

Exploring Rollback Inference for Aspect-based Sentiment Analysis

Anonymous ACL submission

Abstract

With the giant help from pre-trained large language models (LLMs), templated sequence of how to organize the aspect-level elements become the hottest research target while only a few of them move their steps to inference, not to mention utilizing the semantic connection between aspect-level elements during it. We argue that, compared with the high computational cost methods of training language models, considering the inference process can also bring us potential benefits. Motivated by this, we propose *rollback inference* for aspect-based sentiment analysis, which can boost the performance of fine-tuned large language models with a tiny cost, and adapt to various language models. Extensive experiments in three datasets and multiple language models underscore the effectiveness of our proposed rollback inference strategies and the value of the semantic connections in inference.

1 Introduction

Aspect-based sentiment analysis (ABSA) has garnered growing interest in the community, encompassing four subtasks: aspect term extraction, opinion term extraction, aspect term category classification, and aspect-level sentiment classification. The initial two subtasks focus on extracting the aspect term and the opinion term present in the sentence. The objectives of the last two subtasks are to identify the category and sentiment polarity related to the extracted aspect term.

The sentiment quadruple extraction task, which is composed of four subtasks, poses a significant challenge for traditional classification-based models due to its complexity. In response to this challenge, recent studies have adopted a unified generative approach that circumvents the need for explicit modelling of the ABSA problem. This approach treats either the class index (Yan et al., 2021), or the desired sentiment element sequence (Zhang et al., 2021b,a; Bao et al., 2022), as the target output of

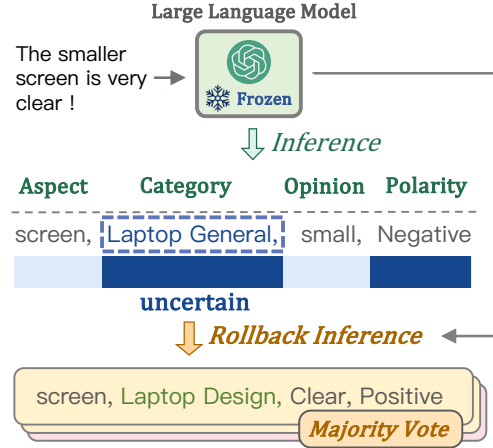


Figure 1: Example of proposed rollback inference framework.

the generative model. By doing so, these studies aim to simplify the overall task and improve its effectiveness.

However, most previous studies have concentrated on enhancing the training phase of generative models for sentiment analysis (Zhang et al., 2021b; Bao et al., 2022; Hu et al., 2022), simply adopting greedy search from left to right and neglecting the significance of the inference stage. As a result, the majority of these models rely on post-processing steps to ensure structural integrity (Bao et al., 2022, 2023b). In addition, these models fail to grasp the correlations among sentiment elements during inference (Hu et al., 2022), thus compromising the comprehensiveness of the sentiment analysis.

In this study, we direct our attention to the inference stage of sentiment generation models. We observe that, once the model is uncertain about one element, it tends to perform a similar attitude on the other elements in the same quadruple that are semantically connected. As shown in Figure 1, uncertainty is reported for both the category of aspect and polarity.

Motivated by this, we introduce a novel self-

consistency rollback inference framework along with a set of rollback strategies to better capture the correlations among sentiment elements during inference and improve its overall effectiveness. As illustrated in Figure 1, we employ an entropy-based mechanism to assess the uncertainty of sentiment elements during inference. When an element is deemed uncertain based on its entropy score, we launch a rollback procedure. This rollback is performed on a specific span determined by our proposed rollback strategies, resampling the span multiple times to get diverse results. Finally, we employ a majority vote mechanism to determine the final results for the rollback span.

The detailed evaluation shows that our model significantly advances the state-of-the-art performance on several benchmark datasets. In addition, the empirical studies also indicate that the proposed rollback inference strategy is more effective than other inference strategies.

2 Related Work

Generative ABSA: Research on ABSA typically follows a progression from addressing individual sub-tasks to dealing with their intricate combinations. The initial focus is often on predicting a single sentiment element (Wang et al., 2021; Hu et al., 2019; Tang et al., 2016; Chen et al., 2022; Liu et al., 2021; Seoh et al., 2021; Zhang et al., 2022). Many studies also delve into exploring the joint extraction of sentiment elements (Xu et al., 2020; Li et al., 2022; Bao et al., 2023a; Zhang and Qian, 2020).

More recently, there are some attempts to tackle ABSA problem in a generative manner (Zhang et al., 2021a), either treating the class index (Yan et al., 2021) or the desired sentiment element sequence (Zhang et al., 2021b) as the target of the generation model. For example, Yan et al. (2021) employed a sequence-to-sequence pre-trained model to generate the sequence of aspect terms and opinion words directly. Meanwhile, Zhang et al. (2021a) proposed a paraphrasing model that utilized the knowledge of the pre-trained model via casting the original task to a paraphrase generation process. In addition, Bao et al. (2022) addressed the importance of correlations among sentiment elements, and proposed an opinion tree generation model to jointly detect all sentiment elements in a tree structure.

Decoding Strategies for LLMs: Multiple decod-

ing strategies for language models have been proposed on general tasks to explicitly promote diversity in the decoding process in the literature, e.g., temperature sampling (Ackley et al., 1985; Ficer and Goldberg, 2017), top-k sampling (Radford et al., 2019; Holtzman et al., 2018; Fan et al., 2018), nucleus sampling (Holtzman et al., 2020). Besides, for improving accuracy, Self-consistency(COT-SC) decoding (Wang et al., 2023) has been proposed to explore multiple different ways of thinking leading to its unique correct answer.

However, the regeneration of the entire sequence in COT-SC is not applicable to the ABSA task as the reasoning associations between quadruples are not as strong as the general reasoning process. Huang et al. (2023) controlled text generation with arbitrary plugins during inference, which however requires to be trained separately. Gou et al. (2023) employed a majority vote decoding over different template orders, treating elements equally without semantic distinction. Hu et al. (2023) somehow proposed marginalized unlikelihood learning to suppress the uncertainty-aware mistake tokens.

Unlike previous works that often require complex pre-processing or post-processing steps, our method does not need such procedures. Instead, it easily integrates with fine-tuned language models, achieving significant improvements with only a minor increase in inference time. This makes our strategy a practical and efficient solution for enhancing sentiment analysis during the inference stage.

3 Aspect-based Sentiment Analysis with Rollback Inference

As shown in Figure 2, we introduce a novel *rollback inference framework* for generative aspect-based sentiment analysis.

To begin, we first fine-tune a large language model and freeze its parameters before entering the inference stage. Next, during inference, we propose an entropy-based mechanism to assess the uncertainty of sentiment elements and resample the uncertain span (detailed in Section 4) multiple times to get diverse results to construct the candidates pool. Finally, we obtain a final self-consistency result for the rollback span with a majority vote mechanism over the candidates.

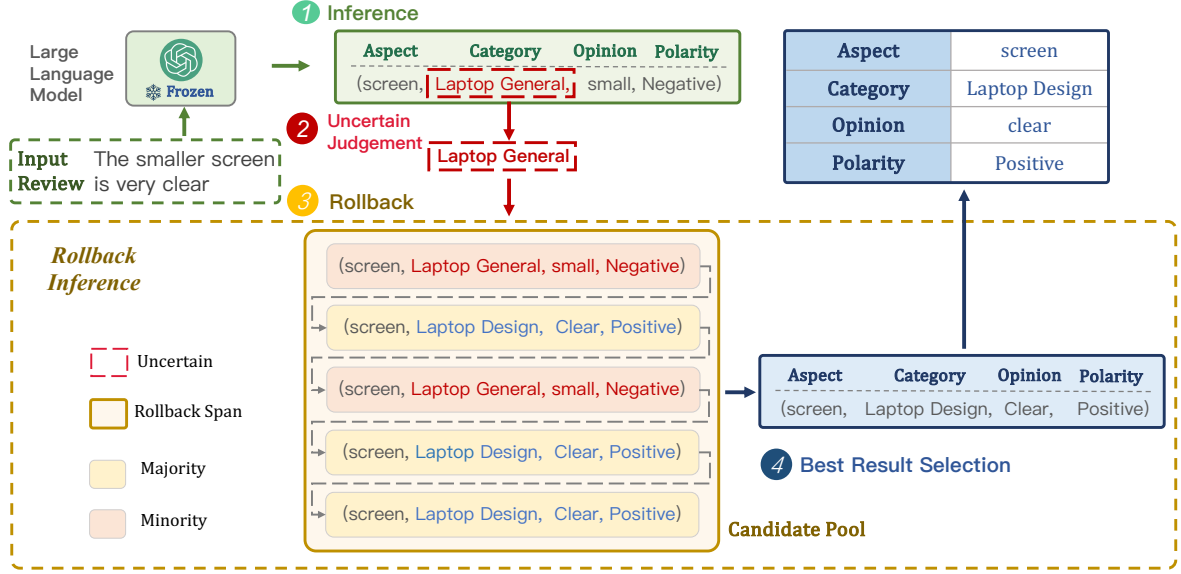


Figure 2: Overview of proposed rollback inference framework.

3.1 Generative Aspect-based Sentiment Analysis

In this study, we fine-tune the pre-trained large language model LLaMA (Touvron et al., 2023) as our foundation. This model receives a review sentence as input and produces sentiment quadruples as output. Each quadruple encompasses critical information regarding the sentiment expressed: the aspect term, opinion term, aspect category, and polarity.

Given the token sequence $x = x_1, \dots, x_{|x|}$ as input, the model outputs the linearized representation $y = y_1, \dots, y_{|y|}$. The decoder predicts the output sequence token-by-token. At the i -th step of generation, the decoder predicts the i -th token y_i in the linearized form, and decoder state h_i^d as:

$$y_i, h_i^d = ([h_1^d, \dots, h_{i-1}^d], y_{i-1}) \quad (1)$$

The conditional probability of the whole output sequence $p(y|x)$ is progressively combined by the probability of each step $p(y_i|y_{<i}, x)$:

$$p(y|x) = \prod_{i=1}^{|y|} p(y_i|y_{<i}, x) \quad (2)$$

where $y_{<i} = y_1 \dots y_{i-1}$, and $p(y_i|y_{<i}, x)$ are the probabilities over target vocabulary V .

The objective function is to maximize the output linearized opinion tree X_T probability given the review sentence X_O . Therefore, we optimize

the negative log-likelihood loss function:

$$\mathcal{L} = -\frac{1}{|\tau|} \sum_{(X_O, X_T) \in \tau} \log p(X_T|X_O; \theta) \quad (3)$$

where θ is the model parameters, and (X_O, X_T) is a $(sentence, tree)$ pair in training set τ , then

$$\begin{aligned} \log p(X_T|X_O; \theta) &= \\ &= \sum_{i=1}^n \log p(x_T^i|x_T^1, x_T^2, \dots, x_T^{i-1}, X_O; \theta) \end{aligned} \quad (4)$$

where $p(x_T^i|x_T^1, x_T^2, \dots, x_T^{i-1}, X_O; \theta)$ is calculated by the decoder.

3.2 Uncertain Element Judgement

During the inference stage of our generative aspect-based sentiment analysis model, we introduce an uncertain judgement mechanism to address elements that need rollback. This mechanism is triggered whenever the model generates a token with low confidence. Instead of accepting this uncertain token, we rollback to a previous state and re-generate the semantically connected span.

To quantify the model's certainty, we adopt information entropy as a metric. Specifically, for each generation step i , we calculate the entropy E_i using the formula:

$$E_i = -\sum_j^M P(x_j) \log(P(x_j)) \quad (5)$$

Here, $P(x_j)$ represents the output probability of the j -th token in the vocabulary, and M denotes the

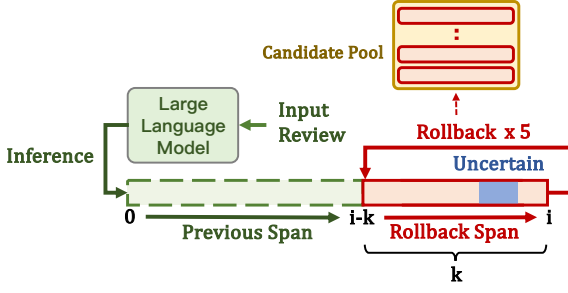


Figure 3: Example of rollback inference.

vocabulary size. A higher entropy E_i indicates that the model is less certain about its choice at step i .

When the entropy exceeds a predefined threshold, we consider the model to be uncertain and initiate the rollback process. This involves revisiting the semantically connected span and potentially generating a new set of candidates. The most confident candidate is then selected as the new output, ensuring that the model’s predictions are both self-consistent and reliable.

3.3 Rollback Inference

When an element is judged to be uncertain during the generation process, we employ a rollback strategy to revisit the corresponding span related to that element as shown in Figure 3. We adopt sampling in rollback inference, which choosing next token randomly with probability distribution instead of greedy search to ensure the diversity of candidates. Since the span of rollback is the key issue of this stage, we will discuss it in the next section.

We first generate sequence normally if there are no elements judged uncertain (green bar in Figure 3). Once an element is judged uncertain as the blue printed, we would like to rollback the corresponding span (printed red) related to it. Assuming we rollback at step i with a length k (determined by specific strategy), we would retreat the steps back to step $i - k$ and resample the following sequence to step i multiple times, the rollbacked span would be served as a candidate.

By rolling back multiple times, we can construct a pool of candidates for the uncertain sub-sequence. This pool provides the model with multiple options to choose from, increasing the chances of finding a more accurate and self-consistent prediction. The final prediction is then selected from this pool based on a predefined criterion, such as the highest confidence score or majority voting.

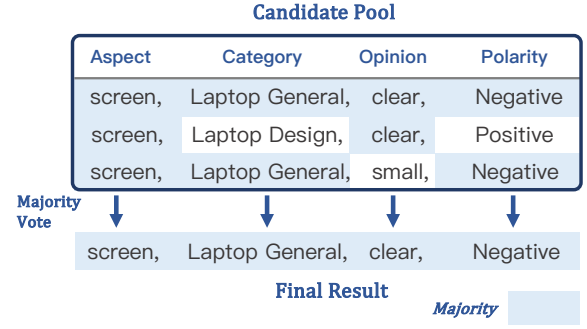


Figure 4: Illustration of the best result selection.

3.4 Best Result Selection

After constructing a pool of candidates for the uncertain sub-sequence, we proceed to select the best result from among these candidates as the final output. As illustrated in Figure 4, our approach involves dividing each candidate into its constituent sentiment elements. We then tally the votes for each element by counting the number of occurrences of its type (e.g., aspect, opinion, and polarity).

The sentiment element with the highest number of votes is subsequently selected as the final result. This majority voting mechanism allows us to leverage the collective wisdom of the model’s predictions, thereby increasing its confidence in the chosen output, especially for uncertain sub-sequences.

4 Rollback Inference Strategies

In this section, we first introduce the utilization of an opinion tree structure as a means to systematically organize and represent various sentiment elements. This tree structure serves as the backbone for rollback inference strategies. We then introduce different rollback inference strategies designed to select suitable candidates for uncertain sub-sequences in it.

4.1 Opinion Tree Construction

As shown in Figure 5, the opinion tree is hierarchically structured, beginning with a root node. The children of this root node are quadruple sub-trees, each rooted at an aspect node. These aspect nodes are then connected to category and opinion nodes, which together form the branches of the sub-tree. Polarity nodes are positioned as the successors of the corresponding opinion nodes, completing the structural representation of sentiment elements.

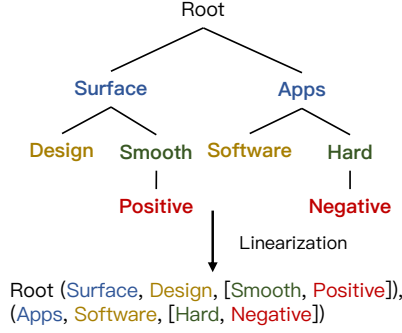


Figure 5: Example of the opinion tree structure.

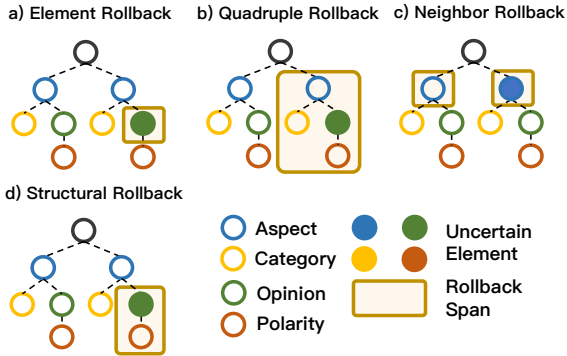


Figure 6: Illustration of proposed rollback inference strategies.

The linearization of this tree structure results in the final target sequence, which preserves the hierarchical relationships and semantic connections among sentiment elements.

4.2 Element Rollback

Element rollback inference (ER) represents a fundamental rollback strategy characterized by its narrow rollback span, which minimizes the additional inference time required.

As illustrated in Figure 6(a), when a token within an element is determined to be uncertain, the element would be regarded as the rollback span and underwent rollback multiple times to construct a pool of candidates.

4.3 Quadruple Rollback

Quadruple rollback inference (QR) is an intuitive strategy that recognizes the natural co-relation among the elements within a quadruple. This approach designs a holistic packaging strategy to address the entire quadruple as a unified entity.

As shown in Figure 6(b), when a token within the sub-sequence of a quadruple is deemed uncertain, the entire quadruple undergoes rollback. This means that instead of focusing solely on the

Domain	Train	Dev.	Test
Restaurant	1,529	171	582
Laptop	2,929	326	816
Phone	4,366	925	913

Table 1: Distribution of three domains.

uncertain token, quadruple rollback considers the broader context provided by the other elements within the quadruple.

4.4 Neighbor Rollback

Neighbor rollback inference (NR) is a strategy tailored to the structural formation of data, operating under the assumption that the neighbors (or sibling nodes) of an uncertain element may be influenced by its uncertainty.

As illustrated in Figure 6(c), when a token within an element of a quadruple is determined to be uncertain, neighbor rollback targets the siblings of this element as the rollback span. This means that instead of rolling back the entire quadruple or just the single uncertain element, neighbor rollback focuses on the immediate vicinity of the uncertain element.

4.5 Structural Rollback

In the context of structural opinion trees, the parent node (also known as the root node of a sub-tree) serves as the semantic foundation for the child nodes that originate from it. The uncertainty associated with a parent node has the potential to propagate throughout the entire sub-tree rooted at that node due to the shared semantic connections.

Recognizing this, we have developed a structural rollback inference strategy (SR) tailored to the inherent properties of the opinion tree. This strategy aims to address uncertainty at its source, the parent node, and mitigate its impact on the broader sub-tree structure.

As shown in Figure 6(d), during the inference process, if a token within a sentiment node of the opinion tree is deemed uncertain, the inference continues uninterrupted until it reaches the terminus of the sub-tree rooted at that sentiment node. Once this point is reached, the entire sub-tree undergoes multiple rollbacks initiated by the inference framework.

5 Experiments

In this section, we introduce the datasets used for evaluation and the baseline methods employed for

Method	Restaurant			Laptop			Phone		
	P	R	F1	P	R	F1	P	R	F1
JET	0.5981	0.2894	0.3901	0.4452	0.1625	0.2381	0.3845	0.2213	0.2809
TAS-BERT	0.2629	0.4629	0.3353	0.4715	0.1922	0.2731	0.3453	0.2207	0.2693
Extract-Classify	0.3854	0.5296	0.4461	0.4556	0.2948	0.3580	0.3128	0.3323	0.3223
GAS	0.6069	0.5852	0.5959	0.4160	0.4275	0.4217	0.5072	0.4815	0.4940
Paraphrase	0.5898	0.5911	0.5904	0.4177	0.4504	0.4334	0.4672	0.4984	0.4832
BARTABSA	0.5662	0.5535	0.5598	0.4165	0.4046	0.4105	0.4448	0.4734	0.4586
TODA	0.5904	0.6029	0.5966	0.4359	0.4367	0.4363	0.4720	0.4916	0.4816
DLO	0.5904	0.6029	0.5966	0.4359	0.4367	0.4363	0.5451	0.5173	0.5308
Seq2Path	0.6029	0.5961	0.5995	0.4448	0.4375	0.4411	0.5263	0.4994	0.5125
OTG	0.6191	0.6085	0.6164	0.4395	0.4383	0.4394	0.5302	0.5659	0.5474
One-ASQP	0.6591	0.5624	0.6069	0.4380	0.3954	0.4156	0.5742	0.5096	0.5400
ChatGPT	0.5014	0.3625	0.4207	0.4492	0.3123	0.3541	0.4514	0.4627	0.4569
LLaMA	0.6213	0.6024	0.6117	0.4334	0.4201	0.4266	0.5314	0.5478	0.5394
Ours	0.6585	0.6197	0.6382	0.4470	0.4417	0.4443	0.5387	0.5709	0.5543

Table 2: Comparison with baselines, we report the performance of our proposed model with structure rollback.

comparison. We then report the experimental results conducted from different perspectives, and analyze the effectiveness of the proposed model with different factors.

5.1 Dataset and Experiment Setting

In this study, we use restaurant and laptop domains in ACOS dataset (Cai et al., 2021) and phone domain in Zhou et al. (2023)’s dataset for our experiments. There are 2,286 sentences in the restaurant domain, 4,076 sentences in the laptop domain and 7,115 sentences in the phone domain. The distribution of these three domains can be found in Table 1.

For our decoder-only opinion tree generation model, we employ LLaMA-2-7B¹ and LoRA fine-tune the adapter parameters. In terms of encoder-decoder model, we employ T5². We tune the parameters of our models by grid searching on the validation dataset. We fine-tune the model with 20 epochs and save the model parameters for inference. The LoRA alpha is set to 128 and LoRA rank is set to 64.

The model parameters are optimized by Adam (Kingma and Ba, 2015), the learning rate of fine-tuning is 5e-5. The batch size is set to 16K with a cut-off length of 1024. The LoRA adapter would be merged with the original LLaMA-2-7B parameters and freeze during the inference process. During inference, we do sampling and set the entropy threshold to 0.6, rollback times to 5, top K

to 2, temperature to 0.95 with beam size 1. Our experiments are carried out with an Nvidia RTX 4090 GPU.

In evaluation, a quadruple is viewed as correct if and only if the four elements, as well as their combination, are exactly the same as those in the gold quadruple. On this basis, we calculate the Precision and Recall, and use F1 score as the final evaluation metric for aspect sentiment quadruple extraction (Cai et al., 2021; Zhang et al., 2021a).

5.2 Main Results

In Table 2, we present a comprehensive comparison of our proposed model with various state-of-the-art baselines. These baselines include both extraction-based methods and generative models, as well as large language models.

Extraction-based methods, such as JET (Xu et al., 2020), TAS-BERT (Wan et al., 2020), and Extract-Classify (Cai et al., 2021), typically rely on identifying relevant spans within the input text to extract sentiment quadruples. On the other hand, generative models, such as GAS (Zhang et al., 2021b), Paraphrase (Zhang et al., 2021a), BARTABSA (Yan et al., 2021), GAS (Zhang et al., 2021b), DLO (Hu et al., 2022), Seq2Path (Mao et al., 2022), OTG (Bao et al., 2022), and One-ASQP (Zhou et al., 2023), aim to generate sentiment quadruples from scratch, potentially allowing for more flexibility and creativity in their outputs. Besides, for large language models which have recently emerged as powerful tools for natural language processing tasks, we include zero-shot ChatGPT (Ouyang et al., 2022) and fine-tuned LLaMA-

¹LLaMA-2-7B-Chat, <https://huggingface.co/meta-llama/Llama-2-7b-chat-hf>

²T5_{base} as previous work did, https://huggingface.co/transformers/model_doc/t5.html

Method	Restaurant	Laptop	Phone
Greedy	0.6157	0.4251	0.5367
Ours-ER	0.6216	0.4382	0.5496
Ours-QR	0.6234	0.4420	0.5516
Ours-NR	0.6325	0.4397	0.5535
Ours-SR	0.6382	0.4443	0.5543

Table 3: Performance of rollback inference strategies.

2-7B (Touvron et al., 2023) in our baselines.

As shown in Table 2, we find that generative models outperform previous classification-based methods and the structural generative method surpasses non-structural methods, this indicates that semantic structure does contribute to quadruple extraction. It also shows that the unified generation architecture can fully utilize the rich label semantics by encoding the natural language label into the target output, and it is very helpful for extracting sentiment elements jointly.

Moreover, our proposed model exhibits significant improvements over all prior studies ($p < 0.05$), demonstrating the efficacy of our rollback inference framework when applied to large language models for sentiment element generation. To the best of our knowledge, this is the first attempt to leverage semantic relations explicitly during the inference process.

5.3 Comparison of Rollback Inference Strategies

Table 3 compares the effectiveness of various rollback inference strategies. The Greedy strategy, serving as a baseline, selects the token with the highest probability as the next token without considering the contextual or structural information. In contrast, the four proposed inference strategies, introduced in Section 4, aim to leverage the correlations among sentiment elements during inference.

The results clearly indicate that all of our proposed strategies surpass the Greedy approach, confirming the validity of our hypothesis that utilizing the relationships between sentiment elements in the inference stage is beneficial. Notably, the structure rollback inference strategy emerges as the most effective among all methods. We attribute this superior performance to the strategy’s ability to exploit structural self-consistency associations between sentiment elements, leading to more accurate and consistent predictions.

Furthermore, case studies in Appendix A are

Method	Manner	Time(s)	Avg. F-score
Sampling	Regular	80.58	0.5259
Greedy		79.80	0.5258
Beam		195.69	0.5276
COT-SC	Rollback	403.22	0.5370
Ours-ER		82.57	0.5364
Ours-QR		163.13	0.5390
Ours-NR		194.93	0.5419
Ours-SR		114.02	0.5456

Table 4: Analysis of inference efficiency, the speed is measured with seconds of generating 100 samples.

given to make more intuitive comparisons between the Greedy and proposed structural rollback.

6 Analysis and Discussion

In this section, we first launch the analysis about the computational efficiency of various inference strategies. We then have an analysis of robustness of our strategies towards various base models and templates.

6.1 Analysis of Inference Efficiency

In Table 4, we analyze the inference efficiency of various inference strategies. The first three strategies follow the conventional inference approach, generating tokens forward until the end of the sequence is reached. Sampling selects the next token based on its output probability, Greedy chooses the token with the highest probability, and Beam represents beam search, which maintains a set of candidate sequences at each step. The next five strategies incorporate rollback inference. In addition to the four rollback inference strategies proposed in this work, we also include COT-SC (Wang et al., 2023) as a baseline, where the rollback span covers the entire target sequence.

As evident from the results, the limited choices offered by Sampling and Greedy lead to their relatively poor performance. Beam search and COT-SC, on the other hand, improve upon these methods by maintaining a set of candidate sequences at each step. However, this comes at the cost of reduced inference speed as Beam search must evaluate multiple candidates at each step.

Within our rollback framework, the element rollback inference strategy stands out for its high speed and competitive performance. By limiting the rollback span to individual sentiment elements, it achieves a speed close to that of Greedy inference while still leveraging contextual information for

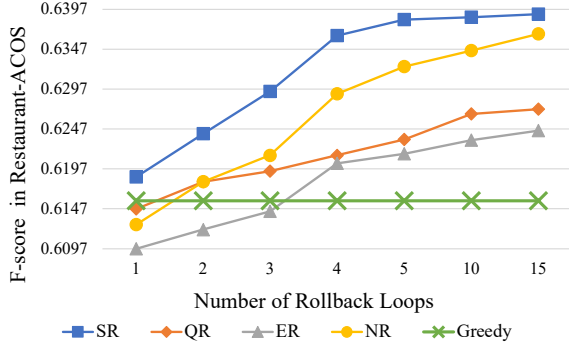


Figure 7: Performance of rollback strategies with different numbers of rollback loop.

improved accuracy. Finally, if we focus solely on performance, the structure rollback inference strategy emerges as the clear winner. It outperforms all other strategies, including COT-SC, while maintaining an acceptable inference speed.

6.2 Impact of Rollback Loops

We conducted a further investigation to assess the impact of rollback loops on our proposed rollback inference strategies. Specifically, we evaluated the performance of four different rollback inference strategies in the Restaurant domain, gradually increasing the number of rollback loops from 1 to 15.

As shown in Figure 7, the performance of all our strategies consistently improved as the number of rollback loops increased. This trend indicates that expanding the pool of candidates through additional rollback iterations enhances the self-consistency of large language models, leading to improved overall performance.

Among the tested strategies, structure rollback inference consistently outperformed the others across all loop counts, aligning with our previous experimental findings. Notably, it was the only strategy capable of surpassing greedy search even with the initial loop count of 1. This finding validates our hypothesis that leveraging the correlations among sentiment elements during inference can provide additional benefits.

To have more detailed investigation, we also have an analysis on the number of candidates in the candidates pool in Appendix B.

6.3 Impact of Language Models

We conducted an investigation to assess the impact of various language models, including LLaMA-2-

Model	Method	Restaurant	Laptop	Phone
LLaMA	Greedy	0.6157	0.4251	0.5367
LLaMA	SR	0.6382	0.4471	0.5543
T5	Greedy	0.6027	0.4129	0.5246
T5	SR	0.6209	0.4389	0.5489
BART	Greedy	0.3956	0.3191	0.3707
BART	SR	0.4177	0.3359	0.3911

Table 5: Results of different language models.

7B, T5-Base, and BART-Base. For each model, we evaluated both the greedy algorithm and structural rollback inference to obtain a comprehensive comparison.

As shown in Table 5, our structural rollback inference strategy proves to be effective across all language models, consistently outperforming the greedy algorithm. This suggests that our strategy is robust and can successfully capture the associations between sentiment elements during the inference stage, regardless of the underlying language model. This is a crucial finding as it highlights the versatility and applicability of our approach to different language models and scenarios.

7 Conclusion

In this study, we move our sight to the inference process of generative ABSA and are motivated to utilize the correlations between sentiment elements during it. We thus propose a self-consistency framework named Rollback Inference Framework along with a set of rollback strategies designed based on the intrinsic characteristics of the connections between sentiment elements in ABSA. Experimental results show that, without requiring complex and expensive training of LLMs, our proposed inference method can achieve state-of-the-art performance in ABSA with fine-tuned LLMs on the trade of a tiny cost in inference time.

The results also validate that, for generative templates that contain semantic connections like ABSA, ignoring utilizing semantic connections during inference could lead to a waste of them.

Limitations

The limitations of our work can be stated from two perspectives. First, we focus on structural rollback inference in ABSA only, more tasks that are close to ABSA like event extraction should be taken into consideration. Secondly, we only adopt the unsupervised entropy-based method to judge rollback

span, more methods of both unsupervised and supervised could be explored.

References

- David H. Ackley, Geoffrey E. Hinton, and Terrence J. Sejnowski. 1985. [A learning algorithm for boltzmann machines](#). *Cognitive Science*, 9(1):147–169.
- Xiaoyi Bao, Xiaotong Jiang, Zhongqing Wang, Yue Zhang, and Guodong Zhou. 2023a. [Opinion tree parsing for aspect-based sentiment analysis](#). In *Findings of the Association for Computational Linguistics: ACL 2023, Toronto, Canada, July 9-14, 2023*, pages 7971–7984. Association for Computational Linguistics.
- Xiaoyi Bao, Zhongqing Wang, Xiaotong Jiang, Rong Xiao, and Shoushan Li. 2022. [Aspect-based sentiment analysis with opinion tree generation](#). In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI 2022, Vienna, Austria, 23-29 July 2022*, pages 4044–4050. ijcai.org.
- Xiaoyi Bao, Zhongqing Wang, and Guodong Zhou. 2023b. [Exploring graph pre-training for aspect-based sentiment analysis](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3623–3634, Singapore. Association for Computational Linguistics.
- 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, Online. Association for Computational Linguistics.
- Chenhua Chen, Zhiyang Teng, Zhongqing Wang, and Yue Zhang. 2022. [Discrete opinion tree induction for aspect-based sentiment analysis](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2051–2064, Dublin, Ireland. Association for Computational Linguistics.
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Jessica Fidler and Yoav Goldberg. 2017. [Controlling linguistic style aspects in neural language generation](#). In *Proceedings of the Workshop on Stylistic Variation*, pages 94–104, Copenhagen, Denmark. Association for Computational Linguistics.
- Zhibin Gou, Qingyan Guo, and Yujiu Yang. 2023. [MvP: Multi-view prompting improves aspect sentiment tuple prediction](#). In *Proceedings of the 61st Annual*

Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 4380–4397, Toronto, Canada. Association for Computational Linguistics.

- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text degeneration](#). In *International Conference on Learning Representations*.
- Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. 2018. [Learning to write with cooperative discriminators](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1638–1649, Melbourne, Australia. Association for Computational Linguistics.
- Mengting Hu, Yinhao Bai, Yike Wu, Zhen Zhang, Liqi Zhang, Hang Gao, Shiwan Zhao, and Minlie Huang. 2023. [Uncertainty-aware unlikelihood learning improves generative aspect sentiment quad prediction](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 13481–13494, Toronto, Canada. Association for Computational Linguistics.
- Mengting Hu, Yike Wu, Hang Gao, Yinhao Bai, and Shiwan Zhao. 2022. [Improving aspect sentiment quad prediction via template-order data augmentation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7889–7900, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Mengting Hu, Shiwan Zhao, Li Zhang, Keke Cai, Zhong Su, Renhong Cheng, and Xiaowei Shen. 2019. [CAN: Constrained attention networks for multi-aspect sentiment analysis](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4601–4610, Hong Kong, China. Association for Computational Linguistics.
- Xuancheng Huang, Zijun Liu, Peng Li, Tao Li, Maosong Sun, and Yang Liu. 2023. [An extensible plug-and-play method for multi-aspect controllable text generation](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15233–15256, Toronto, Canada. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.
- Junjie Li, Jianfei Yu, and Rui Xia. 2022. [Generative cross-domain data augmentation for aspect and opinion co-extraction](#). In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4219–4229, Seattle, United States. Association for Computational Linguistics.

Review text	Method	Output
if it ' s nice outside, request for a table in the balcony	Greedy	(balcony, Ambience General, [NULL, Positive]) ✗
	SR Candidates Pool	(balcony, Ambience General, [nice, Positive]) ✓
		(balcony, Ambience General, [NULL, Positive]) ✗
		(balcony, Ambience General, [nice, Positive]) ✓
		(balcony, Ambience General, [nice, Positive]) ✓
the prior reviews said kid friendly	Greedy	(NULL, Restaurant Miscellaneous, [friendly, Positive]) ✗
	SR Candidates Pool	(NULL, Restaurant Miscellaneous, [friendly, Positive]) ✗
		(NULL, Restaurant Miscellaneous, [friendly, Positive]) ✗
		(NULL, Restaurant Miscellaneous, [NULL, Negative]) ✓
		(NULL, Restaurant Miscellaneous, [NULL, Negative]) ✓
i highly recommend this place to all that want to try indain food for the first time	Greedy	(place, Restaurant Miscellaneous, [recommend, Positive]) ✗
	SR Candidates Pool	(indain food, Food Quality, [recommend, Positive])) ✓
		(indain food, Food Quality, [recommend, Positive]) ✓
		(indain food, Food Quality, [recommend, Positive]) ✓
		(indain food, Food Quality, [recommend, Positive]) ✓
but she is very friendly with certain people , making it even more annoying	Greedy	(NULL, Service General, [friendly, Negative]) ✗
	SR Candidates Pool	(NULL, Service General, [friendly, Negative]) ✗
		(NULL, Service General, [annoying, Negative]) ✓
		(NULL, Service General, [annoying, Negative]) ✓
		(NULL, Service General, [annoying, Negative]) ✓
mercedes restaurant is so tasty, the service is undeniably awesome	Greedy	(mercedes restaurant, Food Quality, [tasty, Positive]) ✗
	SR Candidates Pool	(mercedes restaurant, Food Quality, [tasty, Positive]) ✗
		(NULL, Food Quality, [tasty, Positive])) ✓
		(mercedes restaurant, Food Quality, [tasty, Positive]) ✗
		(NULL, Food Quality, [tasty, Positive])) ✓

Table 6: Cases study, the quadruples in which are organized in (*Aspect, Category, [Opinion, Polarity]*) as introduced in Figure 5.

and the regular Greedy generation of fine-tuned LLaMA-2-7B. We select reviews that are predicted wrongly by Greedy but have been correct through the majority vote of the candidates pool built by SR. The output formation is linearized opinion tree, the quadruples in which are organized as (*Aspect, Category, [Opinion, Polarity]*). As demonstrated in Table 6, these cases are shown in the formation of Greedy output and SR candidates pool, the majority vote would be with a ✓ notation.

The first example: Greedy gives a very typical wrong prediction, it maps “balcony” to “NULL”, neglecting the adjectives “nice” that express clear polarity, while our method operating over majority vote, easily gives a right answer.

The second example: Greedy predicts “friendly”

as the opinion, which is a common adjective yet not an opinion in the review since it was used to describe the unrelated content, leading to the misjudgment of sentiment polarity. Our method roll-backs the span of the sub-tree “ [friendly, Positive]” to a right opinion and the polarity that has a strong semantic connection with it.

The third example: The root uncertain element of the Greedy sequence is “place”, thus our SR roll-backs the entire sub-tree rooted at “place”, which is also the entire quadruple sequence, and gets the correct output on the basis of new sub-trees with semantic connection inside them.

The fourth example: Greedy misunderstands that the “friendly” is used to reinforce the negative sentiment of annoying while SR salvages it with 5 loops

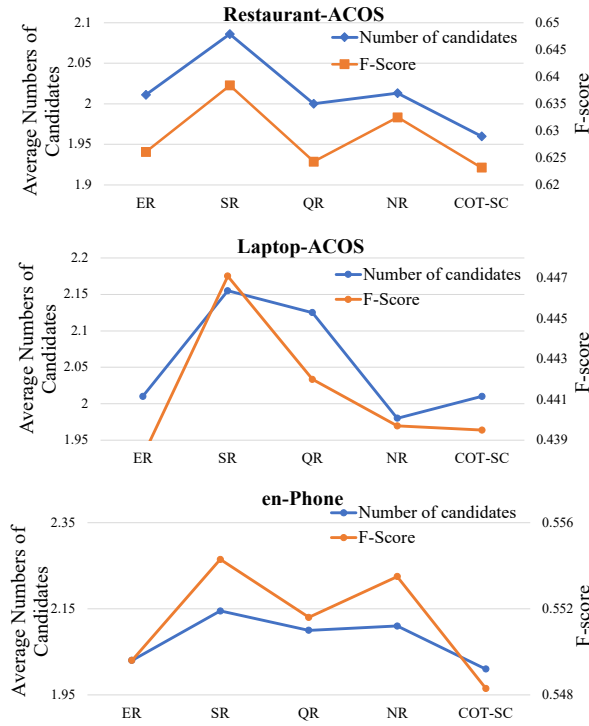


Figure 8: Association between performance and the number of candidates in various strategies.

of rollback.

The fifth example: Based on the entropy threshold, the “mercedes restaurant” is judged uncertain, thus the entire quadruple span would be our rollback span, and the majority vote gives the right answer.

From the cases shown in Table 6, we can find that, with the utilisation of the connection during inference, our method shows significant superiority in improving fine-tuned language models with a tiny cost.

B Impact of Candidates Number

To further investigate the reasons resulting in their different performances, we also compare the inference strategies from the perspective of the association between their performance and the number of candidates in their pools.

Specifically, we compare the strategies from fine to coarse-grained, starting with the rollback of the sentiment element (ER), followed by the rollback of the sub-trees (SR). After that, we perform rollback on the quadruple (QR), the neighbors (NR), and finally the entire sequence (COT-SC). The number of candidates is calculated as the average number of candidates in the candidate pool without duplicates when a rollback occurs.

As shown in Figure 8, the performance generally

follows a similar trend with the average number of candidates after removing duplicates. Among them, the rollback span of the sub-tree, which is our SR, achieves the highest number of candidates and the highest performance. This finding demonstrates that with semantic structure design, SR can offer more options for the framework, which contributes to its highest performance.