# LARA: Linguistic-Adaptive Retrieval-Augmented LLMs for Multi-Turn Intent Classification

Junhua Liu[*1], Yong Keat Tan[*2], and Bin Fu[2](✉)

[1] Forth AI, Singapore
[2] Shopee, Singapore
j@forth.ai, {yongkeat.tan, bin.fu}@shopee.com

**Abstract.** Following the significant achievements of large language models (LLMs), researchers have employed in-context learning for text classification tasks. However, these studies focused on monolingual, single-turn classification tasks. In this paper, we introduce LARA (Linguistic-Adaptive Retrieval-Augmented Language Models), designed to enhance accuracy in multi-turn classification tasks across six languages, accommodating numerous intents in chatbot interactions. Multi-turn intent classification is notably challenging due to the complexity and evolving nature of conversational contexts. LARA tackles these issues by combining a fine-tuned smaller model with a retrieval-augmented mechanism, integrated within the architecture of LLMs. This integration allows LARA to dynamically utilize past dialogues and relevant intents, thereby improving the understanding of the context. Furthermore, our adaptive retrieval techniques bolster the cross-lingual capabilities of LLMs without extensive retraining and fine-tune. Comprehensive experiments demonstrate that LARA achieves state-of-the-art performance on multi-turn intent classification tasks, enhancing the average accuracy by 3.67% compared to existing methods.

**Keywords:** In-Context Learning · LLM · Multi-turn text classification

## 1 Introduction

Chatbots are an essential tool that automatically interacts or converses with customers. It plays a crucial role for international e-commerce platforms due to the rising consumer demand for instant and efficient customer service. Chatbots represent a critical component of dialogue systems [24] that can answer multiple queries simultaneously by classifying intent from the user's utterance to reduce waiting times and operational costs. Naturally, the interaction with users could turn into a multi-turn conversation if they require more detailed information about the query. Developing an intent classification model for a dialogue system isn't trivial, even if it's a typical text classification task. As we must consider

---

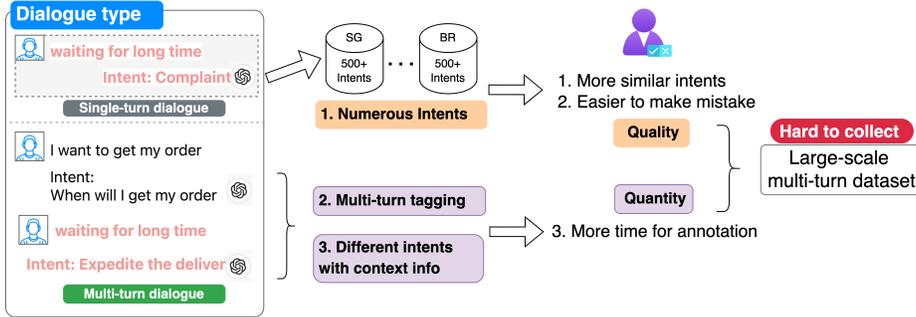[*] These authors contributed equally to this work.

Fig. 1: Annotation challenge of multi-turn intent classification dataset

contextual factors such as historical utterances and intents, failing to understand the session context while recognizing the user intention usually leads to more visible errors, as it would invoke a completely wrong application or provide an unrelated answer [28]. It's not surprising that it faces several challenges in dialogue understanding.

The biggest challenge is that the multi-turn dataset is hard to collect. Some studies have already been done on this problem of dialogue understanding in multi-turn intent classification [18,25,16]. However, they are made under the assumption of the availability of multi-turn training data, which is usually not the case in the real world.

Unlike emotion recognition in conversation (ERC) with only less than 10 classes or topic classification within dialogue state tracking (DST) with tens of topics, there are hundreds of intents within the knowledge base of a chatbot to cover users' specific intents in each market, it increases the complexity of classification tasks and multi-turn data annotation. Annotators can easily make mistakes and spend more time making decisions due to the numerous intents. Combined, these make it a high-cost and time-consuming annotation task, and it's unrealistic to annotate large-scale multi-turn datasets manually. However, the performance will most likely suffer without enough training sample size. This calls for a more efficient method in solving the challenge [13].

To tackle the above challenge, we propose **L**inguistic-**A**daptive **R**etrieval-**A**ugmentation, or LARA, which offers a pipeline of techniques to adopt only single-turn training data to optimize multi-turn dialogue classification. LARA first leverages an XLM-based model trained on single-turn classification datasets for each market, thus simplifying data construction and maintenance. Subsequently, LARA advances the field by selecting plausible candidate intents from user utterances and employing a retriever to gather relevant questions for prompt construction. This process facilitates in-context learning (ICL) with multi-lingual LLMs (MLLMs), significantly enhancing model efficacy without the need for market-specific multi-turn models.

In summary, the contributions of this paper are as follows:

1. We introduce LARA to effectively address multi-turn data collection issue through XLM-based model training and ICL with MLLM.

2. We conduct experiments on our e-commerce multi-turn dataset across six languages, showing that LARA model achieves state-of-the-art results and reduces inference time during ICL with MLLMs.

## 2   Related Work

### 2.1   Intent Classification

Intent classification is a typical text classification task where class labels are intent names. Various neural networks have inspired a wide range of studies aimed at creating neural models for text classification. Various neural model structures such as CNN [5,8], LSTM [29] and GCN [30,11] have proven more effective than conventional methods [9,23] based on statistical features. Additionally, some studies [34,35] utilize label embeddings and train them simultaneously with the input texts. In recent developments, the accomplishments of large-scale pre-training language models (PLM) [7] have generated significant interest in incorporating this pre-training approach [17] into monolingual and multilingual text classification [4]. This has resulted in substantial advancements in few-shot [3] and zero-shot learning [31]. However, most works are single-turn text classification tasks, which is unsuitable for multilingual multi-turn intent classification.

### 2.2   Modeling Multi-turn Dialogue Context

Modelling the multi-turn dialogues is the foundation for dialogue understanding tasks. Previous works adopt bidirectional contextual LSTM [6] to create contextual-aware utterance representation on MultiWOZ intent classification [2]. Recent works use PLM as a sentence encoder [20] on emotion recognition in conversation(ERC). Specifically, [10] used PLM to encode the context and speaker's memory and [15] enhance PLM by integrating multi-turn info from the utterance, context and dialogue structure through fine-tuning. However, all of their tasks adopt the multi-turn dialogue training set, which is hard to collect for an e-commerce chatbot. Our method attempts to combine an XLM-based model trained on the single-turn dataset into an in-context retrieve augmented pipeline with LLM, solving the multi-turn intent classification task in a zero-shot setting.

### 2.3   In-context Retrieval

In-context learning (ICL) with LLM like GPT-3 [1] demonstrates the significant improvement on few-shot/zero-shot NLP tasks. ICL has been successful in tasks like semantic parser [14,21], intent classification [32] and other utterance-level tasks. Some researchers [27,12] also apply ICL for dialogue state tracking, but it didn't perform as well as other methods. Future studies might look into better ways to retrieve dialogues and improve tasks' setup. **Retrieval** part, Most research on in-context learning (ICL) usually deals with single sentences or whole documents, but we are interested in finding and understanding dialogues. Generally, there are two types of systems to find the relevant dialogues: the first

is LM-score based retrieval. They [19,22] check the probability of a language model, like GPT-3, is to decode the right answer based on an example. The second type defines similarity metrics between task results and uses them as the training objective for the retriever. Both K-highest and lowest examples are used as positive and negative samples to help the system learn. The most pertinent research on dialogue retrieval concentrates on areas such as knowledge identification [26] and response selection [33]. Our objectives and settings differ from them.

## 3    Problem Formulation

### 3.1    Single-turn Intent Classification

In the single-turn scenario, the objective is to classify a user's query $q$ into one of the predefined intent classes $\mathcal{I} = \{I_i\}_{i=1}^{k}$, where $k$ often exceeds 200. Queries can range from informational requests to action-oriented dialogues. Recognizing the correct intent enables the dialogue system to provide relevant solutions or perform specific actions. This process is further complicated in transnational applications that must accommodate queries in multiple languages and adapt to localized business operations.

### 3.2    Multi-turn Intent Classification

Multi-turn scenarios involve a series of user queries $\mathcal{Q} = \{q_i\}_{i=1}^{n}$, with the aim of identifying the intent of the final query $q_n$. Unlike single-turn recognition, multi-turn recognition must account for the entire conversational context $\mathcal{C}$, which includes historical queries and their corresponding intents. This context-dependency introduces additional complexity, requiring models to interpret nuanced conversational dynamics and adjust to evolving user intentions over the course of an interaction.

### 3.3    Objective

This work aims to devise a methodology that leverages the readily available single-turn training data to address the inherent challenges of multi-turn intent recognition without the need for extensive multi-turn dataset curation. By proposing a model that dynamically adapts to the conversation's context while minimizing reliance on large-scale annotated multi-turn datasets, we seek to mitigate the significant annotation challenges and enable more effective intent recognition in complex multi-turn dialogues.

## 4    LARA: Linguistic-Adaptive Retrieval-Augmentation

The LARA framework addresses the multi-turn intent recognition challenge through zero-shot in-context learning with single-turn demonstrations, guided
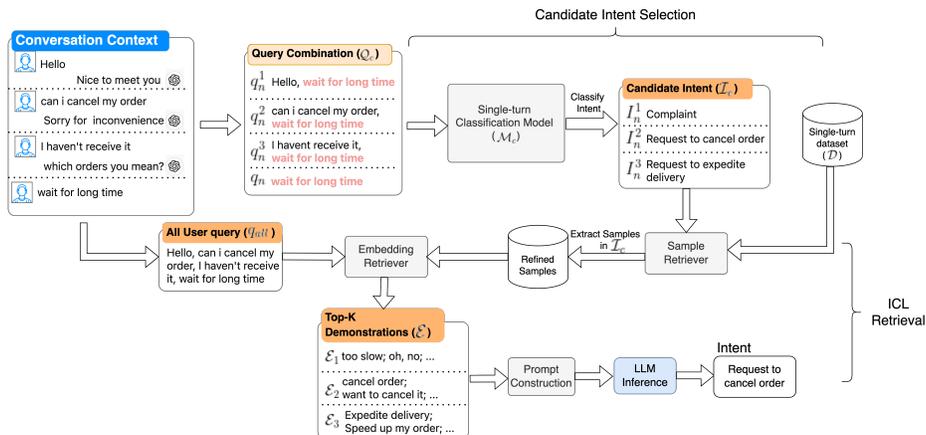
Fig. 2: The pipeline of Linguistic-Adaptive Retrieval-Augmentation

---

**Algorithm 1** Candidate Intent Selection

---

**Require:** $\mathcal{Q}_c = \{q_n, q_n^1, ..., q_n^{n-1}\}$
  **for** each item $q_i$ in $\mathcal{Q}_c$ **do**
    /* Get embedding H of [CLS] token */
    $H_i = \Phi_{XLMR}(q)^{[CLS]} \in \mathbb{R}^d$, $d$ is the hidden dimension

    /* Pass the embedding through a linear layer to get class probability P */
    $P_i = softmax(H \cdot W_c + b_c)$, $W_c \in \mathbb{R}^{d \times |\mathcal{I}|}$ and $b_c \in \mathbb{R}^{|\mathcal{I}|}$

    /* Select the intent with highest probability */
    $I_i = argmax(P_i)$
  **end for**
  **return** $\mathcal{I}_c = \{I_n, I_n^1, ..., I_n^{n-1}\}$

---

by a crafted instruction prompt. First, a single-turn intent classification model $\mathcal{M}_c$ is used to narrow down the intents to be included in the ICL prompt, which hereby referred to as candidate intents. This step is necessary due to the limited LLM context window, and it also helps to filter out extra noises from direct demonstration retrieval. Then, for every candidate intent, in-context demonstrations are selected by retrieving single-turn examples that are semantically-similar to the multi-turn test sample. Finally, an instruction prompt for multi-turn intent recognition is formulated by combining the demonstrations and test user queries.

**Single-turn Intent Recognition Model ($\mathcal{M}_c$)** A text classification model is trained on the annotated single-turn dataset $\mathcal{D}$. Given a query $q$, we adopt the [CLS] token embedding from XLM-RoBERTa-base model with weight $\Phi_{XLMR}$ as the text representation $H$. $\Phi_{XLMR}$ had been further pretrained with contrastive

---

**Algorithm 2** ICL Demonstrations Retrieval

---

**Require:** $\mathcal{I}_c$, $Q$, a positive integer $K$

$\mathcal{I}' = remove\_duplicate(\mathcal{I}_c)$

$q_{all} = text\_concatenate(\mathcal{Q})$

$H_q = \Phi_{XLMR}(q_{all})^{[CLS]}$

**for** each item $I_i$ in $\mathcal{I}'$ **do**
    $X_i = get\_training\_samples\_from\_\mathcal{D}\_for\_intent(I_i)$

    /* get embedding for each training sample */
    $H_{X_i} = \{\Phi_{XLMR}(x_j)^{[CLS]}\}_{j=1}^{|X_i|}, x \in X_i$

    /* calculate text similarity of each training sample with test queries */
    $S_i = \{cosine\_similarity(H_q, h_j)\}_{j=1}^{|H_{X_i}|}, h \in H_{X_i}$

    /* select nearest demonstrations */
    $\mathcal{E}_i \leftarrow$ Top $(K-1)$ $x \in X_i$ based on $S_i$
    Append $r$ of $I_i$ to $\mathcal{E}_i$ /* add representative query to demonstrations of $I_i$ */
**end for**

$\mathcal{E} = \{\mathcal{E}_i\}_{i=1}^{|\mathcal{I}'|}$ /* collect demonstrations of all candidate intents */
$S = \{S_i\}_{i=1}^{|\mathcal{I}'|}$
Sort $\mathcal{E}$ by their scores $S$ in ascending order

**return** $\mathcal{E}$

---

learning to give meaningful representation for [CLS] token. $H$ is then fed into a linear layer along with a softmax function to obtain intent probabilities. The intent with the highest probability is used.

During candidate intents selection, the last query $q_n$ is first combined with each historical query in $\mathcal{C}$ to form a query combination set $\mathcal{Q}_c = \{q_n, q_n^1, ..., q_n^{n-1}\}$, where $q_n^i$ means the text concatenation of $q_i$ with $q_n$ using a comma. $\mathcal{M}_c$'s inference on these combinations yields the candidate intents $\mathcal{I}_c$. The selection process is detailed in Algorithm 1.

**Retrieval Augmentation** $\Phi_{XLMR}$ is also utilized here to gather demonstrations for each intent in $\mathcal{I}_c$ based on their cosine similarity to test queries. Here, the demonstrations refer to a sequence of annotated examples that provide LLM with decision-making evidence and specify an output format for natural language conversion into labels during ICL. The details of ICL demonstration retrieval is included in Algorithm 1.

**Prompt Construction and LLM Inference** The task instruction $\mathcal{T}$, combined with demonstrations $\mathcal{E}$, conversational context $\mathcal{C}$, and the query $q_n$, forms the input prompt $\mathcal{P}$ for the LLM. To accommodate real-time application latency

requirements, two additional methods were explored to constrain the model to generate single-token symbols representing intents, detailed as $\mathcal{P}_{symbolic}$ and $\mathcal{P}_{prepend}$, with examples provided in the appendix. Model outputs are greedily decoded, ensuring efficient and accurate intent recognition.

## 5  Experiments

### 5.1  Dataset

The datasets used in this work consist of user queries in the local languages of eight markets: Brazil, Indonesia, Malaysia, Philippines, Singapore, Thailand, Taiwan, and Vietnam. The queries are related to the E-commerce domain.

| Market | Lang. | Intents | Train(Single-Turn) | Test(Multi-Turn) |
|--------|-------|---------|--------------------|------------------|
| BR | pt | 316 | 66k | 372 |
| ID | id | 481 | 161k | 1145 |
| MY | en,ms | 473 | 74k | 1417 |
| PH | en,fil | 237 | 33k | 189 |
| SG | en | 360 | 76k | 737 |
| TH | th | 359 | 60k | 502 |
| TW | zh-tw | 373 | 31k | 353 |
| VN | vi | 389 | 178k | 525 |

Table 1: The major languages, number of intents, and the number of samples in each market.

Table 1 shows the number of samples we have in each dataset. All the data are collected through the manual annotation by local CS teams of each market. We have the single-turn training data available in abundance over the course of business operation after years. These single-turn samples will serve as the demonstration pool for in-context learning. To evaluate the effectiveness of our methods, we also have the CS teams to manually annotate some real multi-turn online sessions to serve as the test set. Each session queries $\mathcal{Q}$ will only have the last query $q_n$ labelled.

### 5.2  Metrics

We evaluate the accuracy of the methods based only on the label of the last query $q_n$ in each conversation session $\mathcal{Q}$. Other metrics which consider class imbalance are not used as the sampled sessions are expected to reflect the online traffic of each intent, thus better simulates the true online performance.

### 5.3  Baselines

To the best of our knowledge, there is no any existing work that directly addresses the challenge of multi-turn intent recognition with a large number of classes. This is a challenging task due to the lack of labelled data.

In our work, we present two intuitive and realistic approaches as baselines:

**Naive concatenation** All queries in a single session $\mathcal{Q}$ are concatenated using the $\circ$ operation mentioned above, and the concatenation result is fed into the single-turn model $\mathcal{M}_c$ for inference.

**Selective concatenation** In this approach, only one query from $\mathcal{C}_q$ is selected to be concatenated with $q_n$. The intuition is that not all history queries are helpful in understanding the last query, and the excessive use of them might introduce unwanted noise. A concatenation decision model is trained to select the most suitable history query. Depending on the model confidence, there might be cases where no expansion is needed at all.

### 5.4   Implementation Details

The traditional single-turn model, the retriever, and the concatenation decision model used are using backbone initialized with $\Phi_{XLMR}$, a multi-lingual domain specific XLM-RoBERTa-base model continued to be pre-trained with contrastive learning. We use AdamW to finetune the backbone and all other modules with a learning rate of 5e-6 and 1e-3, respectively. In LARA, the LLM used is vicuna-13b-v1.5 on Hugging Face with 13B parameters. All test are run on a single Nvidia V100 GPU card with a 32GB of GPU memory. The number of demonstrations $K$ retrieved for each intent is set at 10 in this experiment. Due to GPU memory constraint, the total number of tokens the in-context learning demonstrations can make up to are limited to 2300 tokens. If exceeded, the number of demonstrations in each candidate intent are pruned equally starting with the ones with the lowest cosine similarity scores to $q_{all}$. During inference time, if the generated intent doesn't match any of the provided options, the intent of $\mathcal{M}_c$ on $q_n$ will be considered as the final result.

## 6   Results and Discussions

Table 2 compares the performance of baselines and LARA. On average, LARA achieves better results than the baselines in any of the prompt variants used. *Naive concatenation* is not always more effective than *Selective concatenation*, showing that naively including all history queries will introduce noises which in turn jeopardizes the performance. However, pseudo-labelling the dataset used to train the concatenation decision model will need to be carefully carried out, and despite the extra steps, it will not necessarily be more effective than the naive method. LARA, on the other hand, can achieve good results on most dataset without any complex pseudo-labelling process. This also highlights the linguistic-adaptivity of the method on broad languages. The only market that it doesn't outperform the baselines is ID, which most probably can be attributed to the

language ability of open-sourced LLMs in handling the local slang and abbreviations in casual conversation. After all, the backbone model used in baselines are pre-trained directly on the in-domain chat log data, while the LLM models are used out-of-the-box.

| Model | Prompt | BR | ID | MY | PH | SG | TH | TW | VN | avg |
|---|---|---|---|---|---|---|---|---|---|---|
| Naive Concat. | - | 50.81% | 60.61% | 57.02% | 47.62% | 60.52% | 56.97% | 65.44% | 76.95% | 60.08% |
| Selective Concat. | - | 52.69% | **63.23%** | 60.20% | 51.32% | 56.99% | 57.77% | 64.02% | 74.10% | 60.97% |
| Vicuna-13B | $\mathcal{P}$ | 52.69% | 61.48% | **65.42%** | <u>54.50%</u> | 65.26% | 60.96% | **67.14%** | <u>77.90%</u> | <u>64.18%</u> |
| Vicuna-13B | $\mathcal{P}_{symbolic}$ | 51.88% | 60.00% | 64.57% | 53.97% | 65.26% | 58.96% | 65.44% | 74.67% | 62.92% |
| Vicuna-13B | $\mathcal{P}_{prepend}$ | <u>54.03%</u> | 61.75% | 64.50% | 53.44% | **65.94%** | <u>61.55%</u> | <u>66.86%</u> | 75.81% | 63.97% |
| Vicuna-13B | $\mathcal{P}_{formatted}$ | **55.65%** | <u>62.88%</u> | <u>64.71%</u> | **55.03%** | <u>65.40%</u> | **61.95%** | 66.86% | **78.10%** | **64.64%** |

Table 2: Performance of LARA compared to baselines, the average here is weighted on the number of test samples in each market. The best performance for each dataset is in boldface, while the second best is underlined.

Replacing the label names with non-related symbols in $\mathcal{P}_{symbolic}$ significantly hurts the performance of in-context learning. On the other hand, minimal changes to label names in $\mathcal{P}_{prepend}$ does not heavily impact the performance. In turn, the inference time is improved by 77%, from 0.75it/s to 1.32it/s on a single V100 card using Hugging Face python library. Interestingly, the model also stopped generating labels which cannot be matched with the options provided in demonstrations, while previously the rate is on average 1.6% using $\mathcal{P}$. Finally, we also tried a new prompt $\mathcal{P}_{formatted}$ based on $\mathcal{P}_{prepend}$. Only a very slight change to the context format is done, but it can outperform the other prompt variants in all dataset, suggesting that giving $\mathcal{C}_{\mathcal{Q}}$ a closer format to $\mathcal{E}$ and the targeted $q_n$ will be more beneficial in the context utilization. Besides, this also hints that the prompt could also be worked on more in the future as it is not extensively tuned in this work.

# 7    Ablation Studies

To validate our motivation and model design, we ablate our model components. The comparison is made on the original $\mathcal{P}$ prompt variant.

## 7.1    Traditional Single-turn Model

We tried not to use $\mathcal{M}_c$ to narrow down the intent candidates before retrieving the demonstrations. All demonstrations are directly retrieved based on their cosine similarity to $q_{all}$. As the original number of retrieved demonstrations for each $\mathcal{Q}$ is dynamic according to number of queries $n$, to ensure fairness, the number of retrieved demonstrations in this variant will also be $K \times n$. From Fig ??, we see that the quality of in-context learning is adversely impacted when the number of intents included in demonstrations is too high.
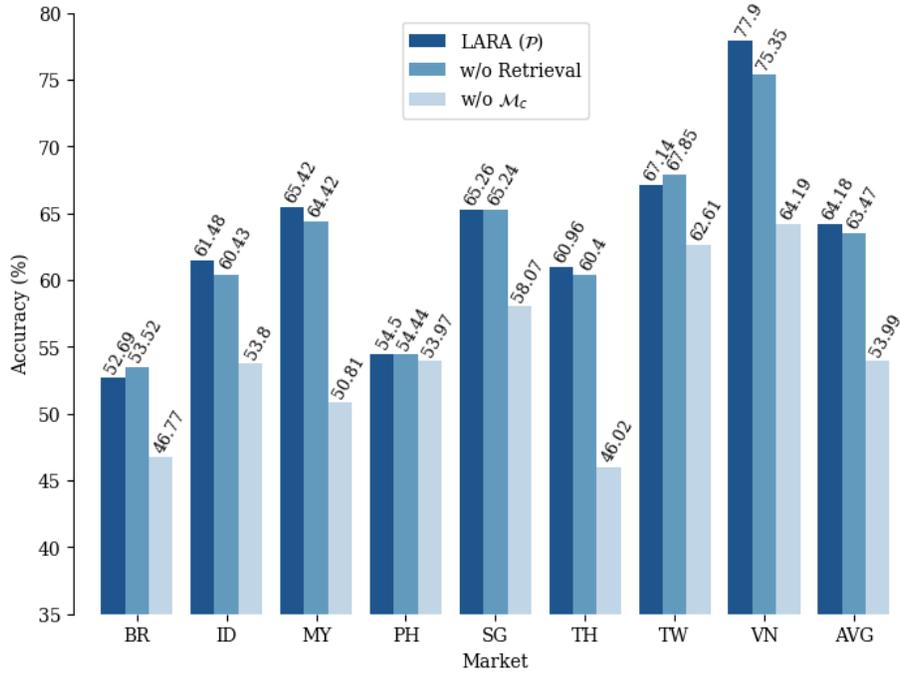
Fig. 3: Ablation on different components of LARA. The last row shows the standard deviations of performance 'w/o Retrieval' over 10 runs. The size of the demonstration pool and the number of intents in each dataset are also included here for the ease of reference.

### 7.2   Retrieval Component

Demonstration selection could also have significant impact on the performance. Thus, we also tried to remove the retrieval component, and randomly sample the demonstrations for each intent. The results are reported with 10 runs on the random sampling. Based on Table **??**, the overall performance will worse if there is no retrieval component, specifically when there are big demonstration pool or high number of intents. Thus, if there is no one good strategy in pruning the demonstration pool, the easier way is to retrieval demonstrations based on their similarity to the task input.

## 8   Conclusion

This paper introduced LARA, a framework that leverages Linguistic-Adaptive Retrieval-Augmentation to address multi-turn intent classification challenges through zero-shot settings across multiple languages. Unlike other supervised Fine-Tuning (SFT) models, which require a multi-turn dialogue set that is hard to collect. Our method only requires a single-turn training set to train a conventional XLM model. It then combines it with an innovative in-context retrieval

augmentation for multi-turn intent classification. LARA demonstrated a notable improvement in accuracy and efficiency, marking a significant advancement in the field of conversational AI.

The empirical results underscore LARA's capability to enhance intent classification accuracy by 3.67% over existing methods while reducing inference time, thus facilitating real-time application adaptability. Its strategic approach to managing extensive intent varieties without exhaustive dataset requirements presents a scalable solution for complex, multi-lingual conversational systems.

# References

1. Brown, T.B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D.M., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., Amodei, D.: Language models are few-shot learners. ArXiv **abs/2005.14165** (2020), `https:// api.semanticscholar.org/CorpusID:218971783`
2. Budzianowski, P., Wen, T.H., Tseng, B.H., Casanueva, I., Ultes, S., Ramadan, O., Gašić, M.: Multiwoz–a large-scale multi-domain wizard-of-oz dataset for task-oriented dialogue modelling. arXiv preprint arXiv:1810.00278 (2018)
3. Cong, X., Yu, B., Liu, T., Cui, S., Tang, H., Wang, B.: Inductive unsupervised domain adaptation for few-shot classification via clustering. In: Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2020, Ghent, Belgium, September 14–18, 2020, Proceedings, Part II. pp. 624–639. Springer (2021)
4. Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, É., Ott, M., Zettlemoyer, L., Stoyanov, V.: Unsupervised cross-lingual representation learning at scale. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics. pp. 8440–8451 (2020)
5. Conneau, A., Schwenk, H., Barrault, L., Lecun, Y.: Very deep convolutional networks for text classification. arXiv preprint arXiv:1606.01781 (2016)
6. Ghosal, D., Majumder, N., Mihalcea, R., Poria, S.: Exploring the role of context in utterance-level emotion, act and intent classification in conversations: An empirical study. In: Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021. pp. 1435–1449 (2021)
7. Kenton, J.D.M.W.C., Toutanova, L.K.: Bert: Pre-training of deep bidirectional transformers for language understanding. In: Proceedings of NAACL-HLT. pp. 4171–4186 (2019)
8. Kim, Y.: Convolutional neural networks for sentence classification. In: Moschitti, A., Pang, B., Daelemans, W. (eds.) Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP). pp. 1746–1751. Association for Computational Linguistics, Doha, Qatar (Oct 2014). https://doi.org/10.3115/v1/D14-1181, `https://aclanthology.org/D14-1181`
9. Kyriakopoulou, A., Kalamboukis, T.: Text classification using clustering. In: Proceedings of the Discovery Challenge Workshop at ECML/PKDD 2006. pp. 28–38. Citeseer (2006)

10. Lee, J., Lee, W.: Compm: Context modeling with speaker's pre-trained memory tracking for emotion recognition in conversation. In: Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. pp. 5669–5679 (2022)

11. Liu, J., Albrethsen, J., Goh, L., Yau, D., Lim, K.H.: Spatial-temporal graph representation learning for tactical networks future state prediction. In: Proceedings of The International Joint Conference on Neural Networks (IJCNN'24) (2024)

12. Madotto, A., Lin, Z., Winata, G.I., Fung, P.: Few-shot bot: Prompt-based learning for dialogue systems. ArXiv **abs/2110.08118** (2021), `https://api.semanticscholar.org/CorpusID:239009514`

13. Mo, F., Nie, J., Huang, K., Mao, K., Zhu, Y., Li, P., Liu, Y.: Learning to relate to previous turns in conversational search. Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (2023), `https://api.semanticscholar.org/CorpusID:259076134`

14. Pasupat, P., Zhang, Y., Guu, K.: Controllable semantic parsing via retrieval augmentation. ArXiv **abs/2110.08458** (2021), `https://api.semanticscholar.org/CorpusID:239016988`

15. Qin, X., Wu, Z., Zhang, T., Li, Y., Luan, J., Wang, B., Wang, L., Cui, J.: Bert-erc: Fine-tuning bert is enough for emotion recognition in conversation. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 13492–13500 (2023)

16. Qu, C., Yang, L., Croft, W.B., Zhang, Y., Trippas, J.R., Qiu, M.: User intent prediction in information-seeking conversations. In: Proceedings of the 2019 Conference on Human Information Interaction and Retrieval. CHIIR '19, ACM (Mar 2019). https://doi.org/10.1145/3295750.3298924, `http://dx.doi.org/10.1145/3295750.3298924`

17. Reimers, N., Gurevych, I.: Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084 (2019)

18. Ren, F., Xue, S.: Intention detection based on siamese neural network with triplet loss. IEEE Access **8**, 82242–82254 (2020). https://doi.org/10.1109/ACCESS.2020.2991484

19. Rubin, O., Herzig, J., Berant, J.: Learning to retrieve prompts for in-context learning. ArXiv **abs/2112.08633** (2021), `https://api.semanticscholar.org/CorpusID:245218561`

20. Shen, W., Wu, S., Yang, Y., Quan, X.: Directed acyclic graph network for conversational emotion recognition. In: Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers). pp. 1551–1560 (2021)

21. Shin, R., Durme, B.V.: Few-shot semantic parsing with language models trained on code. In: North American Chapter of the Association for Computational Linguistics (2021), `https://api.semanticscholar.org/CorpusID:245218525`

22. Shin, R., Lin, C.H., Thomson, S., Chen, C.C., Roy, S., Platanios, E.A., Pauls, A., Klein, D., Eisner, J., Durme, B.V.: Constrained language models yield few-shot semantic parsers. ArXiv **abs/2104.08768** (2021), `https://api.semanticscholar.org/CorpusID:233297024`

23. Wallach, H.M.: Topic modeling: beyond bag-of-words. In: Proceedings of the 23rd international conference on Machine learning. pp. 977–984 (2006)

24. Weld, H., Huang, X., Long, S., Poon, J., Han, S.C.: A survey of joint intent detection and slot-filling models in natural language understanding (2021)

25. Wu, T.W., Su, R., Juang, B.H.: A context-aware hierarchical bert fusion network for multi-turn dialog act detection (2021)

26. Wu, Z., Lu, B.R., Hajishirzi, H., Ostendorf, M.: Dialki: Knowledge identification in conversational systems through dialogue-document contextualization. In: Conference on Empirical Methods in Natural Language Processing (2021), `https://api.semanticscholar.org/CorpusID:237485380`

27. Xie, T., Wu, C.H., Shi, P., Zhong, R., Scholak, T., Yasunaga, M., Wu, C.S., Zhong, M., Yin, P., Wang, S.I., Zhong, V., Wang, B., Li, C., Boyle, C., Ni, A., Yao, Z., Radev, D.R., Xiong, C., Kong, L., Zhang, R., Smith, N.A., Zettlemoyer, L., Yu, T.: Unifiedskg: Unifying and multi-tasking structured knowledge grounding with text-to-text language models. ArXiv **abs/2201.05966** (2022), `https://api.semanticscholar.org/CorpusID:246016124`

28. Xu, P., Sarikaya, R.: Contextual domain classification in spoken language understanding systems using recurrent neural network. In: 2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). pp. 136–140 (2014). https://doi.org/10.1109/ICASSP.2014.6853573

29. Yang, Z., Yang, D., Dyer, C., He, X., Smola, A., Hovy, E.: Hierarchical attention networks for document classification. In: Proceedings of the 2016 conference of the North American chapter of the association for computational linguistics: human language technologies. pp. 1480–1489 (2016)

30. Yao, L., Mao, C., Luo, Y.: Graph convolutional networks for text classification. In: Proceedings of the AAAI conference on artificial intelligence. vol. 33, pp. 7370–7377 (2019)

31. Ye, Z., Geng, Y., Chen, J., Chen, J., Xu, X., Zheng, S., Wang, F., Zhang, J., Chen, H.: Zero-shot text classification via reinforced self-training. In: Proceedings of the 58th annual meeting of the association for computational linguistics. pp. 3014–3024 (2020)

32. Yu, D., He, L., Zhang, Y., Du, X., Pasupat, P., Li, Q.: Few-shot intent classification and slot filling with retrieved examples. In: North American Chapter of the Association for Computational Linguistics (2021), `https://api.semanticscholar.org/CorpusID:233219405`

33. Yuan, C., jie Zhou, W., Li, M., Lv, S., Zhu, F., Han, J., Hu, S.: Multi-hop selector network for multi-turn response selection in retrieval-based chatbots. In: Conference on Empirical Methods in Natural Language Processing (2019), `https://api.semanticscholar.org/CorpusID:202776649`

34. Zhang, H., Xiao, L., Chen, W., Wang, Y., Jin, Y.: Multi-task label embedding for text classification. arXiv preprint arXiv:1710.07210 (2017)

35. Zhang, J., Ye, Y., Zhang, Y., Qiu, L., Fu, B., Li, Y., Yang, Z., Sun, J.: Multi-point semantic representation for intent classification. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 9531–9538 (2020)

# 9    Appendix

## 9.1    Prompt for ICL ($\mathcal{P}$)

Prompt for ICL ($\mathcal{P}$)

```
# Task Description
A chat between a curious user and an artificial intelligence
    assistant. The assistant gives helpful, detailed, and polite
    answers to the user's questions. USER: Determine the intent
```

```
    for the targetted message from the examples, you must use the
    context in the history messages to arrive at the best answer.
# Examples
[Content] Similar Question [Intent] Intent_name_1
[Content] Similar Question [Intent] Intent_name_2
[Content] Similar Question [Intent] Intent_name_3

# Note
DO NOT create new intent on your own, you must strictly use the
    intents in the examples.
DO NOT provide any explanation.
Output ONLY ONE intent for the targgetted message.
Consider the context from previous messages if the targetted
    message is unclear.

# Context
message 1: User's query
message 2: User's query with Entity
[Content] Last user's query

# Output
ASSISTANT: [Intent] <Model generated Intent_name>
```

## 9.2  Prompt for ICL ($\mathcal{P}_{symbolic}$)

In $\mathcal{P}_{symbolic}$ the original label name $l$ of each intent in $\mathcal{E}_{symbolic}$ are replaced with single-token symbols, e.g. 'A', 'B', ..., which bear no meaning to the intents they represented. Explanation will be made in the instruction prompt $\mathcal{T}_{symbolic}$ to link the symbols back to their original intent label $y_j$, and the model is instructed to generated the symbols instead of full label names.

<div align="center">Prompt for ICL ($\mathcal{P}_{symbolic}$)</div>

```
# Task Description
Content is Same as P

# Examples
[Content] Similar Question [Intent] A
[Content] Similar Question [Intent] B
<omitted>
[Content] Similar Question [Intent] B

# Intent options
A is Intent_name_1
B is Intent_name_2
```

```
# Note
Content is Same as P

# Context
Format is same as P

# Output
ASSISTANT: [Intent]
```

## 9.3  $\mathcal{P}_{prepend}$

In $\mathcal{P}_{prepend}$, representative symbols for each intent will be prepend to the original label name $l$, such that they are separated by an extra character as boundary, e.g. label "logistics>how long will it take to receive order?" will be represented as "A>logistics>how long will it take to receive order?". Note that the instruction prompt $\mathcal{T}$ remains the same, the trick is to limit the model generation token count to 1 on API level.

Prompt for ICL ($\mathcal{P}_{prepend}$)

```
# Task Description
Content is Same as P

# Examples
[Content] Similar Question [Intent] A>Intent_name_1
[Content] Similar Question [Intent] B>Intent_name_2
<omitted>
[Content] Similar Question [Intent] B>Intent_name_2

# Note
Content is Same as P

# Context
Format is same as P

# Output
ASSISTANT: [Intent] B
```

## 9.4  $\mathcal{P}_{formatted}$

Prompt for ICL ($\mathcal{P}_{formatted}$)

```
# Task Description
Content is Same as P
```

```
# Examples
```
Format is same as $\mathcal{P}_{prepend}$

```
# Note
```
Content is Same as $\mathcal{P}_{prepend}$

```
# Context
[History message 1] User's query
[History message 2] User's query with Entity
[Content] that is the order id
```

```
# Output
ASSISTANT: [Intent]
```