

Representative Chain-of-Reasoning for Aspect Sentiment Quad Prediction

Anonymous ACL submission

Abstract

Aspect Sentiment Quad Prediction (ASQP) is a crucial sentiment analysis task that has attracted increasing attention. The most recent studies focus on generating complete sentiment quadruples through end-to-end generative models. However, these methods heavily depend on labeled data quality and quantity, performing poorly in low-resource scenarios and less suitable for real-world applications. To address these issues, we propose a novel *Representative Chain-of-Reasoning* framework (RCR), with the aim of providing representative knowledge for large language models (LLMs) and fully activating their reasoning capabilities for ASQP. Specifically, we develop a Chain Prompting (ChaPT) module to decompose the ASQP task into three subtasks using the step-by-step reasoning mechanism. Then, a Representative Demonstration Retriever (RepDR) is introduced to provide ChaPT with representative demonstrations, balancing diversity and similarity, and enhancing the reasoning capabilities of LLMs at each step. Experimental results confirm the superiority of RCR in both zero-shot and few-shot scenarios, significantly surpassing existing counterparts.

1 Introduction

Given a review text, Aspect Sentiment Quad Prediction (ASQP) aims to predict a comprehensive sentiment view in the form of quadruples (Zhang et al., 2022a, 2023b), each consisting of aspect category, aspect term, opinion term, and sentiment polarity, denoted as (c, a, o, s) . For example, given the review sentence, "*The food is great and the environment is even better.*", the ASQP task requires predicting two sentiment quadruples: $(\textit{food quality}, \textit{food}, \textit{great}, \textit{positive})$ and $(\textit{ambiance general}, \textit{environment}, \textit{better}, \textit{positive})$. ASQP is a challenging task due to the complexity of sentence structure and the diversity of sentiment expressions, making it difficult to recognize all sentiment quadruples.

Recently, the end-to-end generative models have been extensively applied to solve the ASQP task by generating sentiment quadruples directly from the review text and achieved promising results (Peper and Wang, 2022). A successful application is to construct sequences in natural language format as generation targets, including annotated sentences (Zhang et al., 2021b), paraphrased sentences (Liu et al., 2021; Zhang et al., 2021a; Hu et al., 2022a), and sentiment element sequences (Zhang et al., 2021c). Furthermore, sentiment clues within sentences have been utilized to promote the quadruple generation (Bao et al., 2022). For example, Mao et al. (2022) introduced a parallel generation framework to capture more sentiment information through beam search. Gou et al. (2023) enhanced the model’s expressive capability by increasing the output views by adjusting the generation order of quadruples. Despite their potential, a notable issue is that these models are less suitable in low-resource scenarios (Hu et al., 2022a; Gou et al., 2023). That is because generative models heavily rely on the scale and quality of the labeled dataset while annotating datasets is costly and time-consuming in practical applications.

With the rise of In-Context Learning (ICL), addressing the ASQP task by generative models in zero-shot and few-shot scenarios becomes feasible (Wang et al., 2023c; Zhang et al., 2023a,c). Sun et al. (2023) propose a multi-LLM negotiation strategy, demonstrating LLMs’ ability to solve sentiment analysis problems involving complex contexts (e.g., clauses and irony) under zero-shot conditions. Moreover, the Three-Hop Reasoning (THOR) CoT framework (Fei et al., 2023) achieves state-of-the-art results in implicit sentiment analysis tasks. However, existing research lacks a discussion on applying LLMs to ASQP tasks, and the reasoning capabilities of LLMs are underutilized.

To this end, we propose a Representative Chain-of-Reasoning framework (RCR) that aims to pro-

vide representative knowledge for LLMs and fully activate their reasoning capabilities for the ASQP task. Inspired by the Chain-of-Thought (CoT) prompting (Wei et al., 2022; Zhou et al., 2022; Zhang et al., 2022b), we first introduce a Chain Prompt (ChaPT) module to decompose the one-step ASQP task into three sub-steps, where each step progressively infers aspect-opinion pairs, category-aspect-opinion triplets, and complete quadruples. Hence, a complete sentiment view is obtained through step-by-step reasoning, effectively reducing the ASQP task’s complexity. Additionally, considering LLM’s reasoning capability is influenced greatly by the quality of demonstrations (Lee et al., 2022; Min et al., 2022; Wang et al., 2023b), we develop a Representative Demonstration Retriever (RepDR) module to provide ChaPT with representative demonstrations, balancing diversity and similarity, and thus enhancing their reasoning capabilities at each step. Specifically, we first paraphrase the sentiment quadruples into natural sentences (Zhang et al., 2021a) and calculate their semantic similarities using SBERT (Reimers and Gurevych, 2019). Based on semantic similarities, the triplet that contains an anchor sentence, a positive sentence, and a negative sentence is picked for further fine-tuning this SBERT model. Hence, this fine-tuned SBERT model is good at retrieving representative demonstrations that possess semantic information of different attributes (Wang et al., 2022a; Shi et al., 2023; Qin et al., 2023).

In summary, the main contributions of this work are as follows:

- We introduce ChaPT, a prompting framework based on the chain-of-reasoning concept, which mitigates task complexity through task decomposition and step-by-step reasoning, fully leveraging the reasoning capabilities of LLMs.
- We use the RepDR module to retrieve demonstrations, providing more representative prior information for model reasoning. To the best of our knowledge, this work is the first to propose retrieving both diversity and similarity samples as demonstrations.
- Experimental results show that our proposed model demonstrates superiority in both zero-shot and few-shot scenarios and greatly surpasses existing counterparts.

2 Methodology

2.1 Problem Definition

The ASQP task is defined as follows: given a sentence X , the model predicts all aspect-based sentiment quadruples (Cai et al., 2021; Zhang et al., 2021a), each formulated as (c, a, o, s) which corresponds to aspect category, aspect term, opinion term, and sentiment polarity, respectively. The aspect category c is part of a predefined category set U_c . The aspect term a is the target of opinion. The opinion term o is the subjective statement. Moreover, the sentiment polarity s belongs to the predefined sentiment set $U_s \in \{positive, neutral, negative\}$. Notably, aspect a is generally within the text scope of sentence X , and if aspect a is not explicitly mentioned, it is represented by the specific tag NULL.

Xie et al. (2021) believe that ICL infers conditional probabilities of the predictive target from the prompt, formulated as $p(y | prompt) = \int_{prompt} p(y | prompt) d(prompt)$, where y represents the prediction target and $d(prompt)$ represents the prompt set. ICL infers the maximum probability of generating y by integrating y over the $prompt$. Therefore, we can use ICL to model the ASQP task as

$$\hat{y} = \operatorname{argmax} p(y | X, prompt) \quad (1)$$

where \hat{y} represents all quadruples, and Zhang et al. (2023c) attempt to construct the following standard *prompt* paradigm as input to the LLMs:

Given the sentence X , tag all the (category, aspect, opinion, sentiment) quadruples.

2.2 Representative Chain-of-Reasoning

To fully activate the reasoning capabilities of LLMs for the ASQP task, we propose a Representative Chain-of-Reasoning framework (RCR) consisting of two sub-modules: Chain Prompt Framework (ChaPT) and Representative Demonstration Retriever (RepDR). The former is designed to decompose the ASQP task into three subtasks, while the latter is responsible for providing representative demonstrations to enhance LLMs’ reasoning capabilities.

2.2.1 Chain Prompt Framework

Inspired by the impressive reasoning capabilities demonstrated by Chain of Thought (CoT) in han-

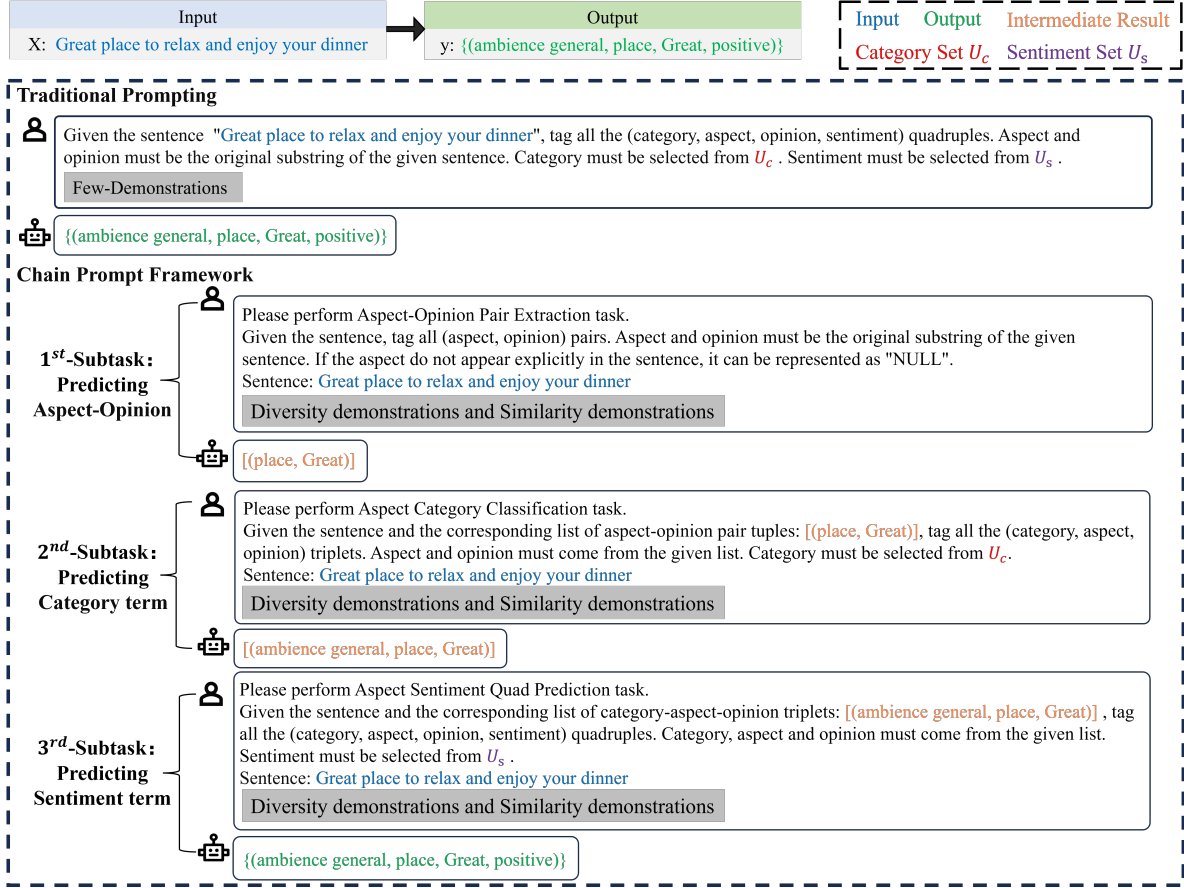


Figure 1: An illustration of our ChaPT framework for Aspect Sentiment Quad Prediction task.

177 dling complex tasks(Wei et al., 2022), we propose
 178 the Chain Prompt framework(ChaPT), shown in
 179 Figure 1, to address the ASQP task by decompos-
 180 ing a one-step ASQP solution into three subtasks.
 181 The details are as follows.

182 Subtask 1. Aspect-Opinion Pair Extraction

183 Empirical studies find that extracting a single as-
 184 pect or opinion alone would ignore their pairwise
 185 relationships, leading to pairing errors (Chen et al.,
 186 2020; Zhao et al., 2020). Therefore, instead of finer-
 187 grained aspect term or opinion term extraction sub-
 188 task, we first consider predicting all aspect-opinion
 189 pairs appearing in the sentence. Mathematically,
 190 this subtask is formulated as:

$$191 \hat{Z}_1 = \operatorname{argmax} p(y | X, \text{prompt}_1) \quad (2)$$

192 where \hat{Z}_1 denotes predicted aspect-opinion
 193 pairs, the template of prompt_1 is defined as:

194 Given the sentence X, tag all the (aspect,
 195 opinion) pairs.

196 Subtask 2. Aspect Category Classification

197 Based on X and the intermediate results \hat{Z}_1 , we
 classify category c from the predefined set U_c and

198 obtain the category-aspect-opinion triplets Z_2 . This
 199 process is represented as:

$$200 \hat{Z}_2 = \operatorname{argmax} p(y | X, Z_1, \text{prompt}_2) \quad (3)$$

201 the template of prompt_2 is as follows.

202 Given the sentence X and the corresponding list
 203 of aspect-opinion pair tuples \hat{Z}_1 , tag all the
 204 (category, aspect, opinion) triplets.

205 Subtask 3. Aspect Sentiment Quad Prediction

206 Based on the intermediate results \hat{Z}_2 , we finally
 predict the complete quadruples y . The final step
 is denoted as:

$$207 \hat{y} = \operatorname{argmax} p(y | X, Z_2, \text{prompt}_3) \quad (4)$$

208 and the template of prompt_3 is as follows.

209 Given the sentence X and the corresponding list of
 210 category-aspect-opinion triplets \hat{Z}_2 , tag all the
 211 (category, aspect, opinion, sentiment) quadruples.

212 2.2.2 Representative Demonstration Retriever

In few-shot scenarios, we propose the Representa-
 tive Demonstration Retriever (RepDR), a demon-

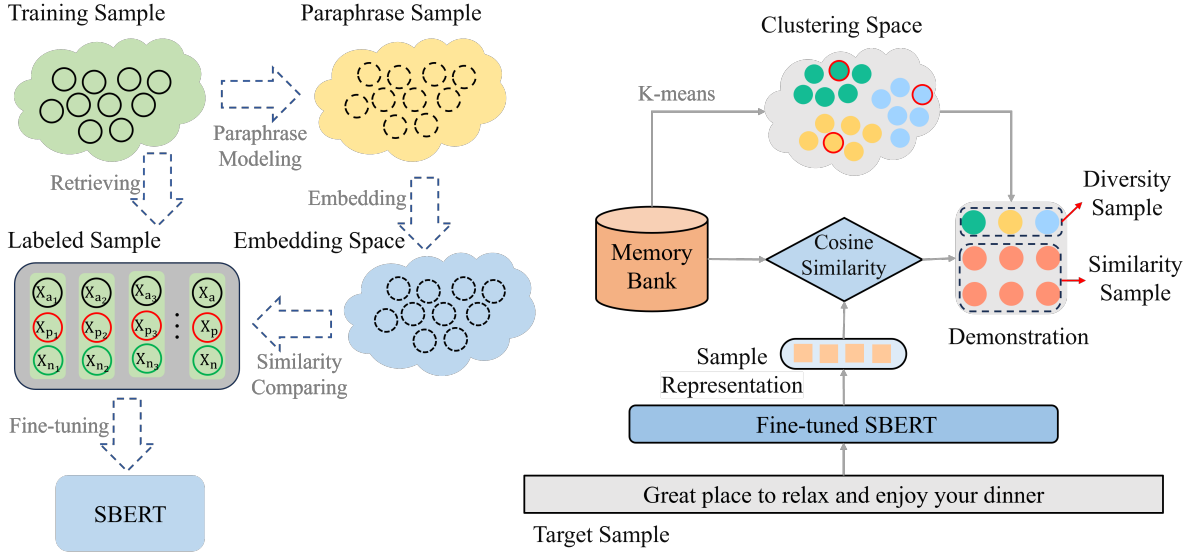


Figure 2: The proposed RepDR module consists of two stages. The first stage generates training samples with pre-trained SBERT, and fine-tunes SBERT using these samples. The second stage generates text embeddings using the fine-tuned model and utilizes clustering and cosine similarity to produce representative demonstrations.

213 stration retriever that balances diversity and simi-
 214 larity. First, we explain the generation of training
 215 sample, then describe training the model with the
 216 labeled sample, and finally show using the trained
 217 model to retrieve representative demonstrations. As
 218 shown in Figure 2.

219 **Generating Training Sample** Since we utilize
 220 clustering and similarity comparison for demon-
 221 stration retrieval, it is vital to train a model that
 222 precisely captures the similarity among sentence
 223 pairs. We choose SBERT (Reimers and Gurevych,
 224 2019), a BERT-based text embedding model, as our
 225 target model. By fine-tuning SBERT with gener-
 226 ated triplet data, we enhance its ability to capture
 227 semantic similarity between sentences. Triplet data
 228 consists of an anchor sentence, a positive sentence,
 229 and a negative sentence without additional labels.
 230 Inspired by paraphrase generation (Zhang et al.,
 231 2021a; Gou et al., 2023; Hu et al., 2022b), we pro-
 232 pose modeling paraphrases for the training set by
 233 linearizing sentiment quadruples (c, a, o, p) into
 234 natural sentences I , as shown in Figure 3.

235 Paraphrase modeling allows us to focus on
 236 the quadruples and ignore unnecessary details in
 237 sentences. We compute their embeddings using
 238 a pre-trained SBERT model for the resulting
 239 paraphrase set U_I . Then, we compare the seman-
 240 tic similarity between these paraphrases by employ-
 241 ing Cosine Similarity. We retrieve the sentence most
 242 similar to the anchor sentence X_a as the positive
 243 sentence X_p and the least similar sentence as

Sentence-1	Our teenage kids love it , too .
Quadruplet-1	(restaurant general, NULL, love , positive)
↓	↓
Paraphrase-1	restaurant general is positive because it is love
Sentence-2	The only thing more wonderful than the food (which is exceptional) is the service .
Quadruplet-2	(food quality, food, exceptional, positive), (service general, service, wonderful, positive)
↓	↓
Paraphrase-2	food quality is positive because food is exceptional and service general is positive because service is wonderful

Figure 3: Two examples of paraphrase modeling. Notably, if the aspect is not explicitly mentioned, it is represented by the implicit pronoun "it". If a sentence contains multiple sentiment quadruples, the paraphrases are concatenated using "and".

244 the negative sentence X_n , thus constructing
 245 triplet training samples (X_a, X_p, X_n) . Notably,
 246 the sentences referred to here are the original
 247 sentences, not the paraphrase sentences, to ensure
 248 the model focuses more attention on quadruples in
 249 the original sentence.

250 **Training Model** We fine-tune the SBERT model
 251 using the typical triplet network. Given triplet data
 252 I , we utilize SBERT (selected mpnet (Song et al.,
 253 2020)) to encode X_a , X_p , and X_n , obtaining em-
 254 beddings O_a , O_p , and O_n . Fine-tuning aims to
 255 minimize the distance between O_a and O_p while
 256 maximizing the distance between O_a and O_n . We
 257 use triplet loss as the loss function, as shown below
 258

in Equation 5:

$$\begin{aligned} \mathcal{L}(O_a, O_p, O_n) \\ = \max(d(O_a, O_p) - d(O_a, O_n) + \alpha, 0) \end{aligned} \quad (5)$$

where $d(O_a, O_p) = \|O_a - O_p\|_2$ represents the Euclidean Distance between embeddings. The hyperparameter α specifies the expected difference between $d(O_a, O_p)$ and $d(O_a, O_n)$.

Retrieving Demonstration Firstly, we use the fine-tuned SBERT to encode all training set sentences into text embeddings and store them in a Memory bank (Wu et al., 2018) to avoid redundant computations. Secondly, we measure the similarity of demonstrations utilizing Cosine Similarity, comparing target samples with the training set to extract the top-k most similar samples. Finally, we propose a diversified demonstration retrieval scheme based on K-means clustering (Arthur et al., 2007). We evaluate the optimal number of clusters by calculating the Silhouette Score and find that the optimal number for both datasets is 3 (see Appendix A). Based on this result, we perform K-means clustering and select the samples closest to the cluster centers as diversity samples that highlight key characteristics.

3 Experiments

3.1 Datasets

We conducted experiments on two public restaurant datasets, Rest15 and Rest16, from the SemEval task (Pontiki et al., 2015, 2016). These datasets, with multiple annotations (Peng et al., 2020; Wan et al., 2020), were aligned by Zhang et al. (2021a) and ultimately served as the standard datasets for the ASQP task. Each sample contains one or more sentiment quadruples. The statistics are shown in Table 1.

3.2 Implementation Details

We utilized two OpenAI models, including ChatGPT (Open, 2022) (gpt-3.5-turbo3) and the newly released GTP-4 (Achiam et al., 2023) (gpt-4o), as the backbone for the ChaPT framework (Section 2.2.1) to evaluate its effectiveness under zero-shot conditions. The temperature for all models was set to 0 to ensure stable predictions.

Moreover, for few-shot scenarios, we employed all-mpnet-base-v1 (Song et al., 2020) as pre-trained SBERT (Section 2.2.2), using a typical triplet network for fine-tuning. During fine-tuning, we used a

	Rest15				Rest16			
	SEN	POS	NUE	NEG	SEN	POS	NUE	NEG
Train	834	1005	34	315	1264	1369	62	558
Dev	209	252	14	81	316	341	23	143
Test	537	453	37	305	544	583	40	176

Table 1: Dataset statistics for Rest15 and Rest16. SEN, POS, NUE and NEG represent the number of sentences, positive, neutral, and negative quadruples, respectively.

batch size of 64, a learning rate of 2e-5, and 5 training epochs. The hyperparameter α of the model was set to 5. Additional implementation details of generative models in low-resource scenarios are provided in Appendix B.

During demonstration retrieval, we used the model at the best checkpoint to re-encode sentences for text similarity comparison and clustering, obtaining demonstrations with similarity and diversity. We only considered three k-shot settings: 1-shot, 5-shot, and 10-shot. For each setting, we maintained a constant number of diversity demonstrations and adjusted the number of similarity demonstrations to achieve k-shot. For example, in the 1-shot scenario, we retrieved 3 diversity samples and the 1 most similar sample.

3.3 Baselines

We employ generative-based models and ICL-based large language models as our comparative baselines. For generative methods, we select the following four models:

- **GAS** (Zhang et al., 2021b) The first attempt to use generative methods to handle aspect-based sentiment analysis, we modify it to use sentiment quadruple sequences as target sequences.
- **Paraphrase** (Zhang et al., 2021a) A paraphrasing modeling framework, using paraphrased sentences as training targets to generate sentiment quadruples end-to-end.
- **DLO/ILO** (Hu et al., 2022a) Selecting the appropriate quadruple generation order as a data augmentation method for paraphrase generation.
- **MVP** (Gou et al., 2023) Enhances the model predictive capability by increasing output views of different quadruple generation orders.

	Rest15				Rest16			
	0-shot	1-shot	5-shot	10-shot	0-shot	1-shot	5-shot	10-shot
Generative-based Baselines								
GAS [†] (Zhang et al., 2021b)	-	4.43	10.65	13.82	-	2.24	16.04	19.03
Paraphrase [†] (Zhang et al., 2021a)	-	7.78	11.53	18.60	-	2.36	12.85	16.34
DLO [†] (Hu et al., 2022a)	-	6.79	13.07	18.92	-	1.90	17.68	28.95
ILO [†] (Hu et al., 2022a)	-	7.25	14.85	20.99	-	2.41	15.71	21.32
MvP [†] (Gou et al., 2023)	-	9.33	18.54	22.82	-	3.62	21.51	29.24
Prompt-based Baselines								
w/ GPT-3.5								
LMMs for SA [†] (Zhang et al., 2023c)	7.77	27.83	26.86	25.74	10.06	28.45	38.63	37.13
THOR [†] (Fei et al., 2023)	10.21	22.13	27.24	23.51	14.11	28.62	37.26	36.04
RCR(Ours)	13.82	28.46	31.44	32.29	19.00	30.60	41.51	42.09
w/ GPT-4								
LMMs for SA [†] (Zhang et al., 2023c)	<u>32.40</u>	<u>35.63</u>	<u>37.65</u>	<u>36.72</u>	35.87	<u>40.56</u>	<u>42.07</u>	<u>40.32</u>
THOR [†] (Fei et al., 2023)	30.57	35.01	36.22	35.81	<u>36.37</u>	38.02	41.58	39.97
RCR(Ours)	33.01	39.78	42.26	44.33	38.07	40.97	48.08	51.23

Table 2: Report the model’s experimental results under zero-shot and few-shot settings. F1 score is used as the evaluation metric. The best and second-best results are indicated in bold and underlined, respectively. The baseline methods, marked with[†], follow the few-shot settings of this work (Zhang et al., 2023c), where k-shot represents sampling k examples for each aspect category.

For ICL methods, we selected the following re-search approaches:

- **LMMs for SA** (Xu et al., 2024) A comprehensive study of sentiment analysis using LLMs, including Flan-T5, FLan-UL2, T5, and GPT-3.5.
- **THOR** (Fei et al., 2023) A Three-hop reasoning (THOR) CoT framework for addressing implicit sentiment analysis issues. In our setup, it serves as one of the benchmarks for ICL methods by modifying the prompt.

Experimental results for these supervised methods are derived from the base pre-trained models (BERT or T5) to ensure a fair comparison.

4 Results and Discussions

4.1 Zero-shot and Few-shot Results

The experimental results are shown in Table 2. Notably, in the zero-shot scenario, the T5-based (Raf-fel et al., 2020) generative model struggled with ASQP tasks and failed to generate effective results. However, the performance gradually improved as the number of samples increased, highlighting the importance of high-quality labeled data for generative models. Compared to the best generative model baseline MVP, the ICL-based LLMs (LLMs for SA) showed significant performance improvements in both zero-shot and few-shot scenarios. For

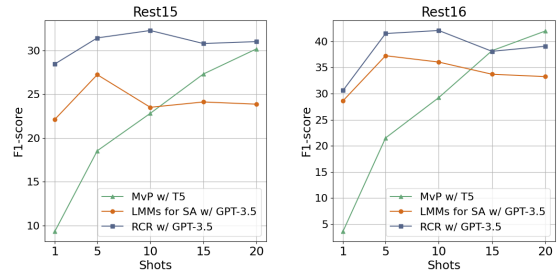


Figure 4: The evaluation curve of the model with varying sample sizes.

GPT-3.5, under few-shot conditions, the average F1 score gained on the Rest15 and Rest16 datasets was 9.91% and 16.61%, respectively. This demonstrates the great potential of LLMs in ASQP tasks. Furthermore, our proposed Representative Chain-of-Reasoning (RCR) framework achieved the best performance with both GPT-3.5 and GPT-4 compared to the original ICL baselines. Specifically, with GPT-3.5, the average F1 score gained on the Rest15 and Rest16 datasets were 4.45% and 4.73%, respectively. With GPT-4, the F1 scores improved by an average of 4.24% and 4.88%. This indicates that the RCR framework provides sufficient prior information for LLMs, fully leveraging their reasoning capabilities in ASQP tasks.

4.2 Ablation Study

We conducted ablation experiments to further validate our RCR framework’s effectiveness. In a

Methods	Rest15			Rest16		
	Pre	Rec	F1	Pre	Rec	F1
RCR	27.99	35.85	31.44	36.82	47.56	41.51
RCR w/o [RepDR]	23.39	34.21	27.78	24.53	32.42	27.92
RCR w/o [ChaPT]	21.83	27.30	24.26	28.16	36.29	31.71
RCR w/o [RepDR,ChaPT]	21.25	26.68	23.66	25.10	30.16	27.40

Table 3: The results of ablation study.

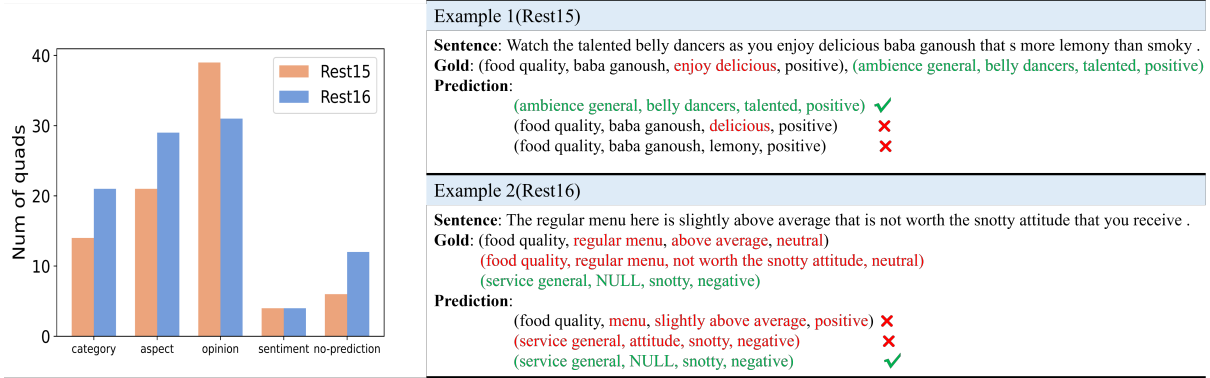


Figure 5: Statistics of error types and two examples of prediction errors. Notably, no-prediction indicates the samples where LLMs made no predictions.

5-shot scenario, we analyzed the impact on the results by removing individual modules, with results shown in Table 3. ChaPT decomposes the ASQP task into subtasks, reducing the complexity of LLM reasoning. RepDR is responsible for providing more accurate prior semantic knowledge to LLMs through demonstration retrieval. The results indicate that removing any module significantly reduces RCR performance, demonstrating the effectiveness of ChaPT and RepDR in stimulating the reasoning capabilities of LLMs.

Furthermore, We observed performance differences across the Rest15 and Rest16 datasets when removing specific modules. For instance, removing the RepDR module resulted in a 13.59% decrease in F1 score for Rest16, but only a 3.66% decrease for Rest15. This indicates that different datasets have varying dependencies on the ChaPT and RepDR modules, reflecting the distinct knowledge support these two components provide to LLMs.

4.3 Influence of Different Sample Sizes

Our preliminary research reveals that the LLMs' reasoning capabilities for ASQP tasks improve significantly with an increased sample size. However, this raises the question of whether this improvement is always directly proportional to the number

of samples. To explore this issue, we further increased the sample size, as shown in Figure 4. We found that the T5-based MvP model's performance steadily improved with more samples, indicating that the generative-based models rely on sufficient high-quality labeled data. Surprisingly, for ICL-based methods, performance tends to decline after reaching a certain sample size threshold. Our analysis suggests two main reasons for this decline. First, a large number of examples provides excessive prior semantic information, causing LLMs to become confused and lose focus on core aspects. Second, lower-ranked samples are poorer in quality and contain more redundancy. Notably, compared to previous ICL methods, the RCR framework mitigates this performance degradation, indicating that RepDR retrieves higher-quality demonstrations and introduces fewer errors.

4.4 Error analysis and Case Study

In order to comprehensively analyze the reasoning errors of our proposed method, we conducted error analysis and case studies. We randomly selected 100 prediction results from each dataset in the 5-shot scenario using GPT-4. The incorrectly predicted quadruples were categorized by error type, as shown in Figure 5. We found that errors were primarily concentrated on the predictions of aspect

and opinion terms in both datasets.

The main reason for this phenomenon is that aspect and opinion terms often appear as text spans rather than individual words. LLMs struggle to match these text spans accurately, as illustrated by Example 1. Another significant cause of errors is the presence of multiple quadruples in the text, which confuses the LLMs. This typically occurs in the first subtask of the ChaPT framework, making it difficult to match each aspect-opinion pair precisely. Example 2 shows an incorrect aspect-opinion pair (attitude, snotty) being generated. Furthermore, errors from the previous subtask can propagate and interfere with predicting aspect categories, leading to cumulative errors. In summary, accurately matching text spans and handling sentences with multiple quadruples are challenging issues that LLMs must address in ASQP problems.

5 Related work

5.1 Aspect Sentiment Quad Prediction

Aspect Sentiment Quad Prediction (ASQP) is a crucial sentiment analysis task that has attracted increasing attention. (Zhang et al., 2022a, 2023b; Zhong et al., 2023). Initially, ASQP was mainly handled using pipeline approaches that combined multiple baseline models (Cai et al., 2021). Further studies have shown that generative models achieve promising results (Zhang et al., 2021b; Bao et al., 2022; Peper and Wang, 2022). For example, Zhang et al. (2021a) introduced a paraphrase model, transforming a quadruple prediction task into a text generation task. Mao et al. (2022) constructed a search tree for the optimal generation path. Bao et al. (2022) developed an opinion tree to jointly detect all sentiment elements. Additionally, many efforts have focused on enhancing generative models through data augmentation. Hu et al. (2022a) first considered selecting the appropriate quadruple generation order as a data augmentation method. Gou et al. (2023) proposed an MVP framework to increase output views. Wang et al. (2023a) suggested generating new data containing quadruples through generation models. However, models trained on specific domain datasets often perform poorly when transferred to other domains.

5.2 In-Context Learning

In-context learning (ICL) refers to the ability of large language models (LLMs) to handle complex tasks with only a few annotated examples with-

out additional training or gradient updates (Zhao et al., 2023). Research on ICL focuses on two main areas. On the one hand, it involves investigating prompting frameworks (Long, 2023; Paranjape et al., 2023; Diao et al., 2023; Li et al., 2024). For example, Wei et al. (2022); Wang et al. (2022b) proposed the Chain of Thought (CoT) to enhance reasoning capabilities. Yao et al. (2024) further refined CoT into the Tree of Thoughts (ToT), maintaining the intermediate thoughts in a search tree and evaluating these thoughts. On the other hand, considerable work studies focus on providing better demonstrations (Li et al., 2022; Min et al., 2022; Li et al., 2023; Wang et al., 2023b). Liu et al. (2022) found that samples closely related to the target data in the embedding space perform better. Building on this idea, Wang et al. (2022a) proposed enhancing inputs by retrieving similar examples. Rubin et al. (2022) introduced a demonstration retriever. Moreover, examples representing diversity can also improve ICL performance (Qin et al., 2023; Xu et al., 2024).

Owing to developments in ICL, some studies have addressed sentiment analysis tasks in zero-shot or few-shot scenarios using ICL, achieving effective results (Wang et al., 2023c). For instance, Zhong et al. (2023) observed that the zero-shot performance of LLMs is comparable to fine-tuned BERT. Sun et al. (2023) proposed a multi-LLM negotiation framework for sentiment analysis. Fei et al. (2023) introduced a THOR framework, significantly enhancing implicit sentiment analysis performance. In light of this, we explore the potential of LLMs for the ASQP problem. To our knowledge, this work is the first to discuss the application of LLMs to ASQP task systematically.

6 Conclusion

In this work, we propose a new RCR framework to solve the ASQP task in low-resource scenarios. To reduce complexity, the chain prompting module (ChaPT) is designed to decompose the ASQP task into three subtasks and enable LLMs to conduct step-by-step reasoning. Furthermore, a representative demonstration retriever (RepDR) is developed to provide ChaPT with demonstrations that balance diversity and similarity, maximizing the reasoning ability of LLMs at each step. Detailed experiments demonstrate the effectiveness of our proposed RCR framework in both zero-shot and few-shot scenarios, enabling GPT-4 to achieve state-of-the-art performance on the ASQP task.

544 Limitations

545 Despite our proposed method achieves state-of-the-
546 art performance in ASQP tasks under low-resource
547 scenarios, our work still has limitations. Firstly,
548 we observe that the performance of RCR improves
549 with the increasing intelligence of the integrated
550 LLM models. Therefore, it is necessary to explore
551 the effects of integrating LLMs of different scales
552 with RCR. Secondly, our proposed ChaPT frame-
553 work requires manually designed prompts, leading
554 to instability in LLM reasoning results as the qual-
555 ity of the prompts varies. Exploring better auto-
556 matic prompt generation strategies could address
557 this issue. Finally, the experiments only validate
558 the improvements of RCR in the ASQP task. Intu-
559 itively, the RCR framework can be easily extended
560 to aspect-based Sentiment analysis subtasks similar
561 to ASQP, such as Aspect Sentiment Triplet Extrac-
562 tion (ASTE), Aspect-Category-Sentiment Detec-
563 tion (ACSD), and Aspect Category Opinion Senti-
564 ment (ACOS).

565 Ethical Statement

566 All our experiments are based on publicly avail-
567 able datasets and code repositories. We maintain
568 impartiality and honesty in our analysis of the
569 experimental results, and our research and work
570 do not harm any individuals or groups. We will
571 open-source our code for further discussion and
572 exploration. Regarding broader impacts, this work
573 may promote further research using large language
574 models(LLMs) for sentiment analysis tasks in low-
575 resource scenarios, contributing to lightweight and
576 automated opinion mining and sentiment analysis
577 in the real world. Additionally, we recognize the ro-
578 bust capabilities and potential risks of LLMs. Thus,
579 we strictly adhere to ethical standards throughout
580 our research to ensure that our work is not misused
581 or causes any negative impact.

582 References

583 Josh Achiam, Steven Adler, Sandhini Agarwal, Lama
584 Ahmad, Ilge Akkaya, Florencia Leoni Aleman,
585 Diogo Almeida, Janko Altenschmidt, Sam Altman,
586 Shyamal Anadkat, et al. 2023. Gpt-4 technical report.
587 *arXiv preprint arXiv:2303.08774*.

588 David Arthur, Sergei Vassilvitskii, et al. 2007. k-
589 means++: The advantages of careful seeding. In
590 *Soda*, volume 7, pages 1027–1035.

591 Xiaoyi Bao, Zhongqing Wang, Xiaotong Jiang, Rong
592 Xiao, and Shoushan Li. 2022. Aspect-based senti-

ment analysis with opinion tree generation. In *IJCAI*,
volume 2022, pages 4044–4050.

Hongjie Cai, Rui Xia, and Jianfei Yu. 2021. Aspect-
category-opinion-sentiment quadruple extraction
with implicit aspects and opinions. 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)*, pages 340–350.

Shaowei Chen, Jie Liu, Yu Wang, Wenzheng Zhang,
and Ziming Chi. 2020. Synchronous double-channel
recurrent network for aspect-opinion pair extraction.
In *Proceedings of the 58th annual meeting of the as-
sociation for computational linguistics*, pages 6515–
6524.

Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong
Zhang. 2023. Active prompting with chain-of-
thought for large language models. *arXiv preprint
arXiv:2302.12246*.

Hao Fei, Bobo Li, Qian Liu, Lidong Bing, Fei Li, and
Tat-Seng Chua. 2023. Reasoning implicit sentiment
with chain-of-thought prompting. In *Proceedings
of the 61st Annual Meeting of the Association for
Computational Linguistics (Volume 2: Short Papers)*,
pages 1171–1182.

Zhibin Gou, Yujiu Yang, et al. 2023. Mvp: Multi-view
prompting improves aspect sentiment tuple predic-
tion. In *The 61st Annual Meeting Of The Association
For Computational Linguistics*.

Mengting Hu, Yike Wu, Hang Gao, Yinhao Bai, and
Shiwan Zhao. 2022a. Improving aspect sentiment
quad prediction via template-order data augmenta-
tion. In *Proceedings of the 2022 Conference on Em-
pirical Methods in Natural Language Processing*,
pages 7889–7900.

Mengting Hu, Yike Wu, Hang Gao, Yinhao Bai, and
Shiwan Zhao. 2022b. Improving aspect sentiment
quad prediction via template-order data augmenta-
tion. In *Proceedings of the 2022 Conference on Em-
pirical Methods in Natural Language Processing*,
pages 7889–7900.

Dong-Ho Lee, Akshen Kadakia, Kangmin Tan, Mahak
Agarwal, Xinyu Feng, Takashi Shibuya, Ryosuke
Mitani, Toshiyuki Sekiya, Jay Pujara, and Xiang Ren.
2022. Good examples make a faster learner: Simple
demonstration-based learning for low-resource ner.
In *Proceedings of the 60th Annual Meeting of the
Association for Computational Linguistics (Volume
1: Long Papers)*, pages 2687–2700.

Kaijian Li, Shansan Gong, and Kenny Q Zhu. 2022.
Few-shot natural language inference generation with
pdd: Prompt and dynamic demonstration. *arXiv
preprint arXiv:2205.10593*.

Mukai Li, Shansan Gong, Jiangtao Feng, Yiheng Xu,
Jun Zhang, Zhiyong Wu, and Lingpeng Kong. 2023.
In-context learning with many demonstration exam-
ples. *arXiv preprint arXiv:2302.04931*.

650	Zekun Li, Baolin Peng, Pengcheng He, Michel Galley,	Maria Pontiki, Dimitris Galanis, Haris Papageorgiou,	705
651	Jianfeng Gao, and Xifeng Yan. 2024. Guiding large	Suresh Manandhar, and Ion Androutsopoulos. 2015.	706
652	language models via directional stimulus prompting.	SemEval-2015 task 12: Aspect based sentiment anal-	707
653	<i>Advances in Neural Information Processing Systems</i> ,	ysis . In <i>Proceedings of the 9th International Work-</i>	708
654	36.	<i>shop on Semantic Evaluation (SemEval 2015)</i> , pages	709
		486–495.	710
655	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan,	Chengwei Qin, Aston Zhang, Anirudh Dagar, and	711
656	Lawrence Carin, and Weizhu Chen. 2022. What	Wenming Ye. 2023. In-context learning with it-	712
657	makes good in-context examples for gpt-3? <i>DeeLIO</i>	erative demonstration selection. <i>arXiv preprint</i>	713
658	2022, page 100.	<i>arXiv:2310.09881</i> .	714
659	Jian Liu, Zhiyang Teng, Leyang Cui, Hanmeng Liu, and	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine	715
660	Yue Zhang. 2021. Solving aspect category sentiment	Lee, Sharan Narang, Michael Matena, Yanqi Zhou,	716
661	analysis as a text generation task. In <i>Proceedings</i>	Wei Li, and Peter J Liu. 2020. Exploring the lim-	717
662	<i>of the 2021 Conference on Empirical Methods in</i>	its of transfer learning with a unified text-to-text	718
663	<i>Natural Language Processing</i> , pages 4406–4416.	transformer. <i>Journal of machine learning research</i> ,	719
		21(140):1–67.	720
664	Jieyi Long. 2023. Large language model guided tree-of-	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	721
665	thought. <i>arXiv preprint arXiv:2305.08291</i> .	Sentence embeddings using siamese bert-networks.	722
666	Yue Mao, Yi Shen, Jingchao Yang, Xiaoying Zhu, and	In <i>Proceedings of the 2019 Conference on Empirical</i>	723
667	Longjun Cai. 2022. Seq2path: Generating sentiment	<i>Methods in Natural Language Processing and the 9th</i>	724
668	tuples as paths of a tree. In <i>Findings of the Associa-</i>	<i>International Joint Conference on Natural Language</i>	725
669	<i>tion for Computational Linguistics: ACL 2022</i> , pages	<i>Processing (EMNLP-IJCNLP)</i> , pages 3982–3992.	726
670	2215–2225.		
671	Sewon Min, Xinxu Lyu, Ari Holtzman, Mikel Artetxe,	Ohad Rubin, Jonathan Herzig, and Jonathan Berant.	727
672	Mike Lewis, Hannaneh Hajishirzi, editor = "Gold-	2022. Learning to retrieve prompts for in-context	728
673	berg Yoav Zettlemoyer, Luke", Zornitsa Kozareva,	learning. In <i>Proceedings of the 2022 Conference</i>	729
674	and Yue Zhang. 2022. Rethinking the role of demon-	<i>of the North American Chapter of the Association</i>	730
675	strations: What makes in-context learning work? In	<i>for Computational Linguistics: Human Language</i>	731
676	<i>Proceedings of the 2022 Conference on Empirical</i>	<i>Technologies</i> , pages 2655–2671.	732
677	<i>Methods in Natural Language Processing</i> .		
678	AI Open. 2022. Introducing chatgpt. open ai.	Zhenmei Shi, Junyi Wei, Zhuoyan Xu, and Yingyu	733
679	Bhargavi Paranjape, Scott Lundberg, Sameer Singh,	Liang. 2023. Why larger language models do in-	734
680	Hannaneh Hajishirzi, Luke Zettlemoyer, and	context learning differently? In <i>RO-FoMo: Robust-</i>	735
681	Marco Tulio Ribeiro. 2023. Art: Automatic multi-	<i>ness of Few-shot and Zero-shot Learning in Large</i>	736
682	step reasoning and tool-use for large language mod-	<i>Foundation Models</i> .	737
683	els. <i>arXiv preprint arXiv:2303.09014</i> .		
684	Haiyun Peng, Lu Xu, Lidong Bing, Fei Huang, Wei Lu,	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-	738
685	and Luo Si. 2020. Knowing what, how and why: A	Yan Liu. 2020. MpNet: Masked and permuted pre-	739
686	near complete solution for aspect-based sentiment	training for language understanding. <i>Advances in</i>	740
687	analysis . <i>Proceedings of the AAAI Conference on</i>	<i>neural information processing systems</i> , 33:16857–	741
688	<i>Artificial Intelligence</i> , page 8600–8607.	16867.	742
689	Joseph Peper and Lu Wang. 2022. Generative aspect-	Xiaofei Sun, Xiaoya Li, Shengyu Zhang, Shuhe Wang,	743
690	based sentiment analysis with contrastive learning	Fei Wu, Jiwei Li, Tianwei Zhang, and Guoyin Wang.	744
691	and expressive structure. In <i>Findings of the Associa-</i>	2023. Sentiment analysis through llm negotiations.	745
692	<i>tion for Computational Linguistics: EMNLP 2022</i> ,	<i>arXiv preprint arXiv:2311.01876</i> .	746
693	pages 6089–6095.		
694	Maria Pontiki, Dimitris Galanis, Haris Papageorgiou,	Hai Wan, Yufei Yang, Jianfeng Du, Yanan Liu, Kunxun	747
695	Ion Androutsopoulos, Suresh Manandhar, Moham-	Qi, and Jeff Z Pan. 2020. Target-aspect-sentiment	748
696	mad AL-Smadi, Mahmoud Al-Ayyoub, Yanyan	joint detection for aspect-based sentiment analysis.	749
697	Zhao, Bing Qin, Orphée De Clercq, Véronique	In <i>Proceedings of the AAAI conference on artificial</i>	750
698	Hoste, Marianna Apidianaki, Xavier Tannier, Na-	<i>intelligence</i> , volume 34, pages 9122–9129.	751
699	talia Loukachevitch, Evgeniy Kotelnikov, Nuria Bel,	An Wang, Junfeng Jiang, Youmi Ma, Ao Liu, and	752
700	Salud María Jiménez-Zafra, and Gülşen Eryiğit.	Naoaki Okazaki. 2023a. Generative data augmen-	753
701	2016. SemEval-2016 task 5: Aspect based sentiment	tation for aspect sentiment quad prediction. In <i>Pro-</i>	754
702	analysis . In <i>Proceedings of the 10th International</i>	<i>ceedings of the 12th Joint Conference on Lexical</i>	755
703	<i>Workshop on Semantic Evaluation (SemEval-2016)</i> ,	<i>and Computational Semantics (* SEM 2023)</i> , pages	756
704	pages 19–30.	128–140.	757
		Shuohang Wang, Yichong Xu, Yuwei Fang, Yang Liu,	758
		Siqi Sun, Ruochen Xu, Chenguang Zhu, and Michael	759
		Zeng. 2022a. Training data is more valuable than you	760

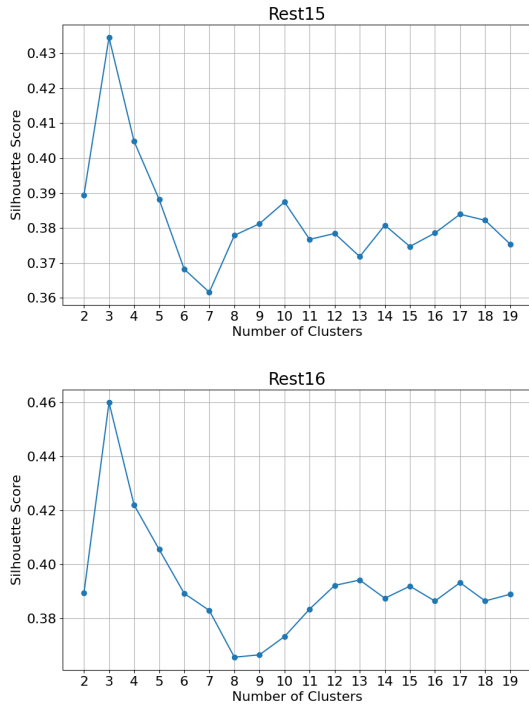


Figure 6: Silhouette Scores for different number of clusters.

867 Rest16. A larger silhouette score indicates better
 868 clustering quality. Therefore, from Figure 6, we
 869 can determine that the optimal number of clusters
 870 for both datasets is 3. To maintain clustering sta-
 871 bility, we first standardize, normalize, and reduce
 872 dimensionality of the sentence embedding.

873 B Implementation Details of Generative 874 Models

875 The few-shot training of generative models fol-
 876 lows the settings proposed by Zhang et al. (2023c),
 877 where k-shot represents sampling k examples for
 878 each aspect category. We set the batch size of all
 879 models to 8, the learning rate to 1e-4, and the train-
 880 ing epochs to 100. All experiments were conducted
 881 using an Nvidia RTX 3090 GPU.