and improve its ability to recognize named entities. Contrastive learning can be naturally and effectively used to reduce the gap among different languages under the cross-lingual scenario with code-switched data. Based on the code-switched corpora, we propose source-codeswitch contrastive learning to minimize the distance between positive samples and maximize the distance between negative samples. Constructing a code-switched sentence for CrossNER involves mixing entity names in both source and target languages to ensure that the code-switched sentence is fluent and accurate in terms of contextual information. For example, a code-switched sentence comprised of English and Chinese entity names can be “**Wo yao qu Bei Jing kan (I want to go to Beijing) the Great Wall (I want to go to Beijing to see the Great Wall)**”. In this sentence, “Bei Jing (Beijing)” is the Chinese location entity while “the Great Wall” is the English translation of the Chinese token “Chang Cheng”. The code-switched sentence containing both source and target entities facilitates entity recognition in CrossNER.

Given a source sentence $X = \{x_1, \dots, x_N\}$ and its code-switched sentence $Y = \{y_1, \dots, y_M\}$ (the target translation is a special case of the code-switched sentence), we consider them a positive sample pair within a batch and the source sentence is paired with other negative samples. The model is encouraged to learn language-agnostic representations for capturing the same meaning across different languages. The cosine similarity is used to measure the distance between sentence pairs. We apply an extra contextual fusion layer to refine the features from the cross-lingual pre-trained model to enhance the contextual information within the sequence. Supposing $H = \{h_1, \dots, h_N\}$ denote the final features of each token in source sentence X and $H' = \{h'_1, \dots, h'_N\}$ denote the features of translated data, we obtain the sentence semantic representation $H_{avg} = f(H)$ and $H'_{avg} = f(H')$, where $f(\cdot)$ denotes the average or pooling operation.

We introduce the contrastive learning between source and code-switched sentences, which helps model to reduce the gap of different language entities with the same meaning as:

$$\mathcal{L}_{sc} = -\log \frac{e^{\text{sim}(H_{avg}, H'_{avg})/\tau}}{\sum_{j=1}^B e^{\text{sim}(H_{avg}, H'_{avg})/\tau}} \quad (3)$$

where $\text{sim}(\cdot)$ means the cosine similarity, B is the number of source sentences in a batch, and τ is the temperature. H_{avg} and H'_{avg} are corresponding representations of the source and code-switched sentence.

Token-to-Token Relation Contrastive Learning

We introduce the token-to-token relation of our cross-lingual NER method and then leverage token-pair relation contrastive learning to cluster the representation of different languages. Token-pair relation is the important feature in the NER task (Shang, Huang, and Mao 2022; Zhu and Li 2022), which can be used to enhance the cross-lingual entity recognition by constraining the cross-lingual token-pair relation.

Given a sentence $X = \{x_1, \dots, x_N\}$, we obtain the hidden representation $H = \{h_1, \dots, h_N\}$ after the encoder layers. The relation representation $\{r_{ij} | (i, j) \in [1, N]\} \in \mathbb{R}^{N \times N \times d_h}$ of token pair (x_i, x_j) is then obtained through the biaffine token-to-token relation layer.

$$\begin{aligned} s_i &= \text{MLP}(h_i), e_j = \text{MLP}(h_j) \\ r_{ij} &= s_i^\top W_1 e_j + W_2 (s_i \oplus e_j) + b \end{aligned} \quad (4)$$

where W_1 , W_2 and b denote the trainable parameters, \oplus means concatenation. MLP is the fully connected layer.

The source or code-switched tokens are fed into the conditional layer normalization module to obtain the relation representation r'_{ij} as:

$$\begin{aligned} r'_{ij} &= \gamma_i \odot \left(\frac{h_{ij} - \mu}{\sigma} \right) + \lambda_i \\ \gamma_i &= W_\alpha h_i + b_\alpha, \lambda_i = W_\beta h_i + b_\beta \\ \mu &= \frac{1}{d_h} \sum_{k=1}^{d_h} h_{jk}, \sigma = \sqrt{\frac{1}{d_h} \sum_{k=1}^{d_h} (h_{jk} - \mu)^2} \end{aligned} \quad (5)$$

where γ_i and λ_i are obtained from h_i . μ and σ are the mean and standard deviation taken across the elements of h_j

respectively. h_{jk} means the k -th dimension of h_j . d_h is the hidden size of h_i . W_α , W_β , b_α , b_β are learned parameters. Finally, we take the token-to-token relation representations r_{ij} and r'_{ij} into an MLP layer to obtain their projection representations z_{ij} and z'_{ij} . The contrastive learning objective of token-to-token relation is defined as:

$$\mathcal{L}_{tc} = -\log \frac{e^{\text{sim}(z_{ij}, z'_{ij})/\tau}}{\sum_{j=1}^{2N} e^{\text{sim}(z_{ij}, z'_{ij})/\tau}} \quad (6)$$

where contrastive learning is applicable to both source and target language data, $\text{sim}(z_{ij}, z'_{ij})$ is the cosine similarity of the source sentence and its counterparts (e.g. code-switched sentence or source sentence).

Self-Training

Due to the scarcity of annotated target language data, we use the self-training to get pseudo labels for target data. First, our cross-lingual NER model Θ_{ner}^{src} is trained on the source labeled dataset $D_{x,r}^{src} = \{(x_i, x_j), r_{ij}\}_{i=1, j=1}^N$:

$$\mathcal{L}_{ce} = -\frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \sum_{r=1}^R r_{ij} \log(P(r_{ij} | (x_i, x_j))) \quad (7)$$

where r_{ij} is the gold relation label of the token pair (x_i, x_j) , $P(r_{ij} | (x_i, x_j))$ is the relation prediction probability of the token pair (x_i, x_j) . The final loss of the model can be composed of multiple contrastive learning losses and the token-to-token relation class loss. The final loss of the cross NER model is accumulated as:

$$\mathcal{L}_{\Theta_{ner}^{src}} = \mathcal{L}_{ce}^{src} + w(\mathcal{L}_{sc}^{src} + \mathcal{L}_{tc}^{src}) \quad (8)$$

Then, the model Θ_{ner}^{src} generates pseudo labels \hat{r}_{ij} of token pairs for unlabeled target language dataset $D_x^{tgt} = \{(x_i)\}_{i=1}^M$. The model is trained on the source and target pseudo dataset. Since pseudo-labels of the target language data may bring extra noise, we only retain the contrastive learning of token pair relation and minimize the mean squared error (Ren et al. 2022) of the source and target model on the prediction distribution. So the training objective of the model Θ_{ner}^{tgt} is defined as:

$$\mathcal{L}_{tgt} = \mathcal{L}_{ce}^{tgt} + w_1 \mathcal{L}_{tc}^{tgt} + w_2 \mathcal{L}_{mse} \quad (9)$$

$$\mathcal{L}_{mse} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N |P(r_{ij} | \Theta_{ner}^{src}), P(r_{ij} | \Theta_{ner}^{tgt})| \quad (10)$$

where the loss \mathcal{L}_{ce}^{tgt} of the target model is similar to Equation 7, \mathcal{L}_{tc}^{tgt} likes Equation 6, $P(r_{ij} | \Theta)$ is the predicted probability of the token pair relation under the model Θ , w_1 and w_2 are the hyperparameters. $|\cdot|$ is the MSE distance.

Experiments

Datasets

XTREME-40 The proposed method is evaluated on the CrossNER dataset from the XTREME benchmark (Hu et al. 2020). Following the previous work (Hu et al. 2020), we use

mBERT																				
Method	af	ar	bg	bn	de	el	es	et	eu	fa	fi	fr	he	hi	hu	id	it	ja	jv	ka
mBERT (Devlin et al. 2018)	76.9	44.5	77.1	68.8	78.8	71.6	74.0	76.3	68.0	48.2	77.2	79.7	56.5	66.9	76.0	46.3	<u>81.1</u>	28.9	<u>66.4</u>	67.7
+Translate Train (Yang et al. 2022a)	74.5	37.6	77.8	73.2	77.2	74.9	69.4	74.1	63.2	43.1	75.9	76.1	55.4	68.1	77.2	<u>48.2</u>	<u>77.2</u>	36.6	55.1	64.4
UniTrans (Wu et al. 2020b)	78.2	47.0	79.5	74.6	79.8	75.6	75.2	76.5	67.2	49.3	75.6	80.1	58.4	72.1	77.9	44.6	78.3	37.6	56.2	69.9
CROP (Yang et al. 2022a)	81.0	48.0	80.8	74.9	80.3	78.7	84.2	<u>78.3</u>	<u>70.6</u>	63.2	<u>79.1</u>	83.5	64.7	<u>77.1</u>	82.5	46.4	79.9	<u>45.3</u>	57.7	74.1
Ours	82.2	50.3	81.1	79.0	82.3	78.9	<u>83.2</u>	79.3	<u>72.2</u>	<u>54.9</u>	79.4	<u>82.3</u>	<u>61.9</u>	78.9	<u>78.1</u>	63.3	82.3	57.2	72.2	70.4
Method	kk	ko	ml	mr	ms	my	nl	pt	ru	sw	ta	te	th	tl	tr	ur	vi	yo	zh	Avg _{all}
mBERT (Devlin et al. 2018)	50.4	60.2	53.7	56.2	61.9	47.6	82.1	79.6	65.2	72.8	50.8	46.8	0.4	71.2	75.5	36.9	69.7	51.7	44.1	61.7
+Translate Train (Yang et al. 2022a)	48.2	61.2	61.0	58.7	67.5	57.3	79.6	78.4	61.2	<u>69.2</u>	62.7	51.2	2.4	72.7	72.6	58.9	69.5	51.1	45.3	62.3
UniTrans (Wu et al. 2020b)	52.5	61.4	<u>63.5</u>	62.3	65.8	59.2	82.4	80.3	64.8	<u>65.2</u>	<u>63.2</u>	56.1	3.1	73.4	<u>77.9</u>	64.1	69.7	50.1	47.4	64.5
CROP (Yang et al. 2022a)	<u>54.9</u>	<u>62.6</u>	72.7	70.6	<u>71.1</u>	<u>61.3</u>	<u>84.6</u>	<u>81.7</u>	<u>69.7</u>	68.3	64.9	61.6	<u>3.9</u>	<u>76.9</u>	80.4	78.0	<u>70.0</u>	<u>51.8</u>	<u>54.4</u>	<u>68.5</u>
Ours	56.0	63.4	62.4	<u>67.9</u>	76.2	64.3	85.6	83.3	71.1	75.3	57.9	60.4	11.2	83.5	74.3	<u>72.4</u>	81.4	65.4	63.5	70.4

Table 1: Experimental results on XTREME-40 initialized by pretrained cross lingual language model mBERT_{base}.

XLM-R																				
Method	af	ar	bg	bn	de	el	es	et	eu	fa	fi	fr	he	hi	hu	id	it	ja	jv	ka
XLM-R (Conneau et al. 2019)	74.6	46.0	78.0	68.3	75.2	75.7	70.2	72.2	59.9	52.0	75.8	76.6	52.4	69.6	78.2	47.4	77.7	21.0	61.8	66.5
+Translate Train (Yang et al. 2022a)	76.2	47.8	79.2	74.3	75.8	67.7	68.4	75.8	61.2	41.0	76.8	76.4	55.0	71.9	76.0	50.6	78.1	35.4	54.7	68.4
UniTrans (Wu et al. 2020b)	78.1	<u>48.1</u>	79.3	74.6	75.2	74.9	73.8	76.9	62.7	49.2	74.6	76.5	53.4	70.4	76.9	48.6	77.3	21.6	<u>62.2</u>	66.8
CROP (Yang et al. 2022a)	80.3	45.2	<u>80.4</u>	<u>75.7</u>	<u>79.6</u>	<u>78.5</u>	83.1	<u>77.2</u>	<u>66.8</u>	65.5	<u>77.9</u>	<u>82.9</u>	<u>63.5</u>	<u>77.4</u>	81.6	46.1	78.8	<u>45.4</u>	63.2	74.0
Ours	<u>79.7</u>	57.0	81.5	79.5	80.2	<u>79.1</u>	<u>79.9</u>	<u>77.7</u>	<u>67.1</u>	<u>55.8</u>	78.1	83.0	64.8	78.1	<u>78.5</u>	51.9	79.9	52.9	61.5	<u>71.4</u>
Method	kk	ko	ml	mr	ms	my	nl	pt	ru	sw	ta	te	th	tl	tr	ur	vi	yo	zh	Avg _{all}
XLM-R (Conneau et al. 2019)	43.2	49.9	62.3	59.6	67.3	53.5	80.2	78.1	64.3	<u>70.3</u>	55.0	50.1	3.0	69.4	<u>78.1</u>	63.6	68.2	47.5	27.7	61.3
+Translate Train (Yang et al. 2022a)	40.1	55.5	60.0	59.8	69.8	61.6	79.6	76.4	60.9	70.0	63.7	50.7	3.4	74.7	72.3	62.7	69.6	46.8	41.2	62.3
UniTrans (Wu et al. 2020b)	46.5	<u>57.2</u>	65.5	64.5	70.2	62.6	81.8	79.4	<u>68.8</u>	68.9	<u>65.1</u>	56.1	4.8	74.8	76.4	71.0	69.8	55.1	44.4	64.2
CROP (Yang et al. 2022a)	<u>50.2</u>	59.8	73.8	71.6	<u>71.8</u>	69.0	83.5	82.3	70.2	69.0	65.6	<u>59.9</u>	3.1	<u>75.5</u>	80.5	80.4	<u>70.1</u>	52.6	<u>50.3</u>	68.2
Ours	53.2	57.0	<u>69.1</u>	<u>67.1</u>	74.4	<u>66.4</u>	<u>83.0</u>	<u>81.2</u>	65.1	71.8	61.8	61.2	19.2	79.3	76.2	80.5	76.7	77.2	60.3	69.7

Table 2: Experimental results on XTREME-40 initialized by pretrained cross lingual language model XLM-R_{base}.

the same split for the train, validation, and test sets, where named entities in Wikipedia are annotated with the LOC, PER, and ORG tags in BOI-2 format. All NER models are trained on the English training data as the source language and then evaluated on other languages.

CoNLL We also conduct experiments on CoNLL-02 and CoNLL-03 datasets (Sang 2002; Sang and Meulder 2003) covering four languages: Spanish (es), Dutch (nl), English (en), and German (de). The cross-lingual NER dataset has 4 entity types, including LOC, ORG, MISC, and PER. We split them into the train, validation, and test sets (Wu et al. 2020b), where English is used as the source language and others are target languages.

Implementation Details and Evaluation

For a fair comparison, we adopt the same structure and model size, which all have 12 layers with an embedding dimension of 768 under the base architecture of both mBERT (Devlin et al. 2018) and XLM-R (Conneau et al. 2019). We set the batch size as 32 samples for XTREME-40 and CoNLL. We use AdamW (Loshchilov and Hutter 2019) for optimization with a learning rate of $1e^{-5}$ for the pre-trained model and $1e^{-3}$ for other extra components. The dimension of the projected representations for contrastive learning is set to 128. More training details can be found in the Appendix. We use average entity-level valid F1 scores of all languages to choose the best checkpoint and report the F1 scores on all test sets. We compare our method with the different strong baselines UniTrans (Wu et al. 2020b), CROP (Yang et al. 2022a), Translate-Train (Yang et al. 2022a), and TSL (Wu et al. 2020a), which are initialized by the cross-lingual pre-trained model by mBERT (Devlin et al. 2018) and XLM-R (Conneau et al. 2019) for model-based transfer. We set a threshold and remove samples that fall below the

threshold to mitigate the noise introduced by pseudo-labeled data. Additionally, we eliminate data that only contains the ‘‘O’’ label and use the continuity of the internal relationship of the entity to remove some discontinuous entity data.

Main Results

XTREME-40 We present the results on the XTREME-40 dataset in Table 1 and 2 by different cross-lingual pre-trained language models including mBERT and XLM-R. Overall, the average F1 score of our method outperforms the previous baselines by a large margin. Compared with the methods (Wu et al. 2020b; Yang et al. 2022a) initialized by mBERT, our work achieves a significant improvement of +1.9 F1 points. For languages id (Indonesian), zh (Chinese), and ja (Japanese) distant from English, our method can further gain +10 points improvement than CROP (Yang et al. 2022a). It is due to the effectiveness of the shared representations across the different languages by multi-view contrastive learning. For the methods initialized by the XLM-R, our model also gets a consistent promotion by +1.5 points compared to the baseline (Yang et al. 2022a) in the average F1 score metric.

CoNLL Table 3 shows the experimental results on CoNLL dataset. To make a fair assessment, we use the mBERT base model. Compared with TSL (Wu et al. 2020a), our approach achieves an average F1 score improvement of +2.3. An averaged +1.3 F1 improvement is gained compared to UniTrans (Wu et al. 2020b) and CROP (Yang et al. 2022a), which means the effectiveness of our proposed multi-view contrastive learning comprised of semantic and token pair relation contrastive. In contrast to the previous models, our method demonstrates favorable enhancements by reducing the gap among different languages in the shared space and clustering cross-lingual entities with the same meaning.

Method	de	es	nl	Avg
X-Lingual Clusters (Täckström et al. 2012)	40.4	59.3	58.4	52.7
+Wikifier (Tsai et al. 2016)	48.1	60.6	61.6	56.8
Inverted Softmax (Smith et al. 2017)	58.5	65.1	65.4	63.0
Cheap Translation (Mayhew et al. 2017)	57.2	64.1	63.4	61.6
BWET (Xie et al. 2018)	57.8	72.4	71.3	67.2
TMP (Jain et al. 2019)	65.2	75.9	74.6	71.9
mBERT (Devlin et al. 2018)	75.0	74.6	77.9	75.8
BERT-f (Wu and Dredze 2019)	71.1	74.5	79.5	75.0
XLNet (Conneau et al. 2019)	73.4	77.4	78.9	76.6
Cross-Augmented (Bari et al. 2020b)	61.5	73.5	69.9	68.3
Meta-learning-based (Wu et al. 2020d)	73.2	76.8	80.4	76.8
TSL (Wu et al. 2020a)	75.3	78.0	81.3	78.2
UniTrans (Wu et al. 2020b)	74.8	79.3	82.9	79.0
MulDA (Liu et al. 2021b)	78.2	77.5	78.4	78.0
+Translate Train (Yang et al. 2022a)	74.2	77.8	79.2	77.1
CROP (Yang et al. 2022a)	<u>80.1</u>	78.1	79.5	<u>79.2</u>
Ours	81.0	79.2	81.3	80.5

Table 3: Experimental results on CoNLL.

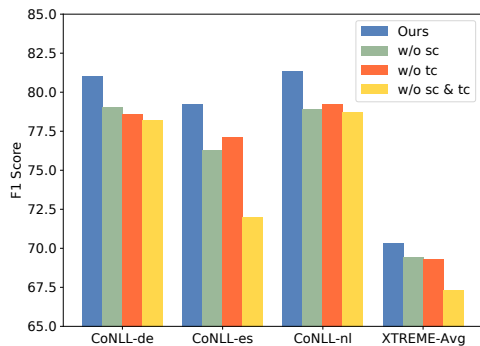


Figure 3: Ablation study of MCL-NER. XTREME-Avg denotes the average F1 scores of 39 languages in XTREME.

Analysis

Ablation Study To verify the effectiveness of MCL-NER, we introduce the following series of ablation experiments: ① Ours, which is the final model with the multi-view contrastive objectives; ② w/o sc, which adopts token-to-token relation contrastive learning; ③ w/o tc, which use the semantic contrastive objective; ④ w/o sc & tc, where semantic and token pair relation contrastive learning are removed. From the ablation experiments in Figure 3, our method outperforms ②, ③, and ④. ② w/o sc, which shows aligned semantic representations in different languages play a pivotal role in CrossNER by enhancing shared semantic representations. ③ w/o tc gets worse performance compared to our method ①, indicating that the contrastive learning of relation between tokens clusters the representations of the similar entities. ④ w/o sc & tc has the worst performance, which verifies that a multi-view of semantic and relation contrastive learning between tokens is the best strategy.

Representation on Token-to-Token Relationships To intuitively understand the effectiveness of token-to-token relation representation contrastive, we take a close look at the token-to-token relations within a sentence. We visualize the distance between the generated representations of the token-

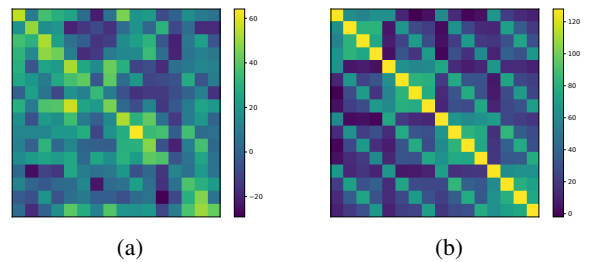


Figure 4: (a) denotes the representation on token-to-token relation without contrastive learning and (b) is the representation on token-to-token relation with contrastive learning.

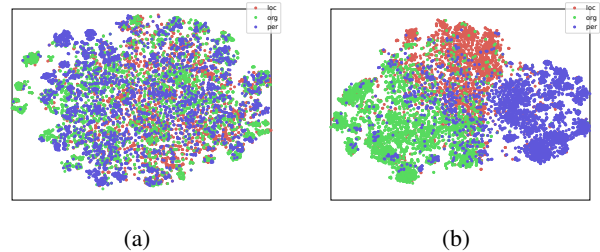


Figure 5: t-SNE visualization of different pre-defined categories (e.g. LOC) in Chinese. (a) and (b) indicate the token-to-token relation representations between entities w and w/o contrastive learning respectively.

to-token relationships without contrastive learning and the counterparts with our method. The distance matrix (cosine similarity) is shown in Figure 4. It can be seen that the similarity between the token-token relationship features of the baseline displays a random pattern. However, two slashes of Figure 4(b) with higher brightness appear in the similarity matrix in a fixed pattern. The positive sample pair has a higher score with brightness, achieving the desired effect.

Distribution of Multilingual Corpora We conduct a t-SNE visualization (van der Maaten and Hinton 2008) of token-to-token relation representations. Figure 5(a) denotes that the model w/o multi-view contrastive learning generates token pair relation representations between different classes. Figure 5(b) shows that MCL-NER w/ multi-view contrastive objectives produce more distinct and distinguishable representations, where entities can be separated into independent regions for better performance of CrossNER.

Effect of Training Data Size We consider the impact from two aspects: the amount of source labeled data and the pseudo-labels target data. The size ratio of the source language data is randomly sampled for analytic experiments according to a certain proportion from 10%, 20%, ..., 100% in Figure 6(a). In Figure (b), we randomly sample the target data with pseudo-labels and put the train data size to $\{1K, 2K, \dots, ALL\}$ sentences to train the model. From Figure 6(a), surprisingly, training with only 10% of the source data brings huge improvements. In Figure 6(b), the overall model performance continues to be improved with increasing pseudo-label target data since the disturbance

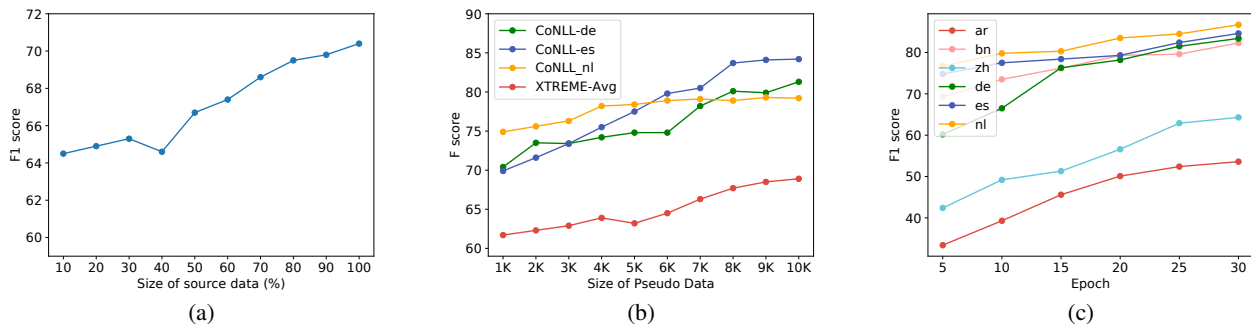


Figure 6: (a) Evaluation results on the validation sets with different source data training sizes. (b) Evaluation results on the validation sets with different pseudo-label target data training sizes. (c) Pseudo label quality for unlabeled target data.

Language	Sentences
German	1981 war Ronald Reagan _{PER} als Präsident der Vereinigten Staaten _{ORG} vereidigt worden .
	TSL: Ronald Reagan _{PER} Präsident der Vereinigten Staaten _{PER}
	UniTrans: Ronald Reagan _{PER} Vereinigten Staaten _{LOC}
	Ours: Ronald Reagan _{PER} Präsident der Vereinigten Staaten _{ORG}
Dutch	Hij stond met Hull City AFC _{ORG} in de finale van de strijd om de FA Cup 2014 _{ORG} , die de ploeg van trainer-coach Steve Bruce _{PER} met 3-2 verloor van Arsenal _{ORG} .
	TSL: Hull City AFC _{ORG} Steve Bruce _{PER} Arsenal _{ORG}
	UniTrans: Hull City AFC _{ORG} FA Cup 2014 _{ORG} trainer-coach Steve Bruce _{PER} Arsenal _{ORG}
	Ours: Hull City AFC _{ORG} FA Cup 2014 _{ORG} Steve Bruce _{PER} Arsenal _{ORG}
Chinese	世界上有多个地方以“剑桥”为地名，分别分布在 英国 _{LOC} 、 美国 _{LOC} 、 加拿大 _{LOC} 、 澳大利亚 _{LOC} 、 新西兰 _{LOC} 等国；此外，“剑桥”还是 剑桥大学 _{ORG} 的简称。
	TSL: 剑桥 _{LOC} 英国 _{LOC} 美国 _{LOC} 加拿大 _{LOC} 澳大利亚 _{LOC} 新西兰 _{LOC} 剑桥大学 _{ORG}
	UniTrans: 剑桥 _{LOC} 英国 _{LOC} 美国 _{LOC} 加拿大 _{LOC} 澳大利亚 _{LOC} 新西兰 _{LOC} 剑桥大学 _{ORG}
	Ours: 英国 _{LOC} 美国 _{LOC} 加拿大 _{LOC} 澳大利亚 _{LOC} 新西兰 _{LOC} 剑桥大学 _{ORG}

Figure 7: Case study on CrossNER. Texts with colorful backgrounds are gold entities and the red fonts are incorrect predictions.

of pseudo-label data noise may cause slight fluctuations in training. When the pseudo-label target data size grows to a certain extent, the improvement of the CrossNER model becomes smaller. The reason is that the model can not learn more useful information from sufficient data unless given new helpful knowledge for CrossNER.

Pseudo Label Quality The quality of pseudo-labels is crucial for the success of self-training cross-lingual NER. To evaluate the effectiveness of our method in improving pseudo-label quality, we use the gold labels of the unlabeled target language data as a reference to measure the F1 score after each epoch. We select representative languages from different language families, including ar, bn, zh, de, es, and nl from XTREME-40. In Figure 6(c), Our experiments show that the quality of pseudo labels about target data is improved with the increasing training epoch, proving the effectiveness of our multi-view contrastive learning cross-lingual NER method. Especially a significant improvement in the F1 score of data zh can be observed, raising nearly 20 points.

Case Study We select three concrete cases of target languages that are similar and distant from the source language in Figure 7. In the German language, TSL and UniTrans have similar mistakes in recognizing the ORG entity “Präsident der Vereinigten Staaten”. TSL incorrectly labeled

the entity type, while UniTrans mislabeled the entity type and wrongly chose the entity span. In the Dutch language, UniTrans has a similar mistake in the span of the PER entity (easily judged as incorrect) and TSL misses an ORG entity. Our method gains the correct prediction by learning the relation of token pairs and achieving more explicit recognition with token-to-token relation contrastive learning. For the Chinese sentence example, TSL and UniTrans recognize the first “Jian Qiao (Cambridge)” as a LOC entity, while our model gets the correct recognition. It can be attributed that our method gets better distinguished semantic meaning than baselines.

Conclusion

In this work, we propose MCL-NER, a multi-view contrastive learning framework for the cross-lingual NER comprising semantic contrastive learning and token-to-token relation contrastive learning. Specifically, we construct the code-switched data by randomly substituting some phrases with the target counterparts for the semantic contrastive learning of the source and the corresponding code-switched sentence. Further, token-to-token contrastive learning enhances the syntactic representation of entities in different languages. Both contrastive learning objectives minimize the semantic gap across different languages even for distant

languages and improve the cross-lingual recognition performance with the more distinct entity representations. Extensive experimental results demonstrate that our approach significantly performs better than strong baselines by a large margin on the XTREME-40 and CoNLL benchmarks.

References

- Ashok, D.; and Lipton, Z. C. 2023. PromptNER: Prompting For Named Entity Recognition. *CoRR*, abs/2305.15444.
- Bari, M. S.; Joty, S. R.; and Jwalapuram, P. 2020. Zero-Resource Cross-Lingual Named Entity Recognition. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, 7415–7423. AAAI Press.
- Chen, Q.; Li, F.-L.; Xu, G.; Yan, M.; Zhang, J.; and Zhang, Y. 2022. Dictbert: Dictionary description knowledge enhanced language model pre-training via contrastive learning. *arXiv preprint arXiv:2208.00635*.
- Chen, T.; Kornblith, S.; Norouzi, M.; and Hinton, G. E. 2020. A Simple Framework for Contrastive Learning of Visual Representations. *CoRR*, abs/2002.05709.
- Chen, W.; Jiang, H.; Wu, Q.; Karlsson, B. F.; and Guan, Y. 2021. AdvPicker: Effectively Leveraging Unlabeled Data via Adversarial Discriminator for Cross-Lingual NER. *CoRR*, abs/2106.02300.
- Chi, Z.; Dong, L.; Wei, F.; Wang, W.; Mao, X.; and Huang, H. 2020. Cross-Lingual Natural Language Generation via Pre-Training. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, 7570–7577. AAAI Press.
- Chuang, C.-Y.; Robinson, J.; Lin, Y.-C.; Torralba, A.; and Jegelka, S. 2020. Debaised contrastive learning. *Advances in neural information processing systems*, 33: 8765–8775.
- Conneau, A.; Khandelwal, K.; Goyal, N.; Chaudhary, V.; Wenzek, G.; Guzmán, F.; Grave, E.; Ott, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. Unsupervised Cross-lingual Representation Learning at Scale. *CoRR*, abs/1911.02116.
- Das, S. S. S.; Katiyar, A.; Passonneau, R. J.; and Zhang, R. 2021. Container: Few-shot named entity recognition via contrastive learning. *arXiv preprint arXiv:2109.07589*.
- Devlin, J.; Chang, M.; Lee, K.; and Toutanova, K. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *CoRR*, abs/1810.04805.
- Gao, T.; Yao, X.; and Chen, D. 2021. Simcse: Simple contrastive learning of sentence embeddings. *arXiv preprint arXiv:2104.08821*.
- Giorgi, J.; Nitski, O.; Wang, B.; and Bader, G. 2020. Declutr: Deep contrastive learning for unsupervised textual representations. *arXiv preprint arXiv:2006.03659*.
- Han, X.; Luo, Y.; Chen, W.; Liu, Z.; Sun, M.; Botong, Z.; Fei, H.; and Zheng, S. 2022. Cross-Lingual Contrastive Learning for Fine-Grained Entity Typing for Low-Resource Languages. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 2241–2250. Dublin, Ireland: Association for Computational Linguistics.
- He, K.; Fan, H.; Wu, Y.; Xie, S.; and Girshick, R. 2020. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 9729–9738.
- Hou, Y.; Che, W.; Lai, Y.; Zhou, Z.; Liu, Y.; Liu, H.; and Liu, T. 2020. Few-shot Slot Tagging with Collapsed Dependency Transfer and Label-enhanced Task-adaptive Projection Network. arXiv:2006.05702.
- Hu, J.; Ruder, S.; Siddhant, A.; Neubig, G.; Firat, O.; and Johnson, M. 2020. XTREME: A Massively Multilingual Multi-task Benchmark for Evaluating Cross-lingual Generalization. arXiv:2003.11080.
- Keung, P.; Lu, Y.; and Bhardwaj, V. 2019. Adversarial Learning with Contextual Embeddings for Zero-resource Cross-lingual Classification and NER. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 1355–1360. Association for Computational Linguistics.
- Khosla, P.; Teterwak, P.; Wang, C.; Sarna, A.; Tian, Y.; Isola, P.; Maschinot, A.; Liu, C.; and Krishnan, D. 2020. Supervised contrastive learning. *Advances in neural information processing systems*, 33: 18661–18673.
- Li, J.; Fei, H.; Liu, J.; Wu, S.; Zhang, M.; Teng, C.; Ji, D.; and Li, F. 2022a. Unified Named Entity Recognition as Word-Word Relation Classification. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, 10965–10973. AAAI Press.
- Li, X.; Bing, L.; Zhang, W.; Li, Z.; and Lam, W. 2021. Unsupervised Cross-lingual Adaptation for Sequence Tagging and Beyond. arXiv:2010.12405.
- Li, Z.; Fu, L.; Wang, X.; Zhang, H.; and Zhou, C. 2022b. RF-BFN: A Relation-First Blank Filling Network for Joint Relational Triple Extraction. In Louvan, S.; Madotto, A.; and Madureira, B., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: Student Research Workshop, ACL 2022, Dublin, Ireland, May 22-27, 2022*, 10–20. Association for Computational Linguistics.
- Liu, L.; Ding, B.; Bing, L.; Joty, S.; Si, L.; and Miao, C. 2021a. 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), 5834–5846. Online: Association for Computational Linguistics.
- Liu, L.; Ding, B.; Bing, L.; Joty, S. R.; Si, L.; and Miao, C. 2021b. MulDA: A Multilingual Data Augmentation Framework for Low-Resource Cross-Lingual NER. In Zong, C.; Xia, F.; Li, W.; and Navigli, R., eds., *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1-6, 2021*, 5834–5846. Association for Computational Linguistics.
- Loshchilov, I.; and Hutter, F. 2019. Decoupled Weight Decay Regularization. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Mayhew, S.; Tsai, C.-T.; and Roth, D. 2017. Cheap translation for cross-lingual named entity recognition. In *Proceedings of the 2017 conference on empirical methods in natural language processing*, 2536–2545.
- Mo, Y.; Tang, H.; Liu, J.; Wang, Q.; Xu, Z.; Wang, J.; Wu, W.; and Li, Z. 2023. Multi-Task Transformer with Relation-Attention and Type-Attention for Named Entity Recognition. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 1–5.
- Ni, J.; Dinu, G.; and Florian, R. 2017. Weakly Supervised Cross-Lingual Named Entity Recognition via Effective Annotation and Representation Projection. In Barzilay, R.; and Kan, M., eds., *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, 1470–1480. Association for Computational Linguistics.
- Ren, J.; Zhang, M.; Yu, C.; and Liu, Z. 2022. Balanced MSE for Imbalanced Visual Regression. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, 7916–7925. IEEE.
- Sang, E. F. T. K. 2002. Introduction to the CoNLL-2002 Shared Task: Language-Independent Named Entity Recognition. In Roth, D.; and van den Bosch, A., eds., *Proceedings of the 6th Conference on Natural Language Learning, CoNLL 2002, Held in cooperation with COLING 2002, Taipei, Taiwan, 2002*. ACL.
- Sang, E. F. T. K.; and Meulder, F. D. 2003. Introduction to the CoNLL-2003 Shared Task: Language-Independent Named Entity Recognition. In Daelemans, W.; and Osborne, M., eds., *Proceedings of the Seventh Conference on Natural Language Learning, CoNLL 2003, Held in cooperation with HLT-NAACL 2003, Edmonton, Canada, May 31 - June 1, 2003*, 142–147. ACL.
- Shang, Y.; Huang, H.; and Mao, X. 2022. OneRel: Joint Entity and Relation Extraction with One Module in One Step. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, 11285–11293. AAAI Press.
- Smith, S. L.; Turban, D. H. P.; Hamblin, S.; and Hammerla, N. Y. 2017. Offline bilingual word vectors, orthogonal transformations and the inverted softmax. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net.
- Tang, W.; Xu, B.; Zhao, Y.; Mao, Z.; Liu, Y.; Liao, Y.; and Xie, H. 2022. UniRel: Unified Representation and Interaction for Joint Relational Triple Extraction. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, 7087–7099. Association for Computational Linguistics.
- van der Maaten, L.; and Hinton, G. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research*, 9(86): 2579–2605.
- Wu, Q.; Lin, Z.; Karlsson, B.; Lou, J.; and Huang, B. 2020a. Single-/Multi-Source Cross-Lingual NER via Teacher-Student Learning on Unlabeled Data in Target Language. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J. R., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, 6505–6514. Association for Computational Linguistics.
- Wu, Q.; Lin, Z.; Karlsson, B. F.; Huang, B.; and Lou, J. 2020b. UniTrans : Unifying Model Transfer and Data Transfer for Cross-Lingual Named Entity Recognition with Unlabeled Data. In Bessiere, C., ed., *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020*, 3926–3932. ijcai.org.
- Wu, Q.; Lin, Z.; Wang, G.; Chen, H.; Karlsson, B. F.; Huang, B.; and Lin, C. 2020c. Enhanced Meta-Learning for Cross-Lingual Named Entity Recognition with Minimal Resources. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, 9274–9281. AAAI Press.
- Wu, Q.; Lin, Z.; Wang, G.; Chen, H.; Karlsson, B. F.; Huang, B.; and Lin, C.-Y. 2020d. Enhanced Meta-Learning for Cross-lingual Named Entity Recognition with Minimal Resources. arXiv:1911.06161.
- Wu, S.; and Dredze, M. 2019. Beto, Bentz, Becas: The Surprising Cross-Lingual Effectiveness of BERT. *CoRR*, abs/1904.09077.
- Xie, J.; Yang, Z.; Neubig, G.; Smith, N. A.; and Carbonell, J. G. 2018. Neural Cross-lingual Named Entity Recognition with Minimal Resources. In Riloff, E.; Chiang, D.; Hockenmaier, J.; and Tsujii, J., eds., *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, 369–379. Association for Computational Linguistics.
- Yang, J.; Huang, S.; Ma, S.; Yin, Y.; Dong, L.; Zhang, D.; Guo, H.; Li, Z.; and Wei, F. 2022a. CROP: Zero-shot Cross-lingual Named Entity Recognition with Multilingual La-

- beled Sequence Translation. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, 486–496. Association for Computational Linguistics.
- Yang, J.; Ma, S.; Huang, H.; Zhang, D.; Dong, L.; Huang, S.; Muzio, A.; Singhal, S.; Hassan, H.; Song, X.; and Wei, F. 2021. Multilingual Machine Translation Systems from Microsoft for WMT21 Shared Task. In Barrault, L.; Bojar, O.; Bougares, F.; Chatterjee, R.; Costa-jussà, M. R.; Federmann, C.; Fishel, M.; Fraser, A.; Freitag, M.; Graham, Y.; Grundkiewicz, R.; Guzman, P.; Haddow, B.; Huck, M.; Jimeno-Yepes, A.; Koehn, P.; Kocmi, T.; Martins, A.; Morishita, M.; and Monz, C., eds., *Proceedings of the Sixth Conference on Machine Translation, WMT@EMNLP 2021, Online Event, November 10-11, 2021*, 446–455. Association for Computational Linguistics.
- Yang, J.; Ma, S.; Zhang, D.; Wu, S.; Li, Z.; and Zhou, M. 2020. Alternating Language Modeling for Cross-Lingual Pre-Training. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, 9386–9393. AAAI Press.
- Yang, J.; Yin, Y.; Ma, S.; Zhang, D.; Li, Z.; and Wei, F. 2022b. High-resource Language-specific Training for Multilingual Neural Machine Translation. In Raedt, L. D., ed., *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, 4461–4467. ijcai.org.
- Yang, J.; Yin, Y.; Ma, S.; Zhang, D.; Wu, S.; Guo, H.; Li, Z.; and Wei, F. 2022c. UM4: Unified Multilingual Multiple Teacher-Student Model for Zero-Resource Neural Machine Translation. In Raedt, L. D., ed., *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, 4454–4460. ijcai.org.
- Ye, D.; Lin, Y.; Li, P.; and Sun, M. 2022. Packed Levitated Marker for Entity and Relation Extraction. In Muresan, S.; Nakov, P.; and Villavicencio, A., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, 4904–4917. Association for Computational Linguistics.
- You, Y.; Chen, T.; Sui, Y.; Chen, T.; Wang, Z.; and Shen, Y. 2020. Graph contrastive learning with augmentations. *Advances in neural information processing systems*, 33: 5812–5823.
- Yu, J.; Bohnet, B.; and Poesio, M. 2020. Named Entity Recognition as Dependency Parsing. arXiv:2005.07150.
- Zhou, R.; Li, X.; Bing, L.; Cambria, E.; and Miao, C. 2023. Improving self-training for cross-lingual named entity recognition with contrastive and prototype learning. *arXiv preprint arXiv:2305.13628*.
- Zhou, R.; Li, X.; Bing, L.; Cambria, E.; Si, L.; and Miao, C. 2022. ConNER: Consistency Training for Cross-lingual Named Entity Recognition. In Goldberg, Y.; Kozareva, Z.; and Zhang, Y., eds., *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022*, 8438–8449. Association for Computational Linguistics.
- Zhu, E.; and Li, J. 2022. Boundary Smoothing for Named Entity Recognition. In Muresan, S.; Nakov, P.; and Villavicencio, A., eds., *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022*, 7096–7108. Association for Computational Linguistics.