

# Holistic Autonomous Driving Understanding by Bird’s-Eye-View Injected Multi-Modal Large Models

Xinpeng Ding<sup>1</sup> Jianhua Han<sup>2</sup> Hang Xu<sup>2\*</sup> Xiaodan Liang<sup>3</sup> Wei Zhang<sup>2</sup> Xiaomeng Li<sup>1\*</sup>

<sup>1</sup>The Hong Kong University of Science and Technology <sup>2</sup>Huawei Noah’s Ark Lab

<sup>3</sup>Sun Yat-Sen University

xdingaf@connect.ust.hk, {hanjianhua4, xu.hang, wz.zhang}@huawei.com  
 xdliang328@gmail.com, eexmli@ust.hk

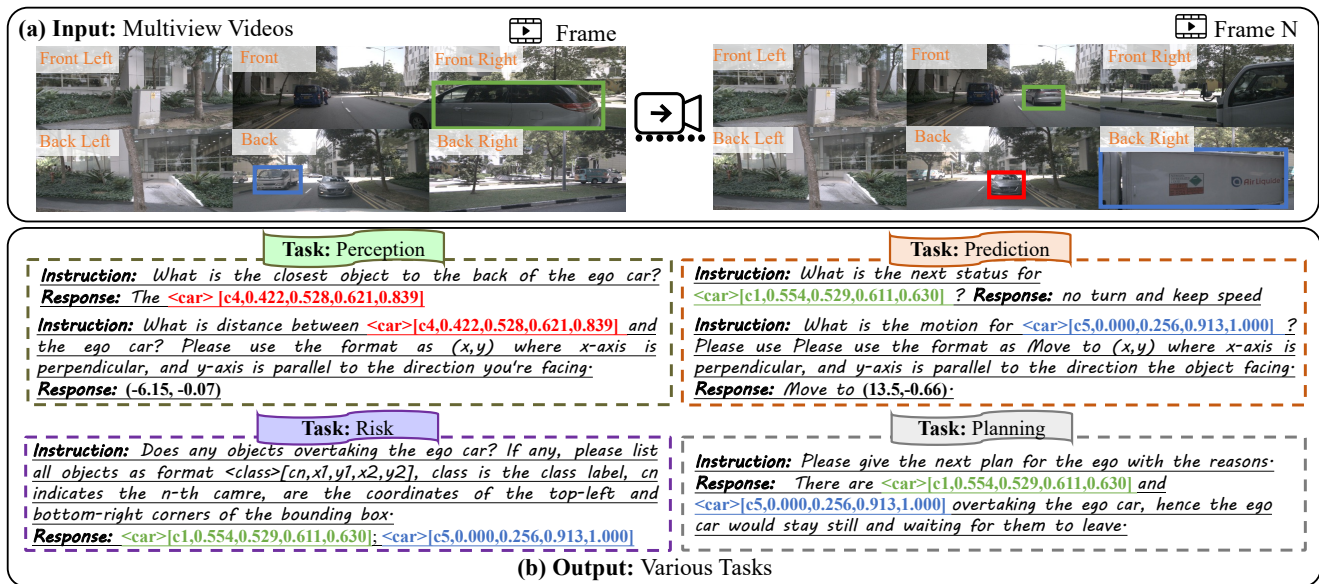


Figure 1. Example of our proposed NuInstruct dataset for holistic language-based autonomous driving. (a) The input are multi-view videos. (b) Various tasks are presented in instruction-response format. There are a total of four tasks, covering 17 subtasks (see in Fig.4 (a)).

## Abstract

The rise of multimodal large language models (MLLMs) has spurred interest in language-based driving tasks. However, existing research typically focuses on limited tasks and often omits key multi-view and temporal information which is crucial for robust autonomous driving. To bridge these gaps, we introduce NuInstruct, a novel dataset with 91K multi-view video-QA pairs across 17 subtasks, where each task demands holistic information ( e.g., temporal, multi-view, and spatial), significantly elevating the challenge level. To obtain NuInstruct, we propose a novel SQL-based method to generate instruction-response pairs automatically, which is inspired by the driving logical progression of humans. We further present BEV-InMLLM,

an end-to-end method for efficiently deriving instruction-aware Bird’s-Eye-View (BEV) features, language-aligned for large language models. BEV-InMLLM integrates multi-view, spatial awareness, and temporal semantics to enhance MLLMs’ capabilities on NuInstruct tasks. Moreover, our proposed BEV injection module is a plug-and-play method for existing MLLMs. Our experiments on NuInstruct demonstrate that BEV-InMLLM significantly outperforms existing MLLMs, e.g. 9% improvement on various tasks. We release our NuInstruct at <https://github.com/xmed-lab/NuInstruct>.

## 1. Introduction

Witnessing the success of multimodal large language models (MLLMs) [3, 5, 11, 13, 20, 22, 34–36, 45, 46], language-

\*Corresponding author

Dataset	Tasks				Information						Scale
	Perception	Prediction	Risk	P w/ R	Multi-view	Temporal	Multi-object	Distance	Position	Road	
BDD-X [16]	✗	✗	✗	✓	✗	✓	✗	✗	✗	✗	20K
Talk2Car [7]	✓	✗	✗	✗	✗	✓	✗	✗	✗	✗	11K
DRAMA [28]	✓	✗	✓	✗	✗	✓	✗	✗	✗	✗	100K
DRAMA-ROLISP [9]	✓	✗	✓	✓	✗	✓	✗	✗	✗	✗	35K
DriveGPT4 [44]	✓	✗	✓	✗	✗	✓	✓	✗	✗	✓	28K
Talk2BEV [8]	✓	✓	✗	✓	✗	✗	✓	✓	✓	✗	20K
Nuscenes-QA [37]	✓	✗	✗	✗	✓	✗	✓	✓	✓	✗	459K
NuPrompt [43]	✓	✗	✗	✗	✓	✓	✓	✗	✗	✗	35K
NuInstruct (Ours)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	91K

Table 1. **Comparison of our NuInstruct with existing language-based driving datasets.** ‘P w/ R’ indicates the planning with reasoning. NuInstruct provides various tasks and comprehensive information ( e.g., including multi-view, temporal, distance, and so on) for comprehensive autonomous driving understanding.

based driving is one of the trends in various autonomous driving tasks [8, 9, 28, 44]. For instance, some researchers ground the instruction prompts to single or multiple objects for 2D or 3D object detection and tracking [7, 41–43]. Nuscenes-QA [37] offers numerous question-answer pairs for multi-view perception tasks in driving scenes. Some advancements, e.g., DRAMA [28] and HiLM-D [9], generating text descriptions for localizing risk objects. Except for perception tasks, DriveGPT4 [44] and GPT-Driver [29] leverage LLMs for interpreting vehicle actions and planning, respectively. Talk2BEV [8] formulate BEV into a JSON file and input it into ChatGPT [30] to conduct autonomous driving understanding.

Although remarkable progress has been achieved, current language-based driving research still exhibits two main shortcomings as shown in Table 1. (i) **Partial tasks.** Existing benchmarks only cover a subset of autonomous driving tasks. However, autonomous driving comprises a series of interdependent tasks, each indispensable to the system’s overall functionality [2]. For instance, it is challenging to make reliable predictions when lacking accurate perception. (ii) **Incomplete information.** The information utilized by existing methods for executing these tasks is often incomplete. Specifically, existing datasets [9, 44] usually consist of single-view-based images, without considering temporal and multi-view information. However, safe driving decisions require a holistic understanding of the environment, e.g., only concerning on the front may neglect an overtaking vehicle in the left [39].

To address the above two problems, we first create **NuInstruct**, a comprehensive language-driving dataset with 91K multi-view video-QA pairs across 17 subtasks (Fig. 4). Our dataset presents more complex tasks than existing benchmarks, demanding extensive information like multi-view, temporal, distance and so on, as shown in Fig. 1 and Table 4. To obtain NuInstruct, we introduce a SQL-based method for the automated generation of instruction-response pairs. This method transforms instruction-response creation into a process utilizing structured query languages (SQLs) [6]

from a database. *Our rationale for the tasks and their corresponding SQL design follows the logical progression of human drivers:* ① *initially observing surrounding objects (Perception),* ② *predicting their behavior (Prediction),* ③ *assessing potential threats such as overtaking vehicles (Risk),* and ④ *ultimately using the previous information to plan a safe route with justified reasoning (Planning with Reasoning).* Finally, to ensure the quality of our NuInstruct, we conduct human or GPT-4 [30] verification to eliminate the erroneous instruction-response pairs. Compared with other data generation methods, e.g., ChatGPT-based [30] or human-based [8], this structured design ensures the generation of instruction-response pairs is both reliable and scalable.

To address the challenging tasks of the proposed NuInstruct, we further extend the current MLLMs to receive more holistic information. Existing MLLMs are constrained by their design for single-view inputs. To overcome this, we provide a Multi-View MLLM (MV-MLLM) with a specialized Multi-view Q-Former capable of processing multi-view video inputs. Although MV-MLLM allows for the capture of Multi-view temporal appearance features, they often miss out on critical information (e.g., distance, spatial) as well as suffer from occlusions. BEV’s feature, a formulation of multi-view inputs, has been widely adopted in traditional autonomous driving models since they can clearly represent object locations and scales (essential for distance/spatial-sensitive tasks) [15]. Leveraging this, we integrate BEV into MV-MLLM to create BEV-InMLLM, enhancing perception and decision-making in autonomous driving by capturing a comprehensive information spectrum. Inspired by this, we integrate BEV into MV-MLLM, obtaining BEV-InMLLM to capture a full spectrum of information for reliable perception and decision-making in autonomous driving. BEV-InMLLM uses a BEV injection module to effectively obtain BEV features aligned with language features for LLMs. This approach is more resource-efficient than training a BEV extractor from scratch with visual-language data like CLIP [38]. Notably, our BEV in-

jection module serves as a plug-and-play solution for existing MLLM

Overall, our contributions can be summarized as follows:

- We curate NuInstruct, a new language-driving dataset with 91K multi-view video-instruction-response pairs across 17 subtasks, using a novel SQL-based method. NuInstruct is currently the most holistic language-driving dataset, to our knowledge. We plan to release our NuInstruct for future research development.
- We propose BEV-InMLMM to integrate instruction-aware BEV features with existing MLLMs, enhancing them with a full suite of information, including temporal, multi-view, and spatial details. Notably, our BEV injection module serves as a plug-and-play solution for existing MLLM.
- Our experiments with NuInstruct demonstrate our proposed methods significantly boost MLLM performance in various tasks, notably outperforming state-of-the-art by 9% on various tasks. Ablation studies show that MV-MLLM enhances multi-view tasks, and BEV-InMLLM is vital for most tasks, emphasizing the importance of spatial information.

## 2. Related Work

**Language-driving datasets and models.** CityScapes-Ref [41], Talk2Car [7] perform language-grounding tasks. ReferKITTI [42] and NuPrompt [37] leverage temporal data for 2D or 3D referring object detection and tracking. Nuscenes-QA [37] offers numerous question-answer pairs for multi-view perception tasks in driving scenes. Some advancements, *e.g.*, DRAMA [28] and HiLM-D [9], generating text descriptions for localizing risk objects. Beyond perception, DriveGPT4 [44] and GPT-Driver [29] leverage LLMs for interpreting vehicle actions and planning, respectively. Talk2BEV [8] formulate BEV into a JSON file and input it into ChatGPT [30] to conduct autonomous driving understanding. Despite these advancements, a common limitation persists: most datasets and models address only part of the autonomous driving tasks with incomplete information. As shown in Table 1 and Fig. 1, in this paper, we propose a challenging dataset containing various tasks that require holistic information, *i.e.*, temporal, multi-view, spatial and so on, to address.

**Multimodal Large Language Models.** Leveraging the capabilities of pre-trained LLMs like LLaMA [40] and Vicuna [4], Multimodal LLMs (MLLMs) are expanding their application spectrum, handling inputs from images [1, 3, 18, 19, 46], videos [20, 45], and 3D data [13] to medical data [17]. In the domain of autonomous driving, DriveGPT4 [44] and Talk2BEV [8] have integrated MLLMs for comprehension. However, these approaches have limitations; DriveGPT4 is confined to single-view inputs, and

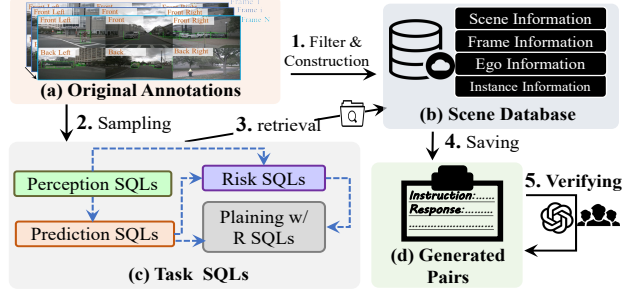


Figure 2. **Procedure of SQL-based data generation.** We formulate the data generation into an SQL-based process, using different task SQLs to retrieve the response from the scene information database. The design of SQLs follows the logical flow of autonomous driving tasks [14], which is represented in blue dashed arrows. ‘Planning w/ R’ indicates the planning with reasoning.

Talk2BEV lacks temporal dynamics and an end-to-end framework. Addressing these gaps, our BEV-InMLLM model assimilates comprehensive temporal, multi-view, and spatial data, for reliable decisions.

## 3. NuInstruct

In this section, we will illustrate the details of our constructed NuInstruct dataset. In Section 3.1, we will discuss the process of data generation. We will then go deeper and provide statistics on our dataset in Section 3.2. Finally, Section 3.3 shows the evaluation metrics to evaluate the performance of different models on the new NuInstruct.

### 3.1. Instruction Data Generation

Our NuInstruct is built on one of the most popular datasets, *i.e.*, Nuscenes [2]. There are six view records for samples of Nuscenes, *i.e.*, Front, Front Left, Front Right, Back Right, Back Left, and Back. These views have some areas of overlap with one another. In Nuscenes, the collected data is annotated with a frequency of 2Hz, and each annotated frame is referred to as a keyframe with annotations.

In our research, we propose an SQL-based approach for the automated generation of four types of instruction-follow data, namely: Perception, Prediction, Risk, and Planning with Reasoning. This methodology aligns with the sequential decision-making stages of human drivers, categorized as follows: **1. Perception:** The initial stage of recognizing surrounding entities. **2. Prediction:** Forecasting the future actions of these entities. **3. Risk:** Identifying imminent dangers, such as vehicles executing overtaking manoeuvres. **4. Planning with Reasoning:** Developing a safe travel plan grounded in logical analysis.

The detailed process is shown in Fig. 2. Specifically, **1).** The **filter & construction step** leverages the (a) original annotations to generate the scene information database (see Fig. 2 (b)). **2).** The **sampling step** first samples

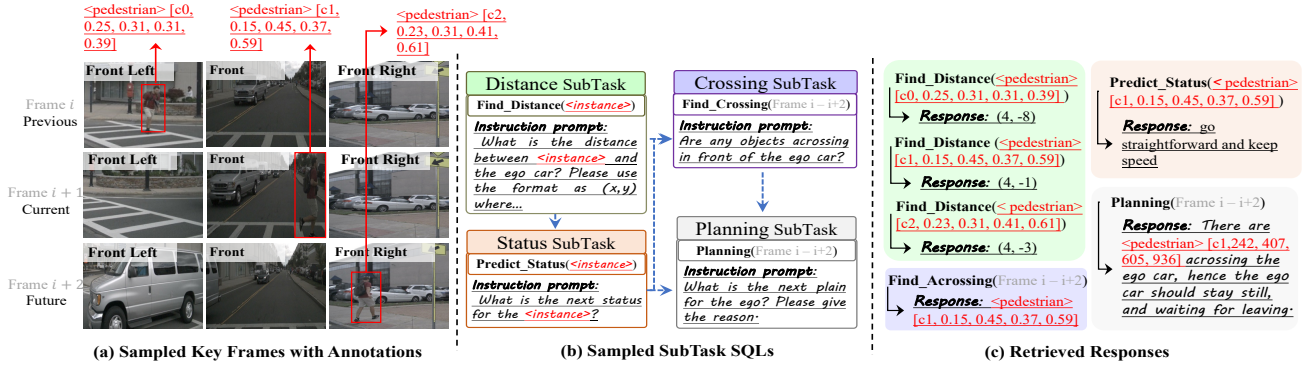


Figure 3. **The illustration of an example for Step 3 retrieval in the data generation process.** (a) **Sampled keyframes with annotations.** Three keyframes with annotations are randomly sampled, and we only select one instance, *i.e.*, the pedestrian (box), in this example for clarity. (b) **Sampled subtask SQLs.** Each subtask SQL consists of two parts, *i.e.*, the subtask function and the instruction prompt. (c) **Retrieved Responses.** The subtask function receives the specific input and retrieves the responses from the scene information database.

three keyframes from the original dataset. Then, as shown in Fig. 2 (c)), we construct a series of pre-defined task SQLs. Each task SQL consists of several subtasks, each of them consisting of a subtask function and an instruction prompt. **3).** The **retrieval step** uses the instruction prompt and the task SQL to retrieve the corresponding response from the scene database. **4).** The **saving step** saves all instruction-response pairs (see Fig. 2 (d)). **5).** The **verifying step** employs either human analysis or LLM-based methods (e.g., GPT-4 [30]) to eliminate erroneous instruction-response pairs, thereby guaranteeing the quality of NuInstruct. Our task SQL design is logically sequenced, and based on the inherent relational flow of autonomous driving tasks, *i.e.*, ‘Perception → Prediction, (Perception, Prediction) → Risk, (Risk, Prediction) → Planing with Reasoning’, where  $a \rightarrow b$  indicates the  $b$  SQL is derived from the  $a$  SQL (blue dashed arrows in Fig. 2 (c)).

We show a more detailed example for **Step 2** and **Step 3** in Fig. 3. Specifically, three keyframes (*i.e.*, from frame  $i$  to frame  $i + 2$ ) with annotations are randomly sampled from the original dataset, and we only select one instance, *i.e.*, the pedestrian (box) for clarity (Fig. 3 (a)). In this case, we choose distance, status, crossing, and planning subtask SQLs from the perception, prediction, risk, and planning with reasoning task SQLs, respectively (shown in Fig. 3 (b)). The status subtask is based on the distance task, since the next status (*e.g.*, speed, direction) for the instance is computed based on the distances of previous, current, and future frames. Each subtask SQL consists of two parts, *i.e.*, the subtask function and the instruction prompt. For example, the distance subtask SQL has `Find_Distance(<instance>)` function and the instruction prompt is “What is the distance between <instance> and the ego car? Please use the format as (x,y) where..”, where <instance> is the input. Finally, as shown in Fig. 3 (c), we use the instance information or frame in-

formation as the input for different subtask functions to retrieve the response from the scene database. Compared with other data generation methods, *e.g.*, ChatGPT-based [44] or human-based [8], this structured design ensures the generation of instruction-response pairs is both reliable and scalable.

We only describe the overview of the data generation in this section. Please refer to the supplementary material for more details about the scene information database, the task SQLs, and the retrieval process.

### 3.2. Data Statistics

To construct our NuInstruct, we sampled a total of 11,850 keyframes from 850 videos within the NuScenes dataset [2]. Subsequent filter yields 55,204 unique instances, which collectively appear 295,828 times across all keyframes. This culminates in an average of approximately 24.96 instances per keyframe. By employing our SQL-based method (Section 3.1), we generated a total of 91,355 instruction-response pairs, encompassing four primary tasks—namely, Perception, Prediction, Risk, and Planning with Reasoning. These tasks are further delineated into 17 sub-tasks. The quantities of task categories are statistically presented in Fig. 4 (a).

Compared with other single-view benchmarks, our dataset covers multi-view information. Hence, we also conduct a statistical analysis of the relations of different views and constructed instruction-response pairs. Fig. 4 (b) shows the distribution of the numbers of the responses based on the views, *i.e.*, for a given view, we record how many responses are derived from information from the view. Through such statistics, we find that to answer the instructions, the system needed to look at multiple views, instead of just a single one. In Fig. 4 (c), for each task, the proportions of responses obtained based on different views are calculated. We find two observations: **(i)** in the case of perception and

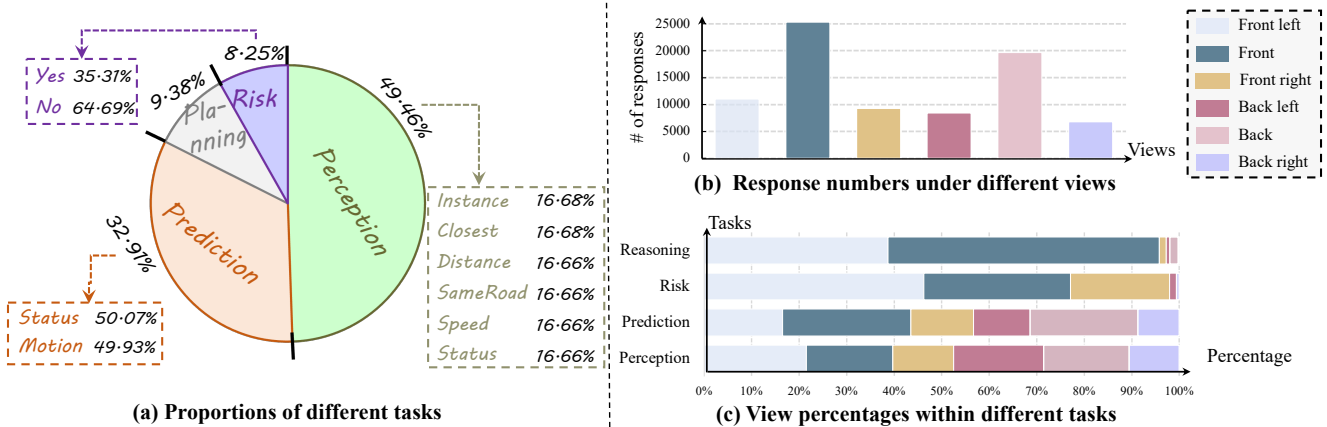


Figure 4. **Statistics of NuInstruct.** (a) **Proportions of different tasks.** The size of the arc represents the proportions of each task, while the same color indicates tasks of the same category. Our task encompasses a diverse range of tasks including perception, prediction, risk, and planning. (b) **Response numbers under different views.** The horizontal axis represents different views, and the vertical axis indicates the number of responses requiring information from the corresponding view. (c) **View percentages within different tasks.** The horizontal and vertical axes represent the proportion of different views and task classes, respectively.

Task	SubTask	Metrics
Perception	Distance, Speeds, Instance Number Closest, Status, Same Road	MAE ↓ Accuracy ↑
Prediction	Motion Ego, Motion Others Status Ego, Status Others	MAE ↓ Accuracy ↑
Risk	All	MAP ↑
Reasoning	All	BLEU [32] ↑

Table 2. **Evaluation metrics for different tasks.** ↓ represents the lower the scores, the better the results, while ↑ means the higher the scores, the better the results. ‘MAE’ indicates the mean absolute error. ‘All’ means the all subtasks.

prediction tasks, the distribution across views is relatively even, showing that our data generation method produces balanced multi-view information; (ii) When it comes to reasoning and risk tasks, the responses predominantly draw on information from the front, left-front, and right-front views. This is reasonable since drivers normally base their ahead or sides to decide the next actions, seldom looking behind.

### 3.3. Evaluation Protocols

**Evaluation Metrics.** Our NuInstruct consists of diverse tasks, making it hard to evaluate the different tasks using one metric. We summarize the evaluation metrics for different tasks in Table 2. For tasks evaluated by MAE, we use the regular expression [10] to obtain values. More detailed computations for different metrics please refer to the supplementary material.

**Data Split.** Our NuInstruct contains a total of 850 videos from NuScenes [2]. We split the all videos into training/validation/testing sets (7.5:1.5:1.5). We train all models in the training set and report the model with the best performance in the validation set on the test set.

## 4. Method

In Section 4.1, we first give a preliminary for our framework, *i.e.*, input, output, task definition and notations. Then, in Section 4.2, we provide Multi-view MLLM (MV-MLLM), a baseline that extends current multimodal large language models (MLLMs) for processing multi-view video inputs. Finally, in Section 4.3, we propose BEV-InMLLM, which injects the bird’s-eye-view (BEV) representations into current MLLMs for better panorama understanding for NuInstruct.

### 4.1. Preliminaries

Different from current MLLMs, the visual inputs for our model are the multi-view videos  $\{\mathbf{V}^i\}_{i=1}^{N_{\text{view}}}$ , where  $N_{\text{view}}$  is the total number of camera views,  $\mathbf{V}^i = \{\mathbf{v}_t^i\}_{t=1}^{N_{\text{frame}}}$  is the  $t$ -th frame in  $\mathbf{V}^i$  and  $N_{\text{frame}}$  is the total number of frames. Instead of the predefined several tasks, we give a specific language instruction, we use a unified model to obtain its corresponding language response, as shown in Fig. 1. For clarity in the following, we use  $\mathbf{L}_{\text{inst}} \in \mathbb{R}^{N_{\text{inst}} \times D_{\text{inst}}}$  and  $\mathbf{L}_{\text{resp}} \in \mathbb{R}^{N_{\text{resp}} \times D_{\text{resp}}}$  to denote the language instruction tokens and the response tokens respectively, which are generated by the language tokenizer [40].  $N_{\text{inst}}/N_{\text{resp}}$  and  $D_{\text{inst}}/D_{\text{resp}}$  are numbers of tokens and dimensions for the instruction/response.

### 4.2. Multi-View MLLM

Existing MLLMs [5, 18, 20, 46] generally consist of three parts: a vision encoder  $f_{\text{vision}}(\cdot)$  to receive the visual input; a connection module (*e.g.*, Q-Former [19])  $f_{\text{connect}}(\cdot)$  to transfer the visual representations to the visual tokens aligned with the language; a large language model (LLM)  $f_{\text{LLM}}(\cdot)$  to receive visual and language instruction tokens to

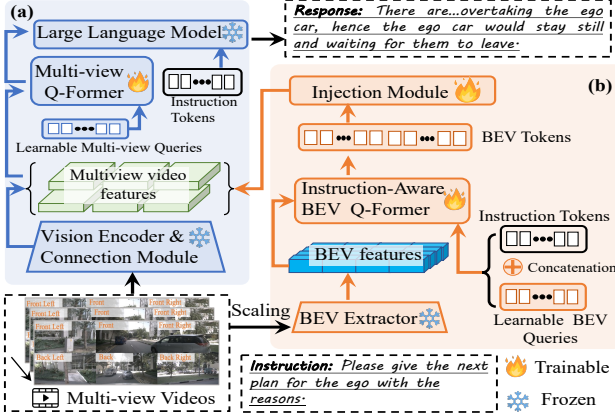


Figure 5. **The overall pipeline of our proposed BEV-InMLLM.** (a) The base multimodal large language model (MLLM) tailored for processing the multi-view videos. (b) The bird’s-eye-view injection module (BEV-In) to inject BEV representations into base MLLM to boost autonomous driving understanding.

generate the response. Since they can only receive single-view input, we propose a baseline model named Multi-view MLLM (MV-MLLM) to enable the current MLLMs to process the multi-view videos, as shown in Fig. 5 (a). Specifically, for a video from a specific view, *i.e.*,  $\mathbf{V}^i$ , we feed it into the vision encoder followed by the connect module to obtain the visual tokens, which can be formulated as:

$$\mathbf{F}_{\text{vis}}^i = f_{\text{connect}}(f_{\text{vision}}(\mathbf{V}^i)) \in \mathbb{R}^{N_{\text{vis}} \times D_{\text{vis}}}, \quad (1)$$

where  $N_{\text{vis}}$  and  $D_{\text{v}}$  are the numbers of the visual tokens and the dimensions respectively. Then, we introduce a multi-view Q-Former (similar to BLIP-2 [18]) to capture the MV visual semantics  $\mathbf{F}_{\text{mv}}$  from  $\{\mathbf{F}_{\text{vis}}^i\}_{i=1}^{N_{\text{view}}}$ . We concatenate  $\{\mathbf{F}_{\text{vis}}^i\}_{i=1}^{N_{\text{view}}}$  along the view dimension to obtain  $\bar{\mathbf{F}}_{\text{mv}} \in \mathbb{R}^{(N_{\text{view}} * N_{\text{vis}}) \times D_{\text{vis}}}$ . The input to the multi-view Q-Former contains a set of  $K_{\text{mv}}$  learnable multi-view queries  $\mathbf{Q}_{\text{mv}} \in \mathbb{R}^{K_{\text{mv}} \times D_{\text{vis}}}$ , which interact with  $\bar{\mathbf{F}}_{\text{mv}}$  through the cross attention, formulated as follows:

$$\mathbf{F}_{\text{mv}} = \text{CrossAttn}(\mathbf{Q}_{\text{mv}}, \bar{\mathbf{F}}_{\text{mv}}) \in \mathbb{R}^{K_{\text{mv}} \times D_{\text{vis}}}. \quad (2)$$

The output  $\mathbf{F}_{\text{mv}}$  then goes through a linear projection (omitted in Fig. 5), and is fed to the LLM. Note that in our MV-MLLM, only the MV Q-Former is trainable and other parameters are frozen to fully retain the knowledge of the pre-trained models.

### 4.3. BEV-Injected MLLM

The BEV approach has become pivotal in autonomous driving for its precise depiction of object positioning, essential for tasks like perception and planning [21, 27]. Integrating BEV into our MV-MLLM offers enhanced visual-spatial analysis for autonomous driving. While BEV can be constructed from the multi-view features  $\{\mathbf{F}_{\text{vis}}^i\}_{i=1}^{N_{\text{view}}}$

by physical transformations such as LSS [33], similar to NuScenes-QA [37], the plain of ViTs in current MLLMs limits their perception capabilities [31]. Replacing the ViT with BEV-specific backbones, *e.g.*, ResNet [12] or Swin Transformer [23], diminishes visual-language alignment [9]. Furthermore, the limited input resolutions generate small feature maps, which are hard to scale up to the high-resolution BEV representation.

To address the above problems, we propose the BEV-injected MLLM (BEV-InMLLM), which uses a BEV injection module (BEV-In) to obtain the BEV information aligned with LLMs in a data-efficient and resource-light way. As shown in Fig. 5 (b), we obtain the high-quality BEV features  $\mathbf{F}_{\text{bev}} \in \mathbb{R}^{W \times H \times D_{\text{bev}}}$  from the pre-trained BEV extractor [15, 33], where  $W$ ,  $H$  and  $D_{\text{bev}}$  denote the width, height and dimensions. The following are two key components of BEV-In, *i.e.*, an instruction-aware BEV Q-Former and an injection module.

**Instruction-aware BEV Q-Former.** We introduce the instruction-aware BEV Q-Former to ignore the redundant and irrelevant to the given instructions from  $\mathbf{F}_{\text{bev}}$ . The input queries for the instruction-aware BEV Q-Former blend two parts: the instruction tokens  $\mathbf{L}_{\text{inst}}$  for instruction-related aspects, and the learnable BEV queries for extracting the useful information pertinent to the instruction from  $\mathbf{F}_{\text{bev}}$ . The process of the instruction-aware BEV Q-Former is defined as:

$$\mathbf{F}_{\text{instbev}} = \text{CrossAttn}(\mathbf{Q}_{\text{bev}} \oplus \mathbf{L}_{\text{inst}}, \mathbf{F}_{\text{bev}}) \quad (3)$$

where  $\oplus$  indicates the concatenation,  $\mathbf{Q}_{\text{bev}} \in \mathbb{R}^{K_{\text{bev}} \times D_{\text{vis}}}$  and  $K_{\text{bev}}$  indicates the BEV queries and their the numbers,  $\mathbf{F}_{\text{instbev}} \in \mathbb{R}^{(K_{\text{bev}} + N_{\text{inst}}) \times D_{\text{vis}}}$  are the instruction-aware BEV tokens.

**Injection module.** Our injection module fuses multi-view features  $\bar{\mathbf{F}}_{\text{mv}}$  with instruction-aware BEV tokens  $\mathbf{F}_{\text{instbev}}$  through cross-attention:

$$\bar{\mathbf{F}}_{\text{mv}} = \bar{\mathbf{F}}_{\text{mv}} + \text{CrossAttn}(\bar{\mathbf{F}}_{\text{mv}}, \mathbf{F}_{\text{instbev}}), \quad (4)$$

where the enhanced  $\bar{\mathbf{F}}_{\text{mv}}$  contains both (i) temporal multi-view cues for scene comprehension and (ii) spatial-aware BEV information for precise perception and planning tasks. We keep our BEV-In module efficient by making only two components trainable: the BEV Q-Former and the injection module, while the BEV feature extractor remains frozen to maintain feature quality.

## 5. Experiments

**Implementation and training details.** We experiment on three base MLLMs, *i.e.* We evaluated our MV-MLLM and BEV-InMLLM on three base MLLMs: BLIP2 [19], Video-LLama [45], and MiniGPT-4 [46]. To adapt BLIP2 and MiniGPT-4, which are image-centric, for video input, we used the spatiotemporal adapter (ST-Adapter)[31], following[9], while preserving their pre-trained parameters.

Method	Perception						Prediction		Risk						Planning with Reasoning $\uparrow$
	Dis $\downarrow$	Sped $\downarrow$	# Ins $\downarrow$	Clos $\uparrow$	Sta $\uparrow$	SameR $\uparrow$	Mot $\downarrow$	Sta $\uparrow$	App $\uparrow$	LaneC $\uparrow$	Onco $\uparrow$	Cro $\uparrow$	Over $\uparrow$	Brak $\uparrow$	
BLIP-2* [18]	29.3	5.6	4.4	18.5	15.9	23.8	8.7	38.7	6.1	10.6	15.4	16.7	3.7	21.5	24.9
MV-MLLM	26.8	5.3	3.9	28.2	17.7	31.5	6.8	43.6	15.2	19.3	18.4	22.7	9.6	22.7	30.1
BEV-InMLLM	<b>23.3</b>	<b>3.6</b>	<b>3.2</b>	<b>33.6</b>	<b>18.9</b>	<b>31.6</b>	<b>3.8</b>	<b>45.2</b>	<b>16.8</b>	<b>21.0</b>	<b>19.7</b>	<b>23.9</b>	<b>10.5</b>	<b>27.5</b>	<b>33.3</b>
MiniGPT-4* [46]	30.2	6.2	6.3	20.2	17.3	24.2	8.7	39.6	7.8	12.5	16.9	18.7	4.8	21.8	26.3
MV-MLLM	28.6	4.7	4.1	27.5	18.5	30.7	7.2	44.2	15.5	18.9	19.1	23.3	8.2	23.1	32.3
BEV-InMLLM	<b>23.6</b>	<b>3.1</b>	<b>3.8</b>	<b>32.9</b>	<b>19.2</b>	<b>31.5</b>	<b>4.2</b>	<b>46.5</b>	<b>17.3</b>	<b>20.5</b>	<b>21.5</b>	<b>24.5</b>	<b>9.4</b>	<b>26.8</b>	<b>35.6</b>
Video-LLama [45]	29.9	6.5	5.4	22.3	16.7	20.9	9.3	39.3	6.2	10.9	16.2	18.4	4.1	21.3	25.3
MV-MLLM	28.9	6.2	4.4	27.9	19.6	30.9	9.3	44.3	16.5	18.7	19.9	23.0	6.5	26.6	31.4
BEV-InMLLM	<b>24.5</b>	<b>3.5</b>	<b>4.2</b>	<b>31.6</b>	<b>19.0</b>	<b>34.6</b>	<b>4.1</b>	<b>44.7</b>	<b>17.7</b>	<b>22.5</b>	<b>21.4</b>	<b>26.1</b>	<b>8.7</b>	<b>27.9</b>	<b>35.2</b>

Table 3. **Performance comparison with state-of-the-arts on NuInstruct.** Optimal scores are highlighted in bold. Note that all models are fine-tuned on the training set of NuInstruct in the same setting. ‘\*’ means we use the spatiotemporal adapter to enable the image-based MLLM to receive the video input. For clarity, we employ abbreviations to denote the names of subtasks instead of their full designations, *i.e.*, Dis = Distance, Sped = Speeds, # Ins = Instance number, Clos = Closest, Sta = Status, SameR = In the same road, Mot = Motion, App = Approach, LaneC = Lane changing, Onco = On coming, Cro = Crossing, Over = Overtaking, Brak = Braking. Best results are reported in **Bold**.

We initialized all MLLMs with their official pre-trained weights, freezing these during training and only training the parameters of ST-Adapters and our additional modules (MV Q-Former, BEV Q-Former, and injection module). We choose LSS [33] and BEVFormer [21] as our BEV extractors, and  $W$  and  $H$  are both set to 200.  $N_{\text{view}}$ ,  $K_{\text{mv}}$ ,  $K_{\text{bev}}$  are set to 6, 32 and 32 respectively. The dimensions  $D_{\text{vis}}$  and  $D_{\text{resp}}$  are both set to 1408, the same as the dimension of the same as EVA\_CLIP hidden feature dim used by BLIP-2. The input is resized and cropped to the spatial size of  $224 \times 224$ , and each video is uniformly sampled 3 frames. We use AdamW [26] as the optimizer and cosine annealing scheduler [25] as the learning rate scheduler with an initial learning rate of  $1e-4$ , and all models are trained in 20 epochs.

### 5.1. State-of-the-art Comparison

We select three advanced MLLMs, *i.e.*, BLIP-2 [19], MiniGPT-4 [46] and Video-LLama [45] as our base models. For each MLLM, we apply our proposed modules to obtain MV-MLLM and BEV-InMLLM. All models are finetuned in the same setting. We report our results on NuInstruct test set in Table 3. To conserve space, we aggregate the reporting of two subtasks, ‘motion ego’ and ‘motion others’. A similar approach is adopted for ‘status ego’ and ‘status others’. From Table 3, we observe that equipped with our proposed modules, there is a significant increase in the evaluation metrics on all tasks, demonstrating its effectiveness. More specifically, the integration of temporal and multi-view information (MV-MLLM) substantially improves risk and planning tasks by 5% and 6%, respectively. Furthermore, injecting BEV into MV-MLLM, *i.e.*, BEV-InMLLM, benefits tasks sensitive to distance and position, *e.g.*, perception and prediction. For perception tasks, a comparison between ours and existing BEV SOTAs without LLMs can be found in the supplementary material.

Task	Temporal	Multi-view	Spatial	Holistic
SubTask	Sped, Mot, Sta	# Ins, SameR	Dis, Mot, Sta, Clos	Risk, Planning

Table 4. **Reclassified tasks.** To better analyze the impact of each module on autonomous driving tasks, we reclassify the sub-tasks into four main tasks based on their dependency on different types of information. ‘#’ indicates the numbers. ‘Holistic’ means those tasks that require all information, *i.e.*, temporal, multi-view, and spatial.

Model	Temporal		Multi-view		Spatial		Whole $\uparrow$
	$\downarrow$	$\uparrow$	$\downarrow$	$\uparrow$	$\downarrow$	$\uparrow$	
(a) Full	3.7	32.8	3.8	31.5	13.9	32.9	22.2
(b) w/o Video	7.4	28.5	4.2	30.4	14.2	30.8	21.0
$\Delta$	<b>-3.7</b>	<b>-4.3</b>	<b>-0.4</b>	<b>-1.1</b>	<b>-0.3</b>	<b>-1.1</b>	<b>-1.2</b>
(c) w/o MV	5.2	31.2	6.0	28.3	14.4	32.5	21.6
$\Delta$	<b>-1.5</b>	<b>-1.6</b>	<b>-2.2</b>	<b>-3.2</b>	<b>-0.5</b>	<b>-0.4</b>	<b>-0.6</b>
(d) w/o BEV	6.0	31.4	4.1	30.7	18.0	30.0	20.1
$\Delta$	<b>-2.3</b>	<b>-1.4</b>	<b>-0.3</b>	<b>-0.8</b>	<b>-4.1</b>	<b>-2.9</b>	<b>-2.1</b>
(e) Base	10.3	25.8	6.7	22.7	20.8	27.6	12.4
$\Delta$	<b>-6.6</b>	<b>-7.0</b>	<b>-2.9</b>	<b>-8.8</b>	<b>-6.9</b>	<b>-5.3</b>	<b>-9.8</b>

Table 5. **The ablation study of different proposed modules.** ‘w/o’ indicates the without the specific module. ‘Video’, ‘MV’ and ‘BEV’ indicate video input, MV Q-Former (Section 4.2) and BEV injection module (Section 4.3) respectively. The results of performance degradation exceeding 2 are reported in **green**.  $\Delta$  is the difference between a specific model with full model, *i.e.*, line (a).

### 5.2. Ablation Study

In this section, we conduct experiments to evaluate the effect of different proposed modules and different input information. Here, we use MiniGPT-4 [46] as our baseline model. To better analyze the impact of each module on autonomous driving tasks, we reclassify the sub-tasks into four main tasks based on their dependency on different types of information. These are categorized as temporal-related (temporal), multi-view-related (multi-view), spatial-related (Spatial), and holistic-related tasks, as shown in Ta-

extractor	<i>Spatial</i> ↓	<i>Whole</i> ↑	$L_{inst}$	<i>Spatial</i> ↓	<i>Whole</i> ↑	$K_{bev}$	<i>Spatial</i> ↓	<i>Whole</i> ↑	feature	<i>Spatial</i> ↓	<i>Whole</i> ↑
LSS [33]	13.9	21.0	w/o	15.3	20.1	16	14.5	21.0	$F_{mv}$	15.1	20.8
BEVDet [15]	13.6	21.3	w/	<b>13.9</b>	<b>21.5</b>	32	<b>13.9</b>	21.5	$\bar{F}_{mv}$	<b>13.9</b>	<b>21.5</b>
BEVFusion [24] <sup>†</sup>	<b>13.2</b>	<b>21.5</b>				64	<b>13.9</b>	<b>21.6</b>			

(a) **BEV extractor.** We select three advanced BEV extractors. Powerful extractors are more effective.  
<sup>†</sup>: using additional lidar-modal data.

(b) **Instruction tokens  $L_{inst}$ .** ‘w/o’ and ‘w/’ indicate without and with.  $L_{inst}$  (Eq. 3) can extract the instruction-aware BEV.

(c) **BEV query number  $K_{bev}$ .** More numbers for BEV queries  $Q_{bev}$  (Eq. 3) benefit the model.

(d) **Injection feature.** BEV features  $F_{insbev}$  can be injected into  $F_{mv}$  or  $\bar{F}_{mv}$  (Eq. 4). The latter is better due to more information.

Table 6. **BEV Injection module ablation experiments** on NuInstruct. Best results and default settings are reported in **Bold** and **gray**.

ble 4. ‘Temporal’ indicates subtasks related to temporal cues, *e.g.*, the vehicle’s status is determined based on its positions at various times, and so do others. Note that some subtasks may be classified into different tasks, *e.g.*, the status task is in both temporal and spatial tasks. We will report the results of different models under the reclassified tasks in the following.

### 5.2.1 Effect of different proposed modules

In our study, we explore different modules to capture different information for autonomous driving tasks: ST-Adapter accepts videos for temporal, MV-Q Former for multi-view, and BEV Injection module for location, distance, and spatial information in BEV features. We use BEV-InMLLM as a full model including comprehensive information types, then sequentially remove each module to derive the following distinct models: (a) The full model, *i.e.*, BEV-InMLLM introduced in Section 4.3. (b) BEV-InMLLM without temporal cues, *i.e.*, the input is image. (c) BEV-InMLLM without multi-view information, *i.e.*, only single-view input. (d) BEV-InMLLM without BEV information, *i.e.*, without BEV injection module. (e) The baseline model, *i.e.*, MiniGPT-4 [46]. We report the results of different models in Table 5. From the table, we observe the following findings: (i) Compared with (a) and (b), without temporal information, the performance of tasks highly dependent on temporal cues would degrade clearly, proving the importance of video input. We can also observe a similar phenomenon when comparing the results with (a) and (c). (ii) Information contained in BEV is very important for most of autonomous driving tasks, since it clearly presents the surroundings of the ego vehicle, thus aiding the model in making informed decisions.

### 5.2.2 Analysis of BEV Injection Module

We ablate our BEV injection module (BEV-In) (Fig. 5 (b)) using the default settings in Table 6 (see caption). Several intriguing properties are observed.

**BEV extractor.** We compare the performance of different BEV extractors in Table 6a. Our results show that more strong extractor, *e.g.*, BEVDet [15], outperforms the weak one LSS [33]. Furthermore, BEVFusion [24] uses RGB and lidar modality for best performance. Here, we use RGB images for efficiency.

**Instruction tokens  $L_{inst}$  in BEV-In.** Table 6b studies the effect of  $L_{inst}$  in Eq. 3. ‘w/o’ indicate only using  $Q_{bev}$  to in-

teract with  $F_{bev}$ . Results show that using instruction tokens can capture more related BEV features, thus improving the performance by 1.4 for both spatial tasks and holistic tasks.

**BEV query number  $K_{bev}$ .** In Table 6c, we study the influence of BEV query numbers, *i.e.*,  $K_{bev}$  in  $Q_{bev}$  (Eq. 3). As the number increases, the performance would be improved, *e.g.*, 0.6 on the spatial performance with  $K_{bev}$  arise from 16 to 32. Considering setting  $K_{bev}$  to 64 only brings a small improvement, we use 32 as the default settings for computation efficiency.

**Injection feature.** The key design of our BEV-InMLLM is to inject instruction-aware BEV features (*i.e.*,  $F_{insbev}$  in Eq. 4) to the MV-MLLM. In Table 6d, we compare the performance of different features to inject with the  $F_{insbev}$ . Specifically, the multi-view visual semantics  $F_{mv}$  (Section 4.2) and the output of multi-view Q-Former  $\bar{F}_{mv}$  (Eq. 2). We find injecting into  $\bar{F}_{mv}$  achieves better, 0.7% improvement over  $F_{mv}$  on the holistic tasks. The reason is that  $F_{mv}$  is the filtered visual tokens, losing much spatial information.

## 6. Conclusion

In this study, we investigate language-based driving for autonomous drivingtasks. We introduce NuInstruct, featuring 91K multi-view video-instruction-response pairs across 17 subtasks, created via a novel SQL-based method. Our proposed BEV-InMLMM integrates instruction-aware BEV features into MLLMs, enhancing temporal, multi-view, and spatial detail processing. BEV-InMLMM, as a plug-and-play enhancement, boosts MLLM performance on autonomous drivingtasks. Our empirical results on NuInstruct confirm our method’s efficacy.

**Limitations.** The current dataset lacks traffic light information and tasks related to 3D object detection, which we plan to address in future work.

## 7. Acknowledgements

This work is partially supported by the National Natural Science Foundation of China under Grant 62306254, the Hong Kong Innovation and Technology Fund under Grant ITS/030/21, a research grant from the Beijing Institute of Collaborative Innovation (BICI) in collaboration with HKUST under Grant HCIC-004, and a grant from the Research Grants Council of the Hong Kong Special Administrative Region under Grant T45-401/22-N.



## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022. **3**
- [2] Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multi-modal dataset for autonomous driving. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11621–11631, 2020. **2, 3, 4, 5**
- [3] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multi-modal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. **1, 3**
- [4] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org> (accessed 14 April 2023), 2023. **3**
- [5] Wenliang Dai, Junnan Li, Dongxu Li, AnthonyMeng Huat, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. 2023. **1, 5**
- [6] Chris J Date. *A Guide to the SQL Standard*. Addison-Wesley Longman Publishing Co., Inc., 1989. **2**
- [7] Thierry Deruyttere, Simon Vandenhende, Dusan Grujicic, Luc Van Gool, and Marie Francine Moens. Talk2car: Taking control of your self-driving car. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2088–2098, 2019. **2, 3**
- [8] Vikrant Dewangan, Tushar Choudhary, Shivam Chandhok, Shubham Priyadarshan, Anushka Jain, Arun K Singh, Siddharth Srivastava, Krishna Murthy Jatavallabhula, and K Madhava Krishna. Talk2bev: Language-enhanced bird’s-eye view maps for autonomous driving. *arXiv preprint arXiv:2310.02251*, 2023. **2, 3, 4**
- [9] Xinpeng Ding, Jianhua Han, Hang Xu, Wei Zhang, and Xiaomeng Li. Hilm-d: Towards high-resolution understanding in multimodal large language models for autonomous driving. *arXiv preprint arXiv:2309.05186*, 2023. **2, 3, 6**
- [10] Jeffrey EF Friedl. *Mastering regular expressions*. ” O’Reilly Media, Inc.”, 2006. **5**
- [11] Jiahui Gao, Renjie Pi, Jipeng Zhang, Jiacheng Ye, Wanjun Zhong, Yufei Wang, Lanqing Hong, Jianhua Han, Hang Xu, Zhenguo Li, and Lingpeng Kong. G-llava: Solving geometric problem with multi-modal large language model, 2023. **1**
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. **6**
- [13] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. *arXiv preprint arXiv:2307.12981*, 2023. **1, 3**
- [14] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, et al. Planning-oriented autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17853–17862, 2023. **3**
- [15] Junjie Huang, Guan Huang, Zheng Zhu, Yun Ye, and Dalong Du. Bevdet: High-performance multi-camera 3d object detection in bird-eye-view. *arXiv preprint arXiv:2112.11790*, 2021. **2, 6, 8**
- [16] Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. Textual explanations for self-driving vehicles. In *Proceedings of the European conference on computer vision (ECCV)*, pages 563–578, 2018. **2**
- [17] Chunyuan Li, Cliff Wong, Sheng Zhang, Naoto Usuyama, Haotian Liu, Jianwei Yang, Tristan Naumann, Hoifung Poon, and Jianfeng Gao. Llava-med: Training a large language-and-vision assistant for biomedicine in one day. *arXiv preprint arXiv:2306.00890*, 2023. **3**
- [18] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *International Conference on Machine Learning*, pages 12888–12900. PMLR, 2022. **3, 5, 6, 7**
- [19] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. **3, 5, 6, 7**
- [20] KunChang Li, Yanan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023. **1, 3, 5**
- [21] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *European conference on computer vision*, pages 1–18. Springer, 2022. **6, 7**
- [22] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. **1**
- [23] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 10012–10022, 2021. **6**
- [24] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela L Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird’s-eye view representation. In *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pages 2774–2781. IEEE, 2023. **8**
- [25] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. **7**

- [26] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 7
- [27] Yuexin Ma, Tai Wang, Xuyang Bai, Huitong Yang, Yuenan Hou, Yaming Wang, Yu Qiao, Ruigang Yang, Dinesh Manocha, and Xinge Zhu. Vision-centric bev perception: A survey. *arXiv preprint arXiv:2208.02797*, 2022. 6
- [28] Srikanth Malla, Chiho Choi, Isht Dwivedi, Joon Hee Choi, and Jiachen Li. Drama: Joint risk localization and captioning in driving. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1043–1052, 2023. 2, 3
- [29] Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. Gpt-driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023. 2, 3
- [30] OpenAI OpenAI. Gpt-4 technical report. 2023. 2, 3, 4
- [31] Junting Pan, Ziyi Lin, Xiatian Zhu, Jing Shao, and Hongsheng Li ST-Adapter. Parameter-efficient image-to-video transfer learning for action recognition. *Preprint at https://arxiv.org/abs/2206.13559*, 2022. 6
- [32] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 5
- [33] Jonah Philion and Sanja Fidler. Lift, splat, shoot: Encoding images from arbitrary camera rigs by implicitly unprojecting to 3d. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XIV 16*, pages 194–210. Springer, 2020. 6, 7, 8
- [34] Renjie Pi, Jiahui Gao, Shizhe Diao, Rui Pan, Hanze Dong, Jipeng Zhang, Lewei Yao, Jianhua Han, Hang Xu, and Lingpeng Kong Tong Zhang. Detgpt: Detect what you need via reasoning. *arXiv preprint arXiv:2305.14167*, 2023. 1
- [35] Renjie Pi, Lewei Yao, Jiahui Gao, Jipeng Zhang, and Tong Zhang. Perceptiongpt: Effectively fusing visual perception into llm, 2023.
- [36] Renjie Pi, Tianyang Han, Wei Xiong, Jipeng Zhang, Runtao Liu, Rui Pan, and Tong Zhang. Strengthening multi-modal large language model with bootstrapped preference optimization, 2024. 1
- [37] Tianwen Qian, Jingjing Chen, Linhai Zhuo, Yang Jiao, and Yu-Gang Jiang. Nuscenes-qa: A multi-modal visual question answering benchmark for autonomous driving scenario. *arXiv preprint arXiv:2305.14836*, 2023. 2, 3, 6
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [39] Bruce Simons-Morton and Johnathon P. Ehsani. Learning to drive safely: Reasonable expectations and future directions for the learner period. *Safety*, 2(4), 2016. 2
- [40] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023. 3, 5
- [41] ArunBalajee Vasudevan, Dengxin Dai, and LucVan Gool. Object referring in videos with language and human gaze. *Cornell University - arXiv, Cornell University - arXiv*, 2018. 2, 3
- [42] Dongming Wu, Wencheng Han, Tiancai Wang, Xingping Dong, Xiangyu Zhang, and Jianbing Shen. Referring multi-object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14633–14642, 2023. 3
- [43] Dongming Wu, Wencheng Han, Tiancai Wang, Yingfei Liu, Xiangyu Zhang, and Jianbing Shen. Language prompt for autonomous driving. *arXiv preprint arXiv:2309.04379*, 2023. 2
- [44] Zhenhua Xu, Yujia Zhang, Enze Xie, Zhen Zhao, Yong Guo, Kenneth KY Wong, Zhenguo Li, and Hengshuang Zhao. Drivegpt4: Interpretable end-to-end autonomous driving via large language model. *arXiv preprint arXiv:2310.01412*, 2023. 2, 3, 4
- [45] Hang Zhang, Xin Li, Lidong Bing, and at al. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023. 1, 3, 6, 7
- [46] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 1, 3, 5, 6, 7, 8