

# DriveLM: Driving with Graph Visual Question Answering

Anonymous CVPR submission

Paper ID NA

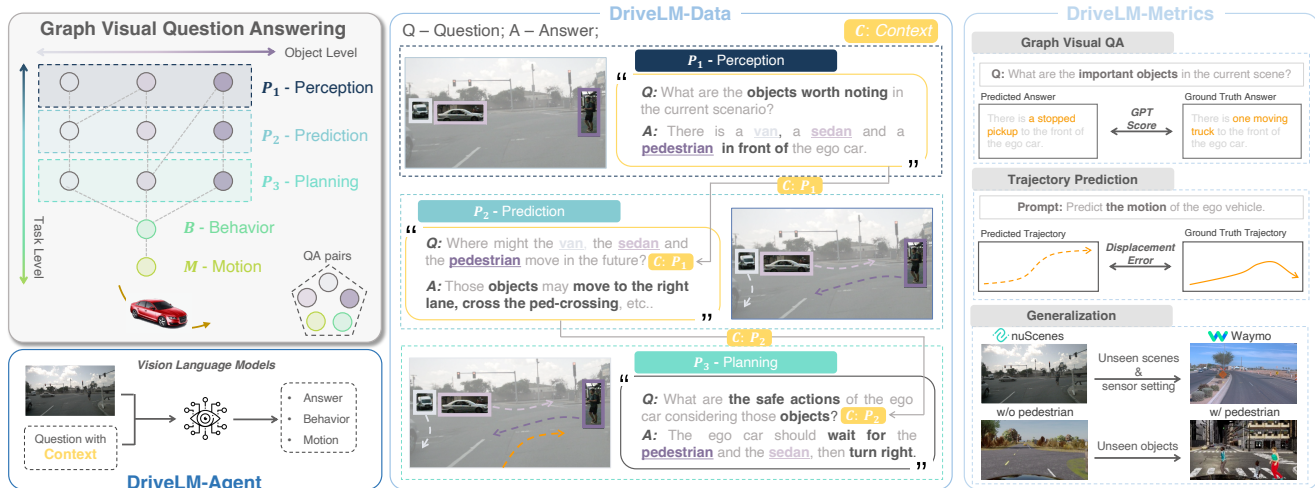


Figure 1. We present **DriveLM**: A new task, dataset, metrics, and baseline for end-to-end autonomous driving. Inspired by [4], DriveLM considers **Graph Visual Question Answering (GVQA)**, where question-answer pairs are interconnected via logical dependencies at the object-level, *i.e.*, interactions between object pairs, and the task-level, *e.g.*, perception  $\rightarrow$  prediction  $\rightarrow$  planning  $\rightarrow$  behavior (discretized action described in natural language)  $\rightarrow$  motion (continuous trajectory). We propose **DriveLM-Data** for training **DriveLM-Agent**, a baseline for GVQA. We validate its effectiveness using the **DriveLM-Metrics** on challenging settings requiring zero-shot generalization.

## Abstract

We study how vision-language models (VLMs) trained on web-scale data can be integrated into end-to-end driving systems to boost generalization and enable interactivity with human users. While recent approaches adapt VLMs to driving via single-round visual question answering (VQA), human drivers reason about decisions in multiple steps. Starting from the localization of key objects, humans estimate object interactions before taking actions. The key insight is that with our proposed task, Graph VQA, where we model graph-structured reasoning through perception, prediction and planning question-answer pairs, we obtain a suitable proxy task to mimic the human reasoning process. We instantiate datasets (DriveLM-Data) built upon nuScenes and CARLA, and propose a VLM-based baseline approach (DriveLM-Agent) for jointly performing Graph VQA and end-to-end driving. The experiments

demonstrate that Graph VQA provides a simple, principled framework for reasoning about a driving scene, and DriveLM-Data provides a challenging benchmark for this task. Our DriveLM-Agent baseline performs end-to-end autonomous driving competitively in comparison to state-of-the-art driving-specific architectures. Notably, its benefits are pronounced when it is evaluated zero-shot on unseen objects or sensor configurations. We hope this work can be the starting point to shed new light on how to apply VLMs for autonomous driving. To facilitate future research, all code, data, and models are available to the public.

## 1. Introduction

Current Autonomous Driving (AD) stacks are still lacking crucial capabilities [4, 5]. One key requirement is generalization, which involves the ability to handle unseen scenarios or unfamiliar objects. A secondary requirement pertains to the interaction of these models with humans, highlighted for example by EU regulations that mandate explainability

\*Equal contribution. †Equal co-advising.

036 in deployment [1]. Furthermore, unlike today’s AD mod- 089  
037 els, humans do not navigate based on geometrically precise 090  
038 bird’s-eye view (BEV) representations [6, 13, 16]. Instead, 091  
039 humans implicitly perform object-centric perception, pre-  
040 diction, and planning (which we refer to as  $P_{1-3}$ ): a rough  
041 identification and localization of key objects, followed by  
042 reasoning about their possible movement and aggregation  
043 of this information into a driving action [22, 27].

044 Simultaneously, another field has been forging ahead:  
045 Vision-Language Models (VLMs) [17, 19, 30, 34]. These  
046 models have several strengths. First, they hold a base un-  
047 derstanding of the world from internet-scale data that could  
048 potentially facilitate generalization for planning in AD. In  
049 fact, this sort of generalization has already been achieved  
050 by VLMs for simpler robotics tasks [9, 35]. Second, the use  
051 of language representations as an input and output offers a  
052 platform for human-friendly interaction with these models,  
053 unlike bounding boxes or trajectories that are more common  
054 to current methods [7, 12, 18, 25]. Finally, VLMs are able  
055 to make decisions in multiple steps linked by logical reason-  
056 ing [2, 8, 31–33, 35]. Importantly, even though they reason  
057 in multiple separate steps, VLMs are end-to-end differen-  
058 tiable architectures, a characteristic that is highly desirable  
059 for autonomous driving [4].

060 Recent work towards enabling the application of VLMs  
061 to AD systems falls into two categories: scene-level or sin-  
062 gle object-level Visual Question Answering (VQA). Scene-  
063 level VQA refers to the task of describing the driving be-  
064 havior by one or two supporting reasons, *e.g.*, “The car  
065 is moving into the right lane because it is safe to do  
066 so.” [14, 15]. Single object-level VQA formulates the un-  
067 derstanding of the ego vehicle’s response to a single ob-  
068 ject by a chain of QAs in the form of “what-which-where-  
069 how-why”, *e.g.*, “The ego vehicle stops because there is a  
070 pedestrian in a white shirt crossing the intersection in front  
071 of the ego vehicle and it does not want to crash into the  
072 pedestrian.” [21, 24, 26]. Unfortunately, neither of these  
073 paradigms provides a suitable proxy task to mimic the  $P_{1-3}$   
074 reasoning process in humans, who consider multiple objects  
075 and reason about each in multiple steps. Therefore, in this  
076 paper, we propose a new task, along with corresponding  
077 datasets and a baseline model architecture (Fig. 1).

078 **Task. Graph Visual Question Answering (GVQA)** in-  
079 volves formulating  $P_{1-3}$  reasoning as a series of question-  
080 answer pairs (QAs) in a directed graph. Its key differ-  
081 ence to the aforementioned VQA tasks for AD is the avail-  
082 ability of logical dependencies between QAs which can be  
083 used to guide the answering process. GVQA also encom-  
084 passes questions regarding behavior and motion planning,  
085 with dedicated metrics (details in Section 2).

086 **Datasets. DriveLM-nuScenes** consist of annotated QAs,  
087 arranged in a graph, linking images with driving behavior  
088 through logical reasoning. In comparison to existing bench-

marks, they provide significantly more text annotations per  
frame (Fig. 2). We pair these training datasets with chal-  
lenging test data for evaluating zero-shot generalization.

**Model. DriveLM-Agent** employs a trajectory tokenizer  
that can be applied to any general VLM [17, 19, 23, 34],  
coupled with a graph prompting scheme that models logi-  
cal dependencies as context inputs for VLMs. The result  
is a simple, elegant methodology to effectively repurpose  
VLMs for end-to-end AD.

Our experiments provide encouraging results. We find  
that GVQA on DriveLM is a challenging task, where cur-  
rent methods obtain moderate scores and better model-  
ing of logical dependencies is likely necessary to achieve  
strong QA performance. Even so, DriveLM-Agent already  
performs competitively to state-of-the-art driving-specific  
models [13] when tested in the open-loop planning setting,  
despite its task-agnostic and generalist architecture. Fur-  
thermore, employing a graph structure improves zero-shot  
generalization, enabling DriveLM-Agent to correctly han-  
dle novel objects unseen during training or deployment on  
the Waymo dataset [28] after training only on nuScenes [3]  
data. From these results, we believe that improving GVQA  
holds great potential towards building autonomous driving  
agents with strong generalization.

## 2. DriveLM: Task, Data, Metrics

Human drivers usually decompose their decision-making  
process into distinct stages that follow a logical progres-  
sion which encompasses the identification and localization  
of key objects, their possible future action and interaction,  
and ego planning based on all this information [10, 20].  
This inspires us to propose the GVQA as the critical ingre-  
dient of DriveLM, which serves as a suitable proxy task to  
mimic the human reasoning process. Within this section, we  
illustrate the formulation of the GVQA task (Section 2.1)  
and introduce DriveLM-Data (Section 2.2) to exemplify the  
instantiation of GVQA using prominent driving datasets.

### 2.1. DriveLM-Task: GVQA

We organize all the Question Answer pairs (QAs) for an im-  
age frame into a graph structure, denoted by  $G = (V, E)$ .  $V$   
stands for the set of vertices, where each vertex represents a  
QA pair  $v = (q, a)$  associated with one or more key objects  
in the scenario. The key difference between GVQA and  
ordinary VQA is that the QAs in GVQA have logical de-  
pendencies, which we formulate as the edges between the  
vertices.  $E \subseteq V \times V$ , is a set of directed edges, where each  
edge  $e = (v_p, v_c)$  connects the parent QA and the child QA.  
We formulate the edge set  $E$  by incorporating two dimen-  
sions: object-level and task-level edges. At the object level,  
we construct the logical edges  $e \in E$  to represent the impact  
of interactions between different objects. For example, the

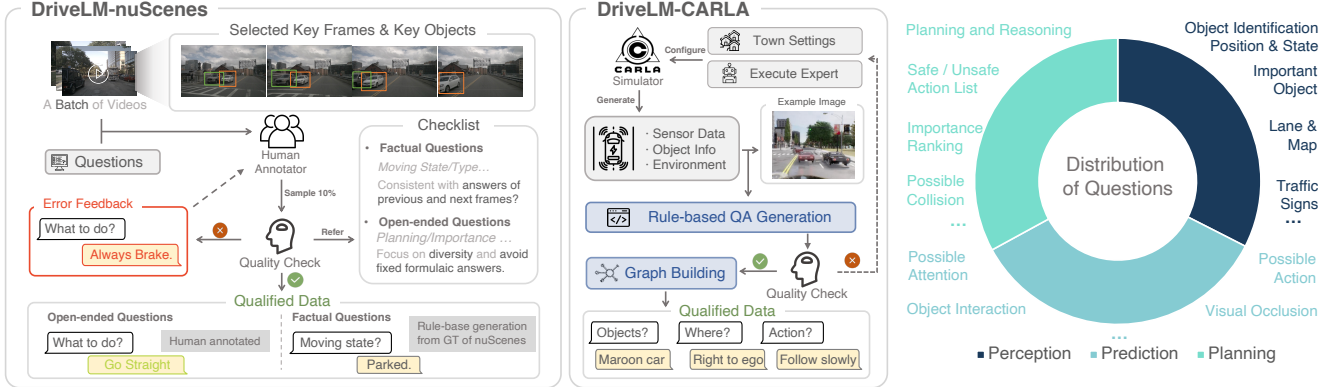


Figure 2. **(Left) Annotation Pipeline:** In DriveLM-nuScenes, we adopt a semi-rule-based QA labeling pipeline, where both the ground truth annotation in nuScenes/OpenLane-V2 and feedback from human annotators are used. A critical part of our pipeline is the multi-round quality check, which guarantees high data quality at reasonable costs. In DriveLM-CARLA, we meet the same standards while exploiting a fully rule-based QA labeling pipeline instead. **(Right) Question Distribution:** The questions in our dataset cover various specific aspects of driving tasks, most of which are annotated by human annotators, making this a suitable proxy for human-like driving reasoning.

139 planning QA node for the sedan is influenced by the per- 171  
 140 ception QA node of the pedestrian in the illustration from 172  
 141 Fig. 1 (center). At the task-level, we establish the logical 173  
 142 edges  $e \in E$  to capture the logical chain of different reason- 174  
 143 ing stages: 175

- 144 • **Perception** ( $P_1$ ): identification, description, and localiza- 176  
 145 tion of key objects in the current scene. 177
- 146 • **Prediction** ( $P_2$ ): estimation of possible action/interaction 178  
 147 of key objects based on perception results. 179
- 148 • **Planning** ( $P_3$ ): possible safe actions of the ego vehicle. 180
- 149 • **Behavior** ( $B$ ): classification of driving decision. 181
- 150 • **Motion** ( $M$ ): waypoints of ego vehicle future trajectory. 182

151 The concepts of perception, prediction, and planning 183  
 152 ( $P_{1-3}$ ) are similar to those in end-to-end AD [4], while the 184  
 153 concepts of motion and behavior are based on the ego ve- 185  
 154 hicle future trajectory. Specifically, we define the motion 186  
 155  $M$  as the ego vehicle future trajectory, which is a set of  $N$  187  
 156 points with coordinates  $(x, y)$  in bird’s-eye view (BEV), de- 188  
 157 noted as  $M = \{(x_0, y_0), (x_1, y_1), \dots, (x_N, y_N)\}$ . Each point 189  
 158 is the offset from the future position and the current po- 190  
 159 sition by a fixed time interval. Then, the distance for  $x, y$  at 191  
 160 each time interval is computed as: 192

$$161 \quad \{x, y\}_{\text{dist}} = \{(\delta_{x,1}, \delta_{y,1}), \dots, (\delta_{x,N}, \delta_{y,N})\}, \quad (1)$$

162 where  $\delta_{x,i} = x_i - x_{i-1}$  and  $\delta_{y,i} = y_i - y_{i-1}$ , for  $i =$  194  
 163  $1, 2, \dots, N$ . The goal of the behavior representation is to 195  
 164 serve as an interface from  $P_{1-3}$  to  $M$ . To obtain a behavior 196  
 165 representation, we map the mean of  $x_{\text{dist}}$  and  $y_{\text{dist}}$  to one of 197  
 166 the predefined bins, where each bin corresponds to a cate- 198  
 167 gory in either speed or steering. These are denoted as  $B_{sp}$  199  
 168 and  $B_{st}$  respectively. In this work, we consider 5 bins: 200

$$169 \quad B_{sp} \in \{\text{fast}_2, \text{fast}_1, \text{moderate}, \text{slow}_1, \text{slow}_2\},$$

$$170 \quad B_{st} \in \{\text{left}_2, \text{left}_1, \text{straight}, \text{right}_1, \text{right}_2\},$$

where the number in the subscript indicates the intensity. 171  
 The combination of the speed and steering categories for 172  
 a trajectory form its behavior category as  $B = (B_{sp}, B_{st})$ . 173  
 While we use a simple definition of  $B$  as a starting point for 174  
 research on driving with VLMs, we note that our formula- 175  
 tion supports the incorporation of more abstract behaviors 176  
 such as a lane changes or overtaking. 177

## 2.2. DriveLM-Data 178

We introduce DriveLM-nuScenes to provide QAs with the 179  
 graph structure defined in Section 2.1, 180

**DriveLM-nuScenes.** We divide the annotation process 181  
 into three steps: selecting key frames from video clips, 182  
 choosing key objects within these key frames, and subse- 183  
 quently annotating the frame-level  $P_{1-3}$  QAs for these key 184  
 objects. A portion of the Perception QAs are generated 185  
 from the nuScenes [3] and OpenLane-V2 [29] ground truth, 186  
 while the remaining QAs are manually annotated. As we 187  
 manually annotate the vast majority of data in DriveLM- 188  
 nuScenes, quality is particularly crucial for this portion. 189  
 When annotating, we conduct multiple rounds of rigorous 190  
 quality checks. In each round, we categorize the data into 191  
 different batches and inspect ten percent of the data in each 192  
 batch. If the qualification rate of manually annotated data in 193  
 this ten percent does not meet expectations, we request the 194  
 annotators to re-label all data in the batch. In Fig. 2 (left), 195  
 we showcase an example of the QA annotation pipeline, 196  
 where all questions undergo quality checks according to our 197  
 standards. As a result, DriveLM-nuScenes stands out from 198  
 previously proposed datasets with its larger scale, greater 199  
 comprehensiveness, and more complex structure. These 200  
 QAs cover various aspects of the driving process, rang- 201  
 ing from perception and prediction to planning, providing 202  
 a comprehensive understanding of autonomous driving sce- 203  
 narios as shown in Fig. 2 (right). 204

### 205 3. Experiments

206 In this section, we present our experimental results that  
207 aim to address the following research questions: (1) How  
208 can VLMs be effectively repurposed for end-to-end au-  
209 tonomous driving? (2) Can VLMs for driving generalize  
210 when evaluated with unseen sensor setups;

211 **Setup.** We now briefly overview the key implementa-  
212 tion details for the two settings used in our experiments  
213 (additional details are provided in the supplementary ma-  
214 terial). All fine-tuning is implemented with LoRA [11].  
215 On DriveLM-nuScenes, we finetune BLIP-2 on the train  
216 split for 10 epochs. We use a batch size of 2 for each GPU,  
217 and the entire training process spans approximately 7 hours  
218 with 8 V100 GPUs.

#### 219 3.1. VLMs for End-to-End Driving

220 In our first experiment, we aim to assess the ability of VLMs  
221 to perform open-loop planning on DriveLM-nuScenes. In  
222 particular, we investigate the impact of the context provided  
223 to the behavior and motion stages. Given sensor data (and  
224 in the case of VLM methods, a text input), the model is  
225 required to predict the ego-vehicle future trajectory in the  
226 form of waypoints.

227 **Baselines.** As a reference for the difficulty of the task,  
228 we provide a simple **Command Mean** baseline. Each  
229 frame in nuScenes is associated with one of 3 commands,  
230 ‘turn left’, ‘turn right’, or ‘go straight’. We output the  
231 mean of all trajectories in the training set whose com-  
232 mand matches the current test frame command. Further,  
233 we compare our approach to the current state-of-the-art on  
234 nuScenes, UniAD [13]. Besides the author-released check-  
235 point, which requires video inputs, we train a single-frame  
236 version (‘UniAD-Single’) for a fair comparison to our  
237 single-frame VLMs. Finally, **BLIP-RT-2** denotes BLIP-  
238 2 [17] fine-tuned on DriveLM-Data with the trajectory to-  
239 kenization scheme. This acts as an indicator for the per-  
240 formance when using an identical network architecture as  
241 DriveLM-Agent, but no context inputs or VQA training  
242 data.

243 **DriveLM-Agent.** We consider 3 variants of DriveLM-  
244 Agent incorporating our proposed changes in steps: (1) a  
245 2-stage version that predicts behavior and then motion (as  
246 described in Section 2.1), but without any  $P_{1-3}$  context  
247 for behavior prediction (‘None’); (2) a ‘Chain’ version that  
248 builds the  $P_{1-3}$  graph, but only passes the final node ( $P_3$ )  
249 to the behavior stage; (3) the full model (‘Graph’) that uses  
250 all QAs from  $P_{1-3}$  as context for  $B$ .

251 **Results.** We show the results for the methods listed above  
252 in Table 1. Among the baselines, BLIP-RT-2 is unable to  
253 match UniAD-Single (though both methods perform well  
254 relative to Command Mean). This shows that the single-  
255 stage approach without any reasoning is unable to compete

Method	Behavior Context	Motion Context	Behavior ( $B$ )			Motion ( $M$ )	
			Acc. $\uparrow$	Speed $\uparrow$	Steer $\uparrow$	ADE $\downarrow$	Col. $\downarrow$
Command Mean	-	-	-	-	-	4.57	5.72
UniAD-Single	-	-	-	-	-	1.80	2.62
BLIP-RT-2	-	-	-	-	-	2.63	2.77
DriveLM-Agent	None	$B$	<b>61.45</b>	<b>72.20</b>	<b>84.73</b>	<b>1.39</b>	<b>1.67</b>
	Chain	$B$	50.43	60.32	75.34	2.07	2.08
	Graph	$B$	57.49	69.89	80.63	1.74	1.89
UniAD [13]	-	-	-	-	-	0.80	0.17

Table 1. **Open-loop Planning on DriveLM-nuScenes.** Using Behavior ( $B$ ) as context for Motion ( $M$ ) enables end-to-end driving with VLMs on par with UniAD-Single, a state-of-the-art driving-specific architecture.

Method	Behavior Context	Motion Context	Behavior ( $B$ )			Motion ( $M$ )	
			Acc. $\uparrow$	Speed $\uparrow$	Steer $\uparrow$	ADE $\downarrow$	FDE $\downarrow$
Command Mean	-	-	-	-	-	7.98	11.41
UniAD-Single	-	-	-	-	-	4.16	9.31
BLIP-RT-2	-	-	-	-	-	2.78	6.47
DriveLM-Agent	None	$B$	35.70	43.90	65.20	2.76	6.59
	Chain	$B$	34.62	41.28	64.55	2.85	6.89
	Graph	$B$	<b>39.73</b>	<b>54.29</b>	<b>70.35</b>	<b>2.63</b>	<b>6.17</b>

Table 2. **Zero-shot Generalization across Sensor Configurations.** Results on 1k randomly sampled frames from the Waymo val set after training on DriveLM-nuScenes. DriveLM-Agent outperforms UniAD-Single and benefits from graph context.

256 with the prior state-of-the-art on nuScenes. However, the  
257 proposed DriveLM-Agent, which predicts behavior as an  
258 intermediate step for motion, provides a significant boost  
259 in performance, surpassing UniAD-Single. This indicates  
260 that with the appropriate prompting, VLMs can be surpris-  
261 ingly competitive for end-to-end driving. Interestingly, in  
262 the experimental setting of Table 1 which does not involve  
263 generalization, the Chain and Graph versions of DriveLM-  
264 Agent do not provide any further advantage over no con-  
265 text. Further, single-frame VLMs fall short in comparison  
266 to the privileged video-based UniAD model, indicating that  
267 VLMs with video inputs may be necessary for this task.

#### 268 3.2. Generalization Across Sensor Configurations

269 As a more challenging setting for evaluating the models  
270 from Section 3.1, we now apply them without any fur-  
271 ther training to a new domain: the Waymo dataset [28].  
272 Waymo’s sensor setup does not include a rear camera, so  
273 we drop this input from UniAD-Single. The VLM methods  
274 only use the front view and do not require any adaptation.

275 **Results.** As shown in Table 2, UniAD-Single does not cope  
276 well with the new sensor configuration, and drops below  
277 BLIP-RT-2 in performance. The multi-stage approach of  
278 DriveLM-Agent provides further improvements. In partic-  
279 ular, the accuracy of speed predictions rises from 43.90 with  
280 no context to 54.29 with the full graph. On the other hand,  
281 the chain approach does not provide sufficient useful infor-  
282 mation, with a speed accuracy of only 41.28.

## References

- 283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339
- [1] Shahin Atakishiyev, Mohammad Salameh, Housam Babiker, and Randy Goebel. Explaining autonomous driving actions with visual question answering. *arXiv preprint arXiv:2307.10408*, 2023. 2
- [2] Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Michal Podstawski, Hubert Niewiadomski, et al. Graph of Thoughts: Solving Elaborate Problems with Large Language Models. *arXiv preprint arXiv:2308.09687*, 2023. 2
- [3] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yuxin Pan, Giancarlo Baldan, and Oscar Beijbom. nuScenes: A multi-modal dataset for autonomous driving. In *CVPR*, 2020. 2, 3
- [4] Li Chen, Penghao Wu, Kashyap Chitta, Bernhard Jaeger, Andreas Geiger, and Hongyang Li. End-to-end autonomous driving: Challenges and frontiers. *arXiv preprint arXiv:2306.16927*, 2023. 1, 2, 3
- [5] Pranav Singh Chib and Pravendra Singh. Recent advancements in end-to-end autonomous driving using deep learning: A survey. *IEEE T-IV*, 2023. 1
- [6] Kashyap Chitta, Aditya Prakash, Bernhard Jaeger, Zehao Yu, Katrin Renz, , and Andreas Geiger. Transfuser: Imitation with transformer-based sensor fusion for autonomous driving. *IEEE T-PAMI*, 2023. 2
- [7] Daniel Dauner, Marcel Hallgarten, Andreas Geiger, and Kashyap Chitta. Parting with misconceptions about learning-based vehicle motion planning. In *CoRL*, 2023. 2
- [8] Ruomeng Ding, Chaoyun Zhang, Lu Wang, Yong Xu, Minghua Ma, Wei Zhang, Si Qin, Saravan Rajmohan, Qingwei Lin, and Dongmei Zhang. Everything of thoughts: Defying the law of penrose triangle for thought generation. *arXiv preprint arXiv:2311.04254*, 2023. 2
- [9] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, et al. PaLM-E: An embodied multimodal language model. In *ICML*, 2023. 2
- [10] John A Groeger. *Understanding driving: Applying cognitive psychology to a complex everyday task*. Routledge, 2013. 2
- [11] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *CoRL*, 2021. 4
- [12] Shengchao Hu, Li Chen, Penghao Wu, Hongyang Li, Junchi Yan, and Dacheng Tao. St-p3: End-to-end vision-based autonomous driving via spatial-temporal feature learning. In *ECCV*, 2022. 2
- [13] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, et al. Planning-oriented autonomous driving. In *CVPR*, 2023. 2, 4
- [14] Jinkyu Kim, Anna Rohrbach, Trevor Darrell, John Canny, and Zeynep Akata. Textual explanations for self-driving vehicles. In *ECCV*, 2018. 2
- [15] Jinkyu Kim, Teruhisa Misu, Yi-Ting Chen, Ashish Tawari, and John Canny. Grounding human-to-vehicle advice for self-driving vehicles. In *CVPR*, 2019. 2
- [16] Hongyang Li, Chonghao Sima, Jifeng Dai, Wenhai Wang, Lewei Lu, Huijie Wang, Jia Zeng, Zhiqi Li, Jiazhi Yang, Hanming Deng, et al. Delving into the devils of bird’s-eye-view perception: A review, evaluation and recipe. *IEEE T-PAMI*, 2023. 2 340 341 342 343 344
- [17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *ICML*, 2023. 2, 4 345 346 347 348
- [18] 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 *ECCV*, 2022. 2 349 350 351 352
- [19] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *NeurIPS*, 2023. 2 353 354
- [20] Charles C Macadam. Understanding and modeling the human driver. *Veh. Syst. Dyn.*, 2003. 2 355 356
- [21] Srikanth Malla, Chiho Choi, Isht Dwivedi, Joon Hee Choi, and Jiachen Li. DRAMA: Joint risk localization and captioning in driving. In *WACV*, 2023. 2 357 358 359
- [22] David Marr. *Vision: A computational investigation into the human representation and processing of visual information*. The MIT Press, 2010. 2 360 361 362
- [23] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023. 2 363 364 365 366
- [24] 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. 2 367 368 369 370
- [25] Katrin Renz, Kashyap Chitta, Otniel-Bogdan Mercea, Almut Sophia Koepke, Zeynep Akata, and Andreas Geiger. Plant: Explainable planning transformers via object-level representations. In *CoRL*, 2022. 2 371 372 373 374
- [26] Enna Sachdeva, Nakul Agarwal, Suhas Chundi, Sean Roelofs, Jiachen Li, Behzad Dariush, Chiho Choi, and Mykel Kochenderfer. Rank2Tell: A multimodal driving dataset for joint importance ranking and reasoning. *arXiv preprint arXiv:2309.06597*, 2023. 2 375 376 377 378 379
- [27] Elizabeth S Spelke and Katherine D Kinzler. Core knowledge. *Dev Sci*, 2007. 2 380 381
- [28] Pei Sun, Henrik Kretschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *CVPR*, 2020. 2, 4 382 383 384 385 386
- [29] Huijie Wang, Tianyu Li, Yang Li, Li Chen, Chonghao Sima, Zhenbo Liu, Bangjun Wang, Peijin Jia, Yuting Wang, Shengyin Jiang, et al. OpenLane-V2: A topology reasoning benchmark for unified 3d HD mapping. In *NeurIPS Datasets and Benchmarks*, 2023. 3 387 388 389 390 391
- [30] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, and Jifeng Dai. VisionLLM: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv preprint arXiv:2305.11175*, 2023. 2 392 393 394 395 396

- 397 [31] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed  
398 Chi, Sharan Narang, Aakanksha Chowdhery, and Denny  
399 Zhou. Self-Consistency improves chain of thought reason-  
400 ing in language models. In *ICLR*, 2023. 2
- 401 [32] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten  
402 Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny  
403 Zhou. Chain-of-thought prompting elicits reasoning in large  
404 language models. In *NeurIPS*, 2022.
- 405 [33] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,  
406 Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan.  
407 Tree of Thoughts: Deliberate problem solving with large lan-  
408 guage models. *arXiv preprint arXiv:2305.10601*, 2023. 2
- 409 [34] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu,  
410 Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao.  
411 LLaMA-Adapter: Efficient fine-tuning of language models  
412 with zero-init attention. *arXiv preprint arXiv:2303.16199*,  
413 2023. 2
- 414 [35] Brianna Zitkovich, Tianhe Yu, Sichun Xu, Peng Xu, Ted  
415 Xiao, Fei Xia, Jialin Wu, Paul Wohlhart, Stefan Welker, et al.  
416 RT-2: Vision-language-action models transfer web knowl-  
417 edge to robotic control. In *CoRL*, 2023. 2