Supplementary Material

FutureVQA

Dataset Creation and Quality control

Our dataset creation aims to provide a benchmark that address VLMs ability in consistant and accurate future reasoning with focus on diverse questions costomized based on different scene. To achieve this we utilize the annotation pipeline operate with both human and AI agent, which to efficiently create the QA.

(1) Human Expert QA Generation and Quality Control: To construct a human-like benchmark dataset with high diversity, we employed five expert annotators to manually generate question-answer pairs based on selected clips from OpenDV-YouTube dataset (Yang et al. 2024a), covering multiple cities with different weathers. Each QA pair was subsequently reviewed and verified by 1–2 annotators to ensure clarity, unambiguity, and answerability based on the given input. Although time-consuming, this process results in a more diverse and naturally phrased QA dataset compared to rule-based or template-driven approaches.

Compared to existing works in the driving domain (see Figure 4), such as nuScenes-OA (Oian et al. 2023) and DRAMA (Malla et al. 2023), which rely on rule-based methods, or OmniDrive (Wang et al. 2024a), which uses GPTgenerated data to construct large-scale datasets, our benchmark prioritizes diversity and human-like reasoning. While DriveLM-ns (Sima et al. 2023) incorporates human annotations for prediction and planning tasks, it still follows a rigid and highly structured question format, regardless of the uniqueness of each video clip. As shown in Table 5, despite being smaller in overall size, our dataset provides over 4× more unique questions, nearly 3× larger vocabulary, and over 400× higher type-token ratio (TTR). Notably, more than 95% of our questions appear fewer than 10 times. In contrast, DriveLM contains over 85% of questions repeated more than 10^2 times, over 20% more than 10^3 times, and over 2% more than 10^4 times, without considering the uniqueness of differences in scene content.

(2) AI Quality Control and Multi-option Generation: To minimize typographical errors, we employ GPT-40 to review all QA pairs generated by human annotators and automatically correct any detected typos. Following this, GPT-40 is further used to generate plausible but incorrect answer options based on the ground-truth answers provided by annotators.

To ensure that the resulting multiple-choice questions remain unambiguous—with only one clearly correct answer—each generated QA pair is manually reviewed by human annotators. This final verification step guarantees the quality and clarity of the multi-option format in our dataset.

Evaluation Protocol

To address the limitations of conventional statistical-based metrics, we adopted an option-based answer format for evaluation, where each question has predefined multiple-choice

Dataset	N. Ques.	N. Uniq. Ques.	Vocab. Size	TTR
DriveLM(Pred.)	123k	15	69	4.1×10^{-5}
DriveLM(Pred.&Percep.)	285k	234	150	4.1×10^{-5}
Ours	2.7k	969	433	$1.8\times\mathbf{10^{-2}}$

Table 5: Comparison of question diversity between our dataset and DriveLM. **N. Ques.** denotes the total number of questions; **N. Uniq. Ques.** represents the number of unique questions after de-duplication; **Vocab. Size** is the number of distinct words used in the questions; and **TTR** (Type-Token Ratio) measures lexical diversity, computed as the ratio of unique words to total words. The results highlight its greater linguistic diversity and reduced reliance on fixed templates.

answers (e.g., A: Yellow). The models were required to provide the corresponding option label (e.g., A) as the answer. See Figure 7 for the prompt and Algorithm 1 for the multi-trials evaluation.

Interestingly, during our experiments, we observed that not all models consistently adhered to this strict answer format. Some models would output answers like A: Yellow or simply Yellow. To account for this, we relaxed the evaluation criteria to accept both answer formats as correct. However, models like Qwen-VL-7B (Bai et al. 2023) still struggled to follow the instructions and produced responses such as "The answer is A", "The answer is A: Yellow", or other variations. Since following instructions is an important part of the evaluation, we did not further relax this restriction, which resulted in lower accuracy for these models, as shown in Table 1.

FutureVQA Prompt

Imagine you are looking at the image **{future_second}**} second after the input frames and answer the following question:

Question: {question}
Options: {options}

Please choose the most appropriate answer from the given options. Respond with the option without any explanation, for example, if the answer is B: Yellow, your answer should be: B

Figure 7: The prompt used to instruct VLMs to predict the future scene and answer the corresponding question.

Algorithm 1: Multi-trial Evaluation for QA Consistency

Require: Model \mathcal{M} , Question Q, Image I, Ground-truth Answer A, Number of Trials N

- 1: for i = 1 to N do
- 2: $Q_i \leftarrow \text{ShuffleOptions}(Q)$
- 3: $P_i \leftarrow \mathcal{M}.\operatorname{predict}(I,Q_i)$
- 4: if $P_i \neq A$ then
- 5: return False
- 6: end if
- 7: **end for**
- 8: return True

Benchmark	Task	T. Size	Cust. Q	Ans. Type	Mul. Trl.	Mul. C.	Pred.	T-Pred.
nuScenes-QA (Qian et al. 2023)	Drive VQA	83.3k**	Х	Mixed	X	√	X	X
BDD-X (Kim et al. 2018)	Drive Action	2.6k	-	Sentence	X	✓	X	X
DRAMA (Malla et al. 2023)	Drive VQA	11.6k	X	Mixed	X	X	X	X
Rank2Tell (Sachdeva et al. 2024)	Drive VQA	-	X	Mixed	X	✓	X	X
OmniDrive (Wang et al. 2024a)	Drive VQA	24k†	✓	Sentence	X	✓	1	X
DriveLM-nS (Sima et al. 2023)	Drive VQA	73k*	X	Sentence	X	✓	1	X
MMBench (Liu et al. 2023c)	Gen. I. QA	1.7k	✓	Options	✓	-	-	-
LngVidBench (Wu et al. 2024)	Gen. V. QA	5.3k	✓	Options	X	-	-	-
Video-MME (Fu et al. 2024)	Gen. V. QA	2.7k	✓	Options	X	-	-	-
MME (Fu et al. 2023a)	Gen. I. QA	2.1k	X	Y/N	X	-	-	-
Ours	Drive VQA	2.8k	✓	Options	✓	✓	✓	√

Table 6: Comparison of existing VLM benchmarks. Key aspects of dataset creation include test size (T. Size), whether questions are customized for different scenarios and video clips (Cust. Q), answer type (Ans. Type), multi-trial evaluation for each question (Mul. Tri), inclusion of multiple cities (Mul. C.), presence of perception tasks (Perc.), inclusion of prediction tasks (Pred.), and whether the dataset challenges VLMs with time-specific prediction (T-Pred.). Our benchmark dataset is the first in the driving domain that does not rely on sentence-based answers, which are subjective and difficult to evaluate. Additionally, it consists of fully human-annotated QA pairs tailored to different scenes, rather than relying on rule-based methods. Furthermore, our dataset challenges VLMs to predict future scenes at specific time intervals, requiring precise temporal reasoning to differentiate between near-future and far-future events.

†: The QA pairs are fully generated by GPT-4. ** Fully rule-based (no human annotators), * Semi-rule-based labeling (with human annotators for certain tasks).

VQA Category

To evaluate the diverse reasoning capabilities of VLMs, we classify VQA tasks into the following categories. These categories are not mutually exclusive, as a single question can belong to multiple categories depending on the type of reasoning required.

- Hallucination: This category evaluates the model's ability to avoid providing incorrect information about objects or features that do not exist in the scene. (e.g., "How many blue cars do you see in this image?") Such questions are especially challenging when an object has just left the scene.
- **General:** General questions involve straightforward scene understanding or recognition tasks that do not require spatial or temporal reasoning. Examples include identifying landmarks, objects, or common scene elements (e.g., "What is the landmark in the middle of the image?").
- **Traffic Understanding:** This category targets trafficrelated reasoning, including understanding road signs, speed limits, or dynamic traffic scenarios. These questions often require knowledge specific to driving environments (e.g., "What is the speed limit here?").
- **Absolute Location:** Absolute location questions focus on the spatial properties of objects in the scene, such as identifying specific positions or attributes relative to the image boundaries (e.g., "What color is the car on the far right of the image?").
- **Relative Position:** Relative position questions require understanding the spatial relationships between multiple objects in the scene. These questions test the model's ability to interpret multiple objects interaction (e.g., "Describe the vehicle in front of the taxi.").

By introducing these categories, we aim to provide a comprehensive evaluation framework for VLMs, covering both basic scene understanding and complex reasoning tasks. See Figure 8 for the examples.

Analysis on Different FutureVQA Categories

To establish a baseline for expected performance in the FutureVQA, we analyze various VLMs on our benchmark dataset across different categories . In this baseline analysis, VLMs perform regular VQA, where the actual images corresponding to the questions are provided as input.

As shown in Figure 9, we evaluate models includes CogVLM (Wang et al. 2023), Yi-VL (Young et al. 2024), LLaVA series (Liu et al. 2023a, 2024) and GPT-40, the results suggest that traffic understanding appears to be a relatively weak area for many existing VQA models. Most models do not exhibit significant differences in their capability to handle absolute or relative position questions. Additionally, for hallucination-related tasks, where models are asked about nonexistent objects, most models perform well when the image is provided, effectively avoiding incorrect predictions. These findings highlight the strengths and weaknesses of current VLMs and provide a foundation for evaluating their potential performance in future image QA tasks.

In Figure 10, we further compare the performance of VLMs across different question categories when asked to predict future scenes. As time progresses, we observe that GPT-4o's performance degrades significantly across all categories, with the most notable decline in questions related to relative and absolute positioning.

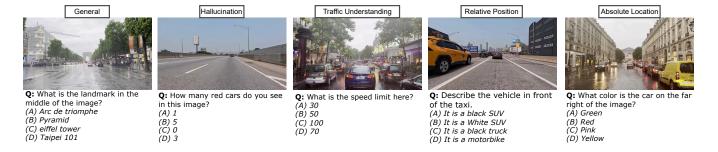


Figure 8: Examples of visual question answering (VQA) tasks categorized into different types: **Hallucination**, **General**, **Traffic Understanding**, **Absolute Location**, and **Relative Position**. Each question is categorized based on the type of reasoning it requires; however, a single question can belong to multiple categories simultaneously, depending on its context and the type of information needed.



Figure 9: Radar plots comparing the performance of various models across five VQA categories: Hallucination, General, Traffic Understanding, Absolute Location, and Relative Position. In this experiment, models perform regular VQA on images, with the actual images provided as input. The plots illustrate the strengths and weaknesses of each model in handling different reasoning tasks, providing a comparative baseline for understanding the capabilities of existing VLMs before extending to future image QA tasks.

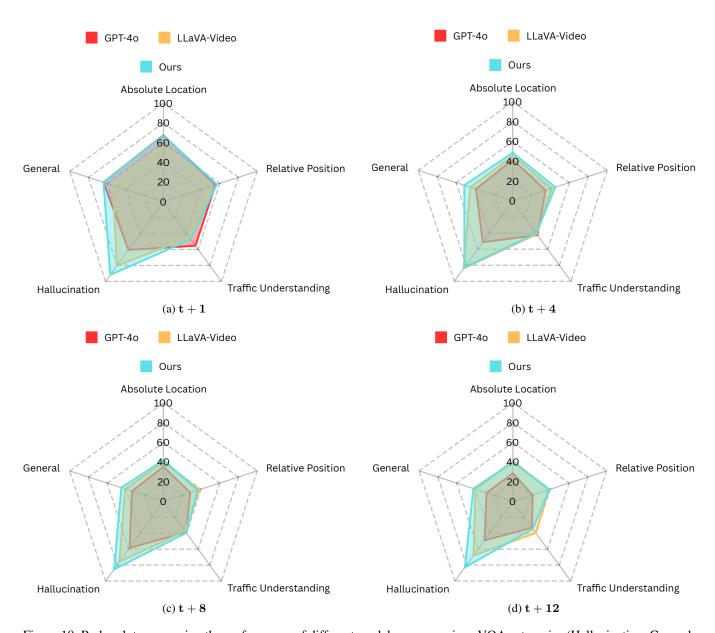


Figure 10: Radar plots comparing the performance of different models across various VQA categories (Hallucination, General, Traffic Understanding, Absolute Location, and Relative Position) at different future time steps: (a) t+1, (b) t+4, (c) t+8, and (d) t+12. The results highlight that while most models maintain robustness in hallucination detection, their performance in other categories, particularly traffic understanding and spatial reasoning, declines as the time offset increases.

Algorithm 2: Temporal Chain-of-Thought Future Scene Reasoning

Require: VLM \mathcal{A} , Observed Frames $I_{t-5:t}$, Target Future Step Δt , Question $Q_{t+t_{\Delta}}$

1: $D_0 \leftarrow$ Initialize empty future description

2: **for** i = 1 to Δt **do**

3: $D_i \leftarrow \mathcal{A}.describe_future(I_{t-5:t}, D_{i-1}, i)$

4: end for

5: $ans \leftarrow \mathcal{A}.answer(Q_{t+t_{\Delta}}, D_T)$

6: return ans

Detail Implementation

Chain-of-Thought

To enhance temporal reasoning, we adopt a Chain-of-Thought (CoT) prompting strategy in which the VLM predicts the future scene progressively, one step at a time. Rather than directly predicting the outcome at a future timestamp, the model is encouraged to reason through each intermediate step—first predicting t = 1, then t = 2, and so on, until the final target time is reached, see Algorithm 2. At each step, the model uses the history frames along with its previous predictions to generate the next future scene description. This design mimics human-like sequential foresight and allows the model to build up an understanding of how the scene may evolve over time. For practical computational efficiency, we limit the maximum number of steps to 4. This step-wise reasoning not only improves temporal consistency but also provides interpretable intermediate predictions that make the model's reasoning process more transparent and grounded in scene dynamics.

Visual Input Encoding

Memory Decay Sampling. Our implementation of the memory decay sampler leverages a transformer-based framework with learnable sampling queries $Q = \{q_1, q_2, \ldots, q_n\}$, where n is the total number of queries set as the initial number of tokens. These queries are initialized at the beginning of training and are optimized to extract temporal information relevant to the task. Let the current time be t_0 , and let the number of tokens provided by the image encoder be n_0 . The decay factor for the frame at time $t_0 - i$ is defined as $\left(\frac{1}{2}\right)^i$. Accordingly, the first $n_0 \cdot \left(\frac{1}{2}\right)^i$ queries are utilized in the cross-attention mechanism to represent the frame at $t_0 - i$.

Adaptive Token Sampling. In our implementation, frame similarity is evaluated by first computing the difference between two consecutive frames, |I(t)-I(t-1)|. To reduce noise introduced by high-frequency details, such as windows on distant buildings in urban environments, a Gaussian filter, G_{σ} , is applied to smooth the difference map while preserving significant changes. Finally, a Sobel operator, S_{xy} , is used to highlight the structural changes between the frames.

During our experiments, we tested multiple Gaussian filter kernel sizes and determined that a kernel size of 13 strikes the best balance between reducing noise and preserving important structural details. The comparison is shown

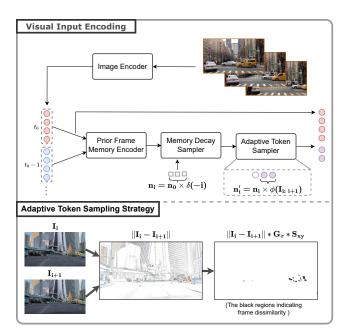


Figure 11: Overview of our visual encoding pipeline. The goal is to minimize the number of tokens while maintaining similar performance. In the context of autonomous driving videos, recent frames typically have greater influence on upcoming events. To reflect this, the **Memory Decay Sampler** assigns fewer queries to older frames, while the **Adaptive Token Sampler** adjusts the number of tokens based on the similarity between adjacent frames. The **Prior Frame Memory Encoder** is a transformer-based module that integrates temporal information from preceding frames.

in Figure 12. After computing the similarity maps, we measure the amount of highlighted area and then scale and cap the values for consistency. On average, the scaling factor is approximately 0.5 across our evaluation dataset.

Additional Evaluation

Ablation Study on Sampling Strategy

Our choice of the number of visual input frames is guided by two main considerations: (1) performance and (2) hardware constraints. The objective is to minimize the number of visual tokens while maintaining competitive performance. Detailed results at specific time steps are provided in Table 7 and Figure 13.

We observe that extending the input range from the past 5 seconds to the past 10 seconds does not lead to significant performance gains, yet results in increased computational cost. On the other hand, reducing the input to only the past 2 seconds leads to a slight drop in performance.

Similarly, ablating either the memory decay sampler or the adaptive token sampler individually does not substantially affect the final accuracy, while reducing visual token usage by approximately 75%. This highlights the efficiency of our sampling strategy in balancing performance and computational cost.



Figure 12: Visualization of frame similarity evaluation using Gaussian smoothing followed by the Sobel operator with different kernel sizes. The input images have a resolution of 1280×720 pixels. The difference between two consecutive frames, |I(t) - I(t-1)|, is computed, smoothed using Gaussian filters with kernel sizes of 5, 9, 11, 13, 15, and 17, and then processed with the Sobel operator, S_{xy} , to highlight changes. For better readability, the colors of the similarity maps are inverted.

Time	Model	Scores ↑					
111111111111111111111111111111111111111		B-3	B-4	R-L	С	M	
+1s	Baseline*	13.4	6.1	25.0	2.4	25.9	
	Ours*	32.2	26.2	40.2	17.7	41.5	
	w/o CoT	28.1	22.7	38.2	15.1	40.2	
	w/o self-sup.	12.1	7.2	24.9	2.7	26.2	
	$t_{0s:-10s}$	32.5	26.5	40.0	17.6	41.3	
	w/o Adpt. Sam.	32.3	26.2	40.4	17.3	41.6	
	w/o Mem. Sam.	31.8	26.2	40.1	17.0	41.2	
	Baseline*	22.5	6.0	24.4	2.6	25.5	
	Ours*	28.6	22.5	37.1	11.8	39.1	
	w/o CoT	25.6	20.1	35.9	11.0	38.4	
+4s	w/o self-sup.	11.6	6.7	24.4	2.0	25.8	
	$t_{0s:-10s}$	32.1	26.2	40.7	17.7	41.5	
	w/o Adpt. Sam.	32.7	26.3	40.6	18.0	41.5	
	w/o Mem. Sam.	32.7	25.6	40.8	18.1	41.5	
	Baseline*	11.2	6.2	25.1	2.2	25.4	
	Ours*	27.5	21.4	36.2	10.1	38.3	
	w/o CoT	24.1	18.6	34.9	9.7	37.6	
+8s	w/o self-sup.	12.0	7.1	24.9	2.3	26.3	
	$t_{0s:-10s}$	32.2	25.8	40.2	17.7	41.5	
	w/o Adpt. Sam.	32.3	26.6	41.5	17.5	41.5	
	w/o Mem. Sam.	32.0	26.4	41.9	17.1	41.5	
+12s	Baseline*	11.3	7.2	23.9	2.2	25.5	
	Ours*	26.7	20.6	35.6	9.4	37.7	
	w/o CoT	25.6	20.1	34.4	8.6	36.9	
	w/o self-sup.	11.5	6.7	24.4	2.0	25.8	
	$t_{0s:-10s}$	31.2	26.0	39.5	17.0	41.5	
	w/o Adpt. Sam.	32.0	25.5	39.6	16.9	41.5	
	w/o Mem. Sam.	32.1	26.2	40.4	17.9	41.5	

Table 7: In this comparison the reference captions are from regular image captioning, while the compared captions are generated by our fine-tuned model which perform future scenes captioning with only previous frames are given. B-3: BLEU-3, B-4: BLEU-4, R-L: ROUGE-L, C: CIDEr, M: METEOR.

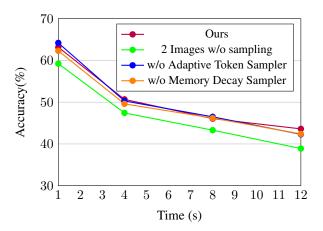


Figure 13: Accuracy over time in the FutureVQA task with different visual input encoding strategies, showing that our sampling approach reduces the number of required tokens while maintaining higher performance.

Hallucination and Mode Collapse

Generating accurate captions is the first and most crucial step in our training methodology. We experimented with various models for this task; however, we observed that not all models are capable of providing objective and accurate captions that comprehensively describe all elements in the scene. For instance, models like LLaVA-V1.5-7B (Liu et al. 2023b) tend to generate repetitive sentences and often produce hallucinations, exaggerating or inaccurately inferring details that are not present in the image. Figure 14 illustrates examples of these issues, showcasing captions that overstate the number of objects in the scene and use overly similar and redundant phrasing. These limitations highlight the need for more robust captioning models to ensure high-quality data generation for downstream tasks.

Prompt Template

In this section, we describe the unified prompt template used for our experiments across three key tasks: captioning evaluation, regular VQA, and FutureVQA. The template, shown in Figure 15, standardizes the model's input format to ensure consistent and fair evaluation.

For captioning evaluation, the model generates captions for a given image, which are subsequently scored by GPT-40 acting as a judge. GPT-40 is instructed to provide a score between 1 and 10 based on the objective aspects of the caption, explicitly disregarding subjective elements such as mood or atmosphere.

In the regular VQA task, the model is provided with the input image and a set of predefined multiple-choice options. It is required to select the most appropriate answer, establishing a baseline for evaluating the model's performance when the image is explicitly available.

LLaVA-1.5-7B





The image features a busy street with multiple cars and trucks. There are **at least 11 cars and 3 trucks visible on the road**. Some cars are positioned closer to the foreground, while others are further back in the scene. The vehicles are driving under a bridge, and there are traffic lights present in the area. The traffic lights are located at different positions along the road, ensuring proper traffic flow and safety.





The image shows a busy street with multiple cars and trucks. **There are at least 11 cars and 2 trucks visible on the road.** Some cars are positioned closer to the foreground, while others are further back in the scene. The vehicles are driving in various directions, creating a sense of movement and activity on the street.

LLaVA-1.6-34B



The image depicts a highway scene with a focus on a road intersection. There is a white SUV in the foreground, stopped at the intersection. The road is marked with white lane markings and a dashed white line indicating a lane change area. A yellow diamond-shaped sign with a black symbol of a person walking is visible on the right side of the road, indicating a pedestrian crossing or area.

Above the intersection, there is a green highway sign with white lettering that reads ""Left Exit 1A,"" suggesting the direction of the exit ramp. The sign also includes a symbol indicating the exit is to the left.

The highway overpass is in the background, with a greenish-gray color. The sky is clear and blue, indicating fair weather conditions. There are no visible people or animals in the image. The style of the image is a standard, non-stylized photograph with no filters or artistic effects applied. The focus is on the road infrastructure and the vehicle, with no additional context or background details provided.



In the image, there is a highway scene with several vehicles and a clear sky. On the left side of the image, there is a black SUV driving on the highway. The SUV is positioned in the left lane, and it appears to be in motion.

In the center of the image, there is a large black van or bus. It is driving in the right lane of the highway. The vehicle is quite large, occupying a significant portion of the image from left to right.

On the right side of the image, there is a black car driving in the right lane of the highway. The car is smaller in comparison to the van and is positioned behind it.

In the background, there are trees and a clear sky, indicating that the highway is likely in a suburban or rural area. There are no people visible in the image. The overall scene suggests a typical day with vehicles traveling on a highway.

Figure 14: Comparison of captions generated by LLaVA-V1.5-7B (Liu et al. 2023b) and LLaVA-V1.6-34B (Liu et al. 2024). While LLaVA-V1.5-7B produces shorter and repetitive captions with occasional hallucination, LLaVA-V1.6-34B generates significantly longer and more detailed descriptions. Additionally, LLaVA-V1.6-34B exhibits a varied response pattern, providing distinct levels of detail and focus when presented with different images.

GPT-40 as Judge for Captioning Evaluation

Please act as an impartial judge and evaluate the quality of the image caption provided by an AI assistant displayed below. Your evaluation should specifically assess the accuracy of object presence and positioning within the image, disregarding any subjective descriptions like vibe, atmosphere, or general impressions. Focus solely on whether the caption correctly reflects the precise positioning and presence of each object mentioned. Begin your evaluation by providing a short explanation. Be as objective as possible. After providing your brief explanation, please rate the response on a scale of 1 to 10 by strictly following this format: '[[rating]]', for example: 'Rating: [[5]]'.

Caption by the AI assistant: {caption}

Regular VQA

Answer the following question based on the image:

Question: **{question}**Options: **{options}**

Please choose the most appropriate answer from the given options. Respond with the option without any explanation, for example, if the answer is B:

Yellow, your answer should be: B

Figure 15: Prompts used for three tasks: GPT-40 as a judge in captioning evaluation and Regular VQA on our annotated evaluation dataset. Each prompt is tailored to the specific requirements of its respective task.

References

- Antol, S.; Agrawal, A.; Lu, J.; Mitchell, M.; Batra, D.; Zitnick, C. L.; and Parikh, D. 2015. VQA: Visual Question Answering. In *International Conference on Computer Vision (ICCV)*.
- Bae, S.; Kyung, D.; Ryu, J.; Cho, E.; Lee, G.; Kweon, S.; Oh, J.; Ji, L.; Chang, E.; Kim, T.; et al. 2024. Ehrxqa: A multi-modal question answering dataset for electronic health records with chest x-ray images. *Advances in Neural Information Processing Systems*, 36.
- Bai, J.; Bai, S.; Yang, S.; Wang, S.; Tan, S.; Wang, P.; Lin, J.; Zhou, C.; and Zhou, J. 2023. Qwen-VL: A Versatile Vision-Language Model for Understanding, Localization, Text Reading, and Beyond. *arXiv preprint arXiv:2308.12966*.
- Chang, C.-P.; Pagani, A.; and Stricker, D. 2024. 3D Spatial Understanding in MLLMs: Disambiguation and Evaluation. arXiv:2412.06613.
- Chang, C.-P.; Wang, S.; Pagani, A.; and Stricker, D. 2024. MiKASA: Multi-key-anchor & scene-aware transformer for 3d visual grounding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14131–14140.
- Chen, K.; Li, Y.; Zhang, W.; Liu, Y.; Li, P.; Gao, R.; Hong, L.; Tian, M.; Zhao, X.; Li, Z.; et al. 2024a. Automated evaluation of large vision-language models on self-driving corner cases. *arXiv* preprint *arXiv*:2404.10595.
- Chen, L.; Sinavski, O.; Hünermann, J.; Karnsund, A.; Willmott, A. J.; Birch, D.; Maund, D.; and Shotton, J. 2024b. Driving with llms: Fusing object-level vector modality for explainable autonomous driving. In 2024 IEEE International Conference on Robotics and Automation (ICRA), 14093–14100. IEEE.
- Deruyttere, T.; Vandenhende, S.; Grujicic, D.; Van Gool, L.; and Moens, M.-F. 2019. Talk2car: Taking control of your self-driving car. *arXiv preprint arXiv:1909.10838*.
- Fatemi, B.; et al. 2024. Test of time: A benchmark for evaluating llms on temporal reasoning. *arXiv*.
- Fu, C.; Chen, P.; Shen, Y.; Qin, Y.; Zhang, M.; Lin, X.; Yang, J.; Zheng, X.; Li, K.; Sun, X.; et al. 2023a. MME: A Comprehensive Evaluation Benchmark for Multimodal Large Language Models. *arXiv preprint arXiv:2306.13394*.
- Fu, C.; Dai, Y.; Luo, Y.; Li, L.; Ren, S.; Zhang, R.; Wang, Z.; Zhou, C.; Shen, Y.; Zhang, M.; et al. 2024. Video-MME: The First-Ever Comprehensive Evaluation Benchmark of Multi-modal LLMs in Video Analysis. *arXiv* preprint arXiv:2405.21075.
- Fu, H.; Zhang, D.; Zhao, Z.; Cui, J.; Liang, D.; Zhang, C.; Zhang, D.; Xie, H.; Wang, B.; and Bai, X. 2025. Orion: A holistic end-to-end autonomous driving framework by vision-language instructed action generation. *arXiv* preprint *arXiv*:2503.19755.
- Fu, J.; Ng, S.-K.; Jiang, Z.; and Liu, P. 2023b. Gptscore: Evaluate as you desire. *arXiv preprint arXiv:2302.04166*.
- Gao, S.; Yang, J.; Chen, L.; Chitta, K.; Qiu, Y.; Geiger, A.; Zhang, J.; and Li, H. 2024. Vista: A Generalizable Driving

- World Model with High Fidelity and Versatile Controllability. In *Advances in Neural Information Processing Systems* (NeurIPS).
- Gopalkrishnan, A.; Greer, R.; and Trivedi, M. 2024. Multi-Frame, Lightweight & Efficient Vision-Language Models for Question Answering in Autonomous Driving. *arXiv* preprint arXiv:2403.19838.
- Goyal, Y.; Khot, T.; Summers-Stay, D.; Batra, D.; and Parikh, D. 2017. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Hassan, M.; Stapf, S.; Rahimi, A.; Rezende, P.; Haghighi, Y.; Brüggemann, D.; Katircioglu, I.; Zhang, L.; Chen, X.; Saha, S.; et al. 2024. GEM: A Generalizable Ego-Vision Multimodal World Model for Fine-Grained Ego-Motion, Object Dynamics, and Scene Composition Control. *arXiv preprint arXiv:2412.11198*.
- Hu, A.; Russell, L.; Yeo, H.; Murez, Z.; Fedoseev, G.; Kendall, A.; Shotton, J.; and Corrado, G. 2023. Gaia-1: A generative world model for autonomous driving. *arXiv* preprint arXiv:2309.17080.
- Hu, X.; Yin, W.; Jia, M.; Deng, J.; Guo, X.; Zhang, Q.; Long, X.; and Tan, P. 2024. Driving World: Constructing World Model for Autonomous Driving via Video GPT. *arXiv* preprint arXiv:2412.19505.
- Huang, Z.; Zhang, J.; and Ohn-Bar, E. 2024. Neural volumetric world models for autonomous driving. In *European Conference on Computer Vision*, 195–213. Springer.
- Hwang, J.-J.; Xu, R.; Lin, H.; Hung, W.-C.; Ji, J.; Choi, K.; Huang, D.; He, T.; Covington, P.; Sapp, B.; et al. 2024. Emma: End-to-end multimodal model for autonomous driving. *arXiv preprint arXiv:2410.23262*.
- Jia, F.; Mao, W.; Liu, Y.; Zhao, Y.; Wen, Y.; Zhang, C.; Zhang, X.; and Wang, T. 2023. Adriver-i: A general world model for autonomous driving. *arXiv preprint arXiv*:2311.13549.
- Jiang, B.; Chen, S.; Liao, B.; Zhang, X.; Yin, W.; Zhang, Q.; Huang, C.; Liu, W.; and Wang, X. 2024. Senna: Bridging large vision-language models and end-to-end autonomous driving. *arXiv preprint arXiv:2410.22313*.
- Kembhavi, A.; Seo, M.; Schwenk, D.; Choi, J.; Farhadi, A.; and Hajishirzi, H. 2017. Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension. In *Proceedings of the IEEE Conference on Computer Vision and Pattern recognition*, 4999–5007.
- Khurana, T.; Hu, P.; Held, D.; and Ramanan, D. 2023. Point cloud forecasting as a proxy for 4d occupancy forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1116–1124.
- Kim, J.; Rohrbach, A.; Darrell, T.; Canny, J.; and Akata, Z. 2018. Textual explanations for self-driving vehicles. In *Proceedings of the European conference on computer vision (ECCV)*, 563–578.
- Lau, J. J.; Gayen, S.; Ben Abacha, A.; and Demner-Fushman, D. 2018. A dataset of clinically generated visual

- questions and answers about radiology images. *Scientific data*, 5(1): 1–10.
- Liao, G.; Li, J.; and Ye, X. 2024. VLM2Scene: Self-Supervised Image-Text-LiDAR Learning with Foundation Models for Autonomous Driving Scene Understanding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, 3351–3359.
- Lin, C.-Y. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.
- Liu, H.; Li, C.; Li, Y.; and Lee, Y. J. 2023a. Improved Baselines with Visual Instruction Tuning. *arXiv:2310.03744*.
- Liu, H.; Li, C.; Li, Y.; Li, B.; Zhang, Y.; Shen, S.; and Lee, Y. J. 2024. LLaVA-NeXT: Improved reasoning, OCR, and world knowledge.
- Liu, H.; Li, C.; Wu, Q.; and Lee, Y. J. 2023b. Visual Instruction Tuning. *NeurIPS*.
- Liu, Y.; Duan, H.; Zhang, Y.; Li, B.; Zhang, S.; Zhao, W.; Yuan, Y.; Wang, J.; He, C.; Liu, Z.; Chen, K.; and Lin, D. 2023c. MMBench: Is Your Multi-modal Model an Allaround Player? *arXiv:2307.06281*.
- Liu, Y.; Iter, D.; Xu, Y.; Wang, S.; Xu, R.; and Zhu, C. 2023d. G-eval: Nlg evaluation using gpt-4 with better human alignment. *arXiv preprint arXiv:2303.16634*.
- Malla, S.; Choi, C.; Dwivedi, I.; Choi, J. H.; and Li, J. 2023. DRAMA: Joint Risk Localization and Captioning in Driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 1043–1052.
- Manivasagam, S.; Wang, S.; Wong, K.; Zeng, W.; Sazanovich, M.; Tan, S.; Yang, B.; Ma, W.-C.; and Urtasun, R. 2020. Lidarsim: Realistic lidar simulation by leveraging the real world. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 11167–11176.
- Nie, M.; Peng, R.; Wang, C.; Cai, X.; Han, J.; Xu, H.; and Zhang, L. 2024. Reason2drive: Towards interpretable and chain-based reasoning for autonomous driving. In *European Conference on Computer Vision*, 292–308. Springer.
- Pan, C.; Yaman, B.; Nesti, T.; Mallik, A.; Allievi, A. G.; Velipasalar, S.; and Ren, L. 2024. VLP: Vision Language Planning for Autonomous Driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14760–14769.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a Method for Automatic Evaluation of Machine Translation. In Isabelle, P.; Charniak, E.; and Lin, D., eds., *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, 311–318. Philadelphia, Pennsylvania, USA: Association for Computational Linguistics.
- Qian, T.; Chen, J.; Zhuo, L.; Jiao, Y.; and Jiang, Y.-G. 2023. NuScenes-QA: A Multi-modal Visual Question Answering Benchmark for Autonomous Driving Scenario. *arXiv* preprint arXiv:2305.14836.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.;

- et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PmLR.
- Renz, K.; Chen, L.; Arani, E.; and Sinavski, O. 2025. Simlingo: Vision-only closed-loop autonomous driving with language-action alignment. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 11993–12003.
- Sachdeva, E.; Agarwal, N.; Chundi, S.; Roelofs, S.; Li, J.; Kochenderfer, M.; Choi, C.; and Dariush, B. 2024. Rank2tell: A multimodal driving dataset for joint importance ranking and reasoning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 7513–7522.
- Sellam, T.; Das, D.; and Parikh, A. P. 2020. BLEURT: Learning robust metrics for text generation. *arXiv* preprint *arXiv*:2004.04696.
- Sima, C.; Renz, K.; Chitta, K.; Chen, L.; Zhang, H.; Xie, C.; Luo, P.; Geiger, A.; and Li, H. 2023. DriveLM: Driving with Graph Visual Question Answering. *arXiv* preprint *arXiv*:2312.14150.
- Vasudevan, A. B.; Dai, D.; and Van Gool, L. 2018. Object referring in videos with language and human gaze. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 4129–4138.
- Wang, S.; Yu, Z.; Jiang, X.; Lan, S.; Shi, M.; Chang, N.; Kautz, J.; Li, Y.; and Alvarez, J. M. 2024a. OmniDrive: A Holistic LLM-Agent Framework for Autonomous Driving with 3D Perception, Reasoning and Planning. *arXiv:2405.01533*.
- Wang, W.; Lv, Q.; Yu, W.; Hong, W.; Qi, J.; Wang, Y.; Ji, J.; Yang, Z.; Zhao, L.; Song, X.; Xu, J.; Xu, B.; Li, J.; Dong, Y.; Ding, M.; and Tang, J. 2023. CogVLM: Visual Expert for Pretrained Language Models. arXiv:2311.03079.
- Wang, X.; Zhu, Z.; Huang, G.; Chen, X.; Zhu, J.; and Lu, J. 2024b. DriveDreamer: Towards Real-World-Drive World Models for Autonomous Driving. In *European Conference on Computer Vision*, 55–72. Springer.
- Wang, Y.; He, J.; Fan, L.; Li, H.; Chen, Y.; and Zhang, Z. 2024c. Driving into the future: Multiview visual forecasting and planning with world model for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14749–14759.
- Wei, J.; Wang, X.; Schuurmans, D.; Bosma, M.; Xia, F.; Chi, E.; Le, Q. V.; Zhou, D.; et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35: 24824–24837.
- Wu, H.; Li, D.; Chen, B.; and Li, J. 2024. Longvideobench: A benchmark for long-context interleaved video-language understanding. *arXiv preprint arXiv:2407.15754*.
- Xu, R.; et al. 2024. Knowledge conflicts for llms: A survey. *arXiv*.
- Yang, J.; Gao, S.; Qiu, Y.; Chen, L.; Li, T.; Dai, B.; Chitta, K.; Wu, P.; Zeng, J.; Luo, P.; Zhang, J.; Geiger, A.; Qiao, Y.; and Li, H. 2024a. Generalized Predictive Model for Autonomous Driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

- Yang, Z.; Chen, L.; Sun, Y.; and Li, H. 2024b. Visual point cloud forecasting enables scalable autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14673–14684.
- You, J.; Shi, H.; Jiang, Z.; Huang, Z.; Gan, R.; Wu, K.; Cheng, X.; Li, X.; and Ran, B. 2024. V2x-vlm: End-to-end v2x cooperative autonomous driving through large vision-language models. *arXiv* preprint arXiv:2408.09251.
- Young, A.; Chen, B.; Li, C.; Huang, C.; Zhang, G.; Zhang, G.; Wang, G.; Li, H.; Zhu, J.; Chen, J.; et al. 2024. Yi: Open foundation models by 01. ai. *arXiv preprint arXiv:2403.04652*.
- Yuan, W.; Neubig, G.; and Liu, P. 2021. Bartscore: Evaluating generated text as text generation. *Advances in Neural Information Processing Systems*, 34: 27263–27277.
- Zhang, H.; Li, X.; and Bing, L. 2023. Video-LLaMA: An Instruction-tuned Audio-Visual Language Model for Video Understanding. *arXiv* preprint arXiv:2306.02858.
- Zhang, P.; Goyal, Y.; Summers-Stay, D.; Batra, D.; and Parikh, D. 2016. Yin and Yang: Balancing and Answering Binary Visual Questions. In *Conference on Computer Vision and Pattern Recognition (CVPR)*.
- Zhang, Y.; Li, B.; Liu, h.; Lee, Y. j.; Gui, L.; Fu, D.; Feng, J.; Liu, Z.; and Li, C. 2024. LLaVA-NeXT: A Strong Zero-shot Video Understanding Model.
- Zhao, G.; Wang, X.; Zhu, Z.; Chen, X.; Huang, G.; Bao, X.; and Wang, X. 2024. Drivedreamer-2: Llm-enhanced world models for diverse driving video generation. *arXiv preprint arXiv:2403.06845*.
- Zheng, L.; Chiang, W.-L.; Sheng, Y.; Zhuang, S.; Wu, Z.; Zhuang, Y.; Lin, Z.; Li, Z.; Li, D.; Xing, E.; et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36: 46595–46623.
- Zhou, X.; Liu, M.; Yurtsever, E.; Zagar, B. L.; Zimmer, W.; Cao, H.; and Knoll, A. C. 2024a. Vision language models in autonomous driving: A survey and outlook. *IEEE Transactions on Intelligent Vehicles*.
- Zhou, Z.; Liu, S.; Han, X.; Liu, H.; Ng, K. W.; Xie, T.; Cong, Y.; Li, H.; Xu, M.; Pérez-Rúa, J.-M.; et al. 2024b. Learning Flow Fields in Attention for Controllable Person Image Generation. *arXiv preprint arXiv:2412.08486*.