

Cyclical Contrastive Learning Based on Geodesic for Zero-shot Cross-lingual Spoken Language Understanding

Anonymous ACL submission

Abstract

Owing to the scarcity of labeled training data, Spoken language understanding (SLU) is still a challenging task in low-resource languages. Therefore, zero-shot cross-lingual SLU attracts more and more attention. Contrastive learning is widely applied to explicitly align representations of similar sentences across different languages. However, the vanilla contrastive learning method may face two problems in zero-shot cross-lingual SLU: (1) the consistency between different languages is neglected; (2) each utterance has two different kinds of SLU labels, i.e. slot and intent, the utterances with one different label are also pushed away without any discrimination, which limits the performance. In this paper, we propose Cyclical Contrastive Learning based on Geodesic (CCLG), which introduces cyclical contrastive learning to achieve the consistency between different languages and leverages geodesic to measure the similarity to construct the positive pairs and negative pairs. Experimental results demonstrate that our proposed framework achieves the new state-of-the-art performance on MultiATIS++ and MTOP datasets, and the model analysis further verifies that CCLG can effectively transfer knowledge between different languages¹.

1 Introduction

Spoken Language Understanding (SLU) holds the central position in the task-oriented dialogue systems (Tur and De Mori, 2011; Qin et al., 2019; Xing and Tsang, 2022; Song et al., 2022). The primary objective of SLU is to comprehend and extract relevant information from user utterances. This capability enables the system to discern the user’s current objective and generate appropriate responses. SLU comprises two critical sub-tasks: intent detection, which focuses on identifying users’ intentions, and slot filling, which entails extracting semantic elements from user queries.

¹Our source code and models will be released after review.

However, the effectiveness of traditional SLU models is intrinsically linked to the availability of extensive annotated data, which poses challenges in scalability. This challenge is particularly evident in the case of low-resource languages, where the lack of substantial labeled datasets exacerbates scalability issues, hindering the seamless deployment and advancement of SLU models. With the demand for language processing solutions extending across diverse linguistic landscapes, the necessity for scalable SLU models that can operate effectively in resource-constrained environments becomes increasingly critical.

To tackle these constraints, the concept of zero-shot cross-lingual SLU generalization has emerged as a central focus of interest and investigation. Recently, mBERT (Devlin et al., 2019) has demonstrated significant advancements in zero-shot cross-lingual SLU. Building upon this work, Liu et al. (2020) introduces an attention-informed mixed-language training approach for cross-lingual SLU. In addition, the exploration of multilingual code-switched settings has been extended by Qin et al. (2020a), which entails aligning a source language with target languages. GL-CL_EF (Qin et al., 2022) employs contrastive learning, leveraging bilingual dictionaries to construct multilingual views of the same utterance, then encouraging their representations to be more similar than those negative example pairs. LAJ-MCL (Liang et al., 2022) proposes to model the utterance-slot-word structure using a multi-level contrastive learning framework to facilitate explicit alignment, further enhancing performance. Although existing zero-shot cross-lingual SLU methods have made promising strides by contrastive learning, we identify two main issues:

(1) **The consistency between different languages is neglected.** Although the code-switching method has been applied to construct positive samples in contrastive learning, we find that the consistency between different languages has not been

effectively established. Specifically, the distances between the corresponding samples in different languages are inconsistent, which affects the transfer of knowledge across different languages.

(2) **The utterances with one different label are also pushed away without discrimination.** Traditional contrastive learning methods utilize code-switching to construct the positive samples and negative samples, bringing tokens with the same label and intent label closer together while pushing other the tokens away. However, this can result in a side effect where tokens with only one different label (slot or intent) can be also indiscriminately pushed away, which undoubtedly hampers the representation modeling of contrastive learning, leading to the suboptimal performance.

In this paper, we propose Cyclical Contrastive Learning based on Geodesic (CCLG) to solve these two problems. For the first problem, we introduce two consistency losses, including the cross-lingual consistency loss and the intra-language consistency loss, aiming to boost consistency between different languages. For the second problem, we abandon the conventional approach of directly employing code-switching to construct positive samples and negative samples in contrastive learning. Instead, we utilize geodesic to reconstruct positive and negative samples and employ geodesic-based similarity instead of the traditional similarity metrics, thereby facilitating the learning of representations.

We conduct experiments on MultiATIS++ (Xu et al., 2020) and MTOP (Li et al., 2021), covering nine and six different languages, respectively. The experimental results show that our framework can outperform previous cross-lingual SLU baselines. The model analysis further indicates that our method can transfer knowledge from high-resource languages to low-resource languages. In summary, our work makes three-fold contributions:

- We use cyclical contrastive learning to achieve consistency between different languages.
- We apply geodesic to construct positive and negative samples in contrastive learning, leading to improved representations of tokens.
- Experiment results show that our framework achieves the new state-of-the-art performance on MultiATIS++ and MTOP datasets.

2 Related Works

The related works are introduced from zero-shot cross-lingual SLU and contrastive learning.

2.1 Zero-shot Cross-lingual SLU

Traditional SLU usually focuses on languages with abundant resources, which limits their widespread use. This limitation has sparked growing interest in a novel approach known as zero-shot cross-lingual SLU. The essence of success in this approach lies in tapping into the linguistic insights present in languages with ample resources. By doing so, it opens up exciting possibilities for overcoming challenges posed by limited data in cross-lingual scenarios. Moreover, it extends the reach of SLU to languages that have been previously overlooked, thereby contributing to a more inclusive and adaptable framework in the field of multilingualism.

In recent years, many cross-lingual embeddings, such as mBERT (Devlin et al., 2019), have shown promising results. Liu et al. (2020) propose code-mixing to construct training sentences containing both the source and target phrases, implicitly fine-tuning mBERT. Building upon it, Qin et al. (2020a) proposes multilingual code-switching data augmentation to better align the source language with all target languages. Additionally, van der Goot et al. (2021) suggests three non-English auxiliary tasks to boost cross-lingual transfer. More recently, SoGo (Zhu et al., 2023) highlights the limitations of the conventional code-switching method and proposes a saliency-based substitution approach for extracting keywords as substitutions. In our method, we use cyclical contrastive learning based on geodesic to further transfer the knowledge from the source language to the target language.

2.2 Contrastive Learning

Contrastive learning aims to learn representations of examples via minimizing the distance between positive pairs and maximizing the distance between negative pairs (Saunshi et al., 2019; Chuang et al., 2020; Liu et al., 2022), a concept initially proposed in the field of computer vision (Chopra et al., 2005; Chen et al., 2020; Wang and Liu, 2021). In natural language processing, contrastive learning is utilized for learning the sentence embeddings (Giorgi et al., 2021; Yan et al., 2021), translation tasks (Pan et al., 2021; Ye et al., 2022), and summarization (Wang et al., 2021; Cao and Wang, 2021). Owing to its strong capability in achieving alignment across different languages, contrastive learning has also been used in zero-shot cross-lingual SLU (Liang et al., 2022; Qin et al., 2022). However, we find two main issues with directly utilizing vanilla conservative

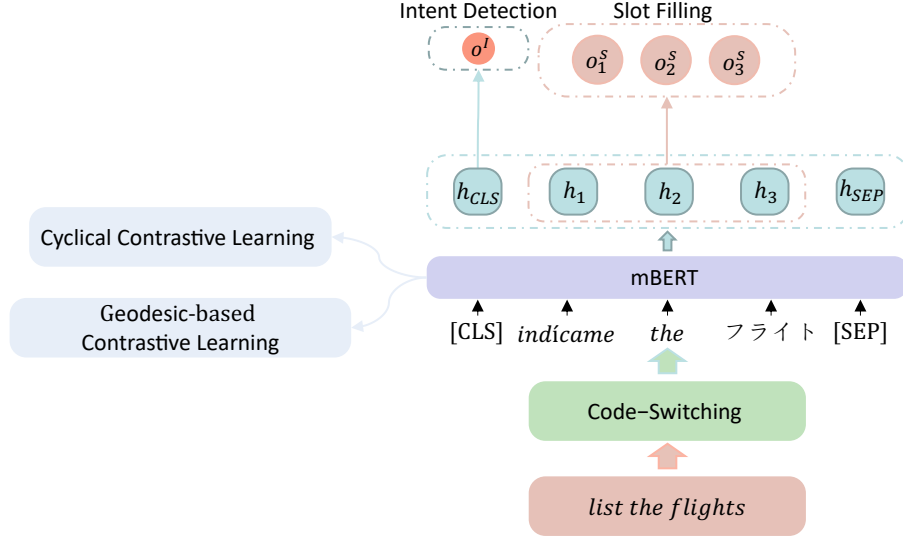


Figure 1: The overview of our approach.

learning in zero-shot cross-lingual SLU. As a result, we propose cyclical contrastive learning based on geodesic to tackle these two issues.

3 Background

SLU comprises two core subtasks, including intent detection and slot filling. Given the input utterance $x = (x_1, x_2, \dots, x_n)$, where n denotes the length of x , intent detection is treated as a classification task, producing the intent label \mathbf{o}^I , and slot filling is a sequence labeling task, mapping each utterance x to a slot output sequence $\mathbf{o}^S = (o_1^S, o_2^S, \dots, o_n^S)$. Due to the intrinsic correlation between intent detection and slot filling, it is common to train a unified SLU model capable of jointly handling both tasks, which is formulated as follows:

$$(\mathbf{o}^I, \mathbf{o}^S) = f(\mathbf{x}) \quad (1)$$

where f denotes the trained model.

Zero-shot cross-lingual SLU task involves training an SLU model on a high-resource source language, such as English, and seamlessly using it on a low-resource target language, such as French. In this scenario, when presented with an instance \mathbf{x}_{target} in the target language, the trained model f can directly generate predictions for both intent and slot values in the target language:

$$(\mathbf{o}_{target}^I, \mathbf{o}_{target}^S) = f(\mathbf{x}_{target}) \quad (2)$$

where $target$ denotes the target language.

4 Method

In this section, we first introduce the Generic SLU Module (Sec. 4.1) and the previous paradigm of

utilizing contrastive learning to enhance zero-shot cross-lingual SLU (Sec. 4.2). Then, we introduce the components of our proposed approach, including Cyclical Contrastive Learning (Sec. 4.3) and Geodesic (Sec. 4.4). At last, we introduce the final Training Objective (Sec. 4.5). The overview of our approach is demonstrated in Figure 1.

4.1 Generic SLU Module

Given the input sentence $x = (x_1, x_2, \dots, x_n)$, the construction of the input sequence is based on each input utterance by incorporating the specific tokens $\mathbf{x} = ([CLS], x_1, x_2, \dots, x_n, [SEP])$ (Devlin et al., 2019). [CLS] serves as the special symbol representing the entire sequence, and [SEP] is employed to separate non-consecutive token sequences. Following Qin et al. (2020a), code-switching is applied to leverage the bilingual dictionaries (Lample et al., 2018) in generating multi-lingual code-switched data as input for the model. The representation of the whole utterance, denoted as $\mathbf{H} = (\mathbf{h}_{CLS}, \mathbf{h}_1, \dots, \mathbf{h}_n, \mathbf{h}_{SEP})$, is obtained by utilizing the pre-trained mBERT (Devlin et al., 2019) model.

For the intent detection task, we utilize the utterance representation \mathbf{h}_{CLS} as input to a classification layer in order to derive the predicted intent:

$$\mathbf{o}^I = \text{softmax}(\mathbf{W}^I \mathbf{h}_{CLS} + \mathbf{b}^I) \quad (3)$$

where \mathbf{W}^I and \mathbf{b}^I are two trainable matrices.

For the slot filling task, we follow the methods proposed in (Wang et al., 2019; Qin et al., 2022), wherein we use the representation of the first sub-token as the whole word representation and lever-

age the hidden states to predict each slot:

$$\mathbf{o}_t^S = \text{softmax}(\mathbf{W}^s \mathbf{h}_t + \mathbf{b}^s) \quad (4)$$

where \mathbf{h}_t is the representation of the first sub-token of word x_t , \mathbf{W}^s and \mathbf{b}^s are two trainable matrices.

4.2 Previous Contrastive Paradigm

Contrastive learning has been applied in zero-shot cross-lingual SLU (Qin et al., 2022; Liang et al., 2022). In general, previous methods aim to bring tokens and the corresponding code-switched tokens (positive pairs) closer together while pushing apart tokens and the non-corresponding tokens (negative pairs). And the previous contrastive loss \mathcal{L}_{CL} can be formulated as follows:

$$\mathcal{L}_{\text{CL}}^I = - \sum_{j=1}^N \log \frac{s(\mathbf{h}_{\text{CLS}}^j, \mathbf{h}_{\text{CLS}}^{j+})}{\sum_{\mathbf{h}_{\text{CLS}}^j \neq \mathbf{h}_{\text{CLS}}^{j'}} s(\mathbf{h}_{\text{CLS}}^j, \mathbf{h}_{\text{CLS}}^{j'})} \quad (5)$$

$$\mathcal{L}_{\text{CL}}^S = - \frac{1}{n} \sum_{j=1}^N \sum_{i=1}^n \log \frac{s(\mathbf{h}_i^j, \mathbf{h}_i^{j+})}{\sum_{\mathbf{h}_i^j \neq \mathbf{h}_i^{j'}} s(\mathbf{h}_i^j, \mathbf{h}_i^{j'})} \quad (6)$$

$$\mathcal{L}_{\text{CL}} = \mathcal{L}_{\text{CL}}^I + \mathcal{L}_{\text{CL}}^S \quad (7)$$

where $s(\cdot)$ denotes the cosine similarity function, $\mathbf{h}_{\text{CLS}}^+$ denotes the positive sample of \mathbf{h}_{CLS} , \mathbf{h}_i^+ denotes the positive sample of \mathbf{h}_i , B denotes the mini-batch of original and code-switched tokens, and N denotes the total number of utterances.

4.3 Cyclical Contrastive Learning

Inspired by previous work (Goel et al., 2022), to improve the consistency between different languages, we introduce two additional consistency losses, including the cross-lingual consistency loss and the intra-language consistency loss.

The cross-lingual consistency loss $\mathcal{L}_{\text{CCL}}^C$ is applied to reduce the discrepancy in similarity scores between the representations of all mismatched pairs of original tokens and code-switched tokens, which can be formulated as follows:

$$\mathcal{L}_{\text{CCL}}^C = \frac{1}{N} \sum_{j=1}^N \sum_{i=1}^N (\langle \mathbf{H}_j, \bar{\mathbf{H}}_i \rangle - \langle \mathbf{H}_i, \bar{\mathbf{H}}_j \rangle)^2 \quad (8)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product function, and $\bar{\mathbf{H}}$ denotes the representation of the corresponding code-switched utterance.

The intra-lingual consistency loss $\mathcal{L}_{\text{CCL}}^I$ is employed to reduce the discrepancy in the similarity scores between the representations of all the original token pairs and corresponding code-switched

token pairs, which can be formulated as follows:

$$\mathcal{L}_{\text{CCL}}^I = \frac{1}{N} \sum_{j=1}^N \sum_{i=1}^N (\langle \mathbf{H}_j, \mathbf{H}_i \rangle - \langle \bar{\mathbf{H}}_i, \bar{\mathbf{H}}_j \rangle)^2 \quad (9)$$

The final cyclical contrastive learning loss \mathcal{L}_{CCL} is the sum of $\mathcal{L}_{\text{CCL}}^C$ and $\mathcal{L}_{\text{CCL}}^I$:

$$\mathcal{L}_{\text{CCL}} = \mathcal{L}_{\text{CCL}}^C + \mathcal{L}_{\text{CCL}}^I \quad (10)$$

4.4 Geodesic

In the previous contrastive paradigm, only the tokens with the same two labels, including intent and slot, are regarded as the positive pairs. Therefore, the tokens with only one different label (slot or intent) are also pushed apart without discrimination, which limits the overall performance. To solve this problem, we use geodesic to discriminate positive pairs in contrastive learning.

The representations of tokens are often embedded within a high-dimensional manifold, and our objective is to gauge the geodesic distance between two points along this manifold. However, calculating the precise geodesic distance proves challenging in the absence of explicit knowledge regarding the manifold’s structure (Kimmel and Sethian, 1998). To address this, we resort to leveraging the K-NN graph (Cover and Hart, 1967) as an approximation to the manifold structure (Surazhsky et al., 2005; Chowdhury et al., 2022). Within this graph, each token \mathbf{h}_i constitutes a node, and connections are established between nodes such that each node links to at most k other nodes.

Specifically, a directed edge is established from the node \mathbf{h}_i to node \mathbf{h}_j if \mathbf{h}_j is one of the k nearest neighbors of \mathbf{h}_i . The weight of each edge $d(\mathbf{h}_i, \mathbf{h}_j)$ is defined utilizing the cosine similarity:

$$d(\mathbf{h}_i, \mathbf{h}_j) = 1 - \mathbf{h}_i \mathbf{h}_j^\top \quad (11)$$

Finally, we employ the shortest path algorithm Dijkstra (Dijkstra, 1959) to compute the length of the shortest path between the two token representations along the obtained weighted directed graph, serving as the final geodesic distance $\mathcal{G}(\mathbf{h}_i, \mathbf{h}_j)$.

For a token \mathbf{h}_i , we define the k tokens with the closest geodesic distance from the code-switched tokens as its positive samples P_i :

$$P_i = \left\{ \mathbf{p}_i^k \right\} = \arg \text{topk}_{\mathcal{G}}(\mathbf{h}_i, \mathbf{h}_j) \quad (12)$$

In vanilla contrastive learning, for negative samples with only one different label and those with

two different labels, the push operation for negative samples is indistinguishable, which clearly undermines the model to learn the correct representations. As a result, we use the geodesic distance to differentially push negative samples away. The similarity $S_G(\mathbf{h}_i, \mathbf{h}_j)$ between different tokens is:

$$S_G(\mathbf{h}_i, \mathbf{h}_j) = \exp(\mathbf{h}_i \mathbf{h}_j^\top \cdot \log \frac{1}{\exp(\mathcal{G}(\mathbf{h}_i, \mathbf{h}_j) + 1)}) \quad (13)$$

By considering the relationships between negative samples while maximizing mutual information, we believe $S_G(\mathbf{h}_i, \mathbf{h}_j)$ is more beneficial than the conventional similarity function. The geodesic-based contrastive learning loss \mathcal{L}_{GCL} are as follows:

$$\mathcal{L}_{\text{GCL}}^I = - \sum_{j=1}^N \log \frac{\sum_{\mathbf{p}_{\text{CLS}}^k \in P_{\text{CLS}}} \exp(\mathbf{h}_{\text{CLS}}^j, \mathbf{p}_{\text{CLS}}^k)}{\sum_{\mathbf{h}_{\text{CLS}}^j \neq \mathbf{h}_{\text{CLS}}^{j'}} S_G(\mathbf{h}_{\text{CLS}}^j, \mathbf{h}_{\text{CLS}}^{j'})} \quad (14)$$

$$\mathcal{L}_{\text{GCL}}^S = - \frac{1}{n} \sum_{j=1}^N \sum_{i=1}^n \log \frac{\sum_{\mathbf{p}_i^k \in P_i} \exp(\mathbf{h}_i^j, \mathbf{p}_i^k)}{\sum_{\mathbf{h}_i^j \neq \mathbf{h}_i^{j'}} S_G(\mathbf{h}_i^j, \mathbf{h}_i^{j'})} \quad (15)$$

$$\mathcal{L}_{\text{GCL}} = \mathcal{L}_{\text{GCL}}^I + \mathcal{L}_{\text{GCL}}^S \quad (16)$$

4.5 Trainig Objective

Following previous work (Qin et al., 2020b, 2022), the intent detection objective \mathcal{L}_I and the slot filling objective \mathcal{L}_S are computed as follows:

$$\mathcal{L}_I = - \sum_{i=1}^{n_I} \hat{y}_i^I \log(\mathbf{o}_i^I) \quad (17)$$

$$\mathcal{L}_S = - \sum_{j=1}^n \sum_{i=1}^{n_S} \hat{y}_j^{i,S} \log(\mathbf{o}_j^{i,S}) \quad (18)$$

where \hat{y}_i^I denotes the gold intent label, $\hat{y}_j^{i,S}$ denotes the gold slot label for the j -th token, n_I denotes the number of gold intent labels, and n_S denotes the number of gold slot labels.

The final training objective \mathcal{L} is as follows:

$$\mathcal{L} = \alpha \mathcal{L}_I + (1 - \alpha) \mathcal{L}_S + \lambda \mathcal{L}_{\text{CCL}} + \gamma \mathcal{L}_{\text{GCL}} \quad (19)$$

5 Experiments

5.1 Datasets and Metrics

We primarily conduct our experiments on two public cross-lingual SLU benchmark datasets, including the MultiATIS++ (Xu et al., 2020) dataset and the MTOP (Li et al., 2021) dataset.

MultiATIS++² dataset is the broadened version of the Multilingual ATIS (Upadhyay et al., 2018) dataset, whose statistics are shown in Table 1. This extension includes human-translated data for an additional six languages: Spanish (es), German (de), Chinese (zh), Japanese (ja), Portuguese (pt), and French (fr), complementing the original languages, Hindi (hi) and Turkish (tr). The dataset comprises 4,478 utterances in the training set, 500 in the validation set, and 893 in the test set, with a total of 18 intents and 84 slots for each language.

Language	Utterances			Intent types	Slot types
	train	valid	test		
hi	1440	160	893	17	75
tr	578	60	715	17	71
others	4488	490	893	18	84

Table 1: Statistics of MultiATIS++ dataset.

MTOP³ is compiled from interactions between humans and assistant systems, with statistics presented in Table 2. MTOP comprises over 100,000 human-translated utterances in six languages (English (en), German (de), Spanish (es), French (fr), Thai (th), Hindi (hi)) across eleven domains. For a fair comparison, we Liang et al. (2022) to use the flat version, divided into 70:10:20 percentage splits for the training set, validation set, and test set.

Number of Total Utterances						Intent types	Slot types
en	de	fr	es	hi	th		
22288	18788	16584	15459	16131	15195	117	78

Table 2: Statistics of MTOP dataset.

Consistent with prior research (Qin et al., 2022; Zhu et al., 2023; Cheng et al., 2023), accuracy serves as the metric for evaluating intent detection, and F1 score is applied to assess slot filling performance. Moreover, overall accuracy is utilized for sentence-level semantic frame parsing evaluation.

5.2 Implementation Details

Following Qin et al. (2022), we utilize the base case of the multilingual BERT (mBERT)⁴(Devlin et al., 2019), featuring $N = 12$ attention heads and $M = 12$ transformer blocks. The learning rate is set to 5×10^{-7} and the total batch size is set to

²<https://github.com/amazon-science/multiatis>

³https://fb.me/mtop_dataset

⁴<https://github.com/google-research/bert/blob/master/multilingual.md>

Intent Accuracy	en	de	es	fr	hi	ja	pt	tr	zh	AVG
ZSJoint [‡] (Chen et al., 2019)	98.54	90.48	93.28	94.51	77.15	76.59	94.62	73.29	84.55	87.00
CoSDA [†] (Qin et al., 2021)	95.74	94.06	92.29	77.04	82.75	73.25	93.05	80.42	78.95	87.32
GL-CLEF* (Qin et al., 2022)	98.77	97.53	97.05	97.72	86.00	82.84	96.08	83.92	87.68	91.95
LAJ-MCL* (Liang et al., 2022)	98.77	98.10	98.10	98.77	84.54	81.86	97.09	85.45	89.03	92.41
DiffSLU* (Mao and Zhang, 2023)	98.86	98.17	98.21	98.93	86.66	82.65	97.21	85.98	89.46	92.90
SoGo* (Zhu et al., 2023)	98.89	98.45	98.15	97.74	83.87	84.75	97.73	85.53	89.10	92.69
FC-MTLF* (Cheng et al., 2023)	98.97	98.21	98.36	99.01	86.72	82.95	97.34	86.02	89.53	93.01
CCLG (ours)	99.35	98.51	98.94	99.43	87.32	85.53	98.79	86.48	89.97	93.81
Slot F1	en	de	es	fr	hi	ja	pt	tr	zh	AVG
ZSJoint [‡] (Chen et al., 2019)	95.20	74.79	76.52	74.25	52.73	70.10	72.56	29.66	66.91	68.08
CoSDA [†] (Qin et al., 2021)	92.29	81.37	76.94	79.36	64.06	66.62	75.05	48.77	77.32	73.47
GL-CLEF* (Qin et al., 2022)	95.39	86.30	85.22	84.31	70.34	73.12	81.83	65.85	77.61	80.00
LAJ-MCL* (Liang et al., 2022)	96.02	86.59	83.03	82.11	61.04	68.52	81.49	65.20	82.00	78.23
DiffSLU* (Mao and Zhang, 2023)	96.16	86.72	85.48	84.26	73.04	74.12	82.52	68.14	83.12	81.51
SoGo* (Zhu et al., 2023)	95.42	87.46	87.01	84.45	74.25	76.69	83.91	67.04	78.53	81.64
FC-MTLF* (Cheng et al., 2023)	96.21	86.87	85.66	84.62	73.18	74.24	82.68	68.22	83.16	81.65
CCLG (ours)	96.83	88.01	87.45	85.22	74.97	77.19	84.17	68.98	83.82	82.96
Overall Accuracy	en	de	es	fr	hi	ja	pt	tr	zh	AVG
ZSJoint [‡] (Chen et al., 2019)	87.23	41.43	44.46	43.67	16.01	33.59	43.90	1.12	30.80	38.02
CoSDA [†] (Qin et al., 2021)	77.04	57.06	46.62	50.06	26.20	28.89	48.77	15.24	46.36	44.03
GL-CLEF* (Qin et al., 2022)	88.02	66.03	59.53	57.02	34.83	41.42	60.43	28.95	50.62	54.09
LAJ-MCL* (Liang et al., 2022)	89.81	67.75	59.13	57.56	23.29	29.34	61.93	28.95	54.76	52.50
DiffSLU* (Mao and Zhang, 2023)	90.06	68.02	59.84	58.08	35.12	43.06	63.04	29.32	55.08	55.74
SoGo* (Zhu et al., 2023)	90.54	72.26	61.05	57.88	39.90	46.95	64.23	29.14	51.31	57.02
FC-MTLF* (Cheng et al., 2023)	91.58	69.54	61.43	59.62	36.86	44.64	64.55	30.86	56.52	57.29
CCLG (ours)	91.97	74.91	62.43	59.99	40.43	47.98	64.95	31.56	57.83	59.12

Table 3: Experiment Results on the MultiATIS++ dataset. We report both individual and average (AVG) results. Results with “*” are obtained from the respective published paper, results with “†” are cited from Qin et al. (2022), and results with “‡” are cited from Liang et al. (2022). The symbol “–” indicates missing results from the published work. Results in **bold** denote our framework significantly outperforms baselines with $p < 0.01$ under t-test.

16. During the training process, the value of label smoothing is set to 0.1, and the dropout rate is set to 0.1. We train the model for 40 epochs, and to avoid overfitting, the training will early-stop if the loss on the development set does not decrease for 10 epochs. We use Adam optimizer (Kingma and Ba, 2015) with $\beta_1 = 0.9$, $\beta_2 = 0.98$, and 4k warm-up updates to optimize parameters. Following the zero-shot setting, we choose the model with the highest overall accuracy based on the English development set and subsequently evaluate on test datasets. For all hyper-parameters, we perform several experiments and select the values with the best performance. α is set to 0.9, λ is set to 0.5, γ is set to 1, and k is set to 5. The experiments are conducted on an NVIDIA A100. Our code is based on PyTorch (Paszke et al., 2019) and Transformers⁵(Wolf et al., 2020) framework.

5.3 Baselines

We compare our proposed approach with the following baselines, including ZSJoint (Chen et al.,

⁵<https://github.com/huggingface/transformers>

Methods	Intent Acc	Slot F1	Overall Acc
ZSJoint [◇]	85.31	67.26	52.15
CoSDA [‡]	90.72	73.34	58.77
CL-CLEF [◇]	88.94	79.86	61.24
LAJ-MCL*	91.04	74.50	60.11
CCLG (ours)	92.42	82.24	64.36

Table 4: Average results of all the languages on MTOP. Results with ‡ are cited from Liang et al. (2022), results with * are from the corresponding published paper, results with ◇ are obtained by our re-implementation, and results in **bold** denote our framework significantly outperforms baselines with $p < 0.01$ under t-test.

2019), CoSDA (Qin et al., 2021), GL-CLEF (Qin et al., 2022), LAJ-MCL (Liang et al., 2022), DiffSLU (Mao and Zhang, 2023), SoGo (Zhu et al., 2023), and FC-MTLF (Cheng et al., 2023), whose details are provided in Appendix A.

5.4 Main Results

The results on MultiATIS++ are shown in Table 3 and the results on MTOP are listed in Table 4. From them, we have the following observations:

Intent Accuracy	en	de	es	fr	hi	ja	pt	tr	zh	AVG
CCLG (ours)	99.35	98.51	98.94	99.43	87.32	85.53	98.79	86.48	89.97	93.81
w/o Cyclical Contrastive Learning	98.21	97.76	97.11	97.74	86.14	84.15	96.01	84.23	88.13	92.16
w/o Geodesic	98.05	97.23	96.54	97.12	85.22	82.05	95.33	83.24	87.42	91.36
Slot F1	en	de	es	fr	hi	ja	pt	tr	zh	AVG
CCLG (ours)	96.83	88.01	87.45	85.22	74.97	77.19	84.17	68.98	83.82	82.96
w/o Cyclical Contrastive Learning	96.13	87.11	86.82	84.75	74.23	76.65	83.76	68.33	83.08	82.32
w/o Geodesic	95.13	86.04	85.03	83.76	69.97	72.44	81.03	64.98	77.01	79.49
Overall Accuracy	en	de	es	fr	hi	ja	pt	tr	zh	AVG
CCLG (ours)	91.97	74.91	62.43	59.99	40.43	47.98	64.95	31.56	57.83	59.12
w/o Cyclical Contrastive Learning	91.13	74.22	62.01	59.56	39.64	47.45	64.33	31.02	56.76	58.46
w/o Geodesic	87.62	65.73	59.14	56.62	34.44	41.02	60.11	28.63	50.14	53.72

Table 5: Ablation study of difference components on the MutliATIS++ dataset.

(1) The methodologies employed in CoSDA, GL-CLEF, LAJ-MCL, and FC-MTLF all incorporate code-switching, and it is evident that they outperform models that do not use this technique, showcasing its effectiveness in enhancing model performance compared to those that do not utilize such strategies. Moreover, our proposed approach goes beyond these established approaches by introducing a novel framework that achieves even greater performance gains. With the relative enhancement of 1.83% in average overall accuracy over the previous state-of-the-art model, our method stands out. This notable improvement can be attributed to our innovative approach based on cyclical contrastive learning based on geodesic.

(2) CCLG obtains notable and consistent advancements across all subtasks, particularly showcasing significant improvements. Its impact is particularly pronounced in low-resource languages compared to high-resource ones. The substantial improvement achieved in these languages surpasses gains observed in other high-resource languages. The success of CCLG in low-resource languages aligns with the original intent of the zero-shot cross-lingual SLU task, which aimed to address challenges in languages with limited training data.

5.5 Ablation Study

To validate the advantages of CCLG from different perspectives, we conduct several ablation studies on the MixATIS++ dataset, the results of which are demonstrated in Table 5.

5.5.1 Effect of Cyclical Contrastive Learning

CCLG makes a pivotal contribution through its innovative cyclical contrastive learning, strategically achieving consistency across different languages.

Methods	Intent Acc	Slot F1	Overall Acc
ChatGPT	73.25	61.57	39.16
Vicuna 1.3 (7B)	72.91	60.40	37.05
LLaMA 2 (7B)	72.86	61.20	37.28
CCLG (ours)	93.81	82.96	59.12

Table 6: Results of LLMs on the MutliATIS++ dataset.

To meticulously evaluate the impact of this module, we conduct an ablation study by excluding \mathcal{L}_{CCL} in Eq. 19, as denoted by "w/o Cyclical Contrastive Learning" in Table 5. A discernible degradation in performance emerges across all metrics for every language when the cyclical contrastive learning module is omitted. We contend that this observed improvement stems from the module’s capability to model the consistency between different languages, particularly beneficial for low-resource languages facing the data scarcity challenges.

5.5.2 Effect of Geodesic

To bolster the effectiveness of geodesic, we conduct an ablation study by excluding \mathcal{L}_{GCL} in Eq. 19. This configuration is denoted as "w/o Geodesic" in Table 5. Significantly, our findings reveal a decline in performance across all metrics for each language, underscoring the importance of geodesic in constructing positive and negative samples in contrastive learning. This ensures a robust and reliable model performance in real-world applications.

5.6 Comparison with Large Language Models

As demonstrated in Table 6, we utilize the evaluation methodology introduced by He and Garner (2023) to assess the performance of ChatGPT (OpenAI, 2023), Vicuna 1.3 (7B) (Zheng et al., 2023), and LLaMA 2(7B) (Touvron et al., 2023). In this

	Text (En):	show	flights	from	burbank	to	st.	louis	on	monday
Ref.	Intent:	atis_flight								
	Slot:	O	O	O	B-fromloc.city_name	O	B-toloc.city_name	I-toloc.city_name	O	B-depart_date.day_name
GL-CL_EF	Intent:	atis_flight								
	Slot:	O	O	O	B-fromloc.city_name	O	O	O	O	B-depart_date.day_name
FC-MTLF	Intent:	atis_flight								
	Slot:	O	O	O	B-fromloc.city_name	O	B-toloc.city_name	O	O	B-depart_date.day_name
CCLG	Intent:	atis_flight								
	Slot:	O	O	O	B-flight_stop	O	O	B-fromloc.city_name	O	B-toloc.city_name
	Text (De):	Zeige	Flüge	von	Burbank	nach	St.	Louis	für	Montag
Ref.	Intent:	atis_flight								
	Slot:	O	O	O	B-fromloc.city_name	O	B-toloc.city_name	I-toloc.city_name	O	B-depart_date.day_name
GL-CL_EF	Intent:	atis_airline								
	Slot:	O	O	O	B-fromloc.city_name	O	O	O	O	O
FC-MTLF	Intent:	atis_airline								
	Slot:	O	O	O	B-fromloc.city_name	O	B-toloc.city_name	O	O	O
CCLG	Intent:	atis_flight								
	Slot:	O	O	O	B-fromloc.city_name	O	B-toloc.city_name	I-toloc.city_name	O	B-depart_date.day_name

Table 7: Case study on MultiATIS++ dataset. Text in red denotes the incorrect predictions.

evaluation, the models are presented with 20 examples each. Despite the impressive performance demonstrated by Large Language Models (LLMs) in few-shot and zero-shot learning tasks, a significant performance gap of approximately 20% persists between these models and CCLG in terms of overall accuracy on the MultiATIS++ dataset. This performance disparity is consistently observed across other datasets as well. The observed performance degradation highlights the persistent challenges that language models encounter in understanding spoken language, despite their advanced few-shot and zero-shot learning capabilities. This underscores the urgent need for dedicated efforts in designing effective zero-shot cross-lingual SLU frameworks. Addressing these challenges is not only crucial but also remains an ongoing and vital task for the NLP community. Further exploration and investigation into innovative approaches are warranted to advance state-of-the-art performance.

5.7 Case Study

To further verify the advancements of our model compared to previous methods in zero-shot cross-lingual SLU, we present a case study across different languages. Specifically, we examine English and German as two representative examples. The results in Table 7 reveal notable distinctions in the performance of GL-CL_EF, FC-MTLF, and CCLG.

In the case of English, all these models correctly predict the intent. However, as the linguistic complexity increases in German, errors become more pronounced in GL-CL_EF and FC-MTLF, while CCLG maintains correct predictions. It exemplifies the robustness and cross-lingual generalizability of

CCLG, outperforming its counterparts in accurately predicting intents across diverse languages, without succumbing to increased linguistic complexity, thereby enhancing overall performance.

In terms of slot filling accuracy, GL-CL_EF and FC-MTLF show some errors in English, whereas CCLG maintains accuracy. Moving to German, the errors in GL-CL_EF and FC-MTLF become more pronounced, while CCLG continues to maintain a high performance. This observed trend highlights the robust nature of CCLG, showcasing its consistent superiority in accurately predicting slots.

6 Conclusion

In this paper, we propose a novel framework CCLG for zero-shot cross-lingual spoken language understanding (SLU), which utilizes cyclical contrastive learning to achieve consistency across different languages and applies geodesic to construct positive samples and negative samples in contrastive learning. Experiments on the MultiATIS++ dataset and the MTOP dataset show that CCLG outperforms the previous best model and achieves a new state-of-the-art performance. Further analysis also demonstrates that our method can indeed transfer knowledge between different languages effectively.

Limitations

While our approach achieves state-of-the-art performance by modifying the traditional contrastive paradigm, we recognize the potential for further enhancements through the incorporation of external knowledge. Given the recent successes observed with LLMs, we anticipate that harnessing LLMs could yield additional improvements in our model’s

485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504

505

506
507
508
509
510
511
512
513
514
515
516
517
518

519
520
521
522
523
524
525
526
527
528
529
530

531

532
533
534
535
536
537
538
539
540
541
542
543

544

545
546
547
548
549
550
551

552	performance. Exploring the integration of LLMs		
553	into our framework represents a promising avenue.		
554	We leave this aspect for future work.		
555	Ethics Statement		
556	We conducted all experiments using publicly avail-		
557	able datasets that are free from offensive content or		
558	information with negative social impact. The main		
559	objective of this paper is to enhance the model’s		
560	capacity for understanding, and our model does		
561	not generate any uncontrollable output. Hence, we		
562	took measures to ensure that our paper adheres to		
563	ethical review guidelines. By prioritizing ethical		
564	considerations, our aim is to contribute responsibly		
565	to the advancement of NLP technology.		
566	References		
567	Shuyang Cao and Lu Wang. 2021. CLIFF: Contrastive		
568	learning for improving faithfulness and factuality in		
569	abstractive summarization. In <i>Proc. of EMNLP</i> .		
570	Qian Chen, Zhu Zhuo, and Wen Wang. 2019. Bert		
571	for joint intent classification and slot filling. <i>arXiv</i>		
572	<i>preprint arXiv:1902.10909</i> .		
573	Ting Chen, Simon Kornblith, Mohammad Norouzi, and		
574	Geoffrey E. Hinton. 2020. A simple framework for		
575	contrastive learning of visual representations. In		
576	<i>Proc. of ICML</i> .		
577	Xuxin Cheng, Wanshi Xu, Ziyu Yao, Zhihong Zhu,		
578	Yaowei Li, Hongxiang Li, and Yuexian Zou. 2023.		
579	Fc-mtlf: a fine-and coarse-grained multi-task learn-		
580	ing framework for cross-lingual spoken language un-		
581	derstanding. In <i>Proc. of Interspeech</i> , volume 2.		
582	Sumit Chopra, Raia Hadsell, and Yann LeCun. 2005.		
583	Learning a similarity metric discriminatively, with		
584	application to face verification. In <i>Proc. of CVPR</i> .		
585	Somnath Basu Roy Chowdhury, Nicholas Monath,		
586	Avinava Dubey, Amr Ahmed, and Snigdha		
587	Chaturvedi. 2022. Unsupervised opinion summa-		
588	rization using approximate geodesics. <i>arXiv preprint</i>		
589	<i>arXiv:2209.07496</i> .		
590	Ching-Yao Chuang, Joshua Robinson, Yen-Chen Lin,		
591	Antonio Torralba, and Stefanie Jegelka. 2020. Debi-		
592	ased contrastive learning. <i>Advances in neural infor-</i>		
593	<i>mation processing systems</i> , 33:8765–8775.		
594	Thomas Cover and Peter Hart. 1967. Nearest neighbor		
595	pattern classification. <i>IEEE transactions on informa-</i>		
596	<i>tion theory</i> , 13(1):21–27.		
597	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and		
598	Kristina Toutanova. 2019. Bert: Pre-training of deep		
599	bidirectional transformers for language understand-		
600	ing. In <i>Proceedings of the 2019 Conference of the</i>		
	<i>North American Chapter of the Association for Com-</i>		601
	<i>putational Linguistics: Human Language Technolo-</i>		602
	<i>gies, Volume 1 (Long and Short Papers)</i> , pages 4171–		603
	4186.		604
	EW Dijkstra. 1959. A note on two problems in connex-		605
	ion with graphs. <i>Numerische Mathematik</i> , 1(1):269–		606
	271.		607
	John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader.		608
	2021. Declutr: Deep contrastive learning for unsuper-		609
	vised textual representations. In <i>Proceedings of the</i>		610
	<i>59th Annual Meeting of the Association for Compu-</i>		611
	<i>tational Linguistics and the 11th International Joint</i>		612
	<i>Conference on Natural Language Processing (Vol-</i>		613
	<i>ume 1: Long Papers)</i> , pages 879–895.		614
	Shashank Goel, Hritik Bansal, Sumit Bhatia, Ryan		615
	Rossi, Vishwa Vinay, and Aditya Grover. 2022. Cy-		616
	clip: Cyclic contrastive language-image pretraining.		617
	<i>Advances in Neural Information Processing Systems</i> ,		618
	35:6704–6719.		619
	Mutian He and Philip N. Garner. 2023. Can ChatGPT		620
	Detect Intent? Evaluating Large Language Models		621
	for Spoken Language Understanding. In <i>Proc. of</i>		622
	<i>Interspeech</i> .		623
	Ron Kimmel and James A Sethian. 1998. Computing		624
	geodesic paths on manifolds. <i>Proceedings of the</i>		625
	<i>national academy of Sciences</i> , 95(15):8431–8435.		626
	Diederik P. Kingma and Jimmy Ba. 2015. Adam: A		627
	method for stochastic optimization. In <i>Proc. of ICLR</i> .		628
	Guillaume Lample, Alexis Conneau, Marc’Aurelio Ran-		629
	zato, Ludovic Denoyer, and Hervé Jégou. 2018.		630
	Word translation without parallel data. In <i>Proc. of</i>		631
	<i>ICLR</i> .		632
	Haoran Li, Abhinav Arora, Shuohui Chen, Anchit		633
	Gupta, Sonal Gupta, and Yashar Mehdad. 2021.		634
	Mtop: A comprehensive multilingual task-oriented		635
	semantic parsing benchmark. In <i>Proceedings of the</i>		636
	<i>16th Conference of the European Chapter of the Asso-</i>		637
	<i>ciation for Computational Linguistics: Main Volume</i> ,		638
	pages 2950–2962.		639
	Shining Liang, Linjun Shou, Jian Pei, Ming Gong,		640
	Wanli Zuo, Xianglin Zuo, and Daxin Jiang. 2022.		641
	Label-aware multi-level contrastive learning for		642
	cross-lingual spoken language understanding. In <i>Pro-</i>		643
	<i>ceedings of the 2022 Conference on Empirical Meth-</i>		644
	<i>ods in Natural Language Processing</i> , pages 9903–		645
	9918.		646
	Risheng Liu, Zhiying Jiang, Shuzhou Yang, and Xin		647
	Fan. 2022. Twin adversarial contrastive learning for		648
	underwater image enhancement and beyond. <i>IEEE</i>		649
	<i>Transactions on Image Processing</i> .		650
	Zihan Liu, Genta Indra Winata, Zhaojiang Lin, Peng		651
	Xu, and Pascale Fung. 2020. Attention-informed		652
	mixed-language training for zero-shot cross-lingual		653
	task-oriented dialogue systems. In <i>Proceedings of</i>		654
	<i>the AAAI Conference on Artificial Intelligence</i> , pages		655
	8433–8440.		656

657	Tianjun Mao and Chenghong Zhang. 2023. DiffSLU: Knowledge Distillation Based Diffusion Model for Cross-Lingual Spoken Language Understanding . In <i>Proc. INTERSPEECH 2023</i> , pages 715–719.	Vitaly Surazhsky, Tatiana Surazhsky, Danil Kirsanov, Steven J Gortler, and Hugues Hoppe. 2005. Fast exact and approximate geodesics on meshes. <i>ACM transactions on graphics (TOG)</i> , 24(3):553–560.	712 713 714 715
661	OpenAI. 2023. ChatGPT (Mar 14 version) [Large language model].	Hugo Touvron, Louis Martin, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. <i>ArXiv preprint</i> .	716 717 718
663	Xiao Pan, Mingxuan Wang, Liwei Wu, and Lei Li. 2021. Contrastive learning for many-to-many multilingual neural machine translation. In <i>Proc. of ACL</i> .	Gokhan Tur and Renato De Mori. 2011. <i>Spoken language understanding: Systems for extracting semantic information from speech</i> . John Wiley & Sons.	719 720 721
666	Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. <i>Advances in neural information processing systems</i> , 32.	Shyam Upadhyay, Manaal Faruqui, Gokhan Tür, Hakkani-Tür Dilek, and Larry Heck. 2018. (almost) zero-shot cross-lingual spoken language understanding. In <i>2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 6034–6038. IEEE.	722 723 724 725 726 727
672	Libo Qin, Wanxiang Che, Yangming Li, Haoyang Wen, and Ting Liu. 2019. A stack-propagation framework with token-level intent detection for spoken language understanding. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 2078–2087.	Rob van der Goot, Ibrahim Sharaf, Aizhan Imankulova, Ahmet Üstün, Marija Stepanović, Alan Ramponi, Siti Oryza Khairunnisa, Mamoru Komachi, and Barbara Plank. 2021. From masked language modeling to translation: Non-english auxiliary tasks improve zero-shot spoken language understanding. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 2479–2497.	728 729 730 731 732 733 734 735 736 737
680	Libo Qin, Qiguang Chen, Tianbao Xie, Qixin Li, Jianguang Lou, Wanxiang Che, and Min-Yen Kan. 2022. Gl-clef: A global-local contrastive learning framework for cross-lingual spoken language understanding. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2677–2686.	Danqing Wang, Jiase Chen, Hao Zhou, Xipeng Qiu, and Lei Li. 2021. Contrastive aligned joint learning for multilingual summarization. In <i>Proc. of ACL Findings</i> .	738 739 740 741
687	Libo Qin, Tailu Liu, Wanxiang Che, Bingbing Kang, Sendong Zhao, and Ting Liu. 2021. A co-interactive transformer for joint slot filling and intent detection. In <i>2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pages 8193–8197. IEEE.	Feng Wang and Huaping Liu. 2021. Understanding the behaviour of contrastive loss. In <i>Proc. of CVPR</i> .	742 743
693	Libo Qin, Minheng Ni, Yue Zhang, and Wanxiang Che. 2020a. Cosda-ml: Multi-lingual code-switching data augmentation for zero-shot cross-lingual NLP. In <i>Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI 2020</i> , pages 3853–3860.	Yuxuan Wang, Wanxiang Che, Jiang Guo, Yijia Liu, and Ting Liu. 2019. Cross-lingual BERT transformation for zero-shot dependency parsing. In <i>Proc. of EMNLP</i> .	744 745 746 747
699	Libo Qin, Xiao Xu, Wanxiang Che, and Ting Liu. 2020b. Agif: An adaptive graph-interactive framework for joint multiple intent detection and slot filling. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1807–1816.	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2020. Transformers: State-of-the-art natural language processing. In <i>Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations</i> , pages 38–45.	748 749 750 751 752 753 754
704	Nikunj Saunshi, Orestis Plevrakis, Sanjeev Arora, Mikhail Khodak, and Hrishikesh Khandeparkar. 2019. A theoretical analysis of contrastive unsupervised representation learning. In <i>Proc. of ICML</i> .	Bowen Xing and Ivor Tsang. 2022. Co-guiding net: Achieving mutual guidances between multiple intent detection and slot filling via heterogeneous semantics-label graphs. In <i>Proc. of EMNLP</i> .	755 756 757 758
708	Mengxiao Song, Bowen Yu, Li Quangang, Wang Yubin, Tingwen Liu, and Hongbo Xu. 2022. Enhancing joint multiple intent detection and slot filling with global intent-slot co-occurrence. In <i>Proc. of EMNLP</i> .	Weijia Xu, Batoool Haider, and Saab Mansour. 2020. End-to-end slot alignment and recognition for cross-lingual nlu. In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5052–5063.	759 760 761 762 763
711		Yuanmeng Yan, Rumei Li, Sirui Wang, Fuzheng Zhang, Wei Wu, and Weiran Xu. 2021. ConSERT: A contrastive framework for self-supervised sentence representation transfer. In <i>Proc. of ACL</i> .	764 765 766 767

768 Rong Ye, Mingxuan Wang, and Lei Li. 2022. Cross-
769 modal contrastive learning for speech translation. In
770 *Proc. of NAACL*.

771 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan
772 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,
773 Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023.
774 Judging llm-as-a-judge with mt-bench and chatbot
775 arena. *arXiv preprint arXiv:2306.05685*.

776 Zhihong Zhu, Xuxin Cheng, Zhiqi Huang, Dongsheng
777 Chen, and Yuexian Zou. 2023. Enhancing code-
778 switching for cross-lingual slu: A unified view of
779 semantic and grammatical coherence. In *Proceed-*
780 *ings of the 2023 Conference on Empirical Methods*
781 *in Natural Language Processing*, pages 7849–7856.

782 A Details of Baselines

783 Here we provide the details of baselines:

784 (1) ZSJoint: We have re-implemented the zero-
785 shot joint model (Chen et al., 2019) (referred to as
786 ZSJoint), trained on the English training set and
787 directly applied to the test sets of target languages.

788 (2) CoSDA: Qin et al. (2021) introduces a dy-
789 namic code-switching method involving random
790 multilingual token-level replacement. For a fair
791 comparison, we utilize both English training data
792 and code-switching data for fine-tuning.

793 (3) GL-CLEF: Qin et al. (2022) proposes a
794 global-local contrastive learning framework for ex-
795 plicit alignment, achieving the different granularity
796 alignments, including sentence-level local intent
797 alignment, token-level local slot alignment, and
798 semantic-level global intent-slot alignment.

799 (4) LAJ-MCL: Liang et al. (2022) introduces
800 a multi-level contrastive learning framework de-
801 signed for zero-shot cross-lingual SLU.

802 (5) DiffSLU: Mao and Zhang (2023) introduces
803 a diffusion model and applies knowledge distil-
804 lation for zero-shot cross-lingual SLU, achieving
805 mutual guidance between intent and slots.

806 (6) SoGo: Zhu et al. (2023) proposes a semantics-
807 coherent and grammar-coherent method to enhance
808 code-switching method for zero-shot cross-lingual
809 SLU, effectively boosting the performance.

810 (7) FC-MTLF: Cheng et al. (2023) introduces a
811 framework for cross-lingual SLU, utilizing code-
812 switching for coarse-grained alignment and ma-
813 chine translation for fine-grained alignment.