# RVISA: Reasoning and Verification for Implicit Sentiment Analysis

Wenna Lai, Haoran Xie, Guandong Xu, Qing Li

*Abstract*—With an increasing social demand for fine-grained sentiment analysis (SA), implicit sentiment analysis (ISA) poses a significant challenge with the absence of salient cue words in expressions. It necessitates reliable reasoning to understand how the sentiment is aroused and thus determine implicit sentiments. In the era of Large Language Models (LLMs), Encoder-Decoder (ED) LLMs have gained popularity to serve as backbone models for SA applications, considering impressive text comprehension and reasoning ability among diverse tasks. On the other hand, Decoder-only (DO) LLMs exhibit superior natural language generation and in-context learning capabilities. However, their responses may contain misleading or inaccurate information. To identify implicit sentiment with reliable reasoning, this study proposes RVISA, a two-stage reasoning framework that harnesses the generation ability of DO LLMs and the reasoning ability of ED LLMs to train an enhanced reasoner. Specifically, we adopt three-hop reasoning prompting to explicitly furnish sentiment elements as cues. The generated rationales are utilized to fine-tune an ED LLM into a skilled reasoner. Additionally, we develop a straightforward yet effective verification mechanism to ensure the reliability of the reasoning learning. We evaluated the proposed method on two benchmark datasets and achieved state-of-the-art results in ISA performance.

*Index Terms*—Implicit sentiment analysis, Large language models, Multi-task learning, Chain-of-Thought.

## I. INTRODUCTION

Sentiment analysis (SA) aims to evoke opinions, sentiments, and emotions through different computational methods [1]. Nowadays, people have demonstrated a stronger willingness to express and share their ideas online about day-to-day activities and global issues. With the increasing demand on social media, SA has gained significant interest considering great commercial value in exploring customer opinions or sentiments from user reviews or other sources of information. Meanwhile, sentiments can assist learning, communication, decision-making, and situation awareness in human-centric environments [2]. Traditionally, SA is classified into three levels, which are document-level, sentence-level, and aspect-level [3]. While the document and sentence level analyze the sentiment towards the overview of a document or sentence,

Wenna Lai is with the Department of Computing, Hong Kong Polytechnic University, Hong Kong, also the School of Computer Science and the Data Science Institute, University of Technology Sydney, Australia.

Haoran Xie is with the School of Data Science, Lingnan University, Hong Kong (email: hrxie@ieee.org).

Guandong Xu is with the School of Computer Science and the Data Science Institute, University of Technology Sydney, Sydney, NSW 2007, and also the Education University of Hong Kong, Hong Kong SAR.

Qing Li is with the Department of Computing, Hong Kong Polytechnic University, Hong Kong SAR.
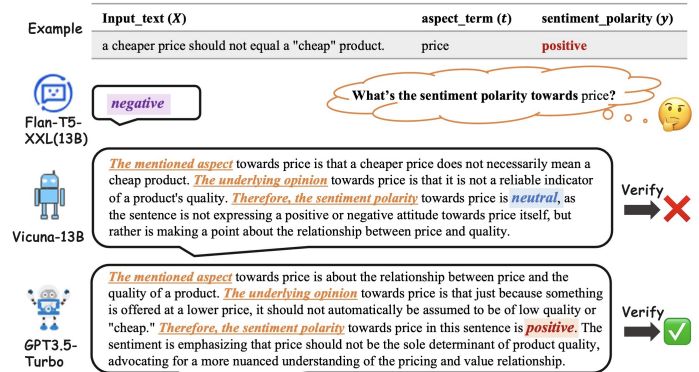
Fig. 1. Different LLMs demonstrate **diverse reasoning abilities** for implicit sentiment analysis. It is conducive to detecting implicit sentiment polarities by explicitly inferring sentiment elements as rationale but **verification is required to ensure reliability**.

aspect-based sentiment analysis (ABSA) is more fine-grained to extract the opinion towards a given aspect or entity. In many cases, there may be multiple aspects in one sentence, making it challenging to pinpoint a specific target and identify the corresponding sentiment.

Considering context information, sentiment analysis can be further classified into implicit sentiment analysis (ISA) and explicit sentiment analysis (ESA), where expressions in ISA contain no explicit polarity markers but still deliver human-aware sentiment polarity [4]. In 2021, [5] split the SemEval-2014 Restaurant and Laptop benchmarks into Explicit Sentiment Expression slice and Implicit Sentiment Expression slice based on the presence of opinion words, drawing attention to ISA in ABSA tasks. [6] conducted pre-experiments on 20 existing sentiment classifiers and investigated that traditional methods performed ineffectively towards the same implicit case. They suggested that majority of traditional classifiers tend to overlook ISA problem and address ISA superficially. Although humans can easily grasp real intent and perceive changes in mood with common sense and reasoning ability, it is more difficult for models to tackle ISA than ESA, due to limited context information and insufficient reasoning skills.

As recent great triumph of large language models (LLMs) has demonstrated impressive complex reasoning with chain-of-thought (CoT) prompting [7], [8] and in-context learning ability [9], more scholars tend to embrace LLMs for downstream applications [10]–[12]. [13] investigated the performance of LLMs in prompt-based inference and observed that for tasks requiring structured sentiment output, like ABSA tasks, both DO LLMs (e.g., GPT-3.5-turbo [14]) and ED LLMs (e.g.,

Flan-T5-XXL [15]) tend to lag behind ED backbone models (e.g., T5-Large [16]) trained with domain-specific data in automatic and human evaluations. The performance can vary significantly with different prompt designs. These indicate that deploying LLMs for ISA directly without training may not fully unleash their reasoning capacity for achieving satisfactory results. [6] first employed CoT fine-tuning on Flan-T5 for ISA and gained improved performance. However, intermediate steps generated by Flan-T5 were most likely to be untrustworthy, with insufficient or duplicate content constrained by weak generation capacity. As illustrated in Figure 1, different LLMs performed diversely in analyzing implicit sentiment towards the aspect term '*price*', given the text '*a cheaper price should not equal a "cheap" product*'. Inferior models, like Flan-T5 in the group of Encoder-Decoder (ED) LLMs, displayed excellent comprehension and reasoning in solving tasks with diverse input information, but limited generation and prompt-based inference capabilities on open-text [17]. They were predisposed to inaccurately predict implicit sentiment in the absence of explicit cues. Conversely, Decoder-only (DO) LLMs with more advanced generation ability, such as Vicuna-13B [18] and GPT-3.5-turbo, demonstrated enhanced proficiency in explicitly deducing sentiment elements pertinent to the context under reasoning prompts, while reliability in achieving accurate or correct responses was not guaranteed. Moreover, LLMs often showcase superior performance with emergent abilities when scaling up at a certain level [19], the direct deployment or fine-tuning of large-scale models (e.g., GPT-3.5-turbo) might be hindered by considerable computational costs. To effectively discern implicit sentiment polarities towards a specified aspect, it is essential to exploit reliable reasoning methods for applicable backbone models.

With this motivation, we attempt to equip ED backbone models with enhanced reasoning ability by explicitly learning from convincing rationales provided by DO LLMs through synchronous verification. Specifically, we follow the sentiment element construction and design corresponding three-hop reasoning prompting to guide DO LLMs in explicitly inferring sentiment elements before determining the final sentiment. Then an ED model is served as the backbone model and fine-tuned based on the generated rationales and golden labels in datasets. To ensure the quality of reasoning learning, we further introduce an answer-based verification mechanism as an additional signal to assess the reliability of the rationale, which promotes dialectical learning to identify and rectify potential inaccuracies.

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

- We propose a novel two-stage learning framework, Reasoning and Verification for Implicit Sentiment Analysis (RVISA), marking the endeavor to improve the proficiency of ED backbone models as adept reasoners in ISA, complemented by the generative strengths of DO LLMs.
- We introduce a straightforward yet efficacious verification mechanism to provide reliable supervision for reasoning learning and improve overall performance.
- The evaluation outcomes on two benchmark datasets underscore the efficacy of our method in achieving state-of-the-art results in ISA performance.

## II. RELATED WORK

In this work, we train a skilled reasoner with the cooperation of LLMs to conduct implicit sentiment analysis, learning fruitful information from rationales generated by reasoning prompting. We draw attention to the existing research on implicit sentiment analysis and methods that learn from reasoning prompting making use of emergent abilities showcased in LLMs.

### A. Implicit Sentiment Analysis

Implicit sentiment analysis has gained considerable attention in the field of sentiment analysis [4], [20]. In the beginning, great efforts have been taken into solving the implicit sentiment detected in sentence level [21], [22]. With the increasing social demand, recent scholars attempted to develop effective paradigms tackling the unique characteristics of implicit sentiment analysis at a more fine-grained level towards the aspect target [5], [6], [23]. To capture the implicit sentiment expression, some research exploited extra knowledge to further improve the learning performance. [5] pre-trained on large-scale sentiment annotated corpora with supervised contrastive learning objectives to align the representation of explicit and implicit sentiment expressions. Instead of making use of external knowledge, [24] generated explicit sentiment augmentation based on the language model itself to enhance implicit classification tendencies. Considering the difficulties of obtaining the full knowledge through additional means, [23] proposed reasoning learning under causal intervention to capture the correlation within the expressions. The relationship within fine-grained sentiment analysis can be summarized into four key sentiment elements involving *target, aspect, opinion, and sentiment polarity*, which are highly close to each other in understanding the underlying sentiment [25]. With the impressive performance of chain-of-thought (CoT) and in-context learning abilities showcased in LLMs, [6] introduced CoT fine-tuning to guide the ED backbone model inferring sentiment elements including implicit sentiment polarities step-by-step in an easy-to-hard manner. Similar to that, our approach makes use of fine-grained sentiment elements as cues for chain-of-thought prompting. But considering the limited generation capabilities of ED LLMs (e.g., Flan-T5 [15]), rather than inferring the sentiment elements from backbone models themselves, we train ED backbone models to become proficient reasoners by leveraging the informative rationale generated from DO LLMs (e.g., GPT-3.5-turbo [14]).

### B. Reasoning Prompting

LLMs have demonstrated impressive complex reasoning abilities with Chain-of-Thought (CoT) prompting [8], [26]. The use of reasoning prompting aims to guide the model in thinking step-by-step and leveraging most of the inference power for task solving. It is discovered effective in boosting the zero-shot or few-shot performance of LLMs [27]–[30]. Figure 2 illustrates various reasoning prompting applying to sentiment analysis. On the left-hand side are commonly used prompting modes including Reasoning and Rationalization:

**Reasoning (RE):**

> Given the sentence [input_text], **what is the sentiment polarity towards [aspect_term], why**?

**Rationalization (RA):**

> Given the sentence [input_text], **the sentiment polarity towards [aspect_term] is [label_polarity], why**?

**Zero-CoT:**

> Given the sentence [input_text], **what is the sentiment polarity towards [aspect_term], why**? Let's think step by step.

**Three-hop Reasoning (TH-RE):**

> Given the sentence [input_text], **what is the sentiment polarity towards [aspect_term], why**?
> Let's think step by step. The mentioned aspect towards [aspect_term] is about ... The underlying opinion towards [aspect_term] is about ... Therefore, the sentiment polarity towards [aspect_term] is ...

**Three-hop Rationalization (TH-RA):**

> Given the sentence [input_text], **the sentiment polarity towards [aspect_term] is [label_polarity], why**?
> Let's think step by step. The mentioned aspect towards [aspect_term] is about ... The underlying opinion towards [aspect_term] is about ... Therefore, the sentiment polarity towards [aspect_term] is ...

Fig. 2. Reasoning promptings applying to sentiment analysis. Left: commonly used prompting modes. Right: three-hop prompting for ISA.

*a) Reasoning (RE):* [31] introduced multi-task learning with reasoning prompting by simultaneously learning the question-answer pairs and question-explanation pairs. The generated rationales for the question will not have the ground truth answer for reference, which prompts the language model to infer the answer according to its step-by-step inference and own judgment. Therefore, the answer showcased in the explanation can be different from the gold label.

*b) Rationalization (RA):* [32] first proposed the idea of rationalization, which attempts to retrieve the explanation for the question by explicitly giving the correct answer. The intuition towards it is to rationalize the question with the golden label and provide the possible reasons behind the question-answer connection.

Besides them, [8] revealed that LLMs are capable of incremental reasoning without exemplars. Simply by incorporating a prompt *"let's think step by step"* (i.e., Zero-CoT in Figure 2), it is universally applicable across tasks. However, the granularity of the reasoning steps generated by Zero-CoT remains unpredictable, hinging on the LLM's inherent knowledge and varying across models. Furthering this exploration, [33] examined the influence of reasoning step length within prompts and suggested that maintaining a certain step size according to the complexity of the task has a critical role in forming the final answer.

Based on these insights, our method ingeniously builds on previous prompting methods by adding a three-hop strategy that uses the construction of sentiment elements to keep critical reasoning steps going, as shown on the right in Figure 2. Considering the heuristic about answer inference, the explanation given by reasoning inference that leads to the correct answer should be more trustworthy for answer prediction. LLMs may falter in complex scenarios where reasoning prompts alone are insufficient, potentially yielding explanations riddled with inaccuracies. To mitigate this issue, [34] employed answer-based filtering to improve rationale quality, with Rationalization serving as the backup option for erroneous explanations. Rather than discarding inaccurate rationales outright or attempting to compensate for them with additional information sources, our approach retains these informative rationales throughout the learning process.

## C. Learning from Rationale

Learning from explanations and empowering the training model with reasoning abilities have been explored in various fields [10]–[12], [34]. LLMs are capable of validating their responses with reasonable intermediate steps [7], [8], rationales can be used as demonstrations [27] or extra fine-tuning data [10], [35], [36] to improve the learning performance. Considering the training cost for LLMs, rationales can also serve as valuable supervised signals for smaller task-specific models, which can be more easily deployed [11], [12], [34], [37]. However, [12] directly kept the answer generated by LLMs as supervision signals, which neglected the potential erroneous occurrence. [34] reorganized the rationale set based on answer-based filtering to mitigate the possibility of error learning, but potential misleading information may still exist even with the guidance of a correct answer. In contrast, our approach augments the overall performance by incorporating an answer-based verification mechanism as an additional layer of supervision for multi-task learning. This innovative strategy not only preserves valuable insights contained within the rationales but also leverages them to refine the learning process with both positive and negative signals.

We compare various prompting methodologies and substantiate the superior performance of our three-hop prompting in the nuanced domain of implicit sentiment analysis through comprehensive experiments. Beyond that, the introduction of the verification mechanism further improves the performance. The integration of three-hop prompting with the verification mechanism effectively navigates complexities inherent in LLM reasoning process.

## III. TWO-STAGE REASONING FRAMEWORK

We propose a novel two-stage framework, RVISA (as shown in Figure 3), aiming to empower ED models with enhanced reasoning ability and incorporate the answer-based verification mechanism for reasoning refinement during model learning. In the initial stage, we leverage DO LLMs to generate insightful rationales and predict labels through our three-hop reasoning prompting approach. Then the verification signals are curated according to the correctness of LLM prediction labels. In the
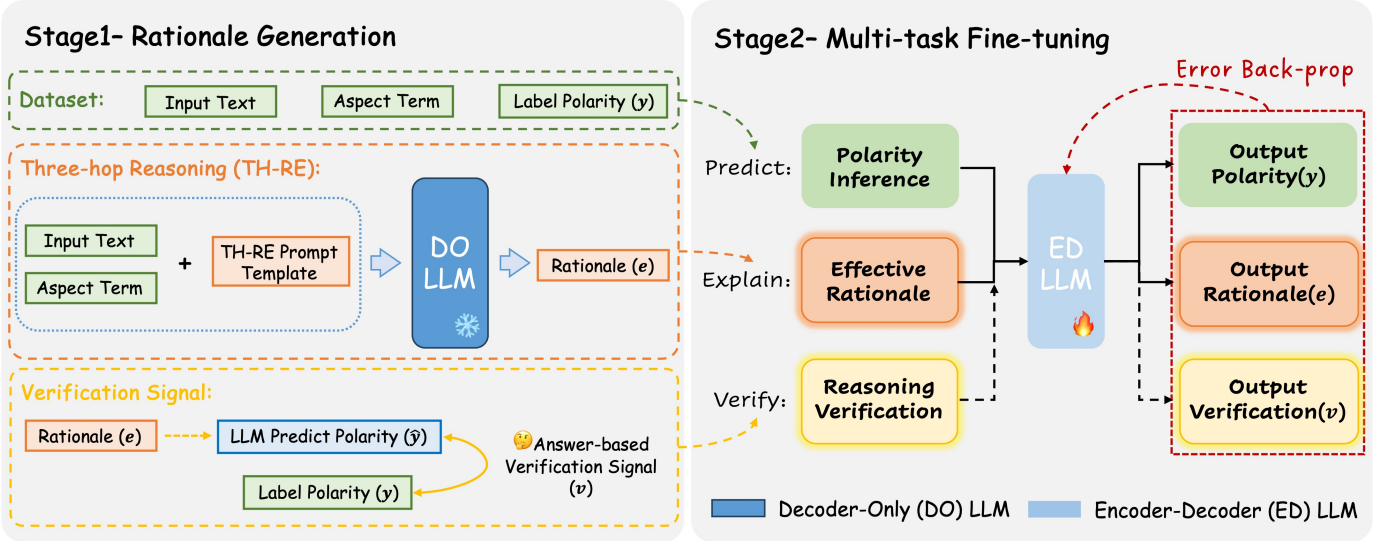
Fig. 3. The overview of proposed two-stage reasoning framework RVISA. Left: rationale generation stage leveraging DO LLM to prepare effective rationales and corresponding answer-based verification signals. Right: multi-task fine-tuning stage to train an ED backbone model as an enhanced reasoner with additional explanation tasks along with verification supervision.

second stage, the generated rationales are employed for multi-task fine-tuning on an ED backbone model. To further ensure the reliability of the generated rationales, we implement a straightforward yet effective verification mechanism with an additional task supervised by the verification signals to guide self-revision during the reasoning learning process. Different tasks are distinguished by task-specific prefixes. It is thought that the model can learn to understand the underlying logic and relationships among sentiment elements that govern implicit sentiment prediction, by training on reasoning rationales and self-verification signals at the same time under the supervision of gold labels that have been annotated.

### A. Problem Definition

In sentiment analysis tasks including explicit sentiment analysis and implicit sentiment analysis, given the dataset $D = (x_i, y_i)^N$, where $1 \leq i \leq N$, $x_i$ represents an input sentence serving as a data example. Within each sentence $x_i$, an aspect term $t_i$ is identified, denoted as $t_i \subset x_i$. The relevant sentiment elements consist of aspect $a_i$, opinion $o_i$, and sentiment polarity $y_i$. The objective of the task is to infer the sentiment polarity $y_i$ towards the aspect term $t_i$, given the input sentence $x_i$ and the specified aspect term $t_i$. In the standard prompting approach for direct fine-tuning, the LLM predicts the sentiment polarity $\hat{y}_i$ solely via $\hat{y}_i = argmax p(y_i|x_i, t_i)$. Without considering the intermediate sentiment elements, this approach potentially limits the ability of models to capture the sentiment nuances present in the text.

### B. Three-hop Reasoning Generation

To improve the generation of informative rationales, we prompt DO LLMs to generate intermediate steps during the inference of implicit sentiments. To understand how the sentiment is aroused, sentiment elements are essential in directing inference process since they contribute to constructing the complete picture of sentiment analysis. Therefore, we further design the three-hop prompting as illustrated in Figure 2, deviating from conventional prompting modes. The objective of this design is to dominate reasoning process by extracting closely associated sentiment elements. Simultaneously, this design is conducive to standardizing the generative structure, facilitating improved learning of patterns and interconnections among rationales. The details of the three-hop reasoning prompting are explained as follows.

*a) Three-hop Reasoning (TH-RE):* Fine-grained sentiment analysis involves dissecting key sentiment elements involving the target, aspect, opinion, and sentiment polarity [25]. Various approaches exist for solving these individual subtasks or their combinations, collectively contributing to a comprehensive sentiment analysis picture. To address the complexity of the task holistically, it is essential to consider the components systematically and tackle them incrementally. [6] first design the prompting based on the CoT strategy by explicitly inferring the sentiment elements and then employ the prompting for three-step generation during model fine-tuning. However, the prompting for each step is inferred separately for a single sentiment element at a time, with the results concatenated as context information for the subsequent step.

In our design, we adopt a structured approach by explicitly presenting sentiment elements in a natural language sequence to construct a three-hop reasoning prompt. This approach underscores the causal relationships among sentiment elements and the final sentiment polarity prediction in a single iteration.

As shown in the template below, we incorporate sentiment elements as cues at the end of *"let's think step-by-step"*, guiding the language model to generate reasoning steps in alignment with the sentiment elements' understanding and finally infer the sentiment polarity. We expect DO LLMs to predict the explanation via $\hat{e}_i = argmax p(e_i|x_i, t_i)$, where $\hat{a}_i, \hat{o}_i, \hat{y}_i \subset \hat{e}_i$.

> Given the sentence $x_i$, what is the sentiment polarity towards $t_i$, why? Let's think step by step.
> The **mentioned aspect** towards $t_i$ is about ... The **underlying opinion** towards $t_i$ is about ... Therefore, the **sentiment polarity** towards $t_i$ is ...

*b) Three-hop Rationalization (TH-RE):* Diverse from TH-RE, we integrate three-hop reasoning with rationalization to establish the Three-hop rationalization (TH-RA) prompting. Specifically, the gold label will be given as the reference, which prompts the LLMs to elucidate the annotated sentiment label through a systematic and step-by-step inference process guided by sentiment elements. We expect DO LLMs to predict the explanation via $\hat{e}_i = argmaxp(e_i|y_i, x_i, t_i)$, where $\hat{a}_i, \hat{o}_i, \hat{y}_i \subset \hat{e}_i$.

> Given the sentence $x_i$, the sentiment polarity towards $t_i$ is $y_i$, why? Let's think step by step.
> The **mentioned aspect** towards $t_i$ is about ... The **underlying opinion** towards $t_i$ is about ... Therefore, the **sentiment polarity** towards $t_i$ is ...

## C. Multi-task Fine-tuning

We employ a multi-task fine-tuning approach to simultaneously learn the rationales generated by the LLM and the annotated labels. Given the dataset $D = \{(x_i, y_i)\}^N$, where $1 \leq i \leq N$, we generate an explanation $e_i$ to serve as a rationale for each input $x_i$ as detailed in Section III-B. Each explanation $e_i$ encompasses a generated label $\hat{y}_i$ from the LLM, denoted as $\hat{y}_i \subset e_i$. Subsequently, we construct a new dataset $D_{exp} = \{(x_i, e_i)\}^N$, where $1 \leq i \leq N$. The objective during the training phase is to effectively utilize the generated content and learn from two distinct tasks: the explanation task utilizing data from $D_{exp}$ and the prediction task utilizing the data from the original dataset $D_{pre} = D = \{(x_i, y_i)\}^N$, where $1 \leq i \leq N$. To further enhance reasoning performance, we introduce the reasoning verification mechanism within the existing multi-task learning framework. This mechanism enhances the learning process by providing verification signals for additional-task learning. The details will be elaborated in the subsequent sections.

*1) Learning with Rationale:* To train the proficient reasoners, we employ the multi-task learning framework and divide the learning task into explanation and prediction, where explanation tends to furnish the rationale based on the input sample and the task objective, while the prediction task focuses solely on inferring sentiment polarity. Through the implementation of multi-task learning, the training phase incorporates the losses associated with both explanation and prediction tasks. The loss function is delineated as follows, where $\mathcal{L}_{exp}$ is the loss for explanation task and $\mathcal{L}_{pre}$ is the loss for prediction task:

$$\mathcal{L}_{loss} = \alpha\mathcal{L}_{exp} + (1 - \alpha)\mathcal{L}_{pre} \tag{1}$$

where the prediction $\mathcal{L}_{pred}$ aims to minimize the cross-entropy loss for label prediction:

$$\mathcal{L}_{pre} = \frac{1}{N}\sum_{N}^{i=1} \ell_{CE}(\hat{y}_i, y_i) \tag{2}$$

while the explanation loss $\mathcal{L}_{exp}$ tends to minimize the generation loss for the rationale, and there exists a subtle distinction between reasoning (RE) and rationalization (RA) scenarios.

$$RE : \mathcal{L}_{exp} = \frac{1}{N}\sum_{N}^{i=1} \ell_{CE}(f(x_i, t_i), \hat{e}_i) \tag{3}$$

$$RA : \mathcal{L}_{exp} = \frac{1}{N}\sum_{N}^{i=1} \ell_{CE}(f(x_i, t_i, y_i), \hat{e}_i) \tag{4}$$

The objective is to equip the model with proficiency in both explanation and prediction, thereby enhancing its reasoning capabilities. However, during the inference phase, only the prediction task is required for evaluation to optimize the inference efficiency and mitigate computational costs.

*2) Reasoning with Verification:* Considering the rationale generated by LLM is directly employed without any post-filtering processes, it might introduce some error patterns that can negatively influence the performance of multi-task fine-tuning. Some research works perform answer-based filtering to improve the rationale quality. [34] directly removed the incorrect rationale given by reasoning prompting based on the final prediction and supplemented it with the rationale generated under rationalization prompting to complete the final rationale set for training. [10] demonstrated that answer-based filtering can also be compensated by a diversity of reasoning paths using diverse reasoning and retaining the rationales leading to the correct answer. In our approach, we preserve the sets of rationales generated by the LLM by introducing a verification signal to facilitate further analysis of rationale quality within the multi-task learning framework. This is achieved by incorporating an additional task for verification.

Specifically, we leverage the rationale set generated by our TH-RE prompting and adopt answer-based verification according to the prediction label $\hat{y}_i$ provided by the LLM and the ground truth annotation $y_i$. Rationales that lead to the correct answer label are deemed to possess higher quality and utility compared to those inferring an incorrect answer label. Based on this premise, we complete the prompting using the following template:

> Given the rationale $e_i$, Please verify whether the above given rationale is reasonable. Return True or False.

To generate the verification signal $v_i$, we validate the reasoning rationales that successfully infer the correct answer and provide a general signal as unreasonable with an *False* label for the other rationales, indicating that they could benefit from further refinement. However, according to our observation, LLMs with larger parameter scales, like GPT-3.5-Turbo, tend to predict ambiguous answers containing dual polarities when

faced with uncertainty in making a final judgment. Therefore, we establish the verification signal $v_i$ based on the following criteria:

$$\begin{cases} \text{if } \hat{y}_i \subset \{\hat{y}_{i(t_1)}, \hat{y}_{i(t_2)}\}, \ \hat{y}_i = \hat{y}_{i(t_1)}, \\ \text{if } \hat{y}_i = y_i, \ v_i \ is \ True \end{cases} \quad (5)$$

In cases where the rationale presents two polarities, $\hat{y}_{i(t_1)}$ and $\hat{y}_{i(t_2)}$, where $t_2 > t_1$, answer-based verification is conducted on $\hat{y}_{i(t_1)}$ based on the First-Fome-First-Served (FCFS) rule, since the label generated earlier is regarded as holding a greater likelihood according to the next token generation. Then the revised loss function incorporating verification signals is formulated as follows:

$$\mathcal{L}_{loss} = \alpha\mathcal{L}_{exp} + \gamma\mathcal{L}_{ver} + (1 - \alpha - \gamma)\mathcal{L}_{pre} \quad (6)$$

where the verification loss concerns the self-validation outcome under the supervision of the verification signal:

$$\mathcal{L}_{ver} = \frac{1}{N}\sum_{N}^{i=1} \ell_{CE}(f(e_i), \hat{v}_i) \quad (7)$$

## IV. Experiments

### A. Setups

In the experiments, we evaluate the results on Restaurant and Laptop datasets in SemEval-2014 [38]. To test the performance for ISA, we follow the prior works utilizing datasets that further labeled with explicit and implicit tags [5]. To generate effective rationales conducive to reasoning learning, we make use of DO LLMs, Vicuna-13B [18] and GPT-3.5-turbo [14] in stage 1 for rationale preparation. Considering the impressive performance of ED style models in understanding input information and comprehension among different tasks, Flan-T5 [15] serves as the backbone LLM during the multi-task fine-tuning stage. We test with different sizes of Flan-T5, scaling from the base model (250M) to the XXL model (13B). For the baseline methods, we compared with the recently reported best results, including seven baseline methods, which are BERT+SPC [39], BERT+ADA [40], BERT+RGAT [41], $BERT_{Asp}$+CEPT [5], $BERT_{Asp}$+SCAPT [5], THOR [6] and ABSA-ESA [24]. Among them, THOR [6] stimulates performance based on CoT prompting with three-step generation. Compared to their method, we utilize a multi-task learning framework during training while directly inferring the final prediction during inference time. To identify the optimal hyperparameters in the training loss, a greedy search is undertaken using the validation set to determine the final values of $\alpha$ and $\gamma$. Without verification supervision, we get the best result when $\alpha = 0.5$ with explanation and prediction tasks only. With the verification supervision, we get the greatest performance when $\alpha = \gamma = 0.3$. The following experiments will follow this hyperparameter setting.

### B. Main Results

*a) Multi-task learning outperforms the baselines:* The main results of baselines and our method, RVISA, are demonstrated in Table I. The evaluation metrics include Accuracy and Macro-F1 score. Notably, as THOR [6] does not provide the accuracy outcome for implicit sentiment, we rerun the results based on the provided source code. It can be seen that RVISA significantly outperforms the baseline methods, irrespective of whether learning is from Vicuna-13B or GPT-3.5-turbo, underscoring the efficacy of learning within the proposed multi-task learning framework.

*b) Strong teachers lead to higher quality learning:* The performance of $RVISA_g$ training under the assistance of GPT-3.5-turbo exhibits enhanced reasoning capabilities in implicit sentiment inference compared to $RVISA_v$ trained by using the rationales generated by Vicuna-13B. This disparity can be attributed to the superior common sense knowledge and reasoning prowess exhibited by GPT-3.5-turbo in producing high-quality rationales, which play a pivotal role in transferring reasoning abilities to the Flan-T5 backbone model. However, the smaller backbone model with a base size (250M) lags behind some of the baseline methods due to its limited generation capacity to derive advantage from rich knowledge through in-context learning.

*c) Explicit rationales learning helps implicit reasoning:* We further compare our method with THOR, which is built upon a chain-of-thought strategy. Instead of eliciting the reasoning ability of the language model through sequential three-step prompting, our method stands out by explicitly giving the rationales as informative resources equipped with a verification mechanism to ensure the learning quality. The comparative results are depicted in Table II. RVISA demonstrates superior performance over THOR in terms of F1 score for implicit sentiment analysis while maintaining competitive results in overall F1 score. This underscores the effectiveness of our method in learning implicit sentiment through reasoning tasks and adeptly capturing implicit relationships among instances. Although THOR claimed the three-step generation during fine-tuning can unleash the reasoning power of the backbone model, the re-run result demonstrated limited improvement in F1 scores. This suggests the vulnerability of THOR to enhance prompt-based inference depending on the backbone model (i.e., Flan-T5 [15]) itself. In contrast, our method prioritizes effective learning from high-quality sentiment information and closely related tasks, offering a more coherent and justifiable approach to achieving high performance in implicit sentiment analysis.

### C. Ablation Study

We conducted an ablation study on the three-hop prompting (TH) and verification mechanism (VE) components, the results of which are summarized in Table III. Our analysis compares the F1 scores in both overall and implicit sentiment scenarios. The findings indicate that the absence of the verification mechanism leads to performance degradation in both cases, with a more significant decline of over one point observed in the implicit sentiment results. This highlights the critical role of verification signals in the context of reasoning learning from LLMs, as the answer-based mechanism aids the backbone model in identifying potential errors or unreasonable attributes during multi-task learning processes.

TABLE I

MAIN RESULTS COMPARED WITH BASELINES ON RESTAURANT AND LAPTOP DATASETS. THE RESULTS WITH † AND ⋆ ARE OBTAINED FROM [5] AND [24], WHILE THE OTHER RESULTS ARE SELF-RERUN OR SELF-IMPLEMENTED. IN OUR METHODS, THE SUBSCRIPTS STAND FOR LEARNING FROM RATIONALES GENERATED BY DIFFERENT MODELS, WHICH ARE VICUNA-13B($v$) AND GPT-3.5-TURBO($g$), RESPECTIVELY. THE SUBSCRIPTS $A$ AND $F$ REPRESENT THE ACCURACY AND MACRO-F1 SCORE.

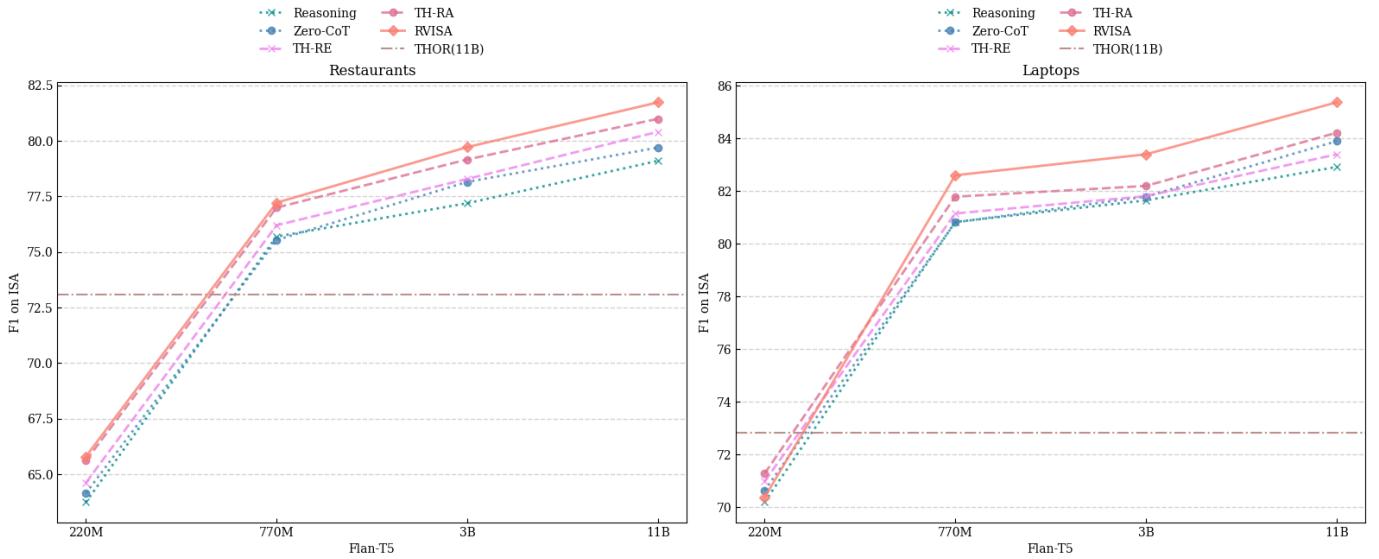| | Restaurant | | | Laptop | | |
|---|---|---|---|---|---|---|
| | $\text{All}_A$ | $\text{All}_F$ | $\text{ISA}_A$ | $\text{All}_A$ | $\text{All}_F$ | $\text{ISA}_A$ |
| *- State-of-the-art baselines* | | | | | | |
| BERT + SPC† (110M) [39] | 83.57 | 77.16 | 65.54 | 78.22 | 73.45 | 69.54 |
| BERT + ADA† (110M) [40] | 87.14 | 80.05 | 65.92 | 78.96 | 74.18 | 70.11 |
| BERT + RGAT† (110M) [41] | 86.60 | 81.35 | 67.79 | 78.21 | 74.07 | 72.99 |
| $\text{BERT}_{Asp}$ + CEPT† (110M) [5] | 87.50 | 82.07 | 67.79 | 81.66 | 78.38 | 75.86 |
| $\text{BERT}_{Asp}$ + SCAPT† (110M) [5] | 89.11 | 83.79 | 72.28 | 82.76 | 79.15 | 77.59 |
| $\text{T5}_{Base}$ + ABSA-ESA⋆ (220M) [24] | 88.29 | 81.74 | 70.78 | 82.44 | 79.34 | 80.00 |
| *- Prompt-based methods* | | | | | | |
| Flan-T5 + prompt (250M) | 86.88 | 79.78 | 65.17 | 81.98 | 77.93 | 73.71 |
| Flan-T5 + prompt (11B) | 89.29 | 83.68 | 75.28 | 81.82 | 77.69 | 75.43 |
| Flan-T5 + THOR (250M) [6] | 87.68 | 81.10 | 68.54 | 81.66 | 77.51 | 74.29 |
| Flan-T5 + THOR (11B) [6] | 88.57 | 82.93 | 73.03 | 82.29 | 78.78 | 76.57 |
| *- Our methods* | | | | | | |
| Flan-T5 + $\text{RVISA}_v$ (250M) | 86.43 | 78.49 | 65.92 | 80.72 | 76.49 | 73.71 |
| Flan-T5 + $\text{RVISA}_g$ (250M) | 86.61 | 78.92 | 66.67 | 81.19 | 77.13 | 75.43 |
| Flan-T5 + $\text{RVISA}_v$ (11B) | 91.25 | 86.57 | 81.65 | 86.52 | 83.28 | 87.43 |
| Flan-T5 + $\text{RVISA}_g$ (11B) | **91.52** | **86.85** | **82.02** | **86.68** | **84.05** | **88.00** |



Fig. 4. The impact of diverse rationales and different model sizes on implicit F1 score. The dashed horizontal line represents the best result of THOR rerun with the Flan-T5-XXL(11B) model on the implicit dataset.

TABLE II

RESULTS COMPARED WITH THOR [6]. THE EVALUATION METRIC IS THE F1 SCORE TRAINED WITH FLAN-T5. THE RESULTS WITH † ARE SELF-RERUN USING THE SOURCE CODE FROM [6].

| | Restaurant | | Laptop | |
|---|---|---|---|---|
| | All | ISA | All | ISA |
| Prompt† (11B) | 83.68 | 74.48 | 77.69 | 72.44 |
| THOR† (11B) | 82.93 | 73.08 | 78.78 | 72.82 |
| $\text{RVISA}_v$ (11B) | 86.57 | 81.73 | 83.26 | 85.36 |
| $\text{RVISA}_g$ (11B) | **86.85** | **82.61** | **84.05** | **86.20** |

TABLE III

ABLATION STUDY OF THREE-HOP PROMPTING (TH) AND VERIFICATION (VE) WITH F1 SCORE METRIC.

| | Restaurant | | Laptop | |
|---|---|---|---|---|
| | All | ISA | All | ISA |
| $\text{RVISA}_v$ | 86.57 | 81.73 | 83.26 | 85.36 |
| - w/o VE | 85.91 | 80.40 | 82.63 | 83.39 |
| - w/o VE and TH | 85.79 | 79.10 | 82.57 | 82.91 |
| $\text{RVISA}_g$ | **86.85** | **82.61** | **84.05** | **86.20** |
| - w/o VE | 86.16 | 80.32 | 82.51 | 83.83 |
| - w/o VE and TH | 85.60 | 79.68 | 82.05 | 83.14 |

In addition, performance goes down even more when we reduce the CoT prompting from three-hop reasoning to reasoning prompting alone, without the sentiment elements to

help guide rationale generation. This happens in both implicit and general scenarios. These observations persist regardless of whether the rationales are generated by Vicuna-13B or GPT-3.5-turbo, indicating that, irrespective of generation quality, the three-hop prompting mechanism plays a pivotal role in steering the correct direction of reasoning for implicit sentiment analysis. Although the impact of performance degradation with three-hop reasoning prompting is less pronounced compared to the absence of the verification mechanism, it is evident that their contributions are mutually reinforcing and indispensable. It is also the essence of multi-task learning, where tasks are strongly related and complement each other.

### D. Further Analysis

*a) The impact of rationale:* In our investigation of the influence of diverse rationales, we conducted training experiments using rationales generated by Vicuna-13B with various prompting methods, including Reasoning, Zero-CoT, Three-hop Reasoning (TH-RE), and Three-hop Rationalization (TH-RA), as depicted in Figure 2. We compared the results with RVISA and THOR, as shown in Figure 4, where RVISA is enhanced by Three-hop Reasoning prompting with the verification mechanism. It can be seen that the model trained with TH-RA demonstrates the second-best results since rationalization prompting can leverage the gold answer as context information to elucidate the underlying logic. This approach facilitates the generation of more reasonable rationales that lead to correct answers. Consequently, TH-RA generally outperforms TH-RE, where TH-RE may produce more problematic responses, resulting in incorrect answers. However, RVISA consistently outperforms both TH-RA and TH-RE, suggesting that the language model, when trained under verification signals, can leverage erroneous or irrational attributes present in TH-RE-generated rationales. This provides a visible solution to utilize LLM-generated labels as an additional verification factor. Furthermore, rationales generated by Reasoning and Zero-CoT methods lag behind Three-hop prompting in most scenarios, underscoring the importance of our designed prompting approach in structuring coherent rationales and eliciting highly relevant sentiment elements within the three-hop prompting.

*b) The impact of model size:* Figure 4 also illustrates the impact of backbone model size on reasoning learning. In the Restaurant dataset, smaller-sized models (i.e., base and large) exhibit marginal performance improvements under the verification mechanism, indicating the limited capabilities of small models to benefit from the prompt-based inference within the multi-task framework. However, as model size increases, the combined benefits of the verification mechanism and three-hop reasoning prompting demonstrate enhanced potential, leading to a widening performance gap compared to the second-best TH-RA method. Notably, with large (770M) size models, RVISA achieves superior performance to the best result of THOR trained with the Flan-T5-XXL (11B) model on both Restaurant and Laptop datasets, showcasing the efficacy of our method in enhancing reasoning abilities for pre-trained models. When it comes to the XXL (11B) size, TH-RE, TH-RA, and RVISA collectively surpass THOR in

implicit sentiment prediction. The Laptop dataset demonstrates similar trends. All prompting methods with XXL size model under the multi-task learning framework surpass the best result of THOR, emphasizing the effectiveness of our proposed framework and the scaling effect influenced by the learning capabilities of the trained model.

## V. Discussion

We propose a two-stage reasoning framework, RVISA, to learn effectively and reliably from the rationales generated by DO LLMs for implicit sentiment analysis. We show that RVISA holds the potential to promote the reasoning and learning ability of ED model under the supervision of verification through extensive experiments. In this section, we discuss the error scenario after fine-tuning using our proposed method and the limitations for further improvements.

*a) Error Analysis:* Our proposed method demonstrates superior performance in implicit sentiment analysis. To further explore the error scene, we calculate the error ratio considering sentiment types, including explicit and implicit, and with the relationship to the corresponding sentiment labels. The result is shown in Figure 5 with the rationales generated from GPT-3.5-turbo. It can be observed that for the Laptop dataset, errors in neutral predictions within the explicit dataset surpass those in the implicit dataset, resulting in the F1 score performance in the implicit dataset exceeding that in all data. Conversely, in the Restaurant dataset, the error ratio associated with neutral predictions in the implicit dataset exceeds that of the original dataset. This observation underscores the significant influence of neutral sentiment distribution on error distribution patterns. Moreover, the ratio of incorrect predictions pertaining to *neutral* polarity exceeds $60\%$. This suggests the nuanced challenges associated with accurately discerning neutral sentiments within sentiment analysis tasks, highlighting the need for further refinement and optimization in model training and inference processes.
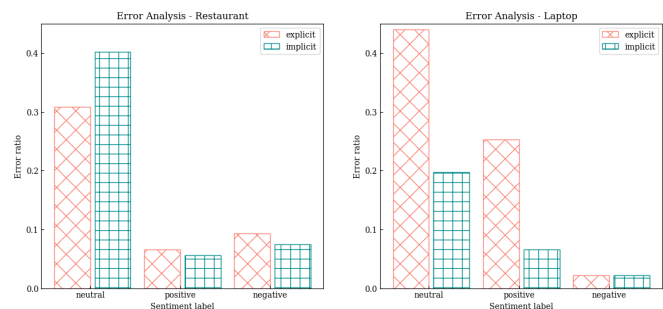


Fig. 5. Error analysis for two datasets with rationales generated by GPT-3.5-turbo. The error ratio here refers to the proportion of the number of error types to the total number of error instances.

*b) Limitations:* In this study, we propose a straightforward yet effective verification mechanism to enhance the overall performance in sentiment analysis. The answer-based verification plays a key role in the RVISA framework, demonstrating its significance in reasoning learning. It is worth noting that while the current answer-based verification signal is effective, there is potential for further enhancement

through the exploration of alternative verification modes or the incorporation of additional pertinent factors. This avenue for future research paves the way for more nuanced and reliable sentiment analysis. On the other hand, the three-hop prompting proves instrumental in generating effective rationales by deducing sentiment elements. It is manually designed with the format drawing on prior works, which poses challenges in further optimization. Given the evolving landscape of advanced techniques focused on optimizing prompts for LLMs, it is unclear whether the prompt can be generated automatically or optimized through the utilization of soft prompts in this study. This raises a feasible direction for further exploration.

## VI. Conclusions

In conclusion, this study sheds light on implicit sentiment analysis in the era of LLMs and proposes a novel two-stage learning framework, RVISA, designed to incorporate reasoning and verification for implicit sentiment analysis. By leveraging the generative prowess of DO LLMs, we empower ED backbone models with enhanced reasoning capabilities. The utilization of three-hop reasoning prompting facilitates the explicit generation of cues guided by sentiment element construction, which is conducive to reasoning learning. Through a straightforward and effective answer-based verification mechanism, we ensure robust and reliable reasoning learning to further improve the proficiency of our ED backbone model in inferring implicit sentiment. The experimental results demonstrate superior performance and achieve state-of-the-art results in ISA on two benchmark datasets.

## Appendix

### Rationale Generation

*a) Vicuna-13B versus GPT-3.5-turbo:* To delve into the quality of rationales generated from Vicuna-13B and GPT-3.5-turbo, the analysis for wrong and ambiguous prediction is conducted as illustrated in Figure 6. In both Restaurant and Laptop datasets, Vicuna-13B exhibited a slightly higher count of incorrect predictions compared to GPT-3.5-turbo. This suggests that stronger models such as GPT-3.5-turbo demonstrate a superior capability to generate higher-quality rationales, leading to more accurate final predictions. However, the percentage of ambiguous predictions originating from GPT-3.5-turbo surpassed that of Vicuna-13B, which indicates that the more powerful model exhibits a greater tendency to generate uncertain expressions rather than provide definitive judgments when deciphering the nuanced sentiment. It also underscores the inherent challenge of capturing subtle nuances in sentiment within constrained contextual information.

*b) Reasoning (Re) versus Three-hop Reasoning (TH-RE):* Table IV shows the results when generating the rationale with GPT-3.5-turbo using diverse prompting. Under the reasoning prompt, the LLM tends to break down the problem into four fine-grained steps, leveraging its own pre-trained knowledge. However, even with more steps of inference, the prediction is the complete opposite of the gold label. On the other hand, with the guidance of three-hop prompting, the rationale follows the generation format as shown in the prompting
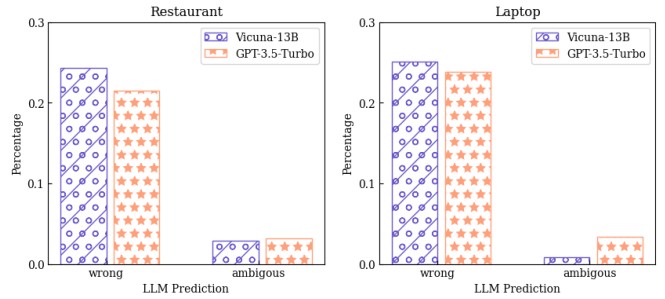


Fig. 6. Wrong prediction and ambiguous prediction analysis for Vicuna-13B and GPT-3.5-turbo, where ambiguous stands for multiple uncertain predictions in the rationale.

template, which is more concise and structured to infer the sentiment elements leading to the final correct prediction.
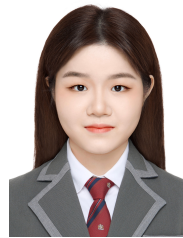
TABLE IV
REASONING(RE) PROMPTING COMPARED WITH THREE-HOP REASONING(TH-RE) PROMPTING FOR GPT-3.5-TURBO MODEL.

| | |
|---|---|
| **Input Text** | a cheaper price should not equal a "cheap" product. |
| **Aspect Term** | price |
| **Sentiment Polarity** | positive |
| **Re Prompt** | Given the sentence {input text}, what's the sentiment polarity towards price, why? |
| **GPT-3.5-turbo** | The sentiment polarity towards price in this sentence is negative. The sentence suggests that a lower price does not necessarily mean that the product is of low quality or "cheap." In other words, the sentence implies that there is a common misconception that cheaper products are of inferior quality, which is not necessarily true. To break it down by step: 1. The sentence contains the word "cheaper," which is a synonym for "lower-priced" or "more affordable." This suggests that the speaker is concerned with the price of the product. 2. However, the sentence also contains the word "cheap," which is often used to describe products of low quality or value. This suggests that the speaker is concerned with the perception of the product's quality. 3. The sentence negates the idea that a cheaper price automatically equals a "cheap" product. This suggests that the speaker believes that there is a common misconception that cheaper products are of inferior quality, which is not necessarily true. 4. Therefore, the sentiment towards price is negative because the sentence implies that people often assume that cheaper products are of low quality, which is not always the case. |
| **TH-RE Prompt** | Given the sentence {input text}, what's the sentiment polarity towards price, why? Let's think step by step. **The mentioned aspect** towards price is about ... **The underlying opinion** towards price is about ... **Therefore, the sentiment polarity** towards price is ... |
| **GPT-3.5-turbo** | **The mentioned aspect** towards price is about the relationship between price and the quality of a product. **The underlying opinion** towards price is that just because something is offered at a lower price, it should not automatically be assumed to be of low quality or "cheap." **Therefore, the sentiment polarity** towards price in this sentence is positive. The sentiment is emphasizing that price should not be the sole determinant of product quality, advocating for a more nuanced understanding of the pricing and value relationship. |

REFERENCES

[1] B. Liu, *Sentiment Analysis and Opinion Mining*, Jan 2012. [Online]. Available: http://dx.doi.org/10.1007/978-3-031-02145-9

[2] R. Das and T. D. Singh, "Multimodal sentiment analysis: A survey of methods, trends, and challenges," *ACM Computing Surveys*, vol. 55, pp. 1 – 38, 2023. [Online]. Available: https://api.semanticscholar.org/CorpusID:257233574

[3] J. Cui, Z. Wang, S.-B. Ho, and E. Cambria, "Survey on sentiment analysis: evolution of research methods and topics," *Artificial Intelligence Review*, p. 8469–8510, Aug 2023. [Online]. Available: http://dx.doi.org/10.1007/s10462-022-10386-z

[4] I. Russo, T. Caselli, and C. Strapparava, "Semeval-2015 task 9: Clipeval implicit polarity of events," in *Proceedings of the 9th International Workshop on Semantic Evaluation (SemEval 2015)*, Jan 2015. [Online]. Available: http://dx.doi.org/10.18653/v1/s15-2077

[5] Z. Li, Y. Zou, C. Zhang, Q. Zhang, and Z. Wei, "Learning implicit sentiment in aspect-based sentiment analysis with supervised contrastive pre-training," in *Conference on Empirical Methods in Natural Language Processing*, 2021. [Online]. Available: https://api.semanticscholar.org/CorpusID:241032941

[6] H. Fei, B. Li, Q. Liu, L. Bing, F. Li, and T.-S. Chua, "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)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 1171–1182. [Online]. Available: https://aclanthology.org/2023.acl-short.101

[7] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting elicits reasoning in large language models," in *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., 2022. [Online]. Available: http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html

[8] T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, "Large language models are zero-shot reasoners," in *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., 2022. [Online]. Available: http://papers.nips.cc/paper_files/paper/2022/hash/8bb0d291acd4acf06ef112099c16f326-Abstract-Conference.html

[9] J. Wei, M. Bosma, V. Y. Zhao, K. Guu, A. W. Yu, B. Lester, N. Du, A. M. Dai, and Q. V. Le, "Finetuned language models are zero-shot learners," in *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022. [Online]. Available: https://openreview.net/forum?id=gEZrGCozdqR

[10] N. Ho, L. Schmid, and S. Yun, "Large language models are reasoning teachers," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 14 852–14 882. [Online]. Available: https://doi.org/10.18653/v1/2023.acl-long.830

[11] N. F. Rajani, B. McCann, C. Xiong, and R. Socher, "Explain yourself! leveraging language models for commonsense reasoning," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Jan 2019. [Online]. Available: http://dx.doi.org/10.18653/v1/p19-1487

[12] C. Hsieh, C. Li, C. Yeh, H. Nakhost, Y. Fujii, A. Ratner, R. Krishna, C. Lee, and T. Pfister, "Distilling step-by-step! outperforming larger language models with less training data and smaller model sizes," in *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, A. Rogers, J. L. Boyd-Graber, and N. Okazaki, Eds. Association for Computational Linguistics, 2023, pp. 8003–8017. [Online]. Available: https://doi.org/10.18653/v1/2023.findings-acl.507

[13] W. Zhang, Y. Deng, B. Liu, S. J. Pan, and L. Bing, "Sentiment analysis in the era of large language models: A reality check," *CoRR*, vol. abs/2305.15005, 2023. [Online]. Available: https://doi.org/10.48550/arXiv.2305.15005

[14] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. F. Christiano, J. Leike, and R. Lowe, "Training language models to follow instructions with human feedback," in *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., 2022. [Online]. Available: http://papers.nips.cc/paper_files/paper/2022/hash/b1efde53be364a73914f58805a001731-Abstract-Conference.html

[15] H. W. Chung, L. Hou, S. Longpre, B. Zoph, Y. Tay, W. Fedus, E. Li, X. Wang, M. Dehghani, S. Brahma, A. Webson, S. S. Gu, Z. Dai, M. Suzgun, X. Chen, A. Chowdhery, S. Narang, G. Mishra, A. Yu, V. Y. Zhao, Y. Huang, A. M. Dai, H. Yu, S. Petrov, E. H. Chi, J. Dean, J. Devlin, A. Roberts, D. Zhou, Q. V. Le, and J. Wei, "Scaling instruction-finetuned language models," *CoRR*, vol. abs/2210.11416, 2022. [Online]. Available: https://doi.org/10.48550/arXiv.2210.11416

[16] C. Raffel, N. M. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *J. Mach. Learn. Res.*, vol. 21, pp. 140:1–140:67, 2019. [Online]. Available: https://api.semanticscholar.org/CorpusID:204838007

[17] Y. Tay, M. Dehghani, V. Q. Tran, X. Garcia, J. Wei, X. Wang, H. W. Chung, D. Bahri, T. Schuster, H. S. Zheng, D. Zhou, N. Houlsby, and D. Metzler, "UL2: unifying language learning paradigms," in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. [Online]. Available: https://openreview.net/pdf?id=6ruVLB727MC

[18] L. Zheng, W.-L. Chiang, Y. Sheng, S. Zhuang, Z. Wu, Y. Zhuang, Z. Lin, Z. Li, D. Li, E. Xing, H. Zhang, J. E. Gonzalez, and I. Stoica, "Judging LLM-as-a-judge with MT-bench and chatbot arena," in *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023. [Online]. Available: https://openreview.net/forum?id=uccHPGDlao

[19] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, E. H. Chi, T. Hashimoto, O. Vinyals, P. Liang, J. Dean, and W. Fedus, "Emergent abilities of large language models," *Trans. Mach. Learn. Res.*, vol. 2022, 2022. [Online]. Available: https://openreview.net/forum?id=yzkSU5zdwD

[20] C. Zong, R. Xia, and J. Zhang, *Sentiment Analysis and Opinion Mining*. Springer Singapore, Jan 2021, p. 163–199. [Online]. Available: http://dx.doi.org/10.1007/978-981-16-0100-2_8

[21] M. Xu, D. Wang, S. Feng, Z. Yang, and Y. Zhang, "KC-ISA: An implicit sentiment analysis model combining knowledge enhancement and context features," in *Proceedings of the 29th International Conference on Computational Linguistics*, N. Calzolari, C.-R. Huang, H. Kim, J. Pustejovsky, L. Wanner, K.-S. Choi, P.-M. Ryu, H.-H. Chen, L. Donatelli, H. Ji, S. Kurohashi, P. Paggio, N. Xue, S. Kim, Y. Hahm, Z. He, T. K. Lee, E. Santus, F. Bond, and S.-H. Na, Eds. Gyeongju, Republic of Korea: International Committee on Computational Linguistics, Oct. 2022, pp. 6906–6915. [Online]. Available: https://aclanthology.org/2022.coling-1.601

[22] D. Zhou, J. Wang, L. Zhang, and Y. He, "Implicit sentiment analysis with event-centered text representation," in *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, Jan 2021. [Online]. Available: http://dx.doi.org/10.18653/v1/2021.emnlp-main.551

[23] S. Wang, J. Zhou, C. Sun, J. Ye, T. Gui, Q. Zhang, and X. Huang, "Causal intervention improves implicit sentiment analysis," in *Proceedings of the 29th International Conference on Computational Linguistics*, N. Calzolari, C.-R. Huang, H. Kim, J. Pustejovsky, L. Wanner, K.-S. Choi, P.-M. Ryu, H.-H. Chen, L. Donatelli, H. Ji, S. Kurohashi, P. Paggio, N. Xue, S. Kim, Y. Hahm, Z. He, T. K. Lee, E. Santus, F. Bond, and S.-H. Na, Eds. Gyeongju, Republic of Korea: International Committee on Computational Linguistics, Oct. 2022, pp. 6966–6977. [Online]. Available: https://aclanthology.org/2022.coling-1.607

[24] J. Ouyang, Z. Yang, S. Liang, B. Wang, Y. Wang, and X. Li, "Aspect-based sentiment analysis with explicit sentiment augmentations," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 17, pp. 18 842–18 850, Mar. 2024. [Online]. Available: https://ojs.aaai.org/index.php/AAAI/article/view/29849

[25] H. Peng, L. Xu, L. Bing, F. Huang, W. Lu, and L. Si, "Knowing what, how and why: A near complete solution for aspect-based sentiment analysis," *Proceedings of the AAAI Conference on Artificial Intelligence*, p. 8600–8607, Jun 2020. [Online]. Available: http://dx.doi.org/10.1609/aaai.v34i05.6383

[26] J. Wei, X. Wang, D. Schuurmans, M. Bosma, B. Ichter, F. Xia, E. H. Chi, Q. V. Le, and D. Zhou, "Chain-of-thought prompting

elicits reasoning in large language models," in *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., 2022. [Online]. Available: http://papers.nips.cc/paper_files/paper/2022/hash/9d5609613524ecf4f15af0f7b31abca4-Abstract-Conference.html

[27] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, "Language models are few-shot learners," in *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., 2020. [Online]. Available: https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfcb4967418bfb8ac142f64a-Abstract.html

[28] X. Wang, J. Wei, D. Schuurmans, Q. V. Le, E. H. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. [Online]. Available: https://openreview.net/pdf?id=1PL1NIMMrw

[29] Y. Fu, H. Peng, A. Sabharwal, P. Clark, and T. Khot, "Complexity-based prompting for multi-step reasoning," in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. [Online]. Available: https://openreview.net/pdf?id=yf1icZHC-l9

[30] Z. Zhang, A. Zhang, M. Li, and A. Smola, "Automatic chain of thought prompting in large language models," in *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. [Online]. Available: https://openreview.net/pdf?id=5NTt8GFjUHkr

[31] P. Hase, S. Zhang, H. Xie, and M. Bansal, "Leakage-adjusted simulatability: Can models generate non-trivial explanations of their behavior in natural language?" in *Findings of the Association for Computational Linguistics: EMNLP 2020*, T. Cohn, Y. He, and Y. Liu, Eds. Online: Association for Computational Linguistics, Nov. 2020, pp. 4351–4367. [Online]. Available: https://aclanthology.org/2020.findings-emnlp.390

[32] O.-M. Camburu, T. Rocktäschel, T. Lukasiewicz, and P. Blunsom, "e-snli: Natural language inference with natural language explanations," in *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31. Curran Associates, Inc., 2018. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2018/file/4c7a167bb329bd92580a99ce422d6fa6-Paper.pdf

[33] M. Jin, Q. Yu, D. Shu, H. Zhao, W. Hua, Y. Meng, Y. Zhang, and M. Du, "The impact of reasoning step length on large language models," *CoRR*, vol. abs/2401.04925, 2024. [Online]. Available: https://doi.org/10.48550/arXiv.2401.04925

[34] S. Li, J. Chen, Y. Shen, Z. Chen, X. Zhang, Z. Li, H. Wang, J. Qian, B. Peng, Y. Mao, W. Chen, and X. Yan, "Explanations from large language models make small reasoners better," *CoRR*, vol. abs/2210.06726, 2022. [Online]. Available: https://doi.org/10.48550/arXiv.2210.06726

[35] J. Huang, S. Gu, L. Hou, Y. Wu, X. Wang, H. Yu, and J. Han, "Large language models can self-improve," in *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, H. Bouamor, J. Pino, and K. Bali, Eds. Singapore: Association for Computational Linguistics, Dec. 2023, pp. 1051–1068. [Online]. Available: https://aclanthology.org/2023.emnlp-main.67

[36] E. Zelikman, Y. Wu, J. Mu, and N. Goodman, "Star: Bootstrapping reasoning with reasoning," in *Advances in Neural Information Processing Systems*, S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, Eds., vol. 35. Curran Associates, Inc., 2022, pp. 15 476–15 488. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2022/file/639a9a172c044fbb64175b5fad42e9a5-Paper-Conference.pdf

[37] L. C. Magister, J. Mallinson, J. Adamek, E. Malmi, and A. Severyn, "Teaching small language models to reason," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, A. Rogers, J. Boyd-Graber, and N. Okazaki, Eds. Toronto, Canada: Association for Computational Linguistics, Jul. 2023, pp. 1773–1781. [Online]. Available: https://aclanthology.org/2023.acl-short.151

[38] M. Pontiki, D. Galanis, J. Pavlopoulos, H. Papageorgiou, I. Androutsopoulos, and S. Manandhar, "Semeval-2014 task 4: Aspect based sentiment analysis," in *Proceedings of the 8th International Workshop on Semantic Evaluation (SemEval 2014)*, Jan 2014. [Online]. Available: http://dx.doi.org/10.3115/v1/s14-2004

[39] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the 2019 Conference of the North*, Jan 2019. [Online]. Available: http://dx.doi.org/10.18653/v1/n19-1423

[40] A. Rietzler, S. Stabinger, P. Opitz, and S. Engl, "Adapt or get left behind: Domain adaptation through bert language model finetuning for aspect-target sentiment classification," in *Proceedings of the 12th Language Resources and Evaluation Conference*. European Language Resources Association, 2020, p. 4933–4941.

[41] K. Wang, W. Shen, Y. Yang, X. Quan, and R. Wang, "Relational graph attention network for aspect-based sentiment analysis," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, D. Jurafsky, J. Chai, N. Schluter, and J. Tetreault, Eds. Online: Association for Computational Linguistics, jul 2020, pp. 3229–3238. [Online]. Available: https://aclanthology.org/2020.acl-main.295

**Wenna Lai** is currently a joint Ph.D. student at the Department of Computing, Hong Kong Polytechnic University, under the supervision of Prof. Qing Li, and the School of Computer Science, University of Technology Sydney, under the supervision of Prof. Guandong Xu. She has been working closely with Prof. Haoran Xie at the School of Data Science, Lingnan University, Hong Kong. Before that, she received her Master's degree in the Department of Electrical and Computer Engineering from the National University of Singapore. Her research interests include Affective Computing and NLP for Social Good.

**Haoran Xie** (Senior Member, IEEE) received a Ph.D. degree in Computer Science from City University of Hong Kong and an Ed.D degree in Digital Learning from the University of Bristol. He is currently the Acting Associate Dean and Associate Professor at the School of Data Science, Lingnan University, Hong Kong. His research interests include artificial intelligence, big data, and educational technology. He has published 411 research publications, including 236 journal articles such as IEEE TPAMI, IEEE TKDE, IEEE TAFFC, and IEEE TCVST. He is the Editor-in-Chief of Natural Language Processing Journal, Computers & Education: Artificial Intelligence and Computers & Education: X Reality. He has been selected as the World's Top 2% Scientists by Stanford University.

**Guandong Xu** (Member, IEEE) received the Ph.D. degree in computer science from Victoria University, Melbourne, VIC, Australia, in 2009. He is currently a Professor and a Program Leader at the School of Computer Science and Data Science Institute, University of Technology Sydney, Sydney, NSW, Australia. His research interests include data science, data analytics, recommender systems, web mining, user modeling, NLP, social network analysis, and social media mining.

**Qing Li** (Fellow, IEEE) received the B.Eng. degree in Computer Science from Hunan Univeristy, Hunan, China, in 1982, and the M.S. and Ph.D. degrees in Computer Science from the University of Southern California, LA, California, USA, in 1985 and 1988, respectively. Qing Li is a Chair Professor and Head at the Department of Computing, The Hong Kong Polytechnic University. His research focuses on data science, web mining, and artificial intelligence. He is a Fellow of IET, a Fellow of IEEE, a member of ACM SIGMOD and IEEE Technical Committee on Data Engineering. He is the chairperson of the Hong Kong Web Society, and is a steering committee member of DASFAA, ICWL, and WISE Society.