

ADS-Edit: A Multimodal Knowledge Editing Dataset for Autonomous Driving Systems

Chenxi Wang^{♡*}, Jizhan Fang^{♡*}, Xiang Chen^{♣*}, Bozhong Tian[♡],
Ziwen Xu[♡], Huajun Chen[♡], Ningyu Zhang^{♡†}

[♡]Zhejiang University

[♣]Nanjing University of Aeronautics and Astronautics
{sunnywxc, zhangningyu}@zju.edu.cn

Abstract

Recent advancements in Large Multimodal Models (LMMs) have shown promise in Autonomous Driving Systems (ADS). However, their direct application to ADS is hindered by challenges such as misunderstanding of traffic knowledge, complex road conditions, and diverse states of vehicle. To address these challenges, we propose the use of Knowledge Editing, which enables targeted modifications to a model’s behavior without the need for full re-training. Meanwhile, we introduce ADS-Edit, a multimodal knowledge editing dataset specifically designed for ADS, which includes various real-world scenarios, multiple data types, and comprehensive evaluation metrics. We conduct comprehensive experiments and derive several interesting conclusions. We hope that our work will contribute to the further advancement of knowledge editing applications in the field of autonomous driving¹.

1 Introduction

The recent Large Multimodal Models (LMMs) (Wang et al., 2024a; Li et al., 2024a; Wu et al., 2024b; Liu et al., 2024c; Team, 2024, 2025) have significantly enhanced capabilities in video understanding and multimodal reasoning. As a key foundation, LMMs have also found initial applications in Autonomous Driving Systems (ADS) (Xing et al., 2024; Sima et al., 2024; Huang et al., 2024b; Yao et al., 2024). However, the direct application of LMMs in ADS yields suboptimal results.

As shown in Figure 1, the reasons for failure can be attributed to the following factors: **(1) Traffic knowledge misunderstanding:** General models misunderstand traffic knowledge leading to suboptimal performance in tasks. **(2) Complex and**

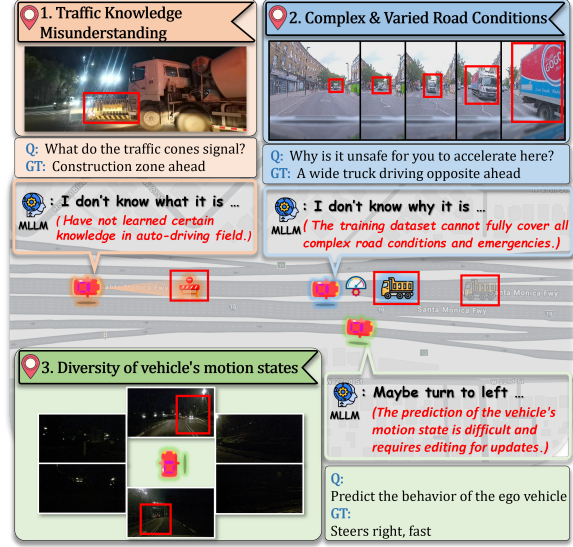


Figure 1: Direct application of LMMs in Autonomous Driving Systems faces several challenges, including the misunderstanding of traffic knowledge, the complex and varied road conditions, and the diverse states of vehicle. Knowledge Editing that enables efficient, continuous, and precise updates to knowledge can effectively address these challenges.

varied road condition: Real-world driving scenarios are highly variable, with training datasets often failing to cover edge cases. **(3) Diversity of vehicle motion states:** Current LMMs struggle to predict unknown and highly dynamic vehicle motion states. These challenges require the model to have the capability of updating knowledge in real-time and continuously.

To address these challenges, we propose the use of Knowledge Editing, which enables targeted modifications to the model’s behavior without the computational burden of full retraining. Unlike traditional fine-tuning, which risks catastrophic forgetting and demands extensive resources, **Knowledge Editing facilitate rapid knowledge updates by selectively altering parameters associated with specific factual or contextual knowledge**

* Equal Contribution.

† Corresponding author.

¹Code and data are available in <https://github.com/zjunlp/EasyEdit>.

(Yao et al., 2023; Gupta et al., 2024; Zhang et al., 2024d; Fang et al., 2024a; Youssef et al., 2025; Jiang et al., 2025; Liu et al., 2024a; Yao et al., 2025). Although Knowledge Editing has been preliminarily explored in the area of LMMs, existing efforts have primarily focused on editing multimodal common knowledge (Cheng et al., 2023; Ma et al., 2024) and multimodal factual knowledge (Zhang et al., 2024a; Huang et al., 2024a). To address this limitation, we pioneer the editing of multimodal domain-specific knowledge in the field of autonomous driving.

Following the principle of leveraging knowledge editing to address the challenges of autonomous driving, we design the **ADS-Edit** benchmark. This benchmark encompasses three real-world scenarios: perception, understanding, and decision making. Additionally, it incorporates three types of data: video, multi-view images, and single image. Furthermore, we establish comprehensive evaluation metrics for knowledge editing in autonomous driving scenarios.

We evaluate four commonly used Knowledge Editing baselines under both single editing and lifelong editing scenarios, such as Prompt, AdaLora (Zhang et al., 2023), GRACE (Hartvigsen et al., 2023), and WISE (Wang et al., 2024b). Thirdly, through our analysis, we obtained a series of interesting findings, including the universality of knowledge editing methods in updating knowledge across various scenarios, their ability to balance editing effectiveness and processing speed, particularly in the context of video data, and the remaining limitations in locality. In general, we conclude our contributions as:

- We are the first to attempt Knowledge Editing on multimodal domain knowledge data, specifically the ADS and effectively addresses the current challenges faced by LMMs when directly applied to ADS.
- We unveil **ADS-Edit**, a Knowledge Editing dataset specifically designed for the ADS. The dataset includes three common types of visual data and encompasses data that evaluates various model capabilities.
- We test four Knowledge Editing baselines under both single editing and lifelong editing settings and analyze several interesting results.

2 Background

LMMs have been applied in autonomous driving, typically after being trained on carefully curated driving datasets. However, when faced with unfamiliar driving scenarios, such as new traffic regulations or predicting driver behavior during traffic congestion, the reliability of LMMs’ decision making can significantly degrade. Furthermore, LMMs struggle to maintain consistent performance when encountering sudden changes in road conditions, such as shifts in weather or unexpected traffic accidents. Finally, for continuously collected driving data, the challenge of updating knowledge in a timely and effective manner remains unresolved for LMMs. Overall, there is **a critical need for a framework that enables LMMs to rapidly and sequentially update knowledges when applied to autonomous driving.**

We propose leveraging knowledge editing techniques to address challenges in autonomous driving scenarios, defined as follows: Let f_θ be a LMM and f_{θ_e} be an edited model. Given the user’s inputs x_e (Includes text t_e and autonomous driving’s multimodal inputs m_e , such as image or videos) and the editing target y_e , the edited model is expected to modify LMM’s outputs within the editing scope to match the editing target, while preserving the original model’s output for the inputs outside the editing scope. The specific definition is as follows:

$$f_{\theta_e}(x) = \begin{cases} y_e & \text{if } x \in I(x_e, y_e) \\ f_\theta(x) & \text{if } x \notin I(x_e, y_e) \end{cases} \quad (1)$$

where $x = (t, m)$ and $I(\cdot)$ means the in-scope of the editing inputs.

3 Benchmark Construction

3.1 Design Principle

To construct a comprehensive benchmark, we propose a tri-axis design principle. This principle organizes evaluation requirements into scenario types: perception, understanding, and decision making, which progressively assess LMMs’ capabilities from basic visual recognition to complex behavioral reasoning. Concurrently, we categorize input modalities into data types: video, multi-view images, and single image. Multiple metrics are designed to evaluate knowledge editing methods, such as reliability, generality and locality.

Scenario Type. *Perception* scenario evaluates LMMs’ basic visual perception capabilities, such

as obstacle detection and vehicle recognition. *Understanding* scenario requires the model to comprehend domain-specific knowledge of autonomous driving, such as traffic rules, beyond basic perception capabilities. *Decision Making* scenario presents a greater challenge, as it requires the model to integrate perception and understanding capabilities to make informed decisions about future driving behaviors. The statistics of scenario types is shown in Figure 2.

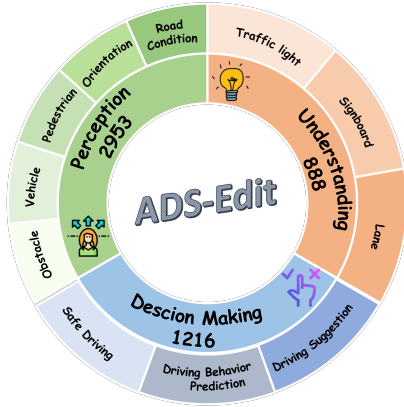


Figure 2: The statistics of scenario types for ADS-Edit.

Data Type. *Video* requires the model to possess the ability to evaluate temporal changes in images. Specifically, the model must be capable of understanding the three dimensions of image width, height, and time. *Multi-view images* are designed to assess the model’s capability when provided with sensor images from multiple viewpoints on the vehicle. *Single image* primarily tests the model’s fundamental perception abilities, such as object recognition and spatial relationship understanding. We will organize these diverse data types to construct a comprehensive benchmark. The statistics of different types of data is shown in Table 1.

	Video	Multi-view	Single	All
Train	1,926	960	1,093	3,979
Test	481	239	358	1,078

Table 1: Statistical information of ADS-Edit data types and dataset splits for training and testing.

Metrics. *Reliability* is to evaluate the success of behavioral modification in the target driving scenarios. *Generality* evaluates the generalization scores when facing similar autonomous driving

scenarios. *Locality* measures whether methods alter unrelated knowledge after updating with domain-specific knowledge of autonomous driving.

3.2 Data Collection

Inspired by knowledge editing applications in unimodal settings, the constructed **ADS-Edit** dataset is composed of visual question answering data. The construction process of the ADS dataset is illustrated in Figure 3.

3.2.1 Reliability Data Construction

We select three well-known datasets of the autonomous driving system: LingoQA (Marcu et al., 2024), DriveLM dataset (Sima et al., 2024), and CODA-LM dataset (Li et al., 2024d) as raw data. Their autonomous driving scenarios consist of video (5 frames), multi-view images (6 views), and single-image, respectively.

Data Preprocess. Notably, the answers in the raw data consist of more tokens (e.g., detailed driving suggestions), distinguishing them from unimodal knowledge editing tasks, such as those represented by triple-based, which have concise token answers. This introduces the following challenges: **1) Suboptimal Editing Performance.** Redundant answers pose difficulties for the implementation of certain editing methods, particularly causal tracing-based approaches such as ROME (Meng et al., 2022). **2) Difficult Evaluation.** Evaluating the accuracy of long text sequences often requires external models to score semantic consistency, which not only suffers from limited accuracy but also incurs high computational costs. Consequently, simplifying the answers in the raw data becomes a necessary step, as shown in Figure 3. We prompt Deepseek-v3 (et al, 2024) to condense the original answers (The prompt template is shown in Appendix C). Meanwhile, we retain the original answers as reserved information for potential use in future research.

Reliability Data Construction. We observe that the data from DriveLM include predictions of vehicle trajectories, which are not feasible for streamlined and generalized data collection. Consequently, we decide to exclude it. For CODA-LM, which contains a significant amount of autonomous driving suggestions and image descriptions, we opt to use Deepseek-v3 to self-generate QA pairs (The prompt template is shown in Appendix C). Subsequently, we shuffle all the data and select a subset

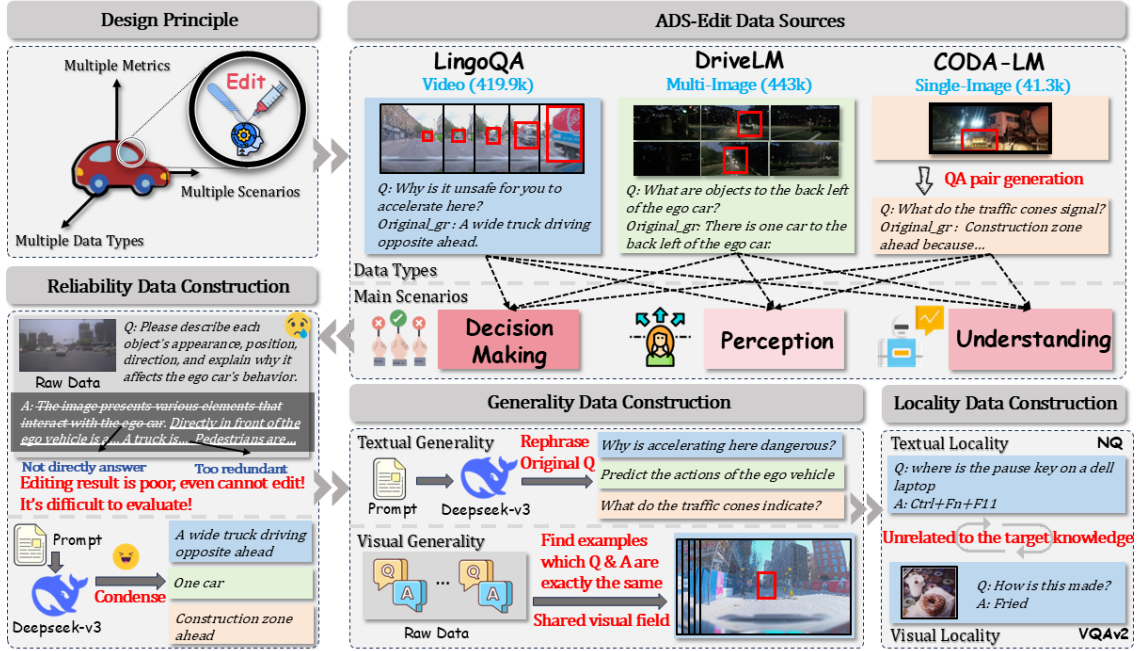


Figure 3: The overview of ADS-Edit construction pipeline.

to serve as reliability data.

3.2.2 Generality Data Construction

Textual Generality Data. Given the same driving scenario, the text query is modified to test whether the model truly understands the underlying context. Following prior works, we prompt Deepseek-v3 to rewrite the original query into a semantically similar query while keeping the answer unchanged (See the prompt in Appendix C).

Multimodal Generality Data. In previous work, text-to-image generation models were employed to create images, or manually curated similar images were used to obtain visually generalized data. However, directly using generative models yields suboptimal performance for the driving domain, which relies exclusively on video data or multiple-view images. The cost of a fully manual curation of similar autonomous driving videos or images is prohibitively high. To address this limitation, we hypothesize that videos corresponding to data with identical questions and answers exhibit similar visual characteristics. Therefore, we only need to match data with identical QA pairs, and then perform non-replacement sampling of two data points per round until exhaustion or only one data point remains. Notably, this process will only occur within the same type of visual data.

3.2.3 Locality Data Construction

To evaluate the capability of editing methods in preserving text locality, we select the Natural Questions dataset (Kwiatkowski et al., 2019), which contains unimodal factual and commonsense knowledge, and randomly sample a subset of it to serve as **Textual Locality data**. For assessing the preservation of multimodal general capabilities, we select a portion of the VQAv2 dataset (Antol et al., 2015), which includes general visual question answers, to construct **Multimodal Locality data**.

3.3 Quality Control

As outlined in the aforementioned process and illustrated in Figure 3, the entire pipeline is fully automated. However, ensuring data quality through manual verification remains a critical step. To this end, we conduct manual verification and calibration for all processes involving AI models, including answer simplification, QA pairs generation from CODA-LM, textual generality data generation, and scenario categorization. Furthermore, for multimodal generality data, we manually sample 20% of the data to assess similarity and filter out those instances with no resemblance. The data quality control tasks are evenly distributed among three annotators with graduate-level education, who independently performed the calibration. In cases of uncertainty, the annotators first independently judge on whether modifications are necessary, and

Method	Reliability \uparrow	T-Generality \uparrow	M-Generality \uparrow	T-Locality \uparrow	M-Locality \uparrow
LLaVA-Onevision					
Prompt	94.25	90.18	95.04	84.86	68.13
AdaLora	78.01	72.76	75.84	85.51	81.12
GRACE	100.00	28.91	28.16	100.00	100.00
WISE	99.10	86.97	95.78	94.18	99.98
Qwen2-VL					
Prompt	90.57	84.98	90.48	89.61	72.44
AdaLora	79.89	75.68	78.76	82.27	69.37
GRACE	100.00	27.01	29.93	100.00	100.00
WISE	94.18	85.20	91.99	94.23	99.85

Table 2: Single edit results on the ADS-Edit. **Reliability** denotes the accuracy of successful editing. **T-Generality**, **M-Generality** represents textual and multimodal generality. **T-Locality**, **M-Locality** refer to the textual and multimodal stability.

then resolve any disagreements through a majority vote. The inter-annotator consistency, measured by Fleiss’s Kappa (κ), demonstrates substantial agreement ($\kappa = 0.912$), based on a randomly selected sample of 200 annotated data points, indicating a high level of consistency within the range of $0.80 \leq \kappa \leq 1.00$.

4 Evaluation

4.1 Experimental settings

We conduct experiments on the most advanced LMMs to date, which are LLaVA-OneVision (Li et al., 2024a) and Qwen2-VL (Wang et al., 2024a). Both models utilize Qwen2-7B as their LLM component. Meanwhile, four classical knowledge editing baselines, including Prompt, AdaLora, GRACE, and WISE are employed to update the knowledge of LMMs. Further details can be found in the Appendix B.

To comprehensively evaluate the performance of various baselines in ADS-Edit, we test in **Single Editing** and **Lifelong Editing** scenarios. Single editing involves updating and evaluating the model immediately after receiving each individual driving data instance. In contrast, lifelong editing represents a more realistic driving scenario, where multiple driving data instances are collected during the vehicle’s moving, requiring continuous integration of this knowledge into the model followed by evaluation. We reasonably assume that the visual processing component (e.g., Vision Transformer) of LMMs provides reliable visual information. Therefore, we focus only on editing the LLM component to modify its understanding of visual input.

4.2 Main Results

4.2.1 Single Editing Results

Memory-based editing methods, such as GRACE and WISE, have proven highly effective in modulating the behavior of LMMs for autonomous driving task prediction. Notably, GRACE achieves a 100% modification rate in both LLaVA-Onevision and Qwen2-VL. Although the ground truth answer is directly provided in the input text, the reliability of the Prompt remains suboptimal. We hypothesize that this is due to the model being affected by interference when processing multimodal driving inputs, which differs from the unimodal scenarios. Under a limited number of training epochs, AdaLora struggles to achieve satisfactory reliability through model parameter updates.

Evaluating the generalization of editing methods to other autonomous driving scenarios and queries is essential. However, GRACE exhibits the worst prediction performance on Generality, with accuracy dropping below 30% in both LMMs, compared to its Reliability. This is attributed to the difficulty of GRACE’s codebook in capturing the representational differences of long-sequence multimodal inputs. The remaining three baselines show similar Generality and Reliability results. Among them, WISE and Prompt achieve approximately 90% on Generality, demonstrating their stable generalization capabilities. AdaLora still struggle to generalize the other driving scenarios. Furthermore, M-Generality of the baselines exceeds T-Generality, suggesting that LMMs tend to focus more on text tokens while making less use of visual tokens. This observation may inspire further research on enhancing the efficient utilization of

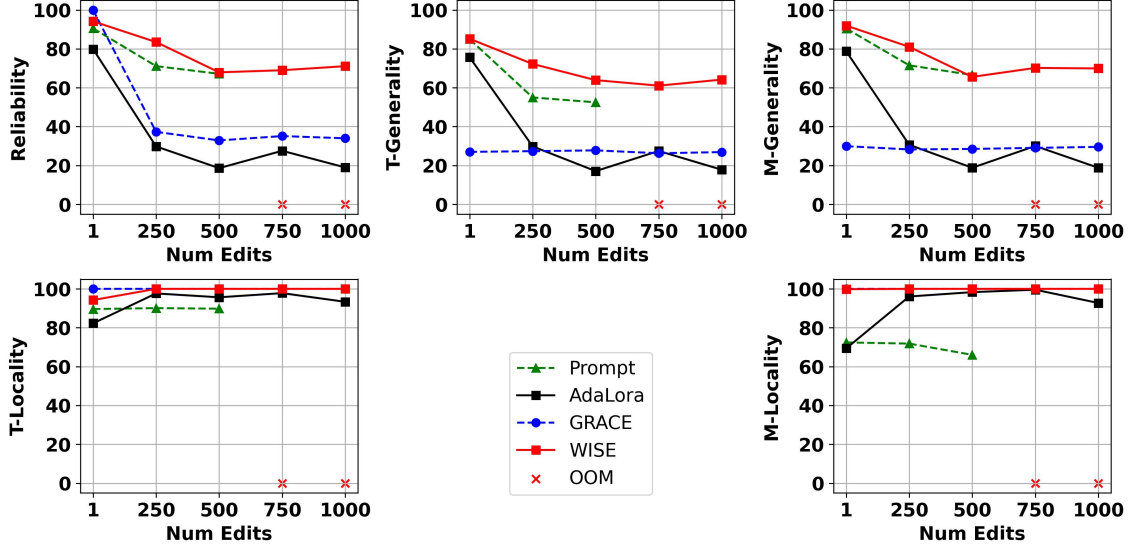


Figure 4: Lifelong Editing results of Qwen2-VL. \times indicates that Prompt triggers an Out-of-Memory (OOM) error at 750 and 1000 editing iterations.

redundant visual representations in LMMs.

While preserving the model’s local capabilities, GRACE demonstrates a significant advantage in Locality. WISE, which employs a dual-memory mechanism, also maintains good Locality, achieving nearly 100%. In contrast, AdaLora and Prompt somewhat disrupt the original model’s Locality. Additionally, we observe that both AdaLora and Prompt exhibit lower M-Locality than T-Locality, which is due to the changes in the LMM’s behavior caused by the multimodal inputs in autonomous driving scenarios.

4.2.2 Lifelong Editing Results

In real-world autonomous driving scenarios, sequential knowledge updates are often required. We sequentially test the effects of lifelong editing 1, 250, 500, 750, and 1000 times across four baseline methods: Prompt, AdaLora, GRACE, and WISE. The lifelong editing results of Qwen2-VL and LLaVA-Onevision are shown in Figure 4 and Figure 8, respectively.

Regarding Reliability and Generality metrics, the performance of the four baselines shows a gradual decline, albeit with some fluctuations. WISE, as a memory-augmented approach, effectively alleviates the knowledge forgetting phenomenon observed in AdaLora, which relies on parameter updates. Although the Prompt method demonstrates also effective lifelong editing performance at 250 and 500 times updates, it leads to an Out-of-Memory (OOM) error after 750 times updates. Due to the

long tokens of multimodal inputs, the codebook struggles to differentiate between distinct representation information, resulting in suboptimal lifelong editing performance for GRACE.

Methods that preserve the original parameters, such as WISE and GRACE, achieved 100% performance on the Locality metric. Due to interference from additional inputs, the Prompt method performs poorly on the Locality metric, with even worse results observed in multimodal locality. AdaLoRA employs low-rank matrix-based parameter updates, which minimally affect the original parameters, resulting in a strong performance on the locality metric.

4.3 Analysis

LLaVA-OneVision vs. Qwen2-VL: Which is Easier to Edit? Although both models use Qwen2 as Large Language Model, LLaVA-Onevision achieve better results across most metrics in terms of reliability and generality. Specifically, on Prompt and WISE methods, Qwen2-VL demonstrates stronger performance in retaining the original predicted answers. Qwen2-VL slightly outperforms LLaVA-OneVision in terms of reliability and generality by using AdaLora. However, locality is significantly compromised compared to LLaVA-OneVision.

Which Scenario Type Knowledges are Easier to Update? We choose the average generality metric, which is the mean of the text and multimodal

generality of Qwen2-VL and LLaVA-OneVision, to assess the performance of the edits. (Since both WISE and GRACE achieve Reliability close to 100%, analyzing the Reliability metric becomes less meaningful. This observation applies to all subsequent analyses in this study.) According to the averaged metrics across three scenarios illustrated in Figure 5, for AdaLora, decision scenario is more challenging to learn compared to perception and understanding scenario, given a fixed number of training epochs. This can be attributed to the higher complexity of decision making data. In contrast, Prompt, GRACE, and WISE show relatively consistent performance metrics across all scenarios, with minimal variation in their generality. **Overall, these knowledge editing methods are broadly applicable to various driving scenarios, effectively updating knowledge in LMMs.**

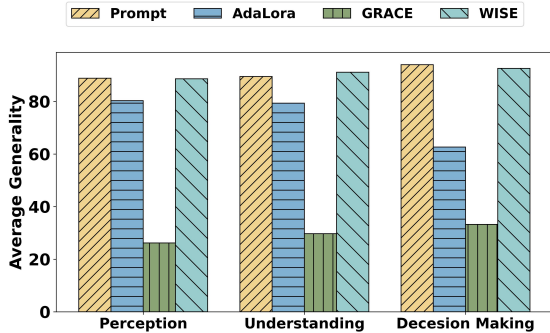


Figure 5: The average generality metric of single editing across different scenarios.

Which Data Type is Easier to Edit? We compute the baseline performance across different types of modality, including video, multi-view images and single image, as shown in Figure 6. WISE demonstrates the highest performance on video data, likely due to its memory mechanism to store and update knowledge in temporal changes in driving scenario. Both Prompt and AdaLoRA exhibit a gradual increase in generality as the length of the visual sequence decreases.

Does Reducing Video Frames Impact the Effectiveness of Knowledge Editing? To assess the impact of video frames on editing performance, we sequentially test videos with 1 to 5 frames as input in single editing setting. The average generality on video data with 1-4 frames is denoted as the effect of the knowledge editing method under low video frame rate conditions, whereas video data with 5 frames is referred to as the maximum frame rate condition. Fewer video frames corre-

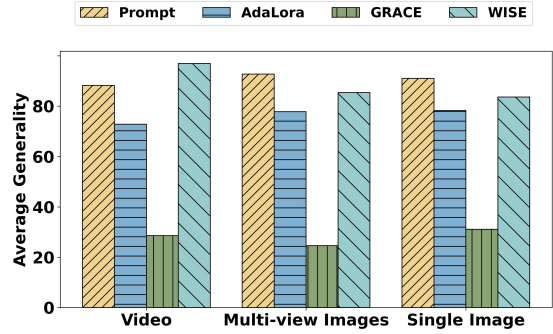


Figure 6: The average generality metric of single editing across different data types.

spond to fewer visual tokens, enabling LMMs to process users’ requests more quickly. Although this reduction in frames results in some loss of information, WISE demonstrates even better performance. Furthermore, the results of both Prompt and AdaLora are not significantly affected under conditions of both low and maximum frame rate. This suggests that **knowledge editing methods in the autonomous driving domain, particularly for video data, can effectively balance processing speed and performance.**

	Low Video Frame rate	Max Video Frame rate
Prompt	90.14	92.61
AdaLora	70.04	74.30
GRACE	34.33	28.54
WISE	99.49	91.37

Table 3: The average generality results from different video frames.

Cases Analysis of Knowledge Editing. From Figure 7, we observe that the edited LMM tends to maintain better locality in the unimodal domain compared to the multimodal one. Interestingly, despite only editing the language model component of the LMM, the model’s visual understanding is impaired. This underscores the distinct differences between the LMM’s capabilities in visual understanding and text-based knowledge storage.

5 Related Work

5.1 Knowledge Editing

Recent advances in knowledge editing have emerged as a pivotal research direction in knowledge updating (Chen, 2024; Liu et al., 2024b; Wang et al., 2024c; Wu et al., 2024a; Xu et al., 2024; Li and Chu, 2024; Nie et al., 2024; Wei et al., 2025; Zhang et al., 2025; Ali et al., 2025; Hwang et al., 2025; Zhao et al., 2025; Markowitz et al.,

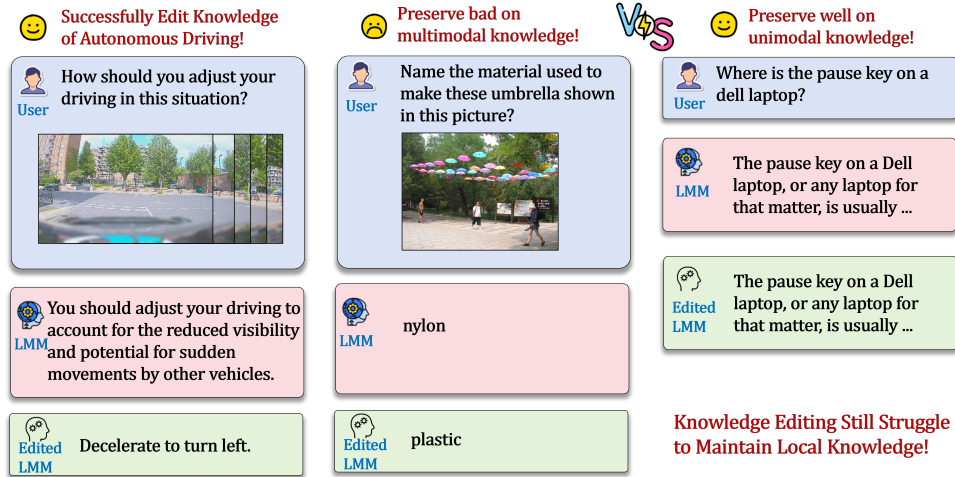


Figure 7: Cases analysis of editing LLaVA-OneVision with WISE.

2025; Yang et al., 2025; Dong et al., 2025; Pan et al., 2025; Wang et al., 2025; Liu et al.; Zhang et al., 2024b; Xu et al., 2025; Bi et al., 2025; He et al., 2025; Li et al., 2025; Shen and Huang, 2025). Following Yao et al. (2023), we systematize current methodologies into two principal paradigms: 1) Parameter Modification Approaches: These methods directly alter the model’s internal representations through targeted interventions (Meng et al., 2022, 2023; Fang et al., 2024b; Cai and Cao, 2024; Mitchell et al., 2022; Zhang et al., 2024c; Tan et al., 2024; Hu et al., 2022; Zhang et al., 2023; Dettmers et al., 2023; Feng et al., 2025). 2) Parameter Preservation Approaches: These strategies maintain original parameters while updating new knowledge (Zheng et al., 2023; Jiang et al., 2024; Wang et al., 2024b; Hartvigsen et al., 2023; Li et al., 2025).

The emergence of Large Multimodal Models (LMMs) has also advanced the field of multimodal knowledge editing. These works mainly focus on the editing of multimodal common knowledge (Cheng et al., 2023; Ma et al., 2024; Pan et al., 2024; Wang et al., 2024d) or multimodal factual knowledge (Li et al., 2024b; Zhang et al., 2024a; Huang et al., 2024a). However, the application of multimodal knowledge editing in the context of domain-specific knowledge, like autonomous driving remains underexplored.

5.2 Large Multi-modal Models for Autonomous Driving

The recent rise of large multimodal models (LMMs) in autonomous driving (Qian et al., 2024; Cui et al., 2024) has spurred the development of ad-

vanced frameworks that integrate perception, prediction, and planning. CODA-LM (Li et al., 2024c) and DriveLM (Sima et al., 2024), have attempted to bridge this gap by integrating hierarchical reasoning and graph-structured visual QA frameworks. CODA-LM evaluates vision-language models on corner cases across perception, prediction, and planning stages, while DriveLM explicitly models logical dependencies between driving phases. Meanwhile, LingoQA (Marcu et al., 2024) expands free-form QA capabilities by integrating action justification and scene understanding. However, these models still face challenges in real-world deployment due to static knowledge representations and modality imbalance.

To address these issues, knowledge editing (Yao et al., 2023) has emerged as a promising direction, enabling targeted modifications to model behavior without full retraining. Yet, existing datasets and methods fail to account for the unique demands of autonomous driving, such as multimodal coherence and lifelong adaptation to evolving scenarios. This gap motivates our work, ADS-Edit, the first multimodal knowledge editing dataset designed explicitly for autonomous driving systems.

6 Conclusion

In this paper, we present to leverage knowledge editing techniques to address the challenges faced by Large Multimodal Models in autonomous driving scenarios. ADS-Edit, a new benchmark specifically is designed for evaluating knowledge editing methods in this domain. Through extensive experimental analysis, we systematically compare

and evaluate the effectiveness of various editing methods, while providing in-depth insights into the underlying reasons for the observed performance variations and failure cases.

Limitations

While this study advances Knowledge Editing for Autonomous Driving System applications, several limitations warrant consideration:

(1) Only VQA Data. In this work, we only test knowledge editing baselines on visual question answering data and don't explore other forms of data related to autonomous driving tasks, such as trajectory prediction. How to successfully update predicted location coordinates as knowledge into LMMs and evaluate their generalization will be explored in future work.

(2) Limited Experimental Baselines. Due to the high cost of LMMs, we did not conduct experiments on larger scale LMMs and discard knowledge editing methods that require more resources, such as MEND and SERAC.

References

- Muhammad Asif Ali, Nawal Daftardar, Mutayyaba Waheed, Jianbin Qin, and Di Wang. 2025. [MQA-KEAL: multi-hop question answering under knowledge editing for arabic language](#). In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 5629–5644. Association for Computational Linguistics.
- Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *International Conference on Computer Vision (ICCV)*.
- Baolong Bi, Shenghua Liu, Yiwei Wang, Yilong Xu, Junfeng Fang, Lingrui Mei, and Xueqi Cheng. 2025. Parameters vs. context: Fine-grained control of knowledge reliance in language models. *arXiv preprint arXiv:2503.15888*.
- Yuchen Cai and Ding Cao. 2024. [O-edit: Orthogonal subspace editing for language model sequential editing](#). *CoRR*, abs/2410.11469.
- Huajun Chen. 2024. Large knowledge model: Perspectives and challenges. *Data Intelligence*, (3).
- Siyuan Cheng, Bozhong Tian, Qingbin Liu, Xi Chen, Yongheng Wang, Huajun Chen, and Ningyu Zhang. 2023. [Can we edit multimodal large language models?](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 13877–13888. Association for Computational Linguistics.
- Can Cui, Yunsheng Ma, Xu Cao, Wenqian Ye, Yang Zhou, Kaizhao Liang, Jintai Chen, Juanwu Lu, Zichong Yang, Kuei-Da Liao, Tianren Gao, Erlong Li, Kun Tang, Zhipeng Cao, Tong Zhou, Ao Liu, Xinrui Yan, Shuqi Mei, Jianguo Cao, Ziran Wang, and Chao Zheng. 2024. [A survey on multimodal large language models for autonomous driving](#). In *IEEE/CVF Winter Conference on Applications of Computer Vision Workshops, WACVW 2024 - Workshops, Waikoloa, HI, USA, January 1-6, 2024*, pages 958–979. IEEE.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. [Qlora: Efficient finetuning of quantized llms](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.
- Zilu Dong, Xiangqing Shen, and Rui Xia. 2025. Memit-merge: Addressing memit's key-value conflicts in same-subject batch editing for llms. *arXiv preprint arXiv:2502.07322*.
- DeepSeek-AI et al. 2024. [Deepseek-v3 technical report](#). Preprint, arXiv:2412.19437.
- Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Xiang Wang, Xiangnan He, and Tat-Seng Chua. 2024a. [Alphaedit: Null-space constrained knowledge editing for language models](#). *CoRR*, abs/2410.02355.
- Junfeng Fang, Houcheng Jiang, Kun Wang, Yunshan Ma, Xiang Wang, Xiangnan He, and Tat-Seng Chua. 2024b. [Alphaedit: Null-space constrained knowledge editing for language models](#). *CoRR*, abs/2410.02355.
- Yujie Feng, Liming Zhan, Zexin Lu, Yongxin Xu, Xu Chu, Yasha Wang, Jiannong Cao, Philip S Yu, and Xiao-Ming Wu. 2025. [Geoedit: Geometric knowledge editing for large language models](#). *arXiv preprint arXiv:2502.19953*.
- Akshat Gupta, Anurag Rao, and Gopala Anumanchipalli. 2024. Model editing at scale leads to gradual and catastrophic forgetting. *arXiv preprint arXiv:2401.07453*.
- Tom Hartvigsen, Swami Sankaranarayanan, Hamid Palangi, Yoon Kim, and Marzyeh Ghassemi. 2023. [Aging with GRACE: lifelong model editing with discrete key-value adaptors](#). In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*.

- Guoxiu He, Xin Song, and Aixin Sun. 2025. Knowledge updating? no more model editing! just selective contextual reasoning. *arXiv preprint arXiv:2503.05212*.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2022. *Lora: Low-rank adaptation of large language models*. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Han Huang, Haitian Zhong, Tao Yu, Qiang Liu, Shu Wu, Liang Wang, and Tieniu Tan. 2024a. *VLKEB: A large vision-language model knowledge editing benchmark*. In *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*.
- Zhijian Huang, Chengjian Feng, Feng Yan, Baihui Xiao, Zequn Jie, Yujie Zhong, Xiaodan Liang, and Lin Ma. 2024b. *Drivemm: All-in-one large multimodal model for autonomous driving*. *CoRR*, abs/2412.07689.
- Seojin Hwang, Yumin Kim, Byeongjeong Kim, and Hwanhee Lee. 2025. Personality editing for language models through relevant knowledge editing. *arXiv preprint arXiv:2502.11789*.
- Houcheng Jiang, Junfeng Fang, Ningyu Zhang, Guojun Ma, Mingyang Wan, Xiang Wang, Xiangnan He, and Tat-seng Chua. 2025. Anyedit: Edit any knowledge encoded in language models. *arXiv preprint arXiv:2502.05628*.
- Yuxin Jiang, Yufei Wang, Chuhan Wu, Wanjuan Zhong, Xingshan Zeng, Jiahui Gao, Liangyou Li, Xin Jiang, Lifeng Shang, Ruiming Tang, Qun Liu, and Wei Wang. 2024. *Learning to edit: Aligning llms with knowledge editing*. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2024, Bangkok, Thailand, August 11-16, 2024*, pages 4689–4705. Association for Computational Linguistics.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur P. Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. *Natural questions: a benchmark for question answering research*. *Trans. Assoc. Comput. Linguistics*, 7:452–466.
- Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. 2024a. *Llava-onevision: Easy visual task transfer*. *CoRR*, abs/2408.03326.
- Jiaqi Li, Miaozeng Du, Chuanyi Zhang, Yongrui Chen, Nan Hu, Guilin Qi, Haiyun Jiang, Siyuan Cheng, and Bozhong Tian. 2024b. *MIKE: A new benchmark for fine-grained multimodal entity knowledge editing*. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 5018–5029. Association for Computational Linguistics.
- Qi Li and Xiaowen Chu. 2024. *Can we continually edit language models? on the knowledge attenuation in sequential model editing*. In *Findings of the Association for Computational Linguistics, ACL 2024, Bangkok, Thailand and virtual meeting, August 11-16, 2024*, pages 5438–5455. Association for Computational Linguistics.
- Shuaike Li, Kai Zhang, Qi Liu, and Enhong Chen. 2025. Mindbridge: Scalable and cross-model knowledge editing via memory-augmented modality. *arXiv preprint arXiv:2503.02701*.
- Yanze Li, Wenhua Zhang, Kai Chen, Yanxin Liu, Pengxiang Li, Ruiyuan Gao, Lanqing Hong, Meng Tian, Xinhai Zhao, Zhenguo Li, Dit-Yan Yeung, Huchuan Lu, and Xu Jia. 2024c. *Automated evaluation of large vision-language models on self-driving corner cases*. *CoRR*, abs/2404.10595.
- Yanze Li, Wenhua Zhang, Kai Chen, Yanxin Liu, Pengxiang Li, Ruiyuan Gao, Lanqing Hong, Meng Tian, Xinhai Zhao, Zhenguo Li, et al. 2024d. *Automated evaluation of large vision-language models on self-driving corner cases*. *arXiv preprint arXiv:2404.10595*.
- Tianci Liu, Ruirui Li, Yunzhe Qi, Hui Liu, Xianfeng Tang, Tianqi Zheng, Qingyu Yin, Monica Xiao Cheng, Jun Huan, Haoyu Wang, et al. *Unlocking efficient, scalable, and continual knowledge editing with basis-level representation fine-tuning*. In *The Thirteenth International Conference on Learning Representations*.
- Zeyu Leo Liu, Shrey Pandit, Xi Ye, Eunsol Choi, and Greg Durrett. 2024a. *Codeupdatearena: Benchmarking knowledge editing on api updates*. *arXiv preprint arXiv:2407.06249*.
- Zeyu Leo Liu, Shrey Pandit, Xi Ye, Eunsol Choi, and Greg Durrett. 2024b. *Codeupdatearena: Benchmarking knowledge editing on API updates*. *CoRR*, abs/2407.06249.
- Zhijian Liu, Ligeng Zhu, Baifeng Shi, Zhuoyang Zhang, Yuming Lou, Shang Yang, Haocheng Xi, Shiyi Cao, Yuxian Gu, Dacheng Li, Xiuyu Li, Yunhao Fang, Yukang Chen, Cheng-Yu Hsieh, De-An Huang, An-Chieh Cheng, Vishwesh Nath, Jinyi Hu, Sifei Liu, Ranjay Krishna, Daguang Xu, Xiaolong Wang, Pavlo Molchanov, Jan Kautz, Hongxu Yin, Song Han, and Yao Lu. 2024c. *Nvila: Efficient frontier visual language models*. *Preprint*, arXiv:2412.04468.
- Yaohui Ma, Xiaopeng Hong, Shizhou Zhang, Huiyun Li, Zhilin Zhu, Wei Luo, and Zhiheng Ma. 2024. *Comprehendedit: A comprehensive dataset and evaluation framework for multimodal knowledge editing*. *CoRR*, abs/2412.12821.

- Ana-Maria Marcu, Long Chen, Jan Hünemann, Alice Karnsund, Benoît Hanotte, Prajwal Chidananda, Saurabh Nair, Vijay Badrinarayanan, Alex Kendall, Jamie Shotton, Elahe Arani, and Oleg Sinavski. 2024. [Lingoqa: Visual question answering for autonomous driving](#). In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LXXVII*, volume 15135 of *Lecture Notes in Computer Science*, pages 252–269. Springer.
- Elan Markowitz, Anil Ramakrishna, Ninareh Mehrabi, Charith Peris, Rahul Gupta, Kai-Wei Chang, and Aram Galstyan. 2025. K-edit: Language model editing with contextual knowledge awareness. *arXiv preprint arXiv:2502.10626*.
- Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. 2022. [Locating and editing factual associations in GPT](#). 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*.
- Kevin Meng, Arnab Sen Sharma, Alex J. Andonian, Yonatan Belinkov, and David Bau. 2023. [Mass-editing memory in a transformer](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D. Manning. 2022. [Fast model editing at scale](#). In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net.
- Ercong Nie, Bo Shao, Zifeng Ding, Mingyang Wang, Helmut Schmid, and Hinrich Schütze. 2024. [BMIKE-53: investigating cross-lingual knowledge editing with in-context learning](#). *CoRR*, abs/2406.17764.
- Haowen Pan, Xiaozhi Wang, Yixin Cao, Zenglin Shi, Xun Yang, Juanzi Li, and Meng Wang. 2025. Precise localization of memories: A fine-grained neuron-level knowledge editing technique for llms. *arXiv preprint arXiv:2503.01090*.
- Kaihang Pan, Zhaoyu Fan, Juncheng Li, Qifan Yu, Hao Fei, Siliang Tang, Richang Hong, Hanwang Zhang, and Qianru Sun. 2024. Towards unified multimodal editing with enhanced knowledge collaboration. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. 2024. [Nuscenes-qa: A multimodal visual question answering benchmark for autonomous driving scenario](#). In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 4542–4550. AAAI Press.
- Ying Shen and Lifu Huang. 2025. Llm braces: Straightening out llm predictions with relevant sub-updates. *arXiv preprint arXiv:2503.16334*.
- Chonghao Sima, Katrin Renz, Kashyap Chitta, Li Chen, Hanxue Zhang, Chengen Xie, Jens Beißwenger, Ping Luo, Andreas Geiger, and Hongyang Li. 2024. [Drivelm: Driving with graph visual question answering](#). In *Computer Vision - ECCV 2024 - 18th European Conference, Milan, Italy, September 29-October 4, 2024, Proceedings, Part LII*, volume 15110 of *Lecture Notes in Computer Science*, pages 256–274. Springer.
- Chenmian Tan, Ge Zhang, and Jie Fu. 2024. [Massive editing for large language models via meta learning](#). In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net.
- Qwen Team. 2024. [Qvq: To see the world with wisdom](#).
- Qwen Team. 2025. [Qwen2.5-vl](#).
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024a. [Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution](#). *CoRR*, abs/2409.12191.
- Peng Wang, Zexi Li, Ningyu Zhang, Ziwen Xu, Yunzhi Yao, Yong Jiang, Pengjun Xie, Fei Huang, and Hua-jun Chen. 2024b. [WISE: Rethinking the knowledge memory for lifelong model editing of large language models](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Pinzheng Wang, Zecheng Tang, Keyan Zhou, Juntao Li, Qiaoming Zhu, and Min Zhang. 2025. Revealing and mitigating over-attention in knowledge editing. *arXiv preprint arXiv:2502.14838*.
- Yiwei Wang, Muhao Chen, Nanyun Peng, and Kai-Wei Chang. 2024c. [Deepedit: Knowledge editing as decoding with constraints](#). *CoRR*, abs/2401.10471.
- Zecheng Wang, Xinye Li, Zhanyue Qin, Chunshan Li, Zhiying Tu, Dianhui Chu, and Dianbo Sui. 2024d. Can we debias multimodal large language models via model editing? In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 3219–3228.
- Zihao Wei, Jingcheng Deng, Liang Pang, Hanxing Ding, Huawei Shen, and Xueqi Cheng. 2025. [Mlake: Multilingual knowledge editing benchmark for large language models](#). In *Proceedings of the 31st International Conference on Computational Linguistics*,

- COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 4457–4473. Association for Computational Linguistics.
- Xiaobao Wu, Liangming Pan, William Yang Wang, and Anh Tuan Luu. 2024a. [AKEW: assessing knowledge editing in the wild](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, EMNLP 2024, Miami, FL, USA, November 12-16, 2024*, pages 15118–15133. Association for Computational Linguistics.
- Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. 2024b. [Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding](#). *Preprint*, arXiv:2412.10302.
- Shuo Xing, Chengyuan Qian, Yuping Wang, Hongyuan Hua, Kexin Tian, Yang Zhou, and Zhengzhong Tu. 2024. [Openemma: Open-source multimodal model for end-to-end autonomous driving](#). *CoRR*, abs/2412.15208.
- Derong Xu, Ziheng Zhang, Zhihong Zhu, Zhenxi Lin, Qidong Liu, Xian Wu, Tong Xu, Wanyu Wang, Yuyang Ye, Xiangyu Zhao, Enhong Chen, and Yefeng Zheng. 2024. [Editing factual knowledge and explanatory ability of medical large language models](#). In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management, CIKM 2024, Boise, ID, USA, October 21-25, 2024*, pages 2660–2670. ACM.
- Hao-Xiang Xu, Jun-Yu Ma, Zhen-Hua Ling, Ningyu Zhang, and Jia-Chen Gu. 2025. [Constraining sequential model editing with editing anchor compression](#). *arXiv preprint arXiv:2503.00035*.
- Wanli Yang, Fei Sun, Jiajun Tan, Xinyu Ma, Qi Cao, Dawei Yin, Huawei Shen, and Xueqi Cheng. 2025. [The mirage of model editing: Revisiting evaluation in the wild](#). *arXiv preprint arXiv:2502.11177*.
- Ruoyu Yao, Yubin Wang, Haichao Liu, Rui Yang, Zengqi Peng, Lei Zhu, and Jun Ma. 2024. [Calmm-drive: Confidence-aware autonomous driving with large multimodal model](#). *CoRR*, abs/2412.04209.
- Yunzhi Yao, Jizhan Fang, Jia-Chen Gu, Ningyu Zhang, Shumin Deng, Huajun Chen, and Nanyun Peng. 2025. [Cake: Circuit-aware editing enables generalizable knowledge learners](#). *arXiv preprint arXiv:2503.16356*.
- Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. 2023. [Editing large language models: Problems, methods, and opportunities](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 10222–10240. Association for Computational Linguistics.
- Paul Youssef, Zhixue Zhao, Daniel Braun, Jörg Schlöterer, and Christin Seifert. 2025. [Position: Editing large language models poses serious safety risks](#). *arXiv preprint arXiv:2502.02958*.
- Junzhe Zhang, Huixuan Zhang, Xunjian Yin, Baizhou Huang, Xu Zhang, Xinyu Hu, and Xiaojun Wan. 2024a. [MC-MKE: A fine-grained multimodal knowledge editing benchmark emphasizing modality consistency](#). *CoRR*, abs/2406.13219.
- Mengqi Zhang, Xiaotian Ye, Qiang Liu, Pengjie Ren, Shu Wu, and Zhumin Chen. 2024b. [Uncovering overfitting in large language model editing](#). *arXiv preprint arXiv:2410.07819*.
- Ningyu Zhang, Bozhong Tian, Siyuan Cheng, Xiaozhuan Liang, Yi Hu, Kouying Xue, Yanjie Gou, Xi Chen, and Huajun Chen. 2024c. [Instructedit: Instruction-based knowledge editing for large language models](#). In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI 2024, Jeju, South Korea, August 3-9, 2024*, pages 6633–6641. ijcai.org.
- Ningyu Zhang, Yunzhi Yao, Bozhong Tian, Peng Wang, Shumin Deng, Mengru Wang, Zekun Xi, Shengyu Mao, Jintian Zhang, Yuansheng Ni, Siyuan Cheng, Ziwen Xu, Xin Xu, Jia-Chen Gu, Yong Jiang, Pengjun Xie, Fei Huang, Lei Liang, Zhiqiang Zhang, Xiaowei Zhu, Jun Zhou, and Huajun Chen. 2024d. [A comprehensive study of knowledge editing for large language models](#). *CoRR*, abs/2401.01286.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. [Adaptive budget allocation for parameter-efficient fine-tuning](#). In *The Eleventh International Conference on Learning Representations*.
- Xue Zhang, Yunlong Liang, Fandong Meng, Songming Zhang, Yufeng Chen, Jinan Xu, and Jie Zhou. 2025. [Multilingual knowledge editing with language-agnostic factual neurons](#). In *Proceedings of the 31st International Conference on Computational Linguistics, COLING 2025, Abu Dhabi, UAE, January 19-24, 2025*, pages 5775–5788. Association for Computational Linguistics.
- Zongkai Zhao, Guozeng Xu, Xiuhua Li, Kaiwen Wei, and Jiang Zhong. 2025. [Fleke: Federated locate-then-edit knowledge editing](#). *arXiv preprint arXiv:2502.15677*.
- Ce Zheng, Lei Li, Qingxiu Dong, Yuxuan Fan, Zhiyong Wu, Jingjing Xu, and Baobao Chang. 2023. [Can we edit factual knowledge by in-context learning?](#) In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 4862–4876. Association for Computational Linguistics.

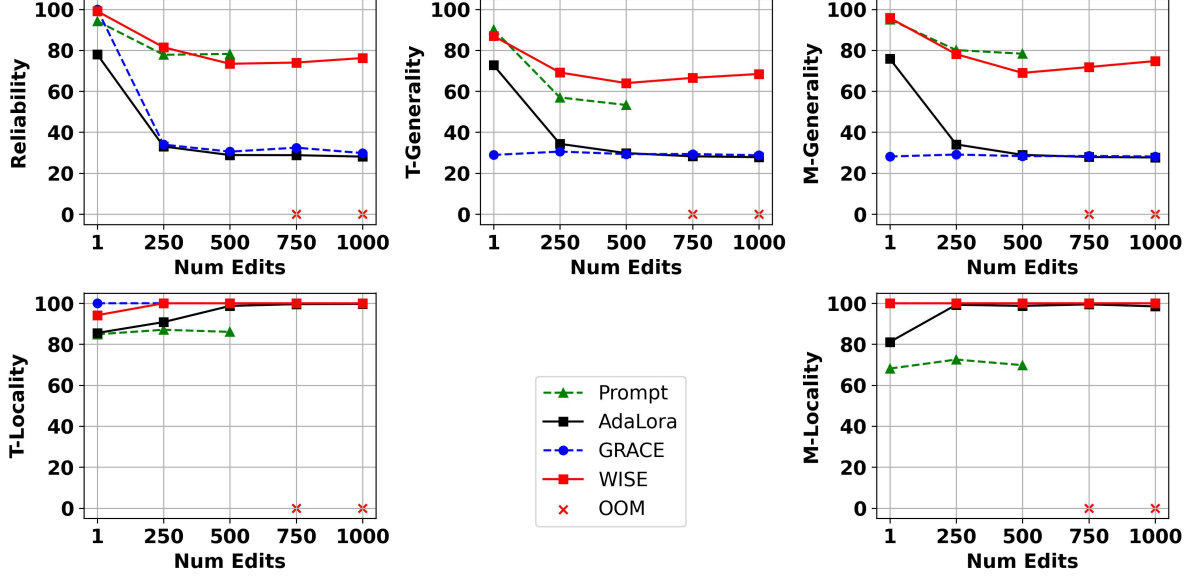


Figure 8: Lifelong Editing results of LLaVA-OneVision. \times indicates that Prompt triggers an Out-of-Memory (OOM) error at 750 and 1000 editing iterations.

A Metrics

Reliability. Reliability is to evaluate the success of behavioral modification in the target driving scenarios. Intuitively, what we need is an updated θ_e with $f(t_e, m_e; \theta_e) = y_e$. To measure the reliability, we use the editing accuracy, as follows:

$$\mathcal{M}_{rel} = \mathbb{E}_{(t_e, m_e, y_e) \sim \mathcal{D}_e} \{f(t_e, m_e; \theta_e) = y_e\} \quad (2)$$

where θ_e refers to the parameters after editing.

Generality. Given similar autonomous driving scenarios, generality assesses whether the edited model f_{θ_e} updates relevant knowledge within (t_e, m_e) to retain the capacity for generalization, so that predict congruent outputs for equivalent inputs (e.g., rephrased textual queries, or similar autonomous driving’s videos). We evaluate text generality and multimodal generality as follows:

$$\mathcal{M}_{gen}^{text} = \mathbb{E}_{t_r \sim \mathcal{N}(t_e)} \{f(t_r, m_e; \theta_e) = y_e\} \quad (3)$$

$$\mathcal{M}_{gen}^{mm} = \mathbb{E}_{m_r \sim \mathcal{N}(m_e)} \{f(t_e, m_r; \theta_e) = y_e\} \quad (4)$$

where m_r presents the rephrased autonomous driving’s image or videos, t_r refers to the rephrased textual queries, and $\mathcal{N}(x)$ denotes to in-scope objects of x .

Locality. After updating the model with domain-specific knowledge in autonomous driving, locality is used to evaluate whether knowledge editing preserves the model’s behavior on unrelated

knowledge (e.g., unimodal factual knowledge and multimodal commonsense knowledge). Following (Cheng et al., 2023), we employ both text locality and multimodal locality to assess the stability of the editing process:

$$\mathcal{M}_{loc}^{text} = \mathbb{E}_{(t, Y) \sim \mathcal{D}_{loc}} \{f(t; \theta_e) = f(t; \theta)\} \quad (5)$$

$$\mathcal{M}_{loc}^{mm} = \mathbb{E}_{(t, m, y) \sim \mathcal{D}_{loc-v}} \{f(t, m; \theta_e) = f(m, t; \theta)\} \quad (6)$$

where \mathcal{D}_{loc} and \mathcal{D}_{loc-v} are distinct from \mathcal{D}_e . The tuples (t, y) and (t, m, y) represent samples drawn from \mathcal{D}_{loc} and \mathcal{D}_{loc-v} , respectively.

B Baselines

Prompt. Directly alter the model’s behavior temporarily through prompts.

AdaLora. LoRA introduces low-rank matrices into Transformer layers, fine-tuning only a small subset of parameters while keeping the majority of the model frozen. Building on LoRA, AdaLora (Zhang et al., 2023) improves LoRA by adaptively allocating the parameter budget based on the importance of weight matrices, using SVD to prune unimportant updates. In our work, we employ AdaLora as the baseline for efficient knowledge editing.

GRACE. GRACE (Hartvigsen et al., 2023) introduces a discrete key-value codebook approach, where edits are cached as latent space mappings without altering the original model weights. This

method enables thousands of sequential edits by storing corrections in a retrieval-based codebook, ensuring minimal interference with unrelated inputs and strong locality . However, GRACE’s reliance on non-parametric representations limits its ability to generalize beyond memorized edits, as it struggles to integrate new knowledge into the model’s reasoning process.

WISE. WISE (Wang et al., 2024b) uses a dual-memory architecture comprising a main memory for pretrained knowledge and a side memory for editing. A routing mechanism decides which memory to use during inference. Additionally, knowledge is partitioned into multiple subspaces through sharding to facilitate conflict-free edits, with subsequent merging of the shards into a unified side memory.

C Benchmark Construction Details

The prompt templates are as shown in Table 4, Table 5, Table 6 and Table 7.

D Case of ADS-Edit benchmark

The case of ADS-Edit benchmark is as shown in Figure 9, Figure 10 and Figure 11

SYSTEM:

You are a helpful assistant.

USER:

Given a question and an answer in the VQA task, you are required to condense the answer into 1 to 5 words. You **MUST FOLLOW THE FOLLOWING RULES:**

1. The condensed answer must be fewer than 5 tokens.
2. The condensed answer must retain the core meaning of the original answer.
3. Return only the condensed answer. **DON'T RESPOND ANYTHING ELSE!**

Here are some examples:

<example 1>

question: What is the current action and its justification? Answer in the form "action, justification".

original answer: The car starts and moves right, because the vehicle in front is pulling away at the green traffic light and then it needs to go around the construction sign on the left of the road.

condensed answer: Moves right.

<example 2>

question: Is there a traffic light? If yes, what color is displayed?

original answer: Yes, a temporary traffic light. It is showing green.

condensed answer: Yes, green.

<example 3>

question: How many pedestrians are present in this image?

original answer: There are no relevant pedestrians, but there is one on the left side-road and one on the right side of the road.

condensed answer: No pedestrians.

question: {question}

original answer: {answer}

condensed answer:

Table 4: Prompt template of answers simplification.

SYSTEM:

You are a helpful assistant.

USER:

Given a question in the Automatic Driving QA task, you are required to rewrite the question in a different way. You MUST FOLLOW THE FOLLOWING RULES:

1. The rephrased question should not have a high degree of overlap with the original wording.
2. The rephrased question should be consistent with the core meaning of the original question.
3. Return only the rephrased question. DON'T RESPOND ANYTHING ELSE!

Here are some examples:

<example 1>

original question: What is the current action and its justification? Answer in the form action, justification:

rephrased question: What is the action being performed at the moment, and what is the reasoning behind it? Please respond with "action, reasoning."

<example 2>

original question: How many pedestrians are in the video?

rephrased question: What is the number of pedestrians depicted in the video?

<example 3>

original question: Is there a traffic light? If yes, what color is displayed?

rephrased question: Does a traffic light exist, and if so, what is its current color?

original question: {question}

rephrased question:

Table 5: Prompt template of question rephrase.

SYSTEM:

You are a helpful assistant.

USER:

Given a description about a road street view in Automatic Driving task, you are required to generate questions and provide corresponding answers from the description.

You MUST FOLLOW THE FOLLOWING RULES:

1. The questions must be related to driving from multiple perspectives, such as traffic object recognition, traffic condition analysis, and so on.
2. The questions cannot be repeated.
3. The corresponding answer must be condensed, such as fewer than 5 tokens.
4. The content of the questions and answers must not include anything not covered in the description.
5. Return only a series of Q&A pairs in the form of a dictionary list. DON'T RESPOND ANYTHING ELSE!

Here are some examples:

<example 1>

description: In the traffic scene observed, there is one vehicle parked on the side of the road. Although it is partially obscured by darkness, it serves as an indication of a potentially active parking area where the ego car should remain vigilant for other vehicles that might enter or exit parking spots. Additionally, there is a traffic sign situated on the right that informs drivers of an upcoming pedestrian crossing. This suggests that the ego car should decrease its speed and be ready to yield to pedestrians who might be in the vicinity.

Q&A: [{"question": "What should the ego car remain vigilant for?","answer": "Other vehicles entering/exiting"}, {"question": "Where is the traffic sign located?","answer": "On the right"}, {"question": "What does the traffic sign indicate?","answer": "Upcoming pedestrian crossing"}, {"question": "What should the ego car do near the pedestrian crossing?","answer": "Decrease speed and yield"}]

<example 2>

description: In the traffic image, there is a small black dog on the left side of the road ahead, appearing to cross from left to right. The dog's presence and potential for unpredictable movement pose a hazard, as animals can change direction suddenly, which could lead to an accident if the autonomous vehicle does not respond correctly. There are no vehicles, traffic signs, traffic lights, traffic cones, barriers, or other objects present in the image that affect driving behavior.

Q&A: [{"question": "Are there any vehicles on the road?","answer": "No"}, {"question": "Are there any traffic signs present?","answer": "No"}, {"question": "Are there any barriers on the road?","answer": "No"}]

<example 3>

description: The traffic scene contains a large cement mixer truck ahead on the road, taking up most of the driving lane. The truck is equipped with rear safety features like warning lights. This vehicle significantly impacts the driving behavior of the ego car, as its size and positioning affect the car's ability to safely overtake and necessitate maintaining a sufficient following distance. Additionally, there are construction barriers and debris on the side of the road to the right. These obstructions narrow the available road space, posing potential hazards that require the ego car to drive cautiously to avoid a collision. There are no vulnerable road users, traffic signs, traffic lights, traffic cones, or other objects present in this image that affect driving conditions.

Q&A: [{"question": "What vehicle is ahead on the road?","answer": "Cement mixer truck"}, {"question": "What safety features does the truck have?","answer": "Warning lights"}, {"question": "How does the truck affect the ego car?","answer": "Impacts overtaking and following distance"}, {"question": "What is on the right side of the road?","answer": "Construction barriers and debris"}, {"question": "How do the obstructions on the right affect driving?","answer": "Narrow the road space"}, {"question": "Are there any vulnerable road users present?","answer": "No"}, {"question": "Are there any traffic signs or lights in the scene?","answer": "No"}]

description: {description}

Q&A:

Table 6: Prompt template of general QA pairs self-generate.

SYSTEM:

You are a helpful assistant.

USER:

Given a driving suggestion about a road street view in Automatic Driving task, you are required to generate questions and provide corresponding answers from the description.

You MUST FOLLOW THE FOLLOWING RULES:

1. The questions must be related to driving from multiple perspectives, such as driving suggestion, traffic condition analysis, and so on.
2. The questions cannot be repeated.
3. The corresponding answer must be condensed, such as fewer than 5 tokens.
4. The content of the questions and answers must not include anything not covered in the description.
5. Return only a series of Q&A pairs in the form of a dictionary list. DON'T RESPOND ANYTHING ELSE!

Here are some examples:

<example 1>

suggestion: The ego car should reduce speed and prepare to stop if necessary, giving ample space for the dog to cross safely. Maintain vigilant observation of the dog's movements and be prepared for sudden changes in its direction. Furthermore, due to the narrowness and environment of the road, the ego car should drive cautiously, being alert for other potential obstacles or road users that may emerge from the side buildings or alleys.

Q&A: [{"question": "What should the ego car do regarding speed?","answer": "Reduce speed"}, {"question": "Why should the ego car drive cautiously?","answer": "Narrow road and potential obstacles"}, {"question": "What other hazards should the driver be alert to?","answer": "Side buildings or alleys"}]

<example 2>

suggestion: Given the proximity to the heavy vehicle ahead and the construction debris to the side, the ego car should maintain a safe following distance, prepare to slow down or stop if the mixer truck's behavior indicates impending stops or turns and refrain from attempting to overtake unless the opposite lane is visibly clear and safe to do so. The ego car should also be prepared for potential hazards from the construction area, such as entering construction vehicles or workers, and remain vigilant for any changes in road width or conditions.

Q&A: [{"question": "What distance should the ego car maintain?","answer": "Safe following distance"}, {"question": "What hazards should the driver be prepared for?","answer": "Construction vehicles or workers"}]

<example 3>

suggestion: Given the wet road conditions, it is recommended that the ego car reduces its speed to ensure a safe stopping distance from the vehicle ahead. It should monitor the green traffic light for changes and be prepared to stop if it switches to yellow or red. The ego car should remain in the current lane, as the presence of vehicles in the left lane may hinder safe lane changing, and the road ahead seems to be leading to an exit which is partially blocked by barriers, making it inaccessible. Adherence to the road signage for navigation should be maintained, but no immediate action is required since there is no indication of an upcoming turn or exit that is accessible. Always be observant for any road users that might appear unexpectedly.

Q&A: [{"question": "What should the ego car do due to wet road conditions?","answer": "Reduce speed"}, {"question": "Which lane should the ego car stay in?","answer": "Current lane"}, {"question": "Why should the ego car avoid changing lanes?","answer": "Left lane vehicles may hinder"}]

suggestion: {suggestion}

Q&A:

Table 7: Prompt template of driving suggestion QA pairs self-generate.



Video Data Example




```
{  
  "data_type": "Obstacle recognition",  
  "image_type": "video",  
  "source": "lingoqa",  
  "src": "Are there any cyclists or motorcyclists on the  
road?",  
  "rephrase": "Is the road occupied by any individuals riding  
bicycles or motorcycles?",  
  "alt": "No cyclists.",  
  "image":  
      
  "image_rephrase":  
      
  "original_gr": "No, there are no cyclists or motorcyclists on  
the road.",  
  "rephrase_images_original_gr": "No, there are no cyclists or  
motorcyclists on the road.",  
  "loc": "nq question: mexican leader who was supported by the  
united states during mexican civil war",  
  "loc_ans": "Benito Juárez",  
  "m_loc":  
      
  "m_loc_q": "What region would you find this type of bear?",  
  "m_loc_a": "northern united states"  
}
```

Figure 9: A video data case of ADS-Edit benchmark.



Multi-views Image Data Example




```
{  
  "data_type": "Driving behavior prediction",  
  "image_type": "multi-image",  
  "source": "drivelm",  
  "src": "Predict the behavior of the ego vehicle.",  
  "rephrase": "What actions is the autonomous vehicle  
expected to take next?",  
  "alt": "Steers right, fast.",  
  "image":  
      
  "image_rephrase":  
      
  "original_gr": "The ego vehicle is slightly steering to the  
right. The ego vehicle is driving fast.",  
  "rephrase_images_original_gr": "The ego vehicle is slightly  
steering to the right. The ego vehicle is driving fast.",  
  "loc": "nq question: who got the first nobel prize in physics",  
  "loc_ans": "Wilhelm Conrad Röntgen",  
  "m_loc":  
      
  "m_loc_q": "Is this surfer regular footed or goofy footed?",  
  "m_loc_a": "regular footed"  
}
```

Figure 10: A multi-views image data case of ADS-Edit benchmark.



Single Image Data Example




```
{  
  "data_type": "Traffic light understanding",  
  "image_type": "single-image",  
  "source": "codalm",  
  "src": "What does the green traffic light indicate?",  
  "rephrase": "What is the meaning or instruction conveyed by  
a green traffic light?",  
  "alt": "Proceed with caution",  
  "image":  
      
    "image_rephrase":  
        
    "original_gr": "Proceed with caution",  
    "rephrase_images_original_gr": "Proceed with caution",  
    "loc": "nq question: mexican leader who was supported by the  
united states during mexican civil war",  
    "loc_ans": "Benito Juárez",  
    "m_loc":  
        
    "m_loc_q": "What century is this?",  
    "m_loc_a": "20th"  
}
```

Figure 11: A single image data case of ADS-Edit benchmark.