

Synergizing Large Language Models and Pre-Trained Smaller Models for Conversational Intent Discovery

Anonymous ACL submission

Abstract

In Conversational Intent Discovery (CID), Small Language Models (SLMs) struggle with overfitting to familiar intents and fail to label newly discovered ones. This issue stems from their limited grasp of semantic nuances and their intrinsically discriminative framework. Therefore, we propose **Synergizing Large Language Models (LLMs) with pre-trained SLMs for CID (SynCID)**. It harnesses the profound semantic comprehension of LLMs alongside the operational agility of SLMs. By utilizing LLMs to refine both utterances and existing intent labels, SynCID significantly enhances the semantic depth, subsequently realigning these enriched descriptors within the SLMs’ feature space to correct cluster distortion and promote robust learning of representations. A key advantage is its capacity for the early identification of new intents, a critical aspect for deploying conversational agents successfully. Additionally, SynCID leverages the in-context learning strengths of LLMs to generate labels for new intents. Thorough evaluations across a wide array of datasets have demonstrated its superior performance over traditional CID methods.¹

1 Introduction

Recognizing user intents within conversational utterances is pivotal for developing intelligent conversational agents (Yilmaz and Toraman, 2020; Shen et al., 2021; Gung et al., 2023). Previous research mainly formulates this problem as a close-world intent classification task (Zhang et al., 2022a; Yehudai et al., 2023). However, in real-world applications, new intents continuously emerge. This spurs increasing interest in the open-world Conversational Intent Discovery (CID) (Zhang et al., 2021c, 2022b; Liang and Liao, 2023), a task that aims to recognize both known and new intents from extensive or even limited amount of user utterances.

¹<https://anonymous.4open.science/r/SynCID-3121>

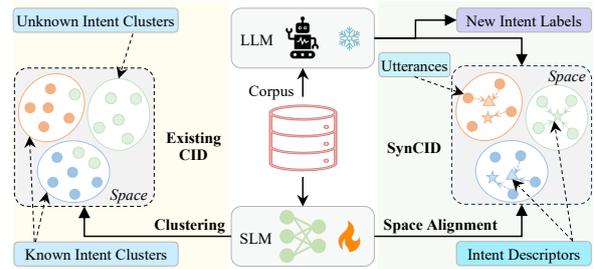


Figure 1: Existing methods rely on SLMs to cluster intents (Left), while our SynCID effectively synergizes LLMs and SLMs via space alignment (Right).

Current attempts at CID primarily rely on pre-trained Small Language Models (SLMs), which fall into two main categories: unsupervised and semi-supervised. Unsupervised methods (Padmasundari and Bangalore, 2018; Shi et al., 2018) firstly train SLMs without using any labeled data to obtain utterance representations, and then cluster them to infer intents. In contrast, semi-supervised methods (Lin et al., 2020; Zhang et al., 2021c; Zhou et al., 2023) leverage the available labeled data for the initial pre-training of SLMs, followed by fine-tuning these models with pseudo supervisory signals on unlabeled utterances for intent recognition. Given the specialized agility of SLMs, these methods can easily fit user utterances and learn discriminative representations for CID.

However, two key challenges persist. The first is *overfitting to known intents*, where these methods struggle to capture the full scope of intents and accurately model known label semantics. This limitation not only biases them towards existing intent categories but also compromises their ability to detect new intents early, a crucial capability for adaptive conversational agents. The second challenge is the *inability to label novel intents*, due to the inherently discriminative architecture of CID models, which falls short in recognizing and labeling emerging intents, marking a critical adaptability gap in current approaches.

069 Recently, Large Language Models (LLMs)
070 (Brown et al., 2020; Chowdhery et al., 2023; Ope-
071 nAI, 2023) have achieved significant breakthroughs
072 in language understanding and generative tasks, in-
073 cluding summarization (Liu et al., 2023) and query
074 rewriting (Anand et al., 2023). Their success in-
075 spires a potential solution for addressing the above
076 challenges by adapting LLMs to enhance intent dis-
077 covery. Yet, the context length limitation of LLMs
078 restricts their direct use in CID, which requires clus-
079 tering thousands of utterances. While integrating
080 user utterances with task-specific prompts to solicit
081 intent labels from LLMs is possible, this prompt-
082 ing method risks generating intent labels without
083 sufficient control, thus leading to unpredictable and
084 uninformative outcomes (Sun et al., 2023).

085 To navigate these challenges while leveraging
086 the strengths of both LLMs and SLMs, we intro-
087 duce SynCID, a framework that synergizes the
088 deep semantic insights of LLMs with the agile,
089 specialized capabilities of SLMs. SynCID employs
090 a dual-prompting mechanism with LLMs to refine
091 both utterances and known intent labels, enhancing
092 the semantic precision of intent descriptors. This
093 refinement process, informed by the nuanced un-
094 derstanding of LLMs, not only clarifies the intent
095 representation but also primes the data for more ef-
096 fective learning. Following this, SLMs are trained
097 through contrastive learning to align the seman-
098 tic spaces of utterances with those of the intent
099 descriptors. This innovative alignment strategy sig-
100 nificantly reduces cluster distortion and improves
101 the system’s ability to detect and label new intents
102 early, addressing the primary limitations of current
103 CID approaches. By selecting a limited number
104 of close-to-center utterances from newly formed
105 intent clusters for in-context learning with LLMs,
106 SynCID achieves precise intent labeling.

107 In summary, our contributions are threefold:

- 108 • We propose SynCID, an effective framework that
109 synergizes powerful LLMs with agile SLMs to
110 identify novel user intents and generate corre-
111 sponding intent labels.
- 112 • We introduce a space alignment schema to align
113 the representation spaces of utterances and the
114 intent descriptors, significantly reducing the risk
115 of overfitting to known intents.
- 116 • Experiments show that SynCID not only outper-
117 forms existing CID methods, but also provides
118 labels for new intent clusters and enables early
119 intent detection.

2 Related Work 120

2.1 Conversational Intent Discovery 121

Unsupervised Methods: Early unsupervised meth-
122 ods (Cheung and Li, 2012; Li et al., 2013) primarily
123 extracted statistical features from unlabeled data to
124 cluster queries with similar intents. Later studies
125 (Xie et al., 2016; Yang et al., 2017; Padmasundari
126 and Bangalore, 2018; Caron et al., 2018; Shi et al.,
127 2018; Hadifar et al., 2019) leveraged deep neural
128 networks to learn robust representations for cluster-
129 ing. More recently, the development of LLMs has
130 facilitated their application in unsupervised intent
131 recognition (De Raedt et al., 2023). Despite the
132 progress, none of these unsupervised CID meth-
133 ods can fully harness supervised signals in learning
134 representations and clustering intents. 135

Semi-supervised Methods: Addressing this limi-
136 tation, semi-supervised methods (Hsu et al., 2018,
137 2019; Han et al., 2019; Lin et al., 2020) focus on
138 integrating limited labeled data with extensive un-
139 labeled data to enhance intent identification. For
140 example, Hsu et al. (2018) transferred prior knowl-
141 edge for clustering via predicting pairwise similar-
142 ities. Further, several semi-supervised CID methods
143 (Zhang et al., 2021b,c; Wei et al., 2022; Zhang
144 et al., 2023; Zhou et al., 2023; Mou et al., 2023)
145 formulated a two-stage schema for CID, which in-
146 volves initially pre-training a base SLM and then
147 iteratively fine-tuning it. This schema significantly
148 enhanced CID by utilizing pseudo supervisory sig-
149 nals from the pre-trained SLM. Yet, it often faces
150 issues related to the quality of these pseudo super-
151 visory signals. Thus, there are also efforts (Mou
152 et al., 2022a,b; Zhang et al., 2022b) refined learning
153 objectives, such as contrastive learning, to learn dis-
154 criminative representations for discerning intents.
155 However, challenges persist in comprehensively
156 grasping the nuanced semantics of both utterances
157 and known intent labels, as well as generating new
158 intent labels, which are addressed by our SynCID
159 by synergizing LLMs and SLMs for CID. 160

2.2 The Synergy Between LLMs and SLMs 161

The emergence of LLMs has recently revolution-
162 ized various NLP tasks (Chowdhery et al., 2023;
163 Black et al., 2022; Touvron et al., 2023), spurring
164 research into their synergy with SLMs for boost-
165 ing performance of small task-specific models. A
166 promising direction in this synergy is using LLMs
167 to create new and high-quality data for training
168 downstream SLMs, enabling them to achieve com-
169

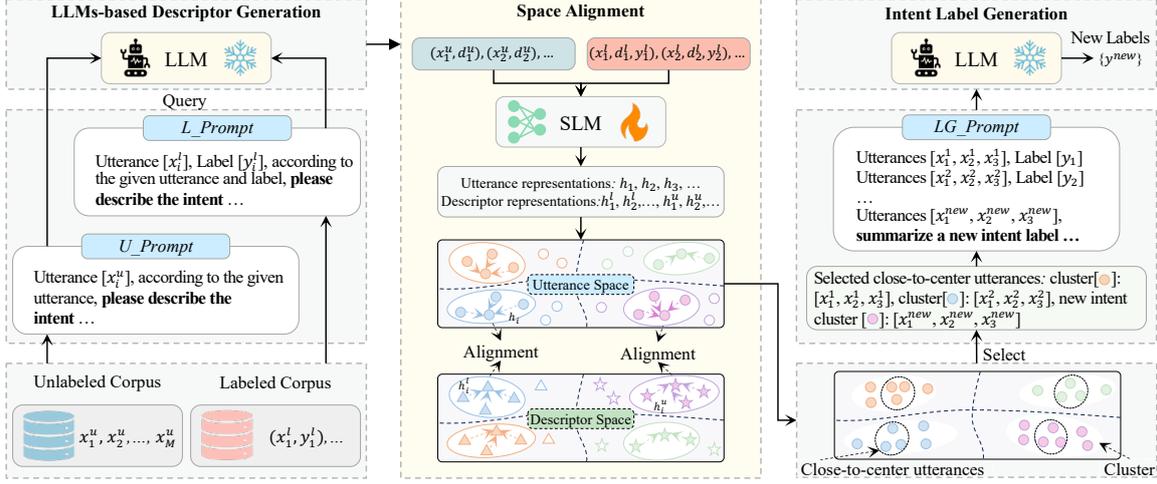


Figure 2: An overview of the proposed SyncID framework. It consists of three stages: LLMs-based Descriptor Generation, Space Alignment, and Intent Label Generation.

petitive performance. (Yang et al., 2020; Ding et al., 2023; Wei and Zou, 2019; Xie et al., 2020). Yet, such a method in CID risks unintentionally altering the semantic meanings of utterances or introducing noise, challenging accurate intent recognition.

Another effective method for synergizing LLMs and SLMs involves distilling task-specific knowledge from LLMs. Wang et al. (2021) showed the potential of GPT-3 as a cost-effective alternative to human labeling. Moreover, researchers like Li et al. (2022), Shridhar et al. (2023), and Hsieh et al. (2023) have utilized LLMs to generate task-specific labels and detailed explanations, facilitating the training of SLMs for reasoning tasks. Nevertheless, all these methods predominantly rely on either using a finite set of labels for annotating data or training generative models for aligning the knowledge from LLMs, which are not applicable in the CID.

In this work, we further the synergy to enhance intent discovery, leveraging a novel space alignment to align LLMs’ comprehensive insights with the agility of SLMs and enabling early detection.

3 The SyncID Framework

3.1 Problem Formulation

Here, we study the CID problem as follows: Let \mathcal{I}_k and \mathcal{I}_{uk} represent the sets of known and unknown intents respectively, where $\{\mathcal{I}_k \cap \mathcal{I}_{uk}\} = \emptyset$ and $|\mathcal{I}_k| + |\mathcal{I}_{uk}| = K$. Here K is the total number of the user intents within the dataset. A typical CID task comprises a set of labeled utterance-intent pairs $\mathcal{D}^l = \{(x_i, y_i)\}_{i=1}^N$, wherein each intent $y_i \in \mathcal{I}_k$, and a set of unlabeled utterances $\mathcal{D}^u = \{(x_i)\}_{i=1}^M$, where the intent of each utterance x_i belongs to

$\{\mathcal{I}_k \cup \mathcal{I}_{uk}\}$. The CID task is to learn a SLM \mathcal{M} to recognize all unknown intents \mathcal{I}_{uk} within \mathcal{D}^u and perform accurate clustering to classify each $x_i \in \{\mathcal{D}^l \cup \mathcal{D}^u\}$ into its corresponding intent category.

3.2 Model Overview

Figure 2 depicts an overview of the proposed SyncID framework for CID. It comprises three stages: **LLMs-based Descriptor Generation** (§3.3) for generating accurate and contextually rich intent descriptors, **Space Alignment** (§3.4) for aligning the representation spaces of utterances and intent descriptors to facilitate the synergy between LLMs and SLMs, and **Intent Label Generation** (§3.5) for producing labels for new intent clusters. We detail these stages in the subsequent subsections.

3.3 LLMs-based Descriptor Generation

This stage aims to leverage LLMs to recapitulate utterances and known intent labels into concise, accurate intent descriptors, eliminating irrelevant content in utterances while enriching the semantics of known intent labels. To achieve this, we develop two prompt templates: U_Prompt and L_Prompt, designed to guide the generation of these descriptors. As illustrated in Figure 2, U_Prompt is constructed as (x_i, p_u) , prompting LLMs to generate descriptors related to the utterances as follows:

$$d_i^u = \text{LLM}(x_i, p_u), \quad (1)$$

where each $x_i \in \{\mathcal{D}^l \cup \mathcal{D}^u\}$ is a user utterance, and p_u denotes the prompt tokens. Similarly, L_Prompt is defined as (x_i, y_i, p_l) for the generation of label-enriched intent descriptors:

$$d_i^l = \text{LLM}(x_i, y_i, p_l), \quad (2)$$

where (x_i, y_i) is an utterance-intent pair in \mathcal{D}^l , and p_l refers to the respective prompt tokens. Crucially, in Equation 2, we integrate each intent label y_i with its corresponding utterance x_i to prompt LLMs for descriptor generation, enhancing the semantics of known user intents. After prompting LLMs to generate corresponding intent descriptors for all utterances and known user intents, we then utilize them to perform space alignment, facilitating the synergy of LLMs and SLMs for recognizing intents. For clarity, we formally redefine the training datasets \mathcal{D}^l and \mathcal{D}^u as follows: $\mathcal{D}^l = \{(x_i, d_i^u, d_i^l, y_i)\}_{i=1}^N$ and $\mathcal{D}^u = \{(x_i, d_i^u)\}_{i=1}^M$. It’s noteworthy that we curate the aforementioned prompt templates without deliberation for better generalization.

3.4 Intent Discovery with Space Alignment

Given the intent descriptors from LLMs, we propose Space Alignment (SA) to synergize LLMs and SLMs for intent recognition. It comprises two sub-strategies: (1) SA with Contrastive Learning, which directly aligns the semantic spaces of utterances and intent descriptors, fostering robust utterance representation learning. (2) SA with Neighbor Filtering, which utilizes intent descriptors to refine neighborhood relationships between utterances, filtering out noise and promoting the formation of compact intent clusters.

SA with Contrastive Learning. Utilizing LLMs’ strength in understanding and generation, we derive intent descriptors that offer more reliable and enriched insights into user intents. To effectively synergize LLMs and SLMs, we align the semantic spaces of utterances and LLM-generated intent descriptors via two contrastive learning objectives. Given the specialized agility of SLMs, this alignment can adeptly fit them into LLMs’ insights, mitigating cluster distortion and enhancing the identification of new intents. Specifically, given a general pre-trained SLMs based CID model \mathcal{M} , we initially extract representations x_i and d_i^u for each utterance x_i and its corresponding intent descriptor d_i^u . Since d_i^u is derived from x_i using LLMs, x_i and d_i^u naturally form a positive pair. Following Gao et al. (2021), we compute an unsupervised contrastive loss between x_i and d_i^u as follows:

$$\mathcal{L}^{ucl} = -\log \frac{e^{\text{sim}(x_i, d_i^u)/\tau_1}}{\sum \mathbb{1}_{[k \neq i]} e^{\text{sim}(x_i, d_k^u)/\tau_1}}, \quad (3)$$

where $\text{sim}(x_i, d_i^u) = \frac{x_i^T d_i^u}{\|x_i\| \|d_i^u\|}$ is the cosine similarity and τ_1 is the temperature. The \mathcal{L}^{ucl} aims

to pull the representation of x_i close to the representation of its associated intent descriptor while maintaining distinction from others.

Additionally, for labeled utterances in \mathcal{D}^l , we further utilize the high-quality supervisory signals to optimize the SynCID. On the one hand, we utilize the supervised contrastive loss to align the extracted representations x_i and d_i^l for utterance x_i and its label-enriched intent descriptor d_i^l , facilitating discriminative representation learning as below:

$$\mathcal{L}^{scl} = -\sum_{j=1}^{\mathcal{Y}_{x_i}} \log \frac{e^{\text{sim}(x_i, d_j^l)/\tau_2}}{\sum \mathbb{1}_{[k \neq i]} e^{\text{sim}(x_i, d_k^l)/\tau_2}}, \quad (4)$$

where τ_2 is the temperature. \mathcal{Y}_{x_i} is the index set of data sharing the same label as x_i .

On the other hand, we compute a standard cross-entropy loss \mathcal{L}^{ce} for the labeled utterances in \mathcal{D}^l to regulate the training of SynCID. It optimizes the model \mathcal{M} to distinguish the target intent classes of utterances from all known intent classes, enhancing the learning of utterance representations. Specifically, we map the utterance representation x_i into a probability distribution using a classifier and maximize the likelihood of its corresponding ground truth class (equation omitted for space). As a result, the overall loss \mathcal{L}_{SACL} is formulated as follows:

$$\mathcal{L}_{SACL} = \mathcal{L}^{ce} + \lambda \mathcal{L}^{ucl} + \eta \mathcal{L}^{scl}, \quad (5)$$

where λ and η denote hyper-parameters that modulate the respective contributions of distinct losses.

SA with Neighbor Filtering. Upon training with \mathcal{L}_{SACL} , SynCID can learn some compact utterance representations for clustering. Yet, these representations are inevitably affected by the utterance noise from either the use of the unsupervised contrastive loss \mathcal{L}^{ucl} or the limited comprehension capability of the model \mathcal{M} . To more effectively synergize LLMs with SLMs and amplify LLMs’ insights for discerning intents, we further enhance SynCID by implementing neighbor utterance filtering, aiming for a more consistent alignment between the semantic spaces of the utterances and the intent descriptors from LLMs. Specifically, for each utterance x_i and its intent descriptor d_i^u , we first retrieve their nearest neighboring utterances \mathcal{N}_{x_i} and intent descriptors $\mathcal{N}_{d_i^u}$ respectively. Owing to the accurate comprehension of LLMs, the intent descriptor d_i^u is anticipated to have cleaner neighbors. Thus, we filter out noisy utterance neighbors by omitting any $x_j \in \mathcal{N}_{x_i}$ where its paired $d_j^u \notin \mathcal{N}_{d_i^u}$, retaining

a purified neighbor set \mathcal{N}'_{x_i} for x_i . During training, we update SynCID via a contrastive learning objective to pull together all filtered neighboring utterances and push apart non-neighbors as follows:

$$\mathcal{L}_{SANF} = - \sum_{j=1}^{\mathcal{N}'_{x_i}} \log \frac{e^{\text{sim}(\mathbf{x}_i, \mathbf{x}_j) / \tau_3}}{\sum_{[p \neq i]} \mathbb{1}_{[p \neq i]} e^{\text{sim}(\mathbf{x}_i, \mathbf{x}_p) / \tau_3}}, \quad (6)$$

where τ_3 is the temperature. Here, we update the neighbor sets \mathcal{N}_{x_i} and $\mathcal{N}_{d_i^u}$ every several epochs for filtering out noisy utterances during training.

3.5 Intent Label Generation

After training models to learn discriminative representations, existing CID methods (Zhang et al., 2022b, 2023) typically utilize clustering algorithms like K -means to group utterances into distinct clusters for inferring intents. Yet, it remains challenging to assign accurate labels for newly identified intent clusters. SynCID addresses this by utilizing the in-context learning capability of LLMs to generate suitable labels for new intent clusters. Specifically, we devise a label generation prompt (LG_Prompt) for extracting labels from LLMs. As illustrated in Figure 2, the LG_Prompt is constructed as:

$$\text{LG_Prompt} = (\text{ICDs}, \text{Center Utterances}, p_c),$$

where $\text{ICDs} = \{(x_1^j, \dots, x_k^j, y_j)\}_{j=1}^n$ is a set of n in-context demonstrations. We can set the number n ranging from 1 to L considering the context size of LLMs. Each demonstration comprises a known intent label y_j and the top- k utterances near the y_j cluster center. *Center Utterances* is a set of utterances (x_1, \dots, x_k) located around the same unknown intent cluster center. p_c is the task description. For each unknown intent cluster, we integrate the top- k utterances allocated to it into the LG_Prompt, prompting LLMs to generate a new intent label y specific to it.

4 Experiments

4.1 Datasets

We conduct experiments on three CID datasets: **BANKING** (Casanueva et al., 2020), **CLINC** (Larson et al., 2019), and **StackOverflow** (Xu et al., 2015). The detailed statistics are reported in Appendix A.1. We keep the same train, development, and test splits as previous work (Zhang et al., 2023). To avoid randomness, we average the experimental results in five random runs. More experimental details are provided in the Appendix A.2.

4.2 Evaluation Metrics

We adopt three standard metrics for evaluating the CID performance: Accuracy (**ACC**) based on the Hungarian algorithm, Adjusted Rand Index (**ARI**), and Normalized Mutual Information (**NMI**). The specific definitions are shown in Appendix A.4. Note that **ACC** is considered as the primary metric, with higher values indicating better performance.

4.3 Baselines

We mainly compare our SynCID with the following SOTA baselines in our experiments:

Unsupervised: (1) **DEC** (Xie et al., 2016), (2) **DCN** (Yang et al., 2017), (3) **SCCL** (Zhang et al., 2021a), (4) **IDAS** (De Raedt et al., 2023).

Semi-supervised: (1) **DTC** (Han et al., 2019), (2) **CDAC+** (Lin et al., 2020), (3) **DeepAligned** (Zhang et al., 2021c), (4) **ProbNID** (Zhou et al., 2023), (5) **DCSC** (Wei et al., 2022), (6) **MTP-CLNN** (Zhang et al., 2022b), (7) **USNID** (Zhang et al., 2023), (8) **CsePL** (Liang and Liao, 2023). More details are provided in Appendix A.5.

4.4 Main Results

4.4.1 CID Performance Comparison

We report the main CID results in Table 1, with the highest performance highlighted in **bold**. We analyze the results as follows:

SynCID consistently outperforms CID baselines by large margins: Table 1 shows that SynCID exceeds all baseline methods in performance across three CID datasets and various KIR settings. For example, SynCID surpasses the top baseline CsePL by averages of 4.35% in ACC, 5.04% in ARI, and 2.07% in NMI on BANKING-25%. Moreover, SynCID shows stronger robustness in relation to the ratio of labeled data available. From BANKING-50% to BANKING-25%, SynCID’s performance merely drops 2.42% in ACC, 2.27% in ARI, and 0.94% in NMI. In contrast, the corresponding metrics for CsePL diminish by 5.88%, 6.30%, and 2.33%, respectively. This suggests that SynCID, leveraging the nuanced understanding from LLMs, learns more robust utterance representations for recognizing intents and effectively alleviates the issue of overfitting to known user intents.

SynCID provides a better way to unleash the power of LLMs for CID. We can observe that our SynCID consistently demonstrates superior performance over previous unsupervised leading baseline IDAS. Specifically, SynCID surpasses IDAS

KIR	Methods	BANKING			CLINC			StackOverflow		
		ACC	ARI	NMI	ACC	ARI	NMI	ACC	ARI	NMI
0%	DEC	38.60	25.32	62.65	48.77	31.71	74.83	59.49	36.23	58.76
	DCN	38.59	25.36	62.72	48.69	31.68	74.77	59.48	36.23	58.75
	SCCL	40.54	26.98	63.89	50.44	38.14	79.35	68.15	34.81	69.11
	USNID	54.83	43.33	75.30	75.87	68.54	91.00	69.28	52.25	72.00
	IDAS	67.43	57.56	82.84	85.48	79.02	93.82	83.82	72.20	81.26
	SynCID	72.89 †	62.42 †	84.20 †	86.80 †	80.85 †	94.23 *	86.90 †	74.42 †	81.95 *
25%	DTC	31.75	19.09	55.59	56.90	41.92	79.35	29.54	17.51	29.96
	CDAC+	48.00	33.74	66.39	66.24	50.02	84.68	51.61	30.99	46.16
	DeepAligned	49.08	37.62	70.50	74.07	64.63	88.97	54.50	37.96	50.86
	ProbNID	55.75	44.25	74.37	71.56	63.25	89.21	54.10	38.10	53.70
	DCSC	60.15	49.75	78.18	79.89	72.68	91.70	-	-	-
	MTP-CLNN	65.06	52.91	80.04	83.26	76.20	93.17	74.70	54.80	73.35
	USNID	65.85	56.53	81.94	83.12	77.95	94.17	75.76	65.45	74.91
	SynCID	75.41 †	65.40 †	85.39 †	87.85 †	82.39 †	94.85 †	87.86 †	76.11 †	82.46 †
50%	DTC	49.85	37.05	69.46	64.39	50.44	83.01	52.92	37.38	49.80
	CDAC+	48.55	34.97	67.30	68.01	54.87	86.00	51.79	30.88	46.21
	DeepAligned	59.38	47.95	76.67	80.70	72.56	91.59	74.52	57.62	68.28
	ProbNID	63.02	50.42	77.95	82.62	75.27	92.72	73.20	62.46	74.54
	DCSC	68.30	56.94	81.19	84.57	78.82	93.75	-	-	-
	MTP-CLNN	70.97	60.17	83.42	86.18	80.17	94.30	80.36	62.24	76.66
	USNID	73.27	63.77	85.05	87.22	82.87	95.45	82.06	71.63	78.77
	SynCID	77.83 †	67.67 †	86.33 †	90.64 †	85.96 †	95.91 *	88.40 †	77.24 †	83.34 †

Table 1: Main performance results on CID across three public datasets. KIR denotes the ratio of known intents. Results are averaged over five random runs. († and * denote p-value<0.01 and p-value<0.05 under t-test respectively.)

by margins of 5.46% in ACC, 4.86% in ARI, and 1.36% in NMI on the BANKING-0%. On the multi-domain CLINC dataset, SynCID records improvements of 1.32% in ACC, 1.83% in ARI, and 0.41% in NMI. It is noteworthy that IDAS utilizes LLMs to refine a frozen pre-trained encoder for discerning intents. Our SynCID, by contrast, dynamically synergizes LLMs and SLMs through the alignment between original utterances and intent descriptors. This observation suggests that our SynCID can effectively unleash LLMs’ nuanced comprehension capability to synergize them with SLMs for CID, guiding the SLMs in learning clarified utterance representations for intent identification.

4.4.2 Generated New Intent Labels

To study the quality of intent labels produced by SynCID, we conduct a comparative analysis between the gold standard labels and SynCID-generated intent labels on the CLINC dataset. Table 2 presents the comparison across different categories of intent labels. We can observe that for those clusters with specific and well-rounded user intent information, SynCID can accurately generate their corresponding intent labels, such as *Book hotel* and *Flight status*. Regarding the clusters that describe general user questions, SynCID can provide intent labels by condensing the user questions into high-level intents. For example, the intents

Gold Intent Label	Generated Intent Label
Book hotel	Book hotel
Flight status	Flight status
Who do you work for	Employer
Do you have pets	Pet ownership
Application status	Credit card application status
Oil change when	Oil change schedule

Table 2: Examples of generated new labels on CLINC.

Who do you work for and *Do you have pets* are succinctly transformed into *Employer* and *Pet ownership*, respectively. As for the clusters with overly general gold labels, i.e., *Application status* and *Oil change when*, the SynCID is able to integrate additional cluster details to construct more specific and accurate intent labels. This analysis indicates that SynCID, leveraging the capabilities of LLMs, can effectively capture the intrinsic intents conveyed within utterances and generate high-quality intent labels for newly identified intents clusters.

4.4.3 Early Detection of New Intents

Effectively identifying new intents at their initial emergence is vital for developing adaptive conversational agents. To meet this practical demand, we evaluate the performance of SynCID in the early discovery of new intents, comparing it with existing top-performing baselines. Table 3 showcases ex-

Shots	Methods	BANKING		
		ACC	ARI	NMI
5	MTP-CLNN	45.72	33.56	69.07
	USNID	43.64	33.00	69.78
	CsePL	47.44	37.34	70.98
	SynCID	56.06	44.40	74.54
10	MTP-CLNN	46.00	35.69	70.54
	USNID	47.29	37.61	72.73
	CsePL	52.31	39.85	73.12
	SynCID	59.01	46.25	75.67
20	MTP-CLNN	50.08	40.15	73.90
	USNID	50.17	40.66	74.77
	CsePL	61.43	49.16	77.33
	SynCID	66.09	54.06	80.22

Table 3: Results of early new intent detection on BANKING-25%. Shots denote the number of utterances within each unknown intent.

470 experimental results in scenarios with a limited number of utterances per unknown intent, specifically at {5, 10, 20} shots. It is observed that existing baselines demonstrate a notable decrease in performance compared to their prior evaluations. In contrast, our SynCID, despite the reduction in performance, consistently surpasses other leading baselines. For example, with 20 utterances per unknown intent on BANKING-25%, SynCID achieves improvements over the baseline CsePL by 4.66% in ACC, 4.90% in ARI, and 2.89% in NMI. Additionally, it is noted that SynCID’s performance gains over existing baselines progressively amplify as the number of utterance shots decreases. With only 5 utterance shots available for each unknown intent, SynCID attains improvements of 8.62% in ACC, 7.06% in ARI, and 3.56% in NMI. We hypothesize this observation can be explained by two main points: (1) Existing methods, which predominantly rely on SLMs, necessitate a sufficient quantity of utterances to cluster intents for reaching competitive performance. (2) In contrast, our SynCID synergizes LLMs and SLMs by aligning the semantic spaces of utterances with intent descriptors, providing a nuanced semantic understanding that compensates for limited data and thus enhancing the early discovery of new intents.

497 4.5 Detailed Analysis

498 4.5.1 Effect of Different LLMs

499 In addition to utilizing *text-davinci-003* as our basic LLM in the experiments, we further explore the use of various distinct LLMs, including the open-sourced *Flan-T5-XXL* (Chung et al., 2022) and the

KIR	Methods	BANKING		
		ACC	ARI	NMI
25%	SynCID- <i>Flan-T5-XXL</i>	73.47	61.90	83.84
	SynCID- <i>gpt-3.5-turbo</i>	74.29	62.86	84.21
	SynCID- <i>davinci-003</i>	75.41	65.40	85.39
	SynCID- <i>gpt-4</i>	77.79	65.95	85.46

Table 4: Effect of different LLMs on BANKING.

KIR	Methods	BANKING		
		ACC	ARI	NMI
0%	SynCID- <i>BERT</i>	72.89	62.42	84.20
	SynCID- <i>E5</i>	74.06	63.91	85.34
25%	SynCID- <i>BERT</i>	75.41	65.40	85.39
	SynCID- <i>E5</i>	77.34	68.16	86.70
50%	SynCID- <i>BERT</i>	77.83	67.67	86.33
	SynCID- <i>E5</i>	79.71	70.27	87.84

Table 5: Effect of different SLMs on BANKING.

close-sourced *gpt-3.5-turbo* and *gpt-4*, for deriving intent descriptors within SynCID. As shown in Table 4, integrating SynCID with different LLMs for intent descriptor generation consistently surpasses the top-performing baseline CsePL, in all three evaluation metrics on the BANKING-25% dataset. Notably, utilizing *gpt-4* for intent descriptor generation yields further enhancements over *text-davinci-003*. We hypothesize that this enhancement is attributable to the superior quality of intent descriptors generated by the more advanced LLM, which are more constructive in accurately fulfilling user intent discovery.

4507 4.5.2 Effect of Different Pre-trained SLMs

4508 The proposed SynCID primarily synergizes the agile responsiveness of the pre-trained SLMs and LLMs’ reliable insights for effectively discovering new intents. We inspect the contribution of different pre-trained SLMs, such as the BERT-based model and the more recent E5 model (Wang et al., 2022), to our SynCID, as detailed in Table 5. We can observe that integrating the E5 model into SynCID leads to further performance enhancements across various known intent rates when compared to the standard SynCID. It suggests that our SynCID framework stands to gain from synergizing LLMs and more advanced pre-trained SLMs.

4509 4.5.3 Effect of Space Alignment

4510 To verify the impact of different contrastive learning objectives within the space alignment on SynCID’s performance, we conduct a comprehensive ablation study on BANKING-25%, with the results detailed in Table 6. Specifically, we selectively re-

KIR	Methods	BANKING		
		ACC	ARI	NMI
25%	SynCID	75.41	65.40	85.39
	- w/o \mathcal{L}^{scl}	72.69	62.42	84.12
	- w/o \mathcal{L}^{ucl}	72.86	62.47	84.11
	- w/o \mathcal{L}_{SANF}	68.20	57.28	81.38

Table 6: Ablation results on BANKING.

536 move three distinct contrastive losses from SynCID
537 for analysis, where *w/o* denotes the model with-
538 out the corresponding loss. Findings from Table 6
539 show a performance decline in CID when any con-
540 trastive loss is excluded. For example, removing
541 \mathcal{L}^{scl} results in SynCID’s performance dropping by
542 2.72% in ACC, 2.98% in ARI, and 1.27% in NMI.
543 Yet, despite these reductions, SynCID variants still
544 maintain competitive performance compared to ex-
545 isting top-performing baselines. This underscores
546 the efficacy of the contrastive learning objectives
547 in the space alignment, highlighting their effective-
548 ness in synergizing the powerful LLMs and the
549 agile SLMs to learn discriminative representations,
550 thereby facilitating the new intent identification.

551 4.5.4 Impact of Descriptor Shots

552 To further validate the efficacy of the intent descrip-
553 tors within the proposed SynCID, we explore the
554 impact of varying intent descriptor shots on Syn-
555 CID’s performance in intent discovery. We con-
556 duct experiments on StackOverflow-25%, where
557 the improvement observed with SynCID is most
558 pronounced, thus providing a solid foundation for
559 this investigation. Figure 3 showcases a compar-
560 ison of the CID performance corresponding to dif-
561 ferent intent descriptor shots within the SynCID. It
562 can be observed that increasing the quantity of the
563 intent descriptors for optimizing the SynCID does
564 not yield substantial improvements in identifying
565 new intents. We hypothesize this can be attributed
566 to the propensity of LLMs to generate similar in-
567 tent descriptors, even when prompted to generate
568 multiple descriptors for a single utterance. These
569 analogous intent descriptors do not provide enough
570 supplementary information for the SynCID while
571 increasing computation costs.

572 4.6 Visualisation of Alleviating Overfitting

573 For a more intuitive analysis of the effect of our
574 SynCID on utterance representation learning, we
575 present the t-SNE visualizations comparing the
576 SynCID framework with the top baseline CsePL, as
577 illustrated in Figure 4. We can observe that the Syn-

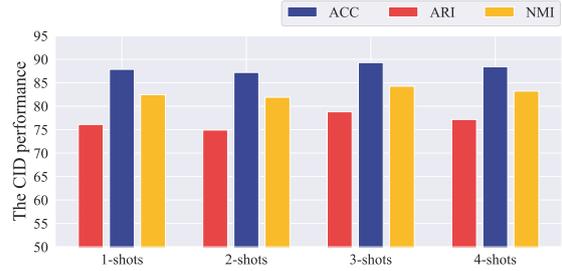


Figure 3: Effect of descriptor shots on StackOverflow.

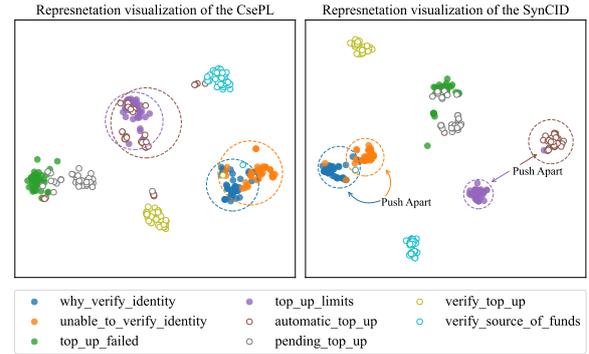


Figure 4: T-SNE visualization. The prefix “UK_” and “K_” denote unknown intents (hollow circles) and known intents (solid circles) respectively.

590 CID framework performs space alignment to align
591 the original utterance semantic space with LLMs’
592 intent descriptor space for representation learning,
593 thereby facilitating the formation of more compact
594 and distinct intent clusters. Additionally, we can
595 notice that SynCID effectively segregates the inter-
596 twined intent clusters, *i.e.*, *UK_automatic_top_up*
597 and *K_top_up_limits*, *K_unable_to_verify_identity*
598 and *K_why_verify_identity*, compared with the
599 CsePL. The visualization of utterance representa-
600 tions demonstrates the proficiency of SynCID in
601 alleviating the issue of overfitting to known intents.
602

603 5 Conclusion

604 In this paper, we introduced SynCID, a novel frame-
work that can effectively synergize LLMs and pre-
trained SLMs for conversational intent discovery.
By aligning LLMs’ reliable insights with the agile
responsiveness of specialized SLMs, SynCID effec-
tively alleviates the risk of overfitting to known
intents in CID. Furthermore, SynCID enables the
LLMs with in-context learning to skillfully produce
labels for newly identified intent clusters. Through
extensive experiments, our findings confirm Syn-
CID’s effectiveness. Deeper analysis reveals that
SynCID not only sets new benchmarks in CID but
also generates appropriate intent labels and enables
early detection of new intents.

605 Limitations

606 Despite the promising results obtained by our Syn-
607 CID, it is important to acknowledge several limita-
608 tions: (1) The SynCID’s reliance on LLMs subjects
609 it to LLMs’ inherent flaws, including biases in the
610 training data and the propensity for hallucinating
611 incorrect information. (2) The financial cost of uti-
612 lizing commercial LLM APIs, such as OpenAI’s,
613 for experiments is significant. In our case, access-
614 ing APIs of LLMs such as *gpt-4*, *gpt-3.5-turbo*,
615 and *davinci-003* for getting all the experimental
616 results incurred a cost of approximately \$510. (3)
617 Our SynCID, similar to existing baselines, assumes
618 a known ground-truth number of intents for clus-
619 tering utterances — a condition that diverges from
620 real-world applications where the exact number of
621 intents remains unknown. To validate SynCID’s
622 effectiveness and robustness, we conduct further
623 experiments with an estimated number of intents
624 and explore the impact of various intent numbers
625 around it on the CID performance of our SynCID.
626 The findings from these additional experiments are
627 detailed in Appendix B.

628 References

629 Avishek Anand, Venkatesh V, Abhijit Anand, and Vinay
630 Setty. 2023. [Query understanding in the age of large
631 language models](#). *CoRR*, abs/2306.16004.

632 Sidney Black, Stella Biderman, Eric Hallahan, Quentin
633 Anthony, Leo Gao, Laurence Golding, Horace
634 He, Connor Leahy, Kyle McDonell, Jason Phang,
635 Michael Pieler, Usvsn Sai Prashanth, Shivanshu
636 Purohit, Laria Reynolds, Jonathan Tow, Ben Wang,
637 and Samuel Weinbach. 2022. [GPT-NeoX-20B: An
638 open-source autoregressive language model](#). In *Big-
639 Science*, pages 95–136.

640 Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
641 Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
642 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
643 Askell, Sandhini Agarwal, Ariel Herbert-Voss,
644 Gretchen Krueger, Tom Henighan, Rewon Child,
645 Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,
646 Clemens Winter, Christopher Hesse, Mark Chen, Eric
647 Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess,
648 Jack Clark, Christopher Berner, Sam McCandlish,
649 Alec Radford, Ilya Sutskever, and Dario Amodei.
650 2020. [Language models are few-shot learners](#). In
651 *NeurIPS*.

652 Mathilde Caron, Piotr Bojanowski, Armand Joulin, and
653 Matthijs Douze. 2018. [Deep clustering for unsuper-
654 vised learning of visual features](#). In *ECCV*, pages
655 139–156.

656 Iñigo Casanueva, Tadas Temčinas, Daniela Gerz,
657 Matthew Henderson, and Ivan Vulić. 2020. [Effi-
658 cient intent detection with dual sentence encoders](#). In
659 *NLP4ConvAI@ACL*, pages 38–45.

660 Jackie Chi Kit Cheung and Xiao Li. 2012. [Sequence
661 clustering and labeling for unsupervised query intent
662 discovery](#). In *WSDM*, pages 383–392.

663 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin,
664 Maarten Bosma, Gaurav Mishra, Adam Roberts,
665 Paul Barham, Hyung Won Chung, Charles Sutton,
666 Sebastian Gehrmann, Parker Schuh, Kensen Shi,
667 Sasha Tsvyashchenko, Joshua Maynez, Abhishek
668 Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vin-
669 odkumar Prabhakaran, Emily Reif, Nan Du, Ben
670 Hutchinson, Reiner Pope, James Bradbury, Jacob
671 Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin,
672 Toju Duke, Anselm Levskaya, Sanjay Ghemawat,
673 Sunipa Dev, Henryk Michalewski, Xavier Garcia,
674 Vedant Misra, Kevin Robinson, Liam Fedus, Denny
675 Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim,
676 Barret Zoph, Alexander Spiridonov, Ryan Sepassi,
677 David Dohan, Shivani Agrawal, Mark Omernick, An-
678 drew M. Dai, Thanumalayan Sankaranarayanan Pilla-
679 i, Marie Pellat, Aitor Lewkowycz, Erica Moreira,
680 Rewon Child, Oleksandr Polozov, Katherine Lee,
681 Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark
682 Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy
683 Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov,
684 and Noah Fiedel. 2023. [Palm: Scaling language mod-
685 eling with pathways](#). *J. Mach. Learn. Res.*, 24:240:1–
686 240:113.

687 Hyung Won Chung, Le Hou, Shayne Longpre, Barret
688 Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi
689 Wang, Mostafa Dehghani, Siddhartha Brahma, et al.
690 2022. [Scaling instruction-finetuned language models](#).
691 *arXiv preprint arXiv:2210.11416*.

692 Marco Cuturi. 2013. [Sinkhorn distances: Lightspeed
693 computation of optimal transport](#). In *NeurIPS*, pages
694 2292–2300.

695 Maarten De Raedt, Frédéric Godin, Thomas De-
696 meester, and Chris Develder. 2023. [IDAS: In-
697 tent discovery with abstractive summarization](#). In
698 *NLP4ConvAI@ACL*, pages 71–88.

699 Bosheng Ding, Chengwei Qin, Linlin Liu, Yew Ken
700 Chia, Boyang Li, Shafiq Joty, and Lidong Bing. 2023. [Is
701 GPT-3 a good data annotator?](#) In *ACL*, pages
702 11173–11195.

703 Tianyu Gao, Xingcheng Yao, and Danqi Chen. 2021. [SimCSE:
704 Simple contrastive learning of sentence
705 embeddings](#). In *EMNLP*, pages 6894–6910.

706 James Gung, Raphael Shu, Emily Moeng, Wesley Rose,
707 Salvatore Romeo, Arshit Gupta, Yassine Benajiba,
708 Saab Mansour, and Yi Zhang. 2023. [Intent induction
709 from conversations for task-oriented dialogue track
710 at DSTC 11](#). In *DSTC-WS*, pages 242–259.

711 Amir Hadifar, Lucas Sterckx, Thomas Demeester, and
712 Chris Develder. 2019. [A self-training approach for](#)

713	short text clustering . In <i>RepL4NLP@ACL</i> , pages	765
714	194–199.	766
715	Kai Han, Andrea Vedaldi, and Andrew Zisserman. 2019.	767
716	Learning to discover novel visual categories via deep	768
717	transfer clustering . In <i>ICCV</i> , pages 8400–8408.	769
718	Cheng-Yu Hsieh, Chun-Liang Li, Chih-Kuan Yeh,	770
719	Hootan Nakhost, Yasuhisa Fujii, Alex Ratner, Ranjay	771
720	Krishna, Chen-Yu Lee, and Tomas Pfister. 2023. Dis-	772
721	tilling step-by-step! outperforming larger language	773
722	models with less training data and smaller model	774
723	sizes . In <i>Findings of ACL</i> , pages 8003–8017.	
724	Yen-Chang Hsu, Zhaoyang Lv, and Zsolt Kira. 2018.	777
725	Learning to cluster in order to transfer across domains	778
726	and tasks . In <i>ICLR</i> .	779
727	Yen-Chang Hsu, Zhaoyang Lv, Joel Schlosser, Phillip	780
728	Odom, and Zsolt Kira. 2019. Multi-class classifica-	781
729	tion without multi-class labels . In <i>ICLR</i> .	782
730	Stefan Larson, Anish Mahendran, Joseph J. Peper,	783
731	Christopher Clarke, Andrew Lee, Parker Hill,	784
732	Jonathan K. Kummerfeld, Kevin Leach, Michael A.	785
733	Laurenzano, Lingjia Tang, and Jason Mars. 2019.	786
734	An evaluation dataset for intent classification and	787
735	out-of-scope prediction . In <i>EMNLP-IJCNLP</i> , pages	
736	1311–1316.	
737	Shiyang Li, Jianshu Chen, Yelong Shen, Zhiyu Chen,	788
738	Xinlu Zhang, Zekun Li, Hong Wang, Jing Qian,	789
739	Baolin Peng, Yi Mao, et al. 2022. Explanations from	790
740	large language models make small reasoners better .	791
741	<i>arXiv preprint arXiv:2210.06726</i> .	
742	Yanen Li, Bo-June Paul Hsu, and ChengXiang Zhai.	792
743	2013. Unsupervised identification of synonymous	793
744	query intent templates for attribute intents . In <i>CIKM</i> ,	794
745	pages 2029–2038.	795
746	Jinggui Liang and Lizi Liao. 2023. ClusterPrompt:	797
747	Cluster semantic enhanced prompt learning for new	798
748	intent discovery . In <i>Findings of EMNLP</i> , pages	799
749	10468–10481.	800
750	Ting-En Lin, Hua Xu, and Hanlei Zhang. 2020. Dis-	801
751	covering new intents via constrained deep adaptive	802
752	clustering with cluster refinement . In <i>AAAI</i> , pages	
753	8360–8367.	
754	Yixin Liu, Alexander R. Fabbri, Pengfei Liu, Dragomir	803
755	Radev, and Arman Cohan. 2023. On learning to	804
756	summarize with large language models as references .	805
757	<i>CoRR</i> , abs/2305.14239.	806
758	Ilya Loshchilov and Frank Hutter. 2019. Decoupled	807
759	weight decay regularization . In <i>ICLR</i> .	
760	Yutao Mou, Keqing He, Pei Wang, Yanan Wu, Jingang	808
761	Wang, Wei Wu, and Weiran Xu. 2022a. Watch the	809
762	neighbors: A unified k-nearest neighbor contrastive	810
763	learning framework for OOD intent discovery . In	811
764	<i>EMNLP</i> , pages 1517–1529.	
	Yutao Mou, Keqing He, Yanan Wu, Zhiyuan Zeng,	812
	Hong Xu, Huixing Jiang, Wei Wu, and Weiran Xu.	813
	2022b. Disentangled knowledge transfer for OOD	814
	intent discovery with unified contrastive learning . In	815
	<i>ACL</i> , pages 46–53.	
	Yutao Mou, Xiaoshuai Song, Keqing He, Chen Zeng,	816
	Pei Wang, Jingang Wang, Yunsen Xian, and Weiran	817
	Xu. 2023. Decoupling pseudo label disambiguation	818
	and representation learning for generalized intent dis-	819
	covery . In <i>ACL</i> , pages 9661–9675.	
	OpenAI. 2023. GPT-4 technical report . <i>CoRR</i> ,	
	abs/2303.08774.	
	Padmasundari and Srinivas Bangalore. 2018. Intent dis-	
	covery through unsupervised semantic text clustering .	
	In <i>INTERSPEECH</i> , pages 606–610.	
	Yilin Shen, Yen-Chang Hsu, Avik Ray, and Hongxia	
	Jin. 2021. Enhancing the generalization for intent	
	classification and out-of-domain detection in SLU .	
	<i>CoRR</i> , abs/2106.14464.	
	Chen Shi, Qi Chen, Lei Sha, Sujian Li, Xu Sun,	
	Houfeng Wang, and Lintao Zhang. 2018. Auto-	
	dialabel: Labeling dialogue data with unsupervised	
	learning . In <i>EMNLP</i> , pages 684–689.	
	Kumar Shridhar, Alessandro Stolfo, and Mrinmaya	
	Sachan. 2023. Distilling reasoning capabilities into	
	smaller language models . In <i>Findings of ACL</i> , pages	
	7059–7073.	
	Jiao Sun, Yufei Tian, Wangchunshu Zhou, Nan Xu, Qian	
	Hu, Rahul Gupta, John Frederick Wieting, Nanyun	
	Peng, and Xuezhe Ma. 2023. Evaluating large lan-	
	guage models on controlled generation tasks . In	
	<i>EMNLP</i> , pages 3155–3168.	
	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier	
	Martinet, Marie-Anne Lachaux, Timothée Lacroix,	
	Baptiste Rozière, Naman Goyal, Eric Hambro,	
	Faisal Azhar, et al. 2023. Llama: Open and effi-	
	cient foundation language models . <i>arXiv preprint</i>	
	<i>arXiv:2302.13971</i> .	
	Liang Wang, Nan Yang, Xiaolong Huang, Binx-	
	ing Jiao, Linjun Yang, Daxin Jiang, Rangan Ma-	
	jumder, and Furu Wei. 2022. Text embeddings by	
	weakly-supervised contrastive pre-training . <i>ArXiv</i> ,	
	abs/2212.03533.	
	Shuohang Wang, Yang Liu, Yichong Xu, Chenguang	
	Zhu, and Michael Zeng. 2021. Want to reduce la-	
	beling cost? gpt-3 can help . In <i>Findings of EMNLP</i> ,	
	pages 4195–4205.	
	Feng Wei, Zhenbo Chen, Zhenghong Hao, Fengxin	
	Yang, Hua Wei, Bing Han, and Sheng Guo. 2022.	
	Semi-supervised clustering with contrastive learning	
	for discovering new intents . <i>CoRR</i> , abs/2201.07604.	
	Jason W. Wei and Kai Zou. 2019. EDA: easy data aug-	
	mentation techniques for boosting performance on	
	text classification tasks . In <i>EMNLP-IJCNLP</i> , pages	
	6381–6387.	

- 820 Junyuan Xie, Ross B. Girshick, and Ali Farhadi. 2016. [Unsupervised deep embedding for clustering analysis](#). In *ICML*, pages 478–487. 874
- 821 [Unsupervised deep embedding for clustering analysis](#). 875
- 822 In *ICML*, pages 478–487. 876
- 823 Qizhe Xie, Zihang Dai, Eduard H. Hovy, Thang Luong,
824 and Quoc Le. 2020. [Unsupervised data augmentation
825 for consistency training](#). In *NeurIPS*.
- 826 Jianming Xu, Peng Wang, Guanhua Tian, Bo Xu, Jun
827 Zhao, Fangyuan Wang, and Hongwei Hao. 2015.
828 [Short text clustering via convolutional neural net-
829 works](#). In *VS@HLT-NAACL*, pages 62–69.
- 830 Bo Yang, Xiao Fu, Nicholas D. Sidiropoulos, and
831 Mingyi Hong. 2017. [Towards k-means-friendly
832 spaces: Simultaneous deep learning and clustering](#).
833 In *ICML*, pages 3861–3870.
- 834 Yiben Yang, Chaitanya Malaviya, Jared Fernandez,
835 Swabha Swayamdipta, Ronan Le Bras, Ji-Ping Wang,
836 Chandra Bhagavatula, Yejin Choi, and Doug Downey.
837 2020. [G-daug: Generative data augmentation for
838 commonsense reasoning](#). In *Findings of EMNLP*,
839 pages 1008–1025.
- 840 Asaf Yehudai, Matan Vetzler, Yosi Mass, Koren Lazar,
841 Doron Cohen, and Boaz Carmeli. 2023. [QAID: ques-
842 tion answering inspired few-shot intent detection](#). In
843 *ICLR*.
- 844 Eyup Halit Yilmaz and Cagri Toraman. 2020. [KLOOS:
845 KL divergence-based out-of-scope intent detection
846 in human-to-machine conversations](#). In *SIGIR*, pages
847 2105–2108.
- 848 Dejian Zhang, Feng Nan, Xiaokai Wei, Shang-Wen
849 Li, Henghui Zhu, Kathleen R. McKeown, Ramesh
850 Nallapati, Andrew O. Arnold, and Bing Xiang. 2021a.
851 [Supporting clustering with contrastive learning](#). In
852 *NAACL-HLT*, pages 5419–5430.
- 853 Hanlei Zhang, Xiaoteng Li, Hua Xu, Panpan Zhang,
854 Kang Zhao, and Kai Gao. 2021b. [TEXTTOIR: An
855 integrated and visualized platform for text open intent
856 recognition](#). In *ACL-IJCNLP*, pages 167–174.
- 857 Hanlei Zhang, Hua Xu, Ting-En Lin, and Rui Lyu.
858 2021c. [Discovering new intents with deep aligned
859 clustering](#). In *AAAI*, pages 14365–14373.
- 860 Hanlei Zhang, Hua Xu, Xin Wang, Fei Long, and Kai
861 Gao. 2023. [USNID: A framework for unsupervised
862 and semi-supervised new intent discovery](#). *CoRR*,
863 abs/2304.07699.
- 864 Haode Zhang, Haowen Liang, Yuwei Zhang, Li-Ming
865 Zhan, Xiao-Ming Wu, Xiaolei Lu, and Albert Y. S.
866 Lam. 2022a. [Fine-tuning pre-trained language mod-
867 els for few-shot intent detection: Supervised pre-
868 training and isotropization](#). In *NAACL-HLT*, pages
869 532–542.
- 870 Yuwei Zhang, Haode Zhang, Li-Ming Zhan, Xiao-Ming
871 Wu, and Albert Lam. 2022b. [New intent discovery
872 with pre-training and contrastive learning](#). In *ACL*,
873 pages 256–269.

A Experimental Details

A.1 Dataset Statistics

Table 7 reports the detailed statistics for the BANKING, CLINC, and StackOverflow datasets. Specifically, the BANKING dataset includes over 13,000 carefully curated customer queries from the banking domain, categorized into 77 unique intents. The CLINC dataset encompasses a diverse collection of 22,500 labeled utterances distributed across 150 intents, covering multiple domains. StackOverflow, sourced from Kaggle.com, is a specialized dataset featuring 20,000 technical questions, organized into 20 distinct categories.

A.2 Implementation Details

For the dataset configuration, we randomly select a portion of intents to be designated as known intents, defining this portion as the known intent rate (KIR) at levels of 0%, 25%, and 50%. The KIR = 0% indicates the unsupervised setting to CID, whereas the KIR > 0% implies the semi-supervised CID setting. From each intent selected as known, we sample 10% of the labeled utterances to create the labeled dataset \mathcal{D}^l . The remaining utterances are considered unlabeled, forming the basis of the unlabeled dataset \mathcal{D}^u .

For the LLMs-based Descriptor Generation and Intent Label Generation, our experiments are conducted with *text-davinci-003* serving as the basic LLM. To ensure deterministic outputs during descriptor generation, the temperature parameter is fixed at 0, and the output is limited to a maximum of 256 tokens. All other parameters are maintained at their default settings.

Within the Space Alignment, we utilize the pre-trained BERT model (*bert-uncased*), featuring a 12-layer transformer architecture, as the foundational SLM for training. The optimization of model parameters is conducted using the AdamW optimizer (Loshchilov and Hutter, 2019). During the SA with Contrastive Learning, the learning rate is set to 5×10^{-5} . The model outputs are projected from a 768-dimensional space to a 128-dimensional space for computing the contrastive loss. The temperatures $\{\tau_1, \tau_2\}$ for Equation 3 and 4 are uniformly set to 0.07. Furthermore, to achieve a balanced integration of \mathcal{L}^{ucl} and \mathcal{L}^{scl} , we apply λ and η values of 1.0. A more detailed analysis of these hyper-parameters is available in Section A.3.

We leverage an early stopping mechanism with a patience setting of 20 epochs on the development

Dataset	Domain	Intents	Utterances
BANKING	banking	77	13,083
CLINC	multi-domain	150	22,500
StackOverflow	question	20	20,000

Table 7: Statistics of datasets used in the experiments.

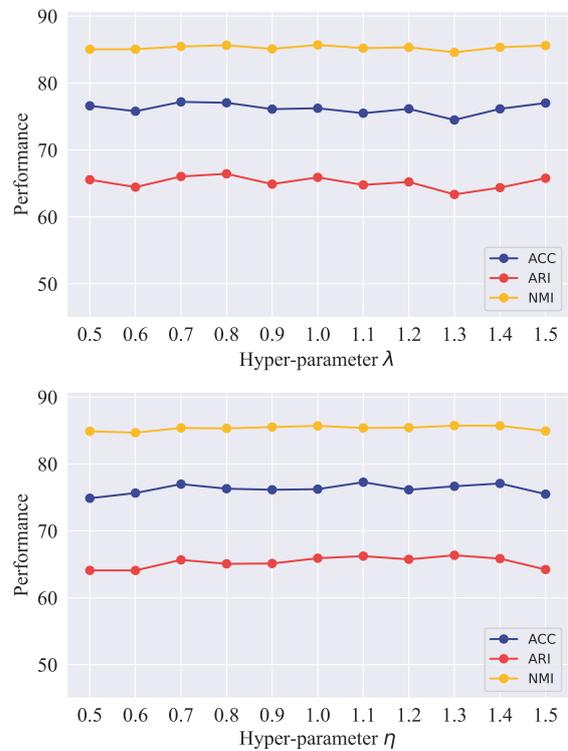


Figure 5: Impact of hyper-parameters λ and η on CID performance.

set to train the model. For the SA with Neighbor Filtering, the learning rate is set to 1×10^{-5} . We set the temperature τ_3 in Equation 6 to 0.07 similarly. Regarding the selection of neighborhood sizes $\{|\mathcal{N}_{x_i}^l|, |\mathcal{N}_{d_i^u}^u|\}$, following Zhang et al. (2022b), we empirically assign the values {100, 50} for the BANKING dataset, {120, 50} for the CLINC dataset, and {1000, 500} for the StackOverflow dataset.

A.3 Hyper-parameter Analysis

We conduct extensive hyper-parameter exploration experiments on BANKING-25% for selecting the proper λ and η to optimize the proposed SynCID. In the experiments, We carefully considered a range of values λ and η , ranging from 0.5 to 1.5. Figure 5 illustrates the effect of different settings of these hyper-parameters on the overall performance of SynCID. It is observed that varying these hyper-

parameters, either by increasing or decreasing their values, does not result in a significant change in the model performance, which demonstrates the robustness and stability of our SynCID.

A.4 Evaluation Metrics

During our experimental analysis, we utilize three metrics for evaluating CID performance: ACC, ARI, and NMI. Specifically, ACC is employed to assess the CID effectiveness by comparing the model’s predicted labels against the actual ground-truth labels. The calculation of ACC is defined as follows:

$$ACC = \frac{\sum_{i=1}^N \mathbb{1}_{y_i = \text{map}(\hat{y}_i)}}{N}$$

where $\{\hat{y}_i, y_i\}$ represent the predicted and ground-truth labels for an input utterance x_i , respectively. The function $\text{map}(\cdot)$ aligns each predicted label \hat{y}_i with its associated ground-truth label y_i , utilizing the Hungarian algorithm for this mapping process.

ARI measures the concordance of the predicted and actual clusters through an assessment of pairwise accuracy within clusters. The formulation of ARI is as follows:

$$ARI = \frac{\sum_{i,j} \binom{n_{i,j}}{2} - [\sum_i \binom{u_i}{2} \sum_j \binom{v_j}{2}] / \binom{N}{2}}{\frac{1}{2} [\sum_i \binom{u_i}{2} + \sum_j \binom{v_j}{2}] - [\sum_i \binom{u_i}{2} \sum_j \binom{v_j}{2}] / \binom{N}{2}}$$

where $u_i = \sum_j n_{i,j}$, and $v_j = \sum_i n_{i,j}$. The total number of samples is given by N , and $n_{i,j}$ indicates the count of sample pairs concurrently classified into the i^{th} predicted and the j^{th} actual cluster.

NMI is calculated to gauge the degree of concordance between the predicted and actual clusters by quantifying the normalized mutual information between them, as delineated below:

$$NMI(\hat{\mathbf{y}}, \mathbf{y}) = \frac{2 \cdot I(\hat{\mathbf{y}}, \mathbf{y})}{H(\hat{\mathbf{y}}) + H(\mathbf{y})}$$

where $\{\hat{\mathbf{y}}, \mathbf{y}\}$ denote the predicted labels and the ground-truth labels respectively. $I(\hat{\mathbf{y}}, \mathbf{y})$ is the mutual information between $\hat{\mathbf{y}}$ and \mathbf{y} . $H(\cdot)$ represents the entropy function.

A.5 Baselines

In this work, we compare the SynCID with the following SOTA baselines in our experiments:

Unsupervised Methods: (1) **DEC** (Xie et al., 2016): An unsupervised intent discovery method that iteratively learns and refines features by optimizing a clustering objective based on an auxiliary distribution. (2) **DCN** (Yang et al., 2017):

A method that combines nonlinear dimensionality reduction with k-means clustering to learn cluster-friendly representations for CID. (3) **SCCL** (Zhang et al., 2021a): An end-to-end clustering framework that jointly optimizes a top-down clustering loss with a bottom-up instance-wise contrastive loss. (4) **IDAS** (De Raedt et al., 2023): An unsupervised method that utilizes LLMs to refine a frozen pre-trained encoder for identifying intents.

Semi-supervised Methods: (1) **DTC** (Han et al., 2019): A semi-supervised deep learning methodology for clustering, featuring an innovative mechanism for estimating the number of intents by leveraging labeled data. (2) **CDAC+** (Lin et al., 2020): An approach based on pseudo-labeling employs pairwise constraints and a target distribution strategy to facilitate the learning process in intent recognition. (3) **DeepAligned** (Zhang et al., 2021c): A semi-supervised technique that addresses inconsistencies in clustering through an alignment strategy, enhancing the learning of utterance embeddings. (4) **ProbNID** (Zhou et al., 2023): A probabilistic framework employs the expectation-maximization technique, considering intent categorizations as potential latent variables. (5) **DCSC** (Wei et al., 2022): An approach for discovering intents through pseudo-labeling incorporates a dual-task mechanism, utilizing the SwAV algorithm alongside the Sinkhorn-Knopp method (Cuturi, 2013) for the assignment of soft clusters. (6) **MTP-CLNN** (Zhang et al., 2022b): A two-stage approach that improves the learning of utterance representations for discovering novel intents by integrating an initial multi-task pre-training with a subsequent nearest neighbor contrastive learning. (7) **USNID** (Zhang et al., 2023): A framework for both unsupervised and semi-supervised intent discovery, featuring a novel strategy for initializing centroids effectively to derive cluster representations using historical clustering information. (8) **CsePL** (Liang and Liao, 2023): A method that employs two-level contrastive learning with label semantic alignment for enhancing the cluster semantics, alongside a soft prompting strategy for identifying new intents.

B Estimate the Intent Number K

Predicting the precise number of intent clusters in conversational intent discovery systems presents a significant challenge in real-world applications. Leveraging the approach presented by Zhang et al. (2021c), our research utilizes the pre-initialized in-

Cluster Num K	BANKING		
	ACC	ARI	NMI
$K = 74$ (predicted)	75.13	63.74	84.74
$K = 77$ (gold)	75.41	65.40	85.39
$K = 71$	73.90	62.43	84.17
$K = 73$	74.94	63.57	84.88
$K = 75$	75.10	63.94	84.78
$K = 79$	75.32	64.73	85.30
$K = 81$	75.39	65.19	85.39

Table 8: Experimental results of different cluster number K under the BANKING-25% setting.

1039 tent features to autonomously ascertain the optimal
1040 number of intent clusters, represented as K . Ini-
1041 tially, we assign a larger estimated number of clus-
1042 ters, K' , and extract feature representations for the
1043 training dataset using a meticulously trained model.
1044 Subsequent clustering via the K-means algorithm
1045 divides these features into distinct groups. From
1046 this division, we distinguish between substantive
1047 intent clusters, characterized by their density and
1048 distinct boundaries, and smaller, less consequential
1049 clusters, which are then disregarded. The criteria
1050 for discerning between these cluster types can be
1051 outlined as follows:

$$1052 \quad K = \sum_{i=1}^{K'} \delta(|S_i| > \rho)$$

1053 where $|S_i|$ is the size of the i^{th} grouped cluster, and
1054 ρ serves as the threshold for filtering. The function
1055 $\delta(\cdot)$ acts as an indicator, yielding a value of 1 when
1056 a specified condition is met.

1057 Results of the experiments are reported in Table
1058 8, where, in addition to the predicted number of
1059 clusters K , we examine the performance across a
1060 range of intent numbers proximal to it. The com-
1061 parative results reveal that SynCID experiences
1062 merely marginal reductions in performance when
1063 confronted with inaccurate numbers of intents, in-
1064 dicated the robustness of SynCID in adapting to
1065 variations in the prediction of intent numbers.