

PICD-Instruct: A Generative Instruction Learning Framework for Few-Shot Multi-Intent Spoken Language Understanding

Anonymous ACL submission

Abstract

Few-shot multi-intent spoken language understanding (SLU) aims to identify users' multiple intents and key slots using a tiny amount of annotated data. Recent advances in large language models (LLMs) have utilized instruction learning frameworks to model intent-slot interdependencies, typically requiring abundant data for effective training. However, in few-shot scenarios, these frameworks face challenges such as mismatches between the number of generated slots and input lengths, relational confusion in multi-intent scenarios and neglect of task-specific variations in intent counts across utterances. To overcome the challenges, we propose PICD-Instruct, a novel generative framework based on Basic Instructions (BI), Pairwise Interaction Instructions (PII) and Contrastive Distinct Instructions (CDI). Specifically, BI directs LLMs to generate entities along with associated words, thereby mitigating mismatches in quantitative correspondences. PII explicitly captures dual-task interdependencies by guiding LLMs to pair each intent with its related entities. CDI enhances understanding of utterances by guiding LLMs to determine whether two utterances share the same intent count. Experimental results on public datasets indicate that PICD-Instruct achieves state-of-the-art performance.

1 Introduction

Spoken Language Understanding (SLU) (Young et al., 2013) is a fundamental component of task-oriented dialogue systems. Among the various aspects of SLU, multi-intent SLU has gained significant attention due to its practical necessity in complex interactive scenarios. This task involves two closely linked subtasks: multi-intent detection and slot filling. Multi-intent detection focuses on identifying the intents embedded within a user utterance, whereas slot filling extracts key semantic information from the utterance. In practical applications, however, obtaining sufficient labeled

data for domain-specific SLU models is often time-intensive and costly. These challenges highlight the critical importance of exploring multi-intent SLU in low-resource settings.

Given the bidirectional relationship between intents and slots, recent models leverage multi-task joint frameworks to capture these interdependencies, achieving strong performance with sufficient training data (Goo et al., 2018; Li et al., 2018; Niu et al., 2019; Liu et al., 2019a; Qin et al., 2020, 2021; Song et al., 2022; Chen et al., 2022; Xing and Tsang, 2022a,b; Mei et al., 2023; Song et al., 2024). Meanwhile, large language models (LLMs) show promise in the zero-shot SLU task (Pan et al., 2023; Zhu et al., 2024) but remain largely designed for single-intent scenarios. For instance, Pan et al. (2023) explored prompt-based zero-shot SLU with ChatGPT, but its slot filling lagged far behind fine-tuned models. Similarly, Zhu et al. (2024) proposed a pseudo-labeling framework to enhance task collaboration but faced error propagation issues. To address these limitations, Xing et al. (2024) first introduced instruction learning into generative multi-intent SLU. Their framework leverages instruction learning and contrastive learning to model intent-slot relationships through mutual prediction of ground-truth labels. By distinguishing task-specific semantics across utterances, this approach enhances SLU reasoning. This raises a key question: Can instruction-guided LLMs achieve superior performance in few-shot multi-intent SLU?

Beyond traditional SLU challenges, LLMs introduce new opportunities by enhancing structured and reliable information extraction (Li et al., 2024). SLU plays a crucial role in intelligent agent-driven task completion, where accurate intent detection ensures effective execution of user commands (Caren Han et al., 2022). Unlike open-ended generation, SLU requires structured output to maintain schema consistency, which is critical for applications in domains such as voice assistants, cus-

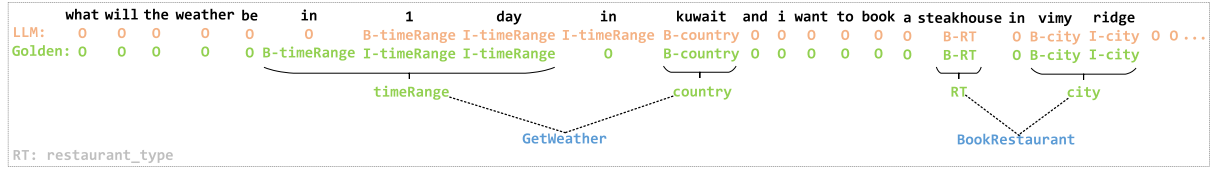


Figure 1: An example from MixSNIPS dataset. Traditional LLMs-generated slot labels are in orange, while golden slot labels and our proposed entity labels are in green. Intent labels are in blue.

tomer service automation, and smart device control (Saxon et al., 2021; Irugalbandara, 2024).

We discover three core challenges in leveraging LLMs for the few-shot multi-intent SLU. Firstly, the uncontrollable nature of LLM-generated outputs poses significant challenges for slot filling, as the number of generated slot often fails to correspond with the input length. This issue is exacerbated in few-shot settings, where limited training data restricts the model’s ability to accurately map slots to tokens. As shown in Fig. 1, the example demonstrates the over-generation and mismatch of slot labels. Secondly, existing generative frameworks fail to effectively capture the semantic dependencies between intents and slots. DC-Instruct (Xing et al., 2024) predicts slot labels based on the provided utterance and intent labels, but it falls short in establishing a one-to-one correspondence between each intent and its associated slots. This leads to confusion in multi-intent scenarios, making it harder for models to learn dual-task interdependencies with limited training data. Thirdly, as an utterance may contain multiple intents, its semantic structure becomes more intricate. Therefore, improving the sensitivity of LLMs to the variations in intent counts across utterances can enhance their understanding of such cases. However, current approaches often overlook this task-specific feature, potentially hindering the models’ ability to effectively comprehend utterances with multiple intents.

To overcome these challenges, we propose PICD-Instruct, a novel generative model based on instruction learning. PICD-Instruct employs three types of instructions: Basic Instructions (BI), Pairwise Interaction Instructions (PII) and Contrastive Distinct Instructions (CDI). BI guides the model in generating intent and slot labels by clearly defining task instructions, providing candidate labels, and specifying output formats. In slot filling, BI utilizes a key-value structure to link entities with specific tokens, effectively avoiding mismatches between generated slots and input lengths observed in the process of directly using LLMs to generate slots.

Considering that each green entity label in Fig. 1 aligns exactly with its associated words, PII incorporates an auxiliary intent-slot pairing task that explicitly models the bidirectional dependencies between intents and slots. By aligning golden intent labels with corresponding entity labels, PII mitigates relational confusions in multi-intent scenarios. CDI enhances the understanding of utterances with multiple intents by introducing a task to determine whether two utterances share the same number of intents. By leveraging positive and negative samples alongside the current utterance, CDI trains the model to distinguish between utterances based on intent counts, thereby improving its comprehension capabilities.

We conduct experiments on two few-shot datasets, FewShotMixATIS and FewShotMixSNIPS (Hua et al., 2024). Experimental results demonstrate that PICD-Instruct significantly outperforms existing baselines, achieving state-of-the-art (SOTA) performance in the few-shot multi-intent SLU task. The ablation study and additional experiments further confirm the robustness and advantages of our model.

In summary, our contributions are three-fold:

(1) We propose PICD-Instruct, a novel generative instruction-learning framework that integrates pairwise interactive instructions and contrastive distinct instructions to overcome challenges in the few-shot multi-intent SLU task.

(2) We advance the explicit modeling of bidirectional dependencies between intents and slots in a low-resource setting, reducing relational confusions in multi-intent scenarios through the application of instruction learning.

(3) PICD-Instruct achieves SOTA performance in the few-shot multi-intent SLU task, as evidenced by extensive experiments and analyses.

2 Related Work

Multi-intent SLU Prevailing models (Kim et al., 2017; Gangadharaiah and Narayanaswamy, 2019) often employ joint modeling to simultaneously learn the two tasks in SLU and capture their rela-

tions. Gangadharaiah and Narayanaswamy (2019) jointly model multiple intent detection and slot filling via a slot-gate mechanism. To better model the two tasks’ interactions, graph neural networks have been widely utilized (Qin et al., 2020, 2021; Xing and Tsang, 2022a,b; Song et al., 2022). The Co-guiding Net (Xing and Tsang, 2022a) pioneers in achieving mutual guidance between the two tasks through a two-stage framework. LCLR (Zhu et al., 2023) proposes to leverage the dual-task correlations in the decoding process. DC-Instruct (Xing et al., 2024) employs instructions for LLMs to predict one subtask’s labels based on the other’s golden labels, effectively capturing the relationships between intents and slots. UGEN (Wu et al., 2022) and PromptSLU (Song et al., 2024) performs multi-intent SLU based on the paradigm of prompt learning.

The above approaches primarily focus on scenarios with abundant training data. However, in few-shot settings, capturing the correlations between the two tasks in SLU becomes significantly more challenging, leading to degraded performance for most models (Hua et al., 2024). While UGEN and DC-Instruct have demonstrated performance in low-resource settings, the few-shot training data they utilize does not align well with real-world application scenarios in terms of sample quantity and distribution. To better simulate practical application scenarios, we employ FewShotMixATIS and FewShotMixSNIPS, two datasets specifically tailored for few-shot scenarios, as the training data for our model. Different from recent works, we propose a novel generative framework incorporating various instructions to ensure the accuracy of LLM outputs. Our approach explicitly captures dual-task interdependencies by reducing relational confusions and effectively harnesses the variations of intent counts across different utterances, enabling improved performance in the few-shot multi-intent SLU task.

Instruction Learning Recently, the rise of LLMs in the natural language processing (NLP) field has positioned instruction learning as a competitive approach across various NLP tasks (Lou et al., 2024; Safa et al., 2024). This paradigm effectively leverages the advanced conversational abilities of LLMs to perform generative tasks, bridging the gap between the pre-training and fine-tuning stages.

In this work, we investigate instruction learning for few-shot multi-intent SLU and propose a

novel model characterized by pairwise interactive instructions and contrastive distinct instructions.

3 Task Definition

As shown in the example in Fig. 1, multi-intent SLU aims to detect all possible intents within an utterance and identify the slot label corresponding to each word. Therefore, multi-intent detection is considered as a multi-label text classification task and slot filling is regarded as a sequence labeling task. The task can be formulated as follows: given an input utterance $X = \{W_1, W_2, \dots, W_n\}$, where n is the length of the utterance. The objective is to predict the correct intents from the candidate intents $I = \{i_1, i_2, \dots, i_m\}$ and indentify the slot label for each word W_i from the candidate slot types $S = \{s_1, s_2, \dots, s_k\}$, where m is the number of intent categories, and k is the number of slot types. However, in real-world scenarios, obtaining sufficient annotated data is often impractical due to data scarcity and annotation costs. This challenge is particularly pronounced in low-resource domains and emerging applications. Therefore, it is crucial to design models capable of handling the multi-intent SLU task effectively in a few-shot setting, where only limited annotated examples are available.

4 Methodology

In this section, we introduce our proposed PICD-Instruct framework. As depicted in Fig. 2, we formulate our instructions in a question-answer (QA) form. The framework includes three types of instructions, each corresponding to a specific task. This approach mitigates the effects of uncontrollable generation by LLMs and more explicitly models the correlations between the two tasks in SLU, reducing relational confusions. In addition, it enhances the model’s ability to understand utterances with multiple intents. The following subsections provide a detailed explanation of our proposed basic instructions (I_1), pairwise interaction instructions (I_2) and contrastive distinct instructions (I_3).

4.1 Basic Instructions

The basic instructions (I_1) aim to guide the model in extracting the intents and named entities expressed in the utterance. The key components of the basic instructions are as follows:

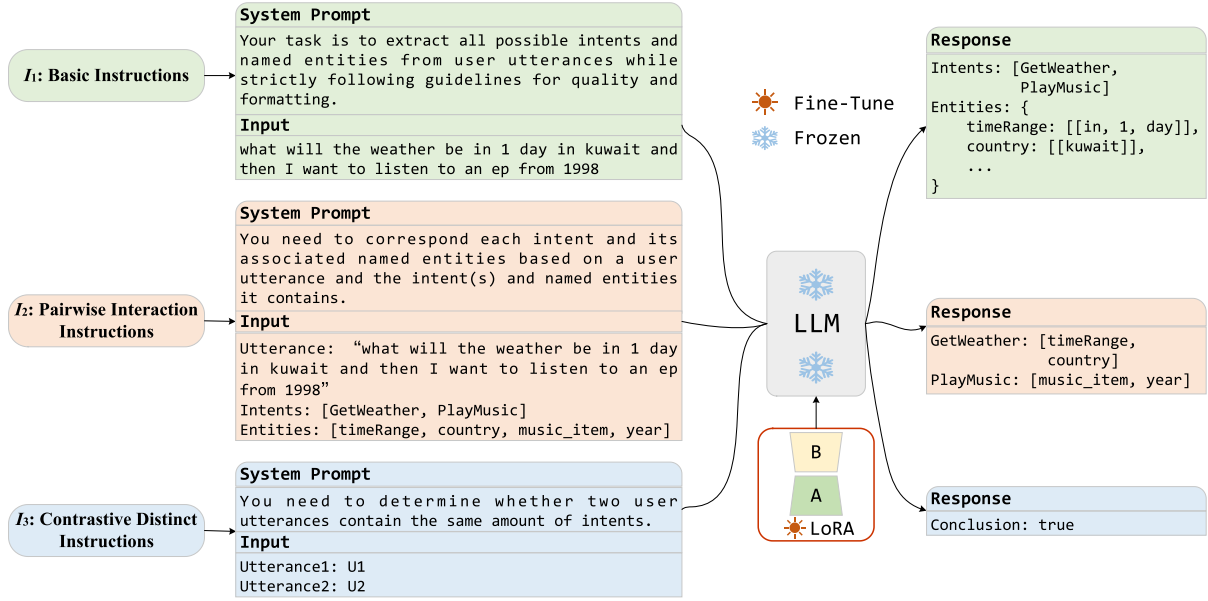


Figure 2: Overview of our framework. Detailed instructions are shown in Appendix A.

[Persona]: You are an expert in multi-intent spoken language understanding. Your task ...

[Instructions]: First, identify the intents in the utterance. The intent options are: {Intent Label Set}. Next, identify the named entities and list each entity with its corresponding words, the entity options are: {Entity Label Set}.

where the persona specifies the model’s role and the tasks to be performed, while the instructions detail the specific steps and requirements. To facilitate result extraction and ensure the controllability of model outputs, the response format for all tasks is standardized to the *JSON* format. It can be formulated as:

$$R = L(SP, I) \quad (1)$$

where SP represents the system prompt, I is the input, L denotes the LLM and R is the response. By converting R into a Python dictionary, we can extract the intents and entities. After obtaining all entities and their corresponding words, inspired by (Wang et al., 2023), we map the words back to their original slot labels using the BIO rule, adhering to the natural left-to-right order of the utterance. This approach allows the LLM to concentrate solely on establishing correspondence between entities and words, disregarding the requirement that the number of final slot labels matches the utterance length. This effectively circumvents the difficulty LLMs face in learning such quantitative correspondences in few-shot scenarios.

4.2 Pairwise Interaction Instructions

To explicitly model dual-task dependencies and reduce relationship confusion, we propose the pairwise interaction instructions (PII). PII is designed to pair each intent with its related entities based on the provided utterance, along with its intent and entity labels. The key components of the PII are as follows:

[Persona]: You are an expert in multi-intent spoken language understanding. You need to ...

[Instructions]: There is a close relationship between each intent and certain named entities. You need to pair them separately.

As shown in Fig. 2, during training, dual-task dependencies are captured by achieving two kinds of alignments. First, in the input part, both the utterance semantics and the labels for the two subtasks are included, achieving a semantic-label alignment for the tasks. Second, dual-task label alignment is established by pairing intent and entity labels in the generation side. With the straightforward mechanism of separate pairing between each intent and its related entities, the mutual dependencies of the two subtasks can be more easily and directly captured by LLMs with their strong few-shot learning capabilities. In addition, it also subtly reduces relational confusions in multi-intent scenarios.

4.3 Contrastive Distinct Instructions

Previous works overlook variations in intent counts among utterances, a factor that aids in

Statistic	FewShotMixATIS					FewShotMixSNIPS				
# K-shot	2-shot	4-shot	6-shot	8-shot	10-shot	2-shot	4-shot	6-shot	8-shot	10-shot
# Original training instances	34	66	100	137	172	14	27	40	54	70
# PICD-Instruct training instances	1,717	6,501	14,950	27,948	44,290	287	1,053	2,380	4,347	7,315
# Testing instances			828					2199		

Table 1: Detail Statistics of FewShotMixATIS and FewShotMixSNIPS.

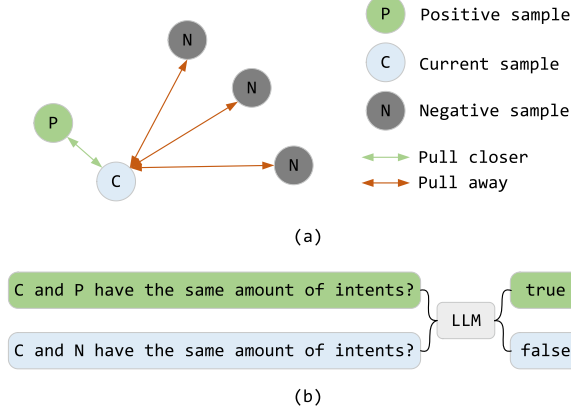


Figure 3: Traditional contrastive learning and our proposed CDI based on instruction learning.

understanding utterances with multiple intents. Inspired by (Xing et al., 2024), such contrastive relationships can be leveraged to enhance the comprehension of utterances and further improve SLU performance. As shown in Fig. 3 (a), traditional contrastive learning aims to optimize representations by pulling similar samples closer in the latent space while pushing dissimilar samples away. To adapt this approach to generative models, we propose straightforward yet effective instructions to implement contrastive learning in the instruction learning paradigm, as shown in Fig. 3 (b). We first sample a positive utterance P and a negative utterance N in relation to the current utterance C. Then we construct instructions to ask the LLM whether C and P, or C and N have the same amount of intents. The expected output is a simple binary response: "true" or "false". The key components of the CDI are as follows:

[Persona]: You are an expert in multi-intent spoken language understanding. You need to ...

[Instructions]: You will be given two user utterances. Each utterance may contain single or multiple intents. You need to judge whether the two utterances contain the same amount of intents.

This approach leverages contrastive relationships to improve the ability of generative LLMs to achieve a deeper understanding of utterances.

4.4 Training and Inference

Training First, an I_3 is constructed for every two samples. Next, an I_1 and an I_2 are created for each sample. To facilitate efficient annotation, GPT-4o¹ is employed to label I_2 . Details of the prompt settings are provided in Appendix B. The shuffled training data is then utilized to train the model in a text-to-text generation form. The training objective is to minimize the negative log-likelihood for each instruction: $\mathcal{L} = -\sum_{n=1}^N \log p(y_n | y_{<n}, I)$. N is the length of the golden output sequence y_1, \dots, y_N and I denotes the current input instruction.

Inference In the inference stage, only I_1 is used to generate predictions for both multiple intent detection and slot filling.

5 Experiments

5.1 Experiment Setup

5.1.1 Dataset

We compare our method with the baselines on two few-shot multi-intent SLU datasets, FewShotMixATIS and FewShotMixSNIPS. They are derived from the MixATIS and MixSNIPS datasets (Qin et al., 2020) using the dynamic sampling algorithm proposed by (Wang et al., 2023). As shown in Table 1, each dataset includes five types of few-shot samples, ranging from 2-shot to 10-shot for training. For testing, we use the test sets of original standard datasets (*i.e.*, MixATIS and MixSNIPS). This setup effectively simulates a realistic application scenario for few-shot multi-intent SLU.

To ensure a balanced number of the three instruction types, oversampling is applied to I_1 and I_2 . The final dataset sizes ranging from 2-shot to 10-shot are presented in the third row of Table 1.

5.1.2 Implementation Details

For PICD-Instruct, we use Qwen2.5-7B² as its backbone model. The model employs AdamW (Loshchilov and Hutter, 2017) as the optimizer with an initial learning rate of 3e-5, a scheduler with a linear warm-up to update and adjust the learning rate. We adopt low-rank adaptation (LoRA)

¹<https://chatgpt.com/>

²<https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

Model	FewShotMixATIS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
BERT	0	57.38	0	4.47	68.37	2.66	12.44	69.54	6.40	25.36	74.23	10.99	36.11	76.66	17.15
RoBERTa	0	48.90	0	0	56.68	0	6.04	65.17	1.33	6.52	68.27	2.17	16.79	70.96	9.18
AGIF+BERT	0	38.28	0	0.60	32.73	0	10.75	48.13	3.02	15.10	38.79	3.50	29.83	56.91	8.94
GL-GIN+BERT	1.21	6.49	0	6.52	21.32	1.57	14.49	32.09	2.90	18.84	33.89	3.26	23.67	49.54	5.56
UGEN	4.47	54.31	1.33	21.98	68.44	6.52	53.50	72.78	15.94	59.30	74.84	19.57	66.67	76.40	22.71
BERT-SIF	30.31	62.51	5.80	37.56	65.74	7.97	58.09	68.20	13.53	61.47	74.90	21.26	62.56	77.61	23.55
ChatGPT	30.07	6.85	0.60	-	-	-	-	-	-	-	-	-	-	-	-
PICD-Instruct	69.57	65.14	18.96	70.29	69.07	21.38	72.71	72.11	24.76	78.86	73.84	27.54	81.28	74.06	27.66

Table 2: Overall results on FewShotMixATIS. I-Acc, S-F1, O-Acc refer to the intent-accuracy, slot F1, and overall accuracy (both intents and slots need to be right), respectively. All models are fine-tuned on the training set of FewShotMixATIS. The version of ChatGPT: gpt-3.5-turbo-16K.

Model	FewShotMixSNIPS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
BERT	4.46	24.84	0.14	3.91	34.59	0	23.78	38.96	0.73	38.06	49.29	3.00	50.34	57.61	4.91
RoBERTa	0.55	8.87	0	1.36	19.04	0	24.51	33.05	0.50	30.38	33.41	0.68	37.79	37.25	0.68
AGIF+BERT	1.27	2.74	0	6.23	7.11	0	17.69	9.12	0.09	21.15	10.03	0.05	14.78	12.53	0.68
GL-GIN+BERT	7.50	0.61	0	14.19	1.48	0	28.06	2.03	0.09	34.20	5.49	0.27	58.21	9.62	0.18
UGEN	2.64	13.10	0	29.65	33.07	0.23	38.84	40.31	1.96	61.57	46.80	4.37	73.08	58.38	7.78
BERT-SIF	37.61	26.29	0.64	56.34	38.32	2.18	64.39	43.34	3.23	65.39	50.18	7.14	74.12	61.75	11.10
ChatGPT	64.48	3.91	0.18	-	-	-	-	-	-	-	-	-	-	-	-
PICD-Instruct	86.45	46.50	5.50	86.77	50.18	7.32	86.99	52.26	8.64	88.18	55.10	10.55	88.09	58.14	11.51

Table 3: Overall results on FewShotMixSNIPS. All models are fine-tuned on the training set of FewShotMixSNIPS.

(Hu et al., 2021) to fine-tune the model with only 55M/28M trainable parameters for FewShotMixATIS/FewShotMixSNIPS. We set the rank to 128/64 for FewShotMixATIS/FewShotMixSNIPS. The batch size is 16 for both datasets. We conduct experiments based on the llamafactory (Zheng et al., 2024) framework to improve the efficiency of implementation. Experiments are conducted on two NVIDIA A5000 GPUs. In multi-intent SLU, accuracy (Acc), F1 score and overall accuracy are used as the metrics for multiple intent detection, slot filling and the SLU semantic frame parsing. Our source code will be released.

5.2 Main Results

We compare our model with BERT (Devlin et al., 2018), RoBERTa (Liu et al., 2019b), ChatGPT, and four other top-performing models. Specifically, AGIF (Qin et al., 2020) presents an adaptive interaction network to achieve fine-grained multiple intent information integration for token-level slot filling. GL-GIN (Qin et al., 2021) introduces a Global-Locally Graph Interaction Network which explores a non-autoregressive model for joint multiple intent detection and slot filling. Wu et al. (2022) proposes a Unified Generative framework (UGEN) based on a prompt-based paradigm and formulates the task as a question-answering problem. BERT-SIF introduces a separate intent-slot interaction framework based on prompt learning to

mitigate relational confusions. The baseline results are sourced from Hua et al. (2024), who implemented the above models using their official code. Due to the limitations of prompt length and costs, the ChatGPT experiment is conducted exclusively in the 2-shot setting. Performance comparisons are presented in Tabel 2 and 3, from which we have the following observations:

(1) *PICD-Instruct achieves new state-of-the-art performance on both datasets.* On the FewShotMixATIS dataset, PICD-Instruct surpasses BERT-SIF in the 2-shot setting by 39.26%, 2.63%, and 13.16% on intent accuracy, slot F1 and overall accuracy, respectively. On the FewShotMixSNIPS dataset, it outperforms BERT-SIF in the 2-shot setting by 48.84%, 20.21% and 4.86% on intent accuracy, slot F1 and overall accuracy. As the amount of training data increases, the performance of our model and all baselines consistently improves across both datasets. This improvement is attributed to our model’s explicit capture of dual-task dependencies via pairwise interaction instructions. The straightforward and effective mechanism significantly reduces training complexity in few-shot scenarios. In addition, our designed contrastive distinct instructions enhance the LLM’s capability to understand utterances with multiple intents. Furthermore, our method of guiding the LLM to generate entities along with their corresponding words effectively mitigates the mismatch between

Model	FewShotMixATIS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
w/o PII, CDI (I_2, I_3)	67.51	64.43	17.75	68.84	68.07	20.65	71.50	70.65	22.83	78.02	72.75	26.69	77.66	73.54	26.81
w/o PII (I_2)	68.24	64.57	17.87	68.96	68.24	20.77	71.62	70.98	22.95	78.26	72.91	26.81	78.14	73.68	27.05
w/o CDI (I_3)	68.48	64.86	18.24	69.20	68.71	21.01	71.98	71.46	23.67	78.50	73.13	27.05	79.23	73.84	27.17
PICD-Instruct	69.57	65.14	18.96	70.29	69.07	21.38	72.71	72.11	24.76	78.86	73.84	27.54	81.28	74.06	27.66
Model	FewShotMixSNIPS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
w/o PII, CDI (I_2, I_3)	84.86	45.14	4.50	85.31	48.27	6.18	86.08	51.16	7.64	86.22	54.31	9.23	86.45	56.25	10.56
w/o PII (I_2)	85.08	45.48	4.64	85.54	48.62	6.41	86.36	51.48	7.82	86.45	54.58	9.64	86.68	56.64	10.83
w/o CDI (I_3)	85.54	46.11	4.96	85.95	49.03	6.82	86.90	51.93	8.05	86.81	54.97	10.14	87.04	57.13	11.28
PICD-Instruct	86.45	46.50	5.50	86.77	50.18	7.32	86.99	52.26	8.64	88.18	55.10	10.55	88.09	58.14	11.51

Table 4: Results of ablation experiments.

the number of slots and the utterance length, a challenge that LLMs typically face when learning quantitative correspondences from a limited amount of annotated data.

(2) *ChatGPT can hardly handle few-shot multi-intent SLU*. The performance of ChatGPT is consistent with recent findings (Pan et al., 2023; Qin et al., 2023). While ChatGPT demonstrates performance comparable to earlier classification-based models in the multiple intent detection task, its performance in slot filling lags far behind other models. We suspect there are two main reasons. First, insufficiently descriptive prompt wording may negatively impact ChatGPT’s performance. We believe advanced in-context learning strategies, such as chain-of-thought prompting, could partially enhance ChatGPT’s performance, while this is beyond the scope of this paper. Second, multi-intent SLU requires task-specific knowledge, which is more effectively acquired through fine-tuning. This finding underscores the need for vertical domain-specific development, particularly for tasks requiring high levels of domain-specific expertise.

5.3 Ablation Study

In this section, we conduct ablation experiments to explore the effect of each component of our PICD-Instruct model. The results are shown in Table 4.

Basic Instructions (BI). Retaining only BI (I_1) still yields significant improvements compared to the previous best-performing model, BERT-SIF, especially in slot filling, where it outperforms ChatGPT. This demonstrates that BI effectively guides the LLM to generate entities along with their corresponding words, simplifying the process of slot filling. Besides, well-crafted instructions fully leverage the few-shot learning capabilities of LLMs, enabling a deeper understanding of the multi-intent SLU task and improving task execution. Detailed instructions are provided in Appendix A.

Pairwise Interaction Instructions (PII). Adding PII (I_2) results in obvious improvements across all metrics and in all few-shot settings. It indicates that PII effectively and explicitly captures the dual-task correlations, leading to substantial performance enhancements. Moreover, PII helps mitigate relational confusions in multi-intent scenarios. The results further verify the fact that a direct and effective interaction mechanism in the instruction learning paradigm is highly beneficial for few-shot learning.

Contrastive Distinct Instructions (CDI). The aim of CDI is to enhance the LLM’s capability to understand utterances with multiple intents. The experimental results reveal that including CDI contributes to improvements in all metrics, verifying its necessity. Besides, combining CDI and PII further enhances the model’s performance. This synergy arises from their individual contributions: CDI and PII excel at their respective tasks, and their integration establishes a strong interdependence. CDI improves the LLM’s initial comprehension of utterances, thereby facilitating multiple intent detection. PII explicitly captures dual-task dependencies, reinforcing the relationship between tasks and enhancing slot filling performance. Therefore, removing any one of CDI and PII leads to performance decreases on all of intent accuracy, slot F1 and overall accuracy.

5.4 Effects of Model Size

To further evaluate the impact of model size on performance, we experiment with 3B, 7B and 14B versions of Qwen2.5 on both datasets. Due to space limitation, we only put results in the 2-shot setting in Table 5, detailed results for other settings are provided in Appendix C. This analysis will help determine whether it is necessary to pursue larger model sizes and understand the trade-offs involved.

As shown in Table 5, the experimental results

Model	FewShotMixATIS			FewShotMixSNIPS		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
Qwen2.5-3B	57.25	57.78	16.55	73.22	36.00	3.32
Qwen2.5-7B	69.57	65.14	18.96	86.45	46.50	5.50
Qwen2.5-14B	71.74	70.04	23.67	88.45	51.12	8.23

Table 5: Results comparison of different model sizes in the 2-shot setting.

Model	FewShotMixATIS			FewShotMixSNIPS		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
w/o PII, CDI (I_2, I_3)	48.03	54.21	8.06	62.26	30.02	0.82
LLaMA3.2-3B	49.52	55.66	9.30	68.62	32.79	2.05
Qwen2.5-3B	57.25	57.78	16.55	73.22	36.00	3.32

Table 6: Results comparison of different model types in the 2-shot setting.

indicate that an increase in Qwen model size leads to improved performance. However, the performance gains in multiple intent detection and slot filling diminish as the model size increases further. For FewShotMixATIS dataset, increasing model parameters from 3B to 7B results in improvements of 12.32% and 7.36% in intent accuracy and slot F1, respectively. However, further increasing parameters from 7B to 14B only yields gains of 2.17% and 4.9% in intent accuracy and slot F1, respectively. A similar trend is observed for the FewShotMixSNIPS dataset, although overall accuracy shows more pronounced improvements when parameters are scaled from 7B to 14B. This suggests that the overall reasoning capability of the LLM improves significantly with increased model size. Consequently, pursuing larger-scale language models may not be essential for achieving substantial performance gains across all metrics in the context of multi-intent SLU.

5.5 Effects of Model Type

To investigate the effectiveness of different model types, we compare the latest versions of LLaMA³ and Qwen. Due to space limitation, only results in the 2-shot setting are presented in Table 6, while detailed results for other settings are included in Appendix D.

As shown in Table 6, the results reveal that Qwen outperforms LLaMA in terms of all metrics. Especially in multiple intent detection, Qwen overpasses LLaMA by 7.73% and 4.6% on FewShotMixATIS and FewShotMixSNIPS, respectively. A possible explanation for this performance gap lies in their foundational capabilities. While LLaMA is primarily trained on English corpora, Qwen excels

³<https://huggingface.co/meta-llama>

Model	FewShotMixATIS	FewShotMixSNIPS
LLaMA3.2-3B	1.33	2.36
Qwen2.5-3B	0.24	0.09

Table 7: Error rate of *JSON* parsing in the 2-shot setting.

in both Chinese and English, potentially allowing it to learn more diverse language patterns during pre-training, which could benefit multi-intent SLU. Another noteworthy observation is the disparity in their *JSON* output format capabilities. As shown in Table 7, Qwen exhibits superior *JSON* output capabilities compared to LLaMA, likely due to its tailored post-training process for generating structured outputs. Specific parsing error analyses are provided in Appendix D. Despite inferior performances of LLaMA, our proposed instructions still demonstrate their effectiveness in few-shot multi-intent SLU. Notably, removing PII and CDI results in significant performance declines across all metrics. This analysis underscores the critical importance of model selection, particularly with respect to capabilities relevant to the task at hand.

6 Conclusion

In this paper, we make in-depth investigations of few-shot multi-intent SLU. We propose PICD-Instruct, a framework designed to address the challenges of generative few-shot multi-intent SLU from three key perspectives. Firstly, we propose basic instructions to tackle the mismatch between generated slot counts and input length. Secondly, we introduce pairwise interaction instructions to explicitly model dual-task dependencies while minimizing relational confusions in multi-intent scenarios. Thirdly, we present contrastive distinct instructions that leverage contrastive relations in intent counts to enhance understanding. Experimental results demonstrate that our proposed model achieves SOTA performance on FewShotMixATIS and FewShotMixSNIPS, thereby highlighting our model’s robust generalization capabilities in a simulated real-world application scenario.

7 Limitations

This paper presents a comprehensive analysis of generative few-shot multi-intent SLU and introduces the PICD-Instruct model, which is based on the paradigm of instruction learning. In fact, detailed descriptions of intent and slot labels could significantly enhance LLMs’ comprehension of multi-intent SLU, as high-quality external knowl-

edge helps mitigate the hallucination issue in LLMs (Wan et al., 2024). In the future, we will explore how to integrate external label knowledge into LLMs to further improve the performance of few-shot multi-intent SLU.

References

Soyeon Caren Han, Siqu Long, Henry Weld, and Josiah Poon. 2022. Spoken language understanding for conversational ai: Recent advances and future direction. *arXiv e-prints*, pages arXiv–2212.

Lisong Chen, Peilin Zhou, and Yuexian Zou. 2022. Joint multiple intent detection and slot filling via self-distillation. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7612–7616. IEEE.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Rashmi Gangadharaiah and Balakrishnan Narayanaswamy. 2019. Joint multiple intent detection and slot labeling for goal-oriented dialog. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 564–569, Minneapolis, Minnesota. Association for Computational Linguistics.

Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu, and Yun-Nung Chen. 2018. Slot-gated modeling for joint slot filling and intent prediction. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 753–757, New Orleans, Louisiana. Association for Computational Linguistics.

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.

Wenbin Hua, Yufan Wang, Rui Fan, Xinhui Tu, and Tingting He. 2024. Unraveling intricacies: A decomposition approach for few-shot multi-intent spoken language understanding. In *2024 IEEE International Conference on Big Data (BigData)*, pages 918–927. IEEE.

Chandra Irugalbandara. 2024. Meaning typed prompting: A technique for efficient, reliable structured output generation. *arXiv preprint arXiv:2410.18146*.

Byeongchang Kim, Seonghan Ryu, and Gary Geunbae Lee. 2017. Two-stage multi-intent detection for spoken language understanding. *Multimedia Tools and Applications*, 76:11377–11390.

Changliang Li, Liang Li, and Ji Qi. 2018. A self-attentive model with gate mechanism for spoken language understanding. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3824–3833.

Yinghao Li, Rampi Ramprasad, and Chao Zhang. 2024. A simple but effective approach to improve structured language model output for information extraction. *arXiv preprint arXiv:2402.13364*.

Yijin Liu, Fandong Meng, Jinchao Zhang, Jie Zhou, Yufeng Chen, and Jinan Xu. 2019a. CM-net: A novel collaborative memory network for spoken language understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1051–1060, Hong Kong, China. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*.

Ilya Loshchilov and Frank Hutter. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*.

Renze Lou, Kai Zhang, and Wenpeng Yin. 2024. Large language model instruction following: A survey of progresses and challenges. *Computational Linguistics*, pages 1–10.

Jie Mei, Yufan Wang, Xinhui Tu, Ming Dong, and Tingting He. 2023. Incorporating bert with probability-aware gate for spoken language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:826–834.

Peiqing Niu, Zhongfu Chen, Meina Song, et al. 2019. A novel bi-directional interrelated model for joint intent detection and slot filling. *arXiv preprint arXiv:1907.00390*.

Wenbo Pan, Qiguang Chen, Xiao Xu, Wanxiang Che, and Libo Qin. 2023. A preliminary evaluation of chatgpt for zero-shot dialogue understanding. *arXiv preprint arXiv:2304.04256*.

Chengwei Qin, Aston Zhang, Zhuosheng Zhang, Jiaao Chen, Michihiro Yasunaga, and Diyi Yang. 2023. Is ChatGPT a general-purpose natural language processing task solver? In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1339–1384, Singapore. Association for Computational Linguistics.

Libo Qin, Fuxuan Wei, Tianbao Xie, Xiao Xu, Wanxiang Che, and Ting Liu. 2021. GL-GIN: Fast and accurate non-autoregressive model for joint multiple intent detection and slot filling. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International*

711	<i>Joint Conference on Natural Language Processing</i>	Bowen Xing and Ivor Tsang. 2022a. Co-guiding net:	766
712	(<i>Volume 1: Long Papers</i>), pages 178–188, Online.	Achieving mutual guidances between multiple intent	767
713	Association for Computational Linguistics.	detection and slot filling via heterogeneous semantics-	768
714	Libo Qin, Xiao Xu, Wanxiang Che, and Ting Liu. 2020.	label graphs . In <i>Proceedings of the 2022 Conference</i>	769
715	AGIF: An adaptive graph-interactive framework for	<i>on Empirical Methods in Natural Language Process-</i>	770
716	joint multiple intent detection and slot filling . In	<i>ing</i> , pages 159–169, Abu Dhabi, United Arab Emi-	771
717	<i>Findings of the Association for Computational Lin-</i>	rates. Association for Computational Linguistics.	772
718	<i>guistics: EMNLP 2020</i> , pages 1807–1816, Online.		
719	Association for Computational Linguistics.	Bowen Xing and Ivor Tsang. 2022b. Group is better	773
720	Abdulfattah Safa, Tamta Kapanadze, Arda Uzunoğlu,	than individual: Exploiting label topologies and label	774
721	and Gözde Gül Şahin. 2024. A systematic survey	relations for joint multiple intent detection and slot	775
722	on instructional text: From representation and down-	filling . In <i>Proceedings of the 2022 Conference on</i>	776
723	stream nlp tasks. <i>arXiv preprint arXiv:2410.18529</i> .	<i>Empirical Methods in Natural Language Processing</i> ,	777
724	Michael Saxon, Samridhi Choudhary, Joseph P	pages 3964–3975, Abu Dhabi, United Arab Emirates.	778
725	McKenna, and Athanasios Mouchtaris. 2021. End-to-	Association for Computational Linguistics.	779
726	end spoken language understanding for generalized		
727	voice assistants. <i>arXiv preprint arXiv:2106.09009</i> .	Steve Young, Milica Gašić, Blaise Thomson, and Ja-	780
728	Feifan Song, Lianzhe Huang, and Houfeng Wang. 2024.	son D. Williams. 2013. Pomdp-based statistical spo-	781
729	A unified framework for multi-intent spoken lan-	ken dialog systems: A review . <i>Proceedings of the</i>	782
730	guage understanding with prompting. In <i>ICASSP</i>	<i>IEEE</i> , 101(5):1160–1179.	783
731	<i>2024-2024 IEEE International Conference on Acous-</i>		
732	<i>tics, Speech and Signal Processing (ICASSP)</i> , pages	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan	784
733	9966–9970. IEEE.	Ye, and Zheyang Luo. 2024. LlamaFactory: Unified	785
734	Mengxiao Song, Bowen Yu, Li Quangang, Wang Yu-	efficient fine-tuning of 100+ language models . In	786
735	bin, Tingwen Liu, and Hongbo Xu. 2022. Enhancing	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	787
736	joint multiple intent detection and slot filling with	<i>sociation for Computational Linguistics (Volume 3:</i>	788
737	global intent-slot co-occurrence. In <i>Proceedings of</i>	<i>System Demonstrations)</i> , pages 400–410, Bangkok,	789
738	<i>the 2022 Conference on Empirical Methods in Natu-</i>	Thailand. Association for Computational Linguistics.	790
739	<i>ral Language Processing</i> , pages 7967–7977.		
740	Fanqi Wan, Xinting Huang, Leyang Cui, Xiaojun Quan,	Zhihong Zhu, Xuxin Cheng, Hao An, Zhichang Wang,	791
741	Wei Bi, and Shuming Shi. 2024. Knowledge verifica-	Dongsheng Chen, and Zhiqi Huang. 2024. Zero-shot	792
742	tion to nip hallucination in the bud . In <i>Proceedings</i>	spoken language understanding via large language	793
743	<i>of the 2024 Conference on Empirical Methods in</i>	models: A preliminary study . In <i>Proceedings of the</i>	794
744	<i>Natural Language Processing</i> , pages 2616–2633, Mi-	<i>2024 Joint International Conference on Computa-</i>	795
745	ami, Florida, USA. Association for Computational	<i>tional Linguistics, Language Resources and Eval-</i>	796
746	Linguistics.	<i>uation (LREC-COLING 2024)</i> , pages 17877–17883,	797
747	Yufan Wang, Jie Mei, Bowei Zou, Rui Fan, Tingting He,	Torino, Italia. ELRA and ICCL.	798
748	and Ai Ti Aw. 2023. Making pre-trained language	Zhihong Zhu, Xuxin Cheng, Zhiqi Huang, Dongsheng	799
749	models better learn few-shot spoken language un-	Chen, and Yuexian Zou. 2023. Towards unified spo-	800
750	derstanding in more practical scenarios . In <i>Findings of</i>	ken language understanding decoding via label-aware	801
751	<i>the Association for Computational Linguistics: ACL</i>	compact linguistics representations . In <i>Findings of</i>	802
752	<i>2023</i> , pages 13508–13523, Toronto, Canada. Associ-	<i>the Association for Computational Linguistics: ACL</i>	803
753	ation for Computational Linguistics.	<i>2023</i> , pages 12523–12531, Toronto, Canada. Associ-	804
754	Yangjun Wu, Han Wang, Dongxiang Zhang, Gang Chen,	ation for Computational Linguistics.	805
755	and Hao Zhang. 2022. Incorporating instructional		
756	prompts into a unified generative framework for joint		
757	multiple intent detection and slot filling. In <i>Proceed-</i>		
758	<i>ings of the 29th International Conference on Compu-</i>		
759	<i>tational Linguistics</i> , pages 7203–7208.		
760	Bowen Xing, Lizi Liao, Minlie Huang, and Ivor Tsang.		
761	2024. Dc-instruct: An effective framework for gen-		
762	erative multi-intent spoken language understanding.		
763	In <i>Proceedings of the 2024 Conference on Empiri-</i>		
764	<i>cal Methods in Natural Language Processing</i> , pages		
765	14520–14534.		

System Prompt
<pre>{ Persona: "You are an expert in multi-intent spoken language understanding. Your task is to extract all possible intents and named entities from user utterances while strictly following guidelines for quality and formatting." Instructions: ["You will be given a user utterance", "Let's think step by step. First, identify the intents in the utterance. The intent options are: {Intent Label Set}." "Next, identify the named entities in the utterance. The named entity options are: {Entity Label Set}." "If an entity appears multiple times in the utterance, list all the words that belong to the entity.", "Make sure not to output any extra content."], OutputFormat: "{Intents: [intent1, intent2], Entities: {entity1: [[word1, word2], [word3, word4]], entity2: [[word5]]}}", Example: "{Utterance: ...}\n{Intents: ..., Entities: ...}" }</pre>
Utterance
<pre>{ Utterance: "what will the weather be in 1 day in kuwait and then I want to listen to an ep from 1998" }</pre>
Response
<pre>{ Intents: [GetWeather, PlayMusic], Entities: { timeRange: [[in, 1, day]], country: [[kuwait]], music_item: [[ep]], year: [[1998]] } }</pre>

I_1 : Basic Instructions

Figure 4: Details of BI (I_1).

System Prompt
<pre>{ Persona: "You are an expert in multi-intent spoken language understanding. You need to correspond each intent and its associated named entities based on a user utterance and the intent(s) and named entities it contains." Instructions: ["You will be given a user utterance with its intents and named entities.", "There is a close relationship between each intent and certain named entities.", "You need to pair them separately in the specified format.", "Make sure not to output any extra content."], OutputFormat: "{Intent1: [entity1], Intent2: [entity2, entity3]}", Example: "{Utterance: ...}\n{Intents: ..., Entities: ...}\n{Intent1: [...], Intent2: [...]}" }</pre>
Utterance
<pre>{ Utterance: "what will the weather be in 1 day in kuwait and then I want to listen to an ep from 1998", Intents: [GetWeather, PlayMusic], Entities: [timeRange, country, music_item, year] }</pre>
Response
<pre>{ GetWeather: [timeRange, country], PlayMusic: [music_item, year] }</pre>

I_2 : Pairwise Interaction Instructions

Figure 5: Details of PII (I_2).

A The Detailed Instructions

This section presents the detailed instructions for BI, PII, and CDI, as illustrated in Figs. 4, 5, and 6, respectively.

B The Prompt Used by GPT-4o

To ensure efficient annotation, we employ GPT-4o to label I_2 , with the corresponding prompt illustrated in Fig. 7. First, we define GPT-4o’s role and provide an example annotation. Next, we introduce a labeling technique designed to improve the quality of the annotations. Finally, we specify the output format.

System Prompt
<pre>{ Persona: "You are an expert in multi-intent spoken language understanding. You need to determine whether two user utterance contain the same amount of intents." Instructions: ["You will be given two user utterances.", "Each utterance may contain single or multiple intents.", "You need to judge whether the two utterances contain the same amount of intents.", "Make sure not to output any extra content."], OutputFormat: "{Conclusion: true}", Example: "{Utterance1: ..., Utterance2: ...}\n{Conclusion: ...}" }</pre>
Utterance
<pre>{ Utterance1: U1, Utterance2: U2 }</pre>
Response
<pre>{ Conclusion: true }</pre>

I_3 : Contrastive Distinct Instructions

Figure 6: Details of CDI (I_3).

你是一个多意图口语理解的专家。现在的任务是根据给定的句子，以及句子包含的意图和实体，将每个意图与它相关联的实体进行配对。下面是一个例子：输入：{ "Utterance": "how much does it cost to rent a car in tacoma and then what's restriction ap68", "Intents": ["atis_ground_fare", "atis_restriction"], "Entities": ["transport_type", "city_name", "restriction_code"] } 输出：{ "atis_ground_fare": ["transport_type", "city_name"], "atis_restriction": ["restriction_code"] } 标注技巧：一条句子如果包含多个意图，那么这条句子可以被逗号或者'and'分隔成多个子句（注意：有的逗号或者'and'可能并不是子句分隔符，需要你按句意判断），每个子句对应一个意图。所以你可以按顺序遍历意图列表，将每一个子句的实体与这个子句的意图联系起来，这样准确率会比较高，注意不要遗漏任何实体。下面每一次我会给你一条case，请你给出标注，不要输出你的思考过程，只输出单行的json代码块结果就行（方便我直接复制）

Figure 7: The prompt used by GPT-4o.

C The Detailed Experimental Results for Model Size

This section presents the detailed experimental results for three parameter sizes across all few-shot settings. As shown in Table 8, performance improves with an increase in model size. Consistent with the findings in Section 5.4, performance gains for most metrics diminish as the model size continues to increase. Therefore, it is crucial to consider both model size and performance together, especially in scenarios with limited computational resources.

D The Detailed Experimental Results for Model Type

This section provides a comprehensive analysis of the experimental results for two model types across all few-shot settings. As shown in Table 9, Qwen surpasses LLaMA significantly on most of the metrics. This disparity can be attributed partly to differences in their foundational capabilities and partly to variations in their ability to handle JSON

Model	FewShotMixATIS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
Qwen2.5-3B	57.25	57.78	16.55	64.86	63.06	17.27	63.89	64.95	20.17	69.57	67.16	21.38	72.83	68.26	21.50
Qwen2.5-7B	69.57	65.14	18.96	70.29	69.07	21.38	72.71	72.11	24.76	78.86	73.84	27.54	81.28	74.06	27.66
Qwen2.5-14B	71.74	70.04	23.67	78.38	70.77	24.76	78.86	72.14	25.36	80.92	75.38	30.68	77.17	76.16	29.71
Model	FewShotMixSNIPS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
Qwen2.5-3B	73.22	36.00	3.32	74.90	40.71	4.14	79.04	41.51	4.73	81.08	45.21	6.23	82.36	46.53	7.23
Qwen2.5-7B	86.45	46.50	5.50	86.77	50.18	7.32	86.99	52.26	8.64	88.18	55.10	10.55	88.09	58.14	11.51
Qwen2.5-14B	88.45	51.12	8.23	86.77	56.65	9.00	88.49	57.50	11.41	91.27	61.58	13.78	90.81	63.02	14.51

Table 8: Results comparison of different model sizes on FewShotMixATIS and FewShotMixSNIPS.

Model	FewShotMixATIS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
LLaMA3.2-3B	49.52	55.66	9.30	56.64	59.68	11.23	58.21	61.55	12.08	56.76	62.98	15.10	67.63	68.92	18.96
Qwen2.5-3B	57.25	57.78	16.55	64.86	63.06	17.27	63.89	64.95	20.17	69.57	67.16	21.38	72.83	68.26	21.50
Model	FewShotMixSNIPS														
	2-shot			4-shot			6-shot			8-shot			10-shot		
	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc	I-Acc	S-F1	O-Acc
LLaMA3.2-3B	68.62	32.79	2.05	69.40	37.31	3.05	68.49	41.08	4.09	77.67	45.12	6.37	81.95	48.54	7.19
Qwen2.5-3B	73.22	36.00	3.32	74.90	40.71	4.14	79.04	41.51	4.73	81.08	45.21	6.23	82.36	46.53	7.23

Table 9: Results comparison of different model types on FewShotMixATIS and FewShotMixSNIPS.

Model	FewShotMixATIS				
	2-shot	4-shot	6-shot	8-shot	10-shot
LLaMA3.2-3B	1.33	0.97	1.33	0.36	0.24
Qwen2.5-3B	0.24	0.12	0.24	0.24	0.24
Model	FewShotMixSNIPS				
	2-shot	4-shot	6-shot	8-shot	10-shot
LLaMA3.2-3B	2.36	1.23	0.68	0.36	0.59
Qwen2.5-3B	0.09	0.27	0.18	0.09	0.05

Table 10: Error rate of *JSON* parsing on FewShotMix-ATIS and FewShotMixSNIPS.

output formats. Table 10 highlights that Qwen exhibits a lower *JSON* parsing error rate compared to LLaMA, a result attributed to its specialized post-training process designed for generating structured outputs, as documented in the official source⁴. Specifically, LLMs frequently generate content such as "*Cutting Knowledge Date: December 2023 Today Date: ...*", where the ellipsis represents the original input, often resulting in errors during *JSON* parsing. This observation underscores that even fine-tuned LLMs can produce unexpected content, emphasizing the critical importance of selecting models with robust controllable generation capabilities.

⁴<https://huggingface.co/Qwen/Qwen2.5-3B-Instruct>