# LMDrive: Closed-Loop End-to-End Driving with Large Language Models

Hao Shao[1,2]    Yuxuan Hu[3]    Letian Wang[4]
Steven L. Waslander[4]    Yu Liu[2,5 ✉]    Hongsheng Li[1,3,5 ✉]

[1]CUHK MMLab    [2]SenseTime Research    [3]CPII under InnoHK
[4]University of Toronto    [5]Shanghai Artificial Intelligence Laboratory

## Abstract

*Despite significant recent progress in the field of autonomous driving, modern methods still struggle and can incur serious accidents when encountering long-tail unforeseen events and challenging urban scenarios. On the one hand, large language models (LLM) have shown impressive reasoning capabilities that approach "Artificial General Intelligence". On the other hand, previous autonomous driving methods tend to rely on limited-format inputs (e.g. sensor data and navigation waypoints), restricting the vehicle's ability to understand language information and interact with humans. To this end, this paper introduces LMDrive, a novel language-guided, end-to-end, closed-loop autonomous driving framework. LMDrive uniquely processes and integrates multi-modal sensor data with natural language instructions, enabling interaction with humans and navigation software in realistic instructional settings. To facilitate further research in language-based closed-loop autonomous driving, we also publicly release the corresponding dataset which includes approximately 64K instruction-following data clips, and the LangAuto benchmark that tests the system's ability to handle complex instructions and challenging driving scenarios. Extensive closed-loop experiments are conducted to demonstrate LMDrive's effectiveness. To the best of our knowledge, we're the very first work to leverage LLMs for closed-loop end-to-end autonomous driving. Codes can be found at our webpage.*

## 1. Introduction

Remarkable progress in autonomous driving has been witnessed in recent years with an increasing number of commercial autonomous vehicles (AVs) deployed on public roads. Generally, state-of-the-art autonomous driving systems can be categorized into two primary approaches: 1) a modular approach where the system is decomposed into several sub-modules such as perception, prediction, and
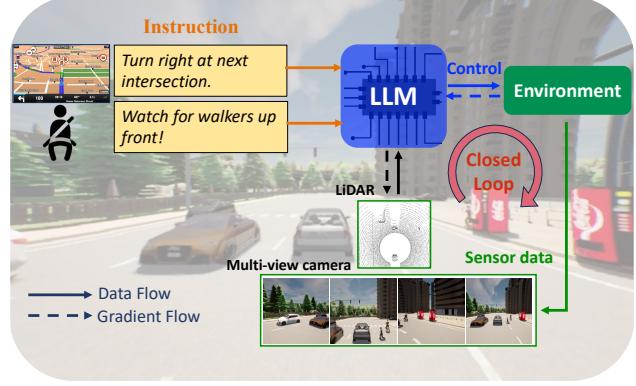


Figure 1. We present LMDrive, the first language-guided closed-loop end-to-end driving framework. LMDrive takes as input the language instruction and multi-modal multi-view sensor data, and outputs control signals in real-time to drive in complex scenarios.

planning, and fixed interfaces are designed to integrate them together [2, 23]; and 2) an end-to-end approach that directly converts sensor data to control signals via a neural network [16, 35]. While both of these approaches are widely adopted and constantly making breakthroughs on challenging benchmarks, both of them share a limitation in that they solely rely on fixed-format inputs such as the sensor data, target waypoints, and action commands, which restricts the agent's ability to comprehend multi-modal information and to interact with humans and the environment. On the other hand, large language models (LLMs) have shown an impressive range of capabilities that approach "Artificial General Intelligence." This encompasses language comprehension, knowledge retrieval, and reasoning. Such capabilities could greatly enhance the safety, controllability, and explanability of autonomous agents. In this work, we seek to answer the question for the first time: *Can we build cognitive autonomous driving systems on top of LLMs, that can interact with human passengers or navigation software simply by natural language?*

Making autonomous systems understand natural language opens profound opportunities for advanced reasoning in complex scenarios and efficient interaction with humans,

✉ Corresponding author.

addressing many previously non-trivial problems. To name a few: 1) in long-tail unforeseen events and challenging urban situations (*e.g.* complex and dense intersections) where modern AV systems typically struggle [41] or even incur serious accidents [37], the language-aware AVs can easily survive by following navigation instructions from passengers or navigation software. 2) AVs can adapt to passengers' sudden notice (*e.g.* small objects that are easily missed by perception systems) simply via natural language, which was previously non-trial and required a large amount of hand-crafted rules.

Toward these appealing properties, many pioneering works have explored the potential of using large language models to enhance the AV system's reasoning abilities, interpretability, and overall performance in open-loop settings. One of the most common strategies [6, 11, 28, 32] is to 1) first use LLMs to transform the scene perception results and navigation commands into textual descriptions; 2) feed these textural descriptions into LLMs to generate textual driving decisions; and then 3) transfer textual driving decisions into executable control commands. While good preliminary results are shown, this type of approach, where different LLMs tackle sub-tasks individually, is hard to be trained in an end-to-end manner, loses the capability to scale with a large amount of data, and is not robust to perception errors and uncertainties. For example, since the LLMs in the latter two stages do not have access to the sensor data, inaccurate or missed detections in the first stage can lead to large accumulative errors in the latter stages. Towards addressing these issues, end-to-end language-based driving methods [46] have been proposed. However, all of these methods undergo training and evaluation in the open-loop setting, where actions are generated and evaluated against the expert actions, but not executed in the actual environments. Notably, when executing navigation instructions such as "turn right", the AV agent should not only generate a sequence of actions, but also consider the changes that the actions bring to the environment. The absence of closed-loop evaluation leads to insufficient consideration of critical issues such as cumulative errors, human-robot interaction, and temporal consistency of actions, which makes the resulting methods difficult to scale beyond a short time horizon making them ineffective in actual systems. To the best of our knowledge, there is no existing paper that leverages LLMs for closed-loop end-to-end autonomous driving.

In this work, we introduce LMDrive, an instruction-following multi-modal LLM model for end-to-end closed-loop autonomous driving. The proposed model can process camera-LiDAR sensor data, comprehend driving instructions in natural language, and directly generate vehicle control signals. A pre-trained LLM model is adopted and kept frozen to maintain its reasoning capability. To adapt the LLM for autonomous driving, multiple camera-LiDAR data encoders and learnable input/output adapters

are integrated. A pre-training strategy specifically designed for the driving task is also introduced for the multi-modal vision encoder. To facilitate the training of LMDrive in a closed-loop setting, we develop a language-guided driving dataset based on the CARLA simulator [12], which simulates the dynamic world and realistic challenging scenarios. To better test the driving model's capability in realistic conversation scenarios, where the instructions might come from humans and navigation software, we 1) consider both navigation instruction and notice instructions; 2) diversify instructions into various phrases via prompt engineering on LLMs; 3) incorporate misleading and unreasonable instructions which are infeasible due to safety concerns or traffic regulations; 4) extend the navigation instructions to include multiple consecutive segments, such as "Take a left here, then another left at the next one". Besides, we provide the corresponding LangAuto evaluation benchmark and the pre-trained LMDrive model for reproducibility, and we hope these can facilitate further research in end-to-end closed-loop language-based autonomous driving.

To summarize, this paper makes the following contributions:

- We propose a novel end-to-end, closed-loop, language-based autonomous driving framework, LMDrive, which interacts with the dynamic environment via multi-modal multi-view sensor data and natural language instructions.
- We provide a dataset with about 64K data clips, where each clip includes one navigation instruction, several notice instructions, a sequence of multi-modal multi-view sensor data, and control signals. The duration of the clip spans from 2 to 20 seconds.
- We present the benchmark LangAuto for evaluating the autonomous agents that take language instructions as navigation inputs, which include misleading/long instructions and challenging adversarial driving scenarios.
- We conduct extensive closed-loop experiments to demonstrate the effectiveness of the proposed framework, and analyze different components of LMDrive to shed light on continuing research along this direction.

## 2. Related Works

### 2.1. End-to-End Autonomous Driving

Much progress [7, 8] has been achieved recently in the field of end-to-end autonomous driving. UniAD [16] devised a framework that incorporates full-stack driving tasks and utilizes query-unified interfaces to communicate between different tasks. ThinkTwice [20] designed a Look Module to retrieve information from critical regions and utilize the features to refine the coarse prediction. ReasonNet [35] exploited temporal and global information of the driving scene to improve perception performance and benefit occlusion detection. InterFuser [34] proposed a transformer-based framework to fully fuse and process information from

multi-modal multi-view sensors for comprehensive scene understanding. TCP [45] proposed an approach that integrates the two branches for trajectory planning and direct control by involving a novel multi-step prediction. LAV [3] introduced some supervisory tasks to learn an viewpoint-invariant representation which can provide a richer supervision signal at training and more information for complex reasoning during inference. Beyond the previously discussed imitation training methods, several approaches have sought to incorporate reinforcement learning strategies. Latent DRL [38] trained in a supervised way to get a latent representation of the environment observation and conducts reinforcement learning using the representation as the input. Roach [48] employed a reinforcement learning agent with privileged access to environmental information and distills a model as the final agent. ASAPRL [44] and TaEcRL [50] exploited abstracted skills to effectively improve reinforcement learning efficiency and the final performance by facilitating effective exploration and reward signaling. However, these end-to-end methods lack the ability to verbally or textually interact with humans (passengers), and usually have low explanatility in the decision-making process.

## 2.2. LLMs in Driving Tasks

Emerging advancements in large language models (LLMs) [10, 17, 19, 39, 40] have been witnessed over the last few months. Furthermore, vision large language models (VLLM) further introduce vision encoders and open the doors for LLMs to interpret not only textual data but also images and data in other modalities [5, 25, 51]. In the field of autonomous driving (AD), recent research has integrated LLMs into the AD system for better explainability and natural interaction with humans. Some studies adopt the visual language model approach, which can handle multi-modal input data and provide the textual description as well as the control signal for the driving scenarios. For example, DRIVEGPT4 [46] proposed a multimodal LLM framework, which takes a sequence of frames as input, then generates responses to human inquiries and predicts control signals for the next step. However, since the framework lacks an input command, the predicted control can not follow the specific navigation command, which denotes that the framework is hard to deploy in real scenarios. Meanwhile, more researchers focus on transforming the driving situations into textual descriptions as the input for the LLM, for directly interpreting and reasoning about comprehensive driving situations In this thread of works, GPT-Driver [28] reformulated motion planning as the task of natural language modeling by converting heterogeneous scene input to language tokens. LanguageMPC [32] leveraged a LLM to reason the complex scenarios and output high-level driving decisions. Then the method tunes a parameter matrix to convert the decision into the low-level control signals. LLM-Driver [6] utilized the numeric vector



Navigation instruction: Just move to the left and get ready to leave the highway.



Navigation instruction: Turn left at the next T-junction
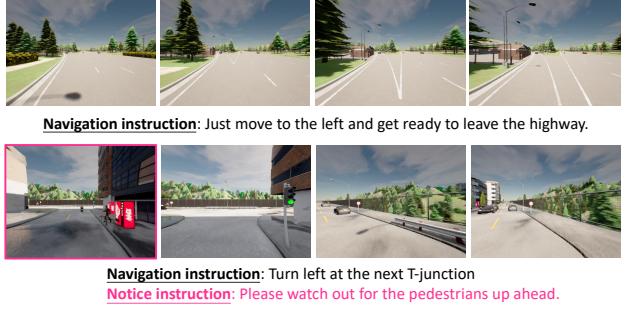Notice instruction: Please watch out for the pedestrians up ahead.

Figure 2. Two examples of the collected data with corresponding labeled navigation instructions and optional notice instructions.
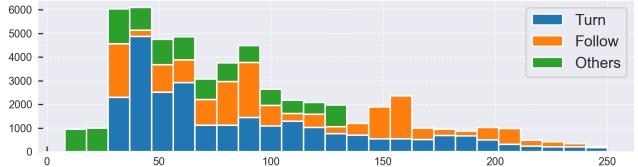


Figure 3. Distribution of parsed clips in terms of clip length and the corresponding navigation instruction type.

as the input modality and fused vectorized object-level 2D scene representation to enable the LLM to answer the questions based on the current environment.

However, this line of work only considered the driving problem in the open-loop settings, and ignored questions such as cumulative error, temporal action consistency, and end-to-end trainability, which are critical for bringing the models into actual closed-loop driving tasks. To the best of our knowledge, we are the first language-based end-to-end autonomous driving method in the closed-loop setting. Relevant datasets, benchmarks, and trained models are also open-sourced to facilitate further research in the community.

## 3. Dataset generation

We aim to develop an intelligent driving agent that can generate driving actions based on three sources of input: 1) sensor data (multi-view camera and LiDAR), so that the agent can generate actions that are aware of and compliant with the current scene; 2) navigation instructions (*e.g.* lane changing, turning), so that the agent can drive to meet the requirement in natural language (instruction from humans or navigation software); and 3) human notice instruction, so that the agent can interact with humans and adapt to human's suggestions and preferences (*e.g.* pay attention to adversarial events, deal with long-tail events, *etc*). In this section, we describe how to generate the multi-modal dataset needed to train the agent, and the prompt design for the navigation instruction and human notice instruction. Specifically, we choose the CARLA [12] as the simulator, because it can simulate a realistic dynamic closed-loop world and it is widely adopted in the field of end-to-end autonomous

| Type | Three randomly chosen instructions of each instruction type |
|------|-------------------------------------------------------------|
| **Follow** | Maintain your current course until the upcoming intersection. <br> In [x] meters, switch to left lane. <br> Ease on to the left and get set to join the highway. |
| **Turn** | After [x] meters, take a left. <br> At the next intersection, just keep heading straight, no turn. <br> You'll be turning left at the next T-junction, alright? |
| **Others** | Feel free to start driving. <br> Slow down now. <br> Head to the point, next one's [x] meters ahead, [y] meters left/right. |
| **Notice** | Watch for walkers up front. <br> Just a heads up, there's a bike ahead. <br> Please be aware of the red traffic signal directly in front of you. |

Table 1. Examples of considered navigation instructions (follow, turn, others) and notice instructions. $[x]$ and $[y]$ represent the float number for a specific distance.

driving. The data collection consists of two stages: 1) collecting sensor data and control signals with an expert agent; and 2) parsing and labeling collected data with instructions.

**Sensor and control data collection.** We utilize a rule-based expert agent [34] to create a dataset including about 3M driving frames. Since the expert agent can access the privileged information in the CARLA, this dataset will include camera data, LiDAR data, and control actions for each frame. To enhance the diversity of the collected dataset, the agent runs on 2.5k routes, 8 towns, and 21 kinds of environmental conditions (*e.g.* weather, time of the day). We use four RGB cameras (left, front, right, rear) and one LiDAR. The side cameras are angled at $60°$. Besides, we center-crop the front image as an additional focus-view image to capture the status of the distant traffic light. The LiDAR has 64 channels and generates 600K points per second.

**Parsing and language annotation.** In the second stage, we parse the collected data into clips, and label each clip with proper navigation instructions and optional notice instructions. The parsing process takes a sequence of frames as input, and segments these frames into clips, where each clip corresponds to one navigation instruction. For instance, if the agent started to turn left at frame $T_0$ and ended at frame $T_n$, we will label $(T_0, T_n)$ as a new clip with the instruction "*Hang a left at the next crossroads*". Besides, if an adversarial event[1] occurs at time $T_a$, we will add one notice instruction into this clip, simulating a real-life scenario where a passenger or a side assistance system would communicate with the driver when an emergency happens. As shown in Figure 2, each clip includes sensor data, control signals, the corresponding navigation instruction, and optional notice instructions. The distribution of the parsed clips in terms of clip length and corresponding instruction is shown in Figure 3. In our dataset, we collect 64K parsed clips and 464K notice instructions.

---

[1]The adversarial events include bad road conditions, the front vehicle's sudden brake, unexpected entities rushing into the road from occluded regions, vehicles running a red traffic light, *etc.*

**Instruction design.** We consider three types of navigation instructions (follow, turn, and others) along with one type of notice instruction, consisting of a total of 56 different instructions. Table 1 shows some examples and the full list can be found in supplementary material. To enable the agent to drive in realistic instructional settings where the instructions come from navigation software or humans, we

- *Diversifying the instructions*: Considering the inherent richness of natural language, for each type of instruction, we utilized ChatGPT API to generate eight different variants, each carrying the same semantic meaning but varying in phrasing. This enables more comprehensive coverage and flexibility in language interpretation, accommodating the diverse ways the same instruction can be conveyed.
- *Incorporating misleading instructions*: in real-world cases, the navigation software or passengers may give misleading instructions to the AV that violate traffic rules or raise safety concerns. For example, on a single-lane road, following an instruction "*Change to left lane*" is dangerous. To improve the robustness of our model against misleading instructions, we simulate these scenarios and add them to our dataset.
- *Connecting multiple instructions:* In many cases, the instructions may consist of two to three consecutive instructions, such as "*Turn right at this intersection, then go straight to the next intersection and turn right again.*" We also construct some consecutive complex instruction data to simulate real navigation-based driving scenarios.

## 4. LMDrive methodology

In this work, we propose LMDrive, a framework that can understand and follow high-level driving instructions through natural language. As illustrated in Figure 4, LM-Drive is composed of two major components: 1) a vision encoder that processes multi-view multi-modal sensor data (camera and LiDAR) for scene understanding and generating visual tokens; 2) a large language model and its associated component (tokenizer, Q-Former, and adapters) that takes in the visual tokens and language instruction, to predict the control signal and whether the given instruction is completed. We will introduce the vision encoder in Section 4.1, and the language model with its associated components in Section 4.2. Finally, we describe the training details in Section 4.3.

### 4.1. Vision encoder

In the visual language community [5, 25, 51], the most common way to align vision and language could be using pre-trained CLIP models [31] to encode image features. However, the large flops and parameter size of the CLIP models increase the difficulty of its deployment in AV systems. Also, AV perception systems are usually in 3D to include LiDAR input. Hence, inspired by InterFuser [34]
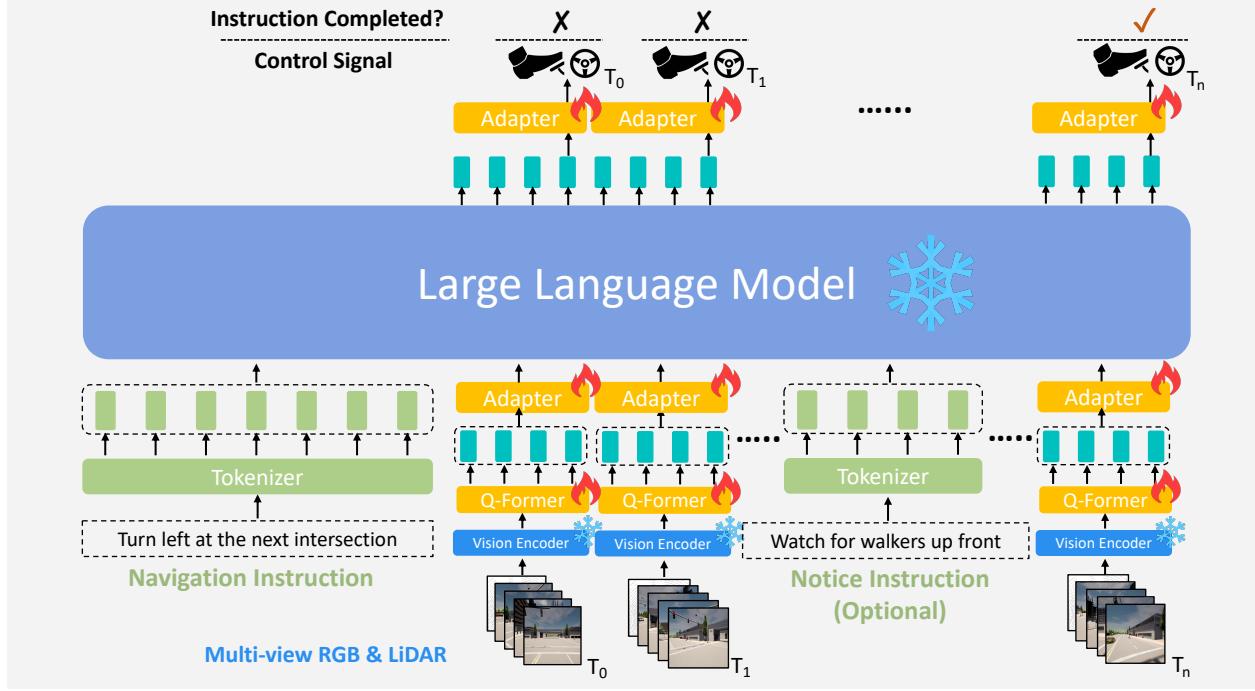
Figure 4. The structure of the proposed LMDrive model, which consists of two major components: 1) a vision encoder that processes multi-view multi-modal sensor data (camera and LiDAR) for scene understanding and generating visual tokens; 2) a large language model and its associated component (tokenizer, Q-Former, and adapters) that processes all the historic visual tokens and the language instructions (navigation instruction and optional notice instruction), to predict the control signal and whether the given instruction is completed.

and TF++ [18], we design a multi-view multi-modality vision encoder to encode/fuse the sensor data. As shown in Figure 5, the vision encoder consists of the sensor encoding part which encodes image and LiDAR input respectively, and a BEV decoder that fuses image and point cloud features to generate visual tokens which are then passed to the language model. Notably, the vision encoder is pre-trained on perception tasks by adding additional prediction heads, then the encoder is frozen for later use by the large language model.

**Sensor encoding.** For each image input, we apply a 2D backbone ResNet [15] to extract the image feature map. The feature map is then flattened into one-dimensional tokens. For a comprehensive understanding of the global context from multiple viewpoints, tokens from different views will be fused by a standard $K_{enc}$-layer transformer encoder, each layer containing Multi-Headed Self-Attention [42], MLP blocks and layer normalization [1]. For the LiDAR input, we adopt a 3D backbone PointPillars [21] to process the raw point cloud data into ego-centered LiDAR features, where each pillar encompasses points within a 0.25m × 0.25m area. PointNet [30] is then used to aggregate features and downsample the feature map to $C \times H \times W$, which subsequently serve as BEV queries.

**BEV decoder.** The encoded sensor features above are then passed into the BEV decoder to generate visual tokens.
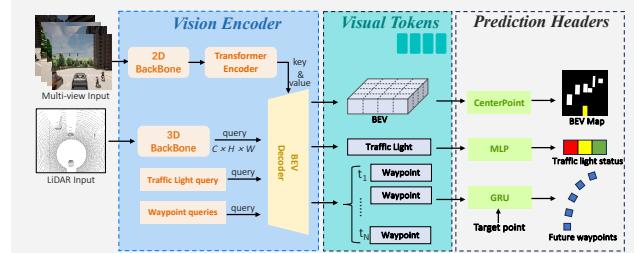


Figure 5. The detailed structure of the vision encoder, which takes as input the multi-view multi-modality sensor data. In the pre-training stage, the vision encoder is appended with prediction headers to perform pre-training tasks (object detection, traffic light status classification, and future waypoint prediction). In the instruction-finetuning stage and inference stage, the prediction headers are discarded, and the vision encoder is frozen to generate visual tokens to feed into the LLM.

Specifically, the BEV decoder is designed as a standard transformer with $K_{dec}$ layers. The BEV point cloud features are fed into the BEV decoder as $H \times W$ queries to attend to the multi-view image features and generate BEV tokens. We also feed $N$ learnable queries and 1 learnable query into the BEV decoder to generate $N$ waypoint point tokens and 1 traffic light token respectively. Thus the three types of visual tokens (BEV, waypoint, and traffic light) will contain rich scene information and will be then presented to the large language model.

5

**Pre-training with prediction headers.** we consider three vision encoder pre-training tasks: object detection, future waypoint prediction, and traffic light status classification. For object detection, the BEV tokens will pass through a one-stage CenterPoint [47] to predict the bounding boxes and velocity of the objects in an $Hm \times Wm$ area. For waypoint prediction, we pass the $N$ waypoint tokens along with the navigation waypoint into the GRU network [9] sequentially to predict $N$ future waypoint. For the traffic light status classification, a 2-layer MLP is applied to the traffic light token. Three corresponding loss terms are considered: 1) the detection loss as in InterFuser [35]; 2) the $l_1$ waypoint loss; and 3) the cross-entropy traffic light state loss. Note that these prediction headers are only used in the pre-training of the vision encoder, and will be discarded in the training of the LLMs and inference of the whole model.

### 4.2. LLM for instruction-following auto driving

As illustrated in Figure 4, in our framework, the LLM functions as the "brain" throughout the entire driving procedure, processing sensor tokens generated by the frozen vision encoder for each frame, comprehending natural language instructions, generating the necessary control signal, and predicting whether the given instruction is completed. Specifically, we choose LLaMA [39] as the language backbone, which has been widely adopted in many language [14, 49] and vision [25, 51] instruction-tuning models. We also have three associated components to bridge LLM with instruction, visual information input, and action prediction: 1) a tokenizer, 2) a Q-Former, 3) two adapters.

**Instruction and visual tokenization.** Given the navigation instruction and optional notice instruction, we apply the LLaMA tokenizer [39] to convert the instruction into textual tokens. Note that the duration of executing one instruction would span from a few seconds to a few minutes, and our model is deployed in the closed-loop setting. Thus at each frame, we utilize all historic sensor information (with maximum limit $T_{max}$) to depress cumulative error and improve the temporal consistency of the model. Specifically, for each frame's multi-view multi-modality sensor input, we utilize the vision encoder pre-trained in the previous section to generate visual tokens ($H \times W$ BEV tokens, $N$ waypoint tokens, and one traffic light token). However, the number of visual tokens (*e.g.* 406 tokens for each frame) quickly grows too large for the LLM because usually hundreds of frames are needed to complete one instruction. To overcome this, we follow BLIP-2 [22] to use the Q-Former to reduce the number of visual tokens. Specifically, for each frame, we employ $M$ learnable queries to attend to the visual tokens via cross-attention layers, which can reduce each frame's visual token number to $M$. Subsequently, we use a 2-layer MLP adapter to convert the tokens extracted by the Q-Former to share the same dimension as the language token, which can then be fed into the LLM.

**Action prediction.** After receiving a sequence of instructional and visual tokens, the LLM predicts the action tokens. One another 2-layer MLP adapter is then applied to predict future waypoints, as well as a flag to indicate whether the given instruction has been completed. Note that to enhance the supervision signal, we will also conduct prediction for every historic frame during training, and only the prediction of the latest frame will be executed at inference time. To get the final control signal, which includes braking, throttling, and steering, following LBC [4], we use two PID controllers for latitudinal and longitudinal control to track the heading and velocity of predicted waypoints respectively.

**Training objectives.** When finetuning the LLM and its associated components, we consider two loss terms: 1) the $l_1$ waypoint loss; 2) the classification loss (cross-entropy), which determines if the current frame finishes the given instruction.

### 4.3. Training details

LMDrive's training consists of two stages: 1) the vision encoder pre-training stage; and 2) the instruction-finetuning stage, to align the instruction/vision and control signal.

**Vision encoder pre-training stage.** The vision encoder takes a single frame's sensor data as input, and we use the dataset collected in Section 3 for training. Specifically, since the instruction annotation process will drop some frames, we use the raw dataset before the instruction annotation for the vision encoder pre-training, which includes data of around 3M frames. Only the vision encoder is pre-trained with perception tasks for scene understanding.

**Instruction-finetuning stage.** The entire system is trained for end-to-end autonomous driving under the guidance of the instruction, where the Q-Former and Adapters are trainable and the other components are frozen. While our LM-Drive takes a sequence of frames as input, during training we set a fixed sequence length $T_{max}$ for building up batch data. The training utilizes the instruction-following data generated in Section 3. To enable the model to reject misleading instructions, we label the corresponding data as 'completed' after the misleading instruction is given for about 1 second. Since the dataset is collected at a high frequency (∼10Hz), the data in adjacent frames are highly similar. To encourage efficient training, following video prediction methods [13, 33, 43], we sample the training frames in a fixed interval, and apply temporal augmentation which randomly shifts the training frames either forward or backward, with the random shift less than the fixed interval.

## 5. LangAuto Benchmark

We propose the LangAuto (Language-guided Autonomous Driving) CARLA benchmark, the first benchmark that evaluates closed-loop driving performance under language instructions. Compared with the previous CARLA bench-

| LLM Backbone | LangAuto | | | LangAuto-Short | | | LangAuto-Tiny | | |
|---|---|---|---|---|---|---|---|---|---|
| | DS ↑ | RC ↑ | IS ↑ | DS ↑ | RC ↑ | IS ↑ | DS ↑ | RC ↑ | IS ↑ |
| Random Init. | 10.7±3.8 | 16.2±4.9 | 0.63±0.04 | 14.2±4.4 | 20.1±4.4 | 0.72±0.04 | 20.1±4.1 | 24.7±5.1 | 0.75±0.03 |
| LLaMA [39] | 31.3±1.5 | 37.1±1.6 | 0.82±0.01 | 42.8±7.2 | 49.1±8.5 | **0.87±0.03** | 52.2±5.3 | 57.8±8.0 | **0.91±0.05** |
| LLaMA2 [40] | 32.8±2.1 | 40.1±2.2 | 0.81±0.02 | 44.8±6.2 | 53.5±5.5 | 0.84±0.02 | 56.1±4.1 | 64.2±4.7 | 0.87±0.04 |
| Vicuna [49] | 33.5±1.9 | 39.3±1.9 | 0.83±0.02 | 45.3±4.9 | 54.3±3.9 | 0.83±0.03 | 55.5±3.9 | 63.1±4.2 | 0.88±0.04 |
| Vicuna-v1.5 [49] | 34.0±3.8 | 39.0±3.3 | **0.85±0.06** | 47.0±4.3 | 56.5±2.4 | 0.83±0.04 | 59.0±2.6 | 69.9±2.3 | 0.84±0.02 |
| LLaVA-v1.5 [24] | **36.2±2.3** | **46.5±4.3** | 0.81±0.03 | **50.6±1.7** | **60.0±3.4** | 0.84±0.04 | **66.5±3.6** | **77.9±2.3** | 0.85±0.02 |

Table 2. Performance comparison of 6 LLM backbones on the LangAuto benchmark. We report the metrics for 3 evaluation runs.

| Module design | DS ↑ | RC ↑ | IS ↑ |
|---|---|---|---|
| Baseline (LLaVA-v1.5) | **36.2±2.3** | **46.5±4.3** | **0.81±0.03** |
| w/o Q-Former | 31.7±3.5 | 41.2±4.4 | 0.79±0.02 |
| w/o using BEV tokens | 33.9±3.9 | 45.9±5.1 | 0.72±0.03 |
| w/o visual pre-training | 16.9±5.1 | 24.1±4.7 | 0.70±0.04 |

Table 3. Ablation study on the module design.

mark, Town05 [29] and Longest6 [8] that navigate the agent with discrete driving commands or target waypoints, our benchmark only provides the AV with a navigation instruction and optional notice instructions in natural language.

Specifically, the LangAuto benchmark covers all 8 publicly available towns in CARLA to include various scenarios (*e.g.* highways, intersections, roundabouts). We also consider 16 kinds of environmental conditions, encompassing the combinations of 7 weather conditions (Clear, Cloudy, Wet, MidRain, WetCloudy, HardRain, SoftRain) and 3 daylight conditions (Night, Noon, Sunset). Besides, LangAuto consists of three tracks to fully test the agent's instruction-following abilities:

- LangAuto track: For each route, navigation instructions are given and updated to the agent based on the agent's current position. We also divide this track into three sub-tracks with different route lengths, to better distinguish the performance. LangAuto where the routes are longer than 500 meters, LangAuto-Short where the route length is between 150 and 500 meters, and LangAuto-Tiny where the route length is shorter than 150 meters.
- LangAuto-Notice track: Based on the LangAuto track, we additional add notice instructions to the agent. This setting simulates real cases where passengers or other assistance systems can give real-time notice in long-trail complex or adversarial scenarios, which is usually hard for the AV system to handle by itself. Ideally, the agent that can comprehend and leverage instruction can achieve better performance.
- LangAuto-Sequential track: Based on the LangAuto track, we merge 10% of consecutive 2 to 3 instructions into a single long instruction. This setting mimics the realistic scenarios where the multi-sentence instructions come from the passengers or navigation software.

Note that misleading instructions will be randomly (∼ 5%) and intermittently given to the driving agent, which lasts for a certain duration (1-2 seconds). The driving agent is

expected to reject these misleading instructions and execute safe actions that are compliant to the current scene, until the next correct instruction is spawned.

**Metrics.** We consider three major metrics introduced by the CARLA LeaderBoard [36]: route completion (RC), infraction score (IS), and driving score (DS). The route completion refers to the percentage of the total route length that has been completed. It only takes into account the distance traveled along the predetermined route, where each segment of the predetermined route corresponds to a navigation instruction. If the agent deviates too far from the route, the agent is regarded as violating the instruction, and this episode is marked as a failure and terminated. The infraction score measures infractions triggered by the agent. When collisions or traffic rule violations occur, the infraction score is decayed by a corresponding discount factor. The driving score is the product of the route completion ratio and the infraction score, describing both driving progress and safety. It is generally recognized as the primary ranking metric.

## 6. Experiments

### 6.1. Experiment Setup

We implement and evaluate our approach on the open-source CARLA simulator of version 0.9.10.1 [12]. For the 2D backbone and 3D backbone of the vision encoder, ResNet-50 [15] backbone is pre-trained on ImageNet, and PointPillars [21] is trained from scratch. $C, H, W, N$ are set as 256, 50, 50, 5 respectively. For the Q-former token number $M$, we empirically found that $M = 4$ tokens achieve decent performance. During the instruction-finetuning stage, we sample the training frames in a fixed interval of 2. Because each parsed clip typically contains multiple notice instructions, to avoid overfitting to notice following and generating over-conservative behaviors, during training we randomly removed notices for 75% of the clips. For the remaining clips, we ensure that a maximum of one notice is included. We refer readers to the appendix for more details.

### 6.2. Quantitative Results

**LLM backbones.** In Table 2, we evaluate our method with different pretrained or randomly initialized LLM models of 7B parameters: LLaMA [39] and LLaMA2 [40] were pretrained with large public language datasets; Vicuna [49]

| LLM Backbone | Benchmark Type | Infraction Score ↑ | Vehicle Collisions ↓ | Pedestrian Collisions ↓ | Layout Collisions ↓ | Red light Violations ↓ | Offroad Infractions ↓ | Blocked Infractions ↓ |
|---|---|---|---|---|---|---|---|---|
| LLaVA-v1.5 | LangAuto | 0.81 | 0.33 | 0.03 | 0.50 | 0.92 | 0.36 | 0.22 |
| | LangAuto-Notice | 0.87 | 0.17 | 0.02 | 0.31 | 0.50 | 0.17 | 0.27 |
| Vicuna-v1.5 | LangAuto | 0.85 | 0.30 | 0.03 | 0.43 | 1.18 | 0.24 | 0.19 |
| | LangAuto-Notice | 0.91 | 0.15 | 0.01 | 0.28 | 0.56 | 0.26 | 0.19 |

Table 4. Performance comparison on LangAuto and LangAuto-Notice benchmarks. The metrics (except infraction score) are normalized by the driven distance (km).

| LLM Backbone | Benchmark Type | DS ↑ | RC ↑ | IS ↑ |
|---|---|---|---|---|
| LLaVA-v1.5 | LangAuto | 36.2 | 46.5 | 0.81 |
| | LangAuto-Sequential | 34.0 | 43.7 | 0.81 |
| Vicuna-v1.5 | LangAuto | 34.0 | 39.0 | 0.85 |
| | LangAuto-Sequential | 31.9 | 37.1 | 0.84 |

Table 5. Performance comparison on LangAuto and LangAuto-Sequential benchmarks.

and Vicuna-v1.5 [49] are additionally finetuned with human conversation data; LLaVA-v1.5 [24] incorporate multi-modal data (*e.g.* language, images) for training. We observe that LLaVA-v1.5 [24] surpasses the other LLM models, which demonstrates the importance of adopting pretrained multi-modal LLMs for instruction-following autonomous driving. Besides, the models finetuned with instruction data perform better than others (Vicuna-v1.5 > LLaMA2 ≈ Vicuna > LLaMA). We also test a randomly initialized 7B LLM model, which struggles to drive properly with the same amount of training data, which demonstrates the necessity of finetuning pretrained LLM for instruction-following driving.

**Module Design.** In Table 3, we conduct ablation studies on different components. First, instead of reducing the number of BEV tokens with the Q-Former, we directly downsample the BEV features to 4 × 4, then feed them into the LLM (denoted as "w/o Q-Former"). The average driving score dropped from 36.2 to 31.7. Second, excluding the BEV tokens input into LLM decoder (denoted as "w/o using BEV tokens" results in a decreased infraction score (0.81 to 0.72). One explanation is that the BEV tokens are important for detecting and reasoning the surrounding obstacles and road structures. Third, we remove the pre-training stage of the vision encoder and directly train it from scratch in the instruction-finetuning stage. The driving score drops to 16.9 (denoted as "w/o visual pre-training"), demonstrating the importance of our proposed vision encoder pre-training.

**LangAuto-Notice Benchmark.** The LangAuto-Notice benchmark provides some notice instructions to the agent when adversarial events happen. As shown in Table 4, our agents can effectively leverage the information of the notice in real-time, resulting in a significant decrease in both collisions and traffic rule violations.

**LangAuto-Sequential Benchmark.** In Table 5, we demonstrate the effect of the LangAuto-Sequential benchmark
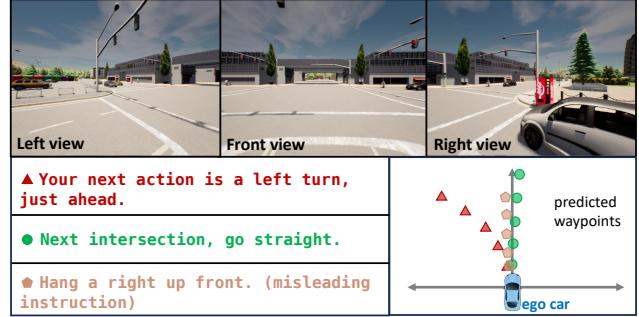


Figure 6. An example of how our LMDrive predicts future waypoints, given sensor inputs and varied navigational instructions. Under the first two instructions, our model predicts different waypoints accordingly. The third instruction is a misleading one (turn right on a left-only lane). The model appropriately rejects the incorrect instruction, generating a slower speed and a safe path that is compliant with the scenario.

where some consecutive 2 to 3 navigation instructions are merged together to a long and complex instruction. This setting additionally requires the agent to be able to be temporally aware of which instruction has been completed and which has not. Our agents based on LLaVA and Vicuna both have a drop in the driving score and route completion ratio.

## 7. Conclusion

In this paper, we introduced LMDrive, a language-guided, end-to-end, closed-loop autonomous driving framework. LMDrive incorporates natural language instructions along with multi-modal sensor data, enabling human-like interaction and navigation in complex driving scenarios. We also propose the language-guided driving dataset, comprising around 64K multi-modal data clips along with corresponding navigation instructions. We established the LangAuto benchmark for evaluating autonomous driving systems considering natural language instructions. The effectiveness of LMDrive was demonstrated through extensive closed-loop experiments, underlining the potential of improving the interaction of autonomous vehicles with humans and the environment. Our work serves as an encouraging starting point for further exploration and developments in the field of language-based closed-loop end-to-end autonomous driving.

# References

[1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016. 5

[2] Sergio Casas, Cole Gulino, Simon Suo, Katie Luo, Renjie Liao, and Raquel Urtasun. Implicit latent variable model for scene-consistent motion forecasting. In *Computer Vision– ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 624– 641. Springer, 2020. 1

[3] Dian Chen and Philipp Krähenbühl. Learning from all vehicles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17222–17231, 2022. 3, 1

[4] Dian Chen, Brady Zhou, Vladlen Koltun, and Philipp Krähenbühl. Learning by cheating. In *Conference on Robot Learning*, pages 66–75. PMLR, 2020. 6

[5] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm's referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 3, 4

[6] Long Chen, Oleg Sinavski, Jan Hünermann, Alice Karnsund, Andrew James Willmott, Danny Birch, Daniel Maund, and Jamie Shotton. Driving with llms: Fusing object-level vector modality for explainable autonomous driving. *arXiv preprint arXiv:2310.01957*, 2023. 2, 3

[7] 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. 2

[8] Kashyap Chitta, Aditya Prakash, Bernhard Jaeger, Zehao Yu, Katrin Renz, and Andreas Geiger. Transfuser: Imitation with transformer-based sensor fusion for autonomous driving. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2, 7

[9] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014. 6

[10] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022. 3

[11] DriveLM Contributors. Drivelm: Drive on language. https://github.com/OpenDriveLab/DriveLM, 2023. 2

[12] Alexey Dosovitskiy, German Ros, Felipe Codevilla, Antonio Lopez, and Vladlen Koltun. Carla: An open urban driving simulator. In *Conference on robot learning*, pages 1–16. PMLR, 2017. 2, 3, 7

[13] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 6

[14] Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and Dawn Song. Koala: A dialogue model for academic research. *Blog post, April*, 1, 2023. 6

[15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5, 7, 1

[16] 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. 1, 2

[17] Bernhard Jaeger, Kashyap Chitta, and Andreas Geiger. Hidden biases of end-to-end driving models. 2023. 3

[18] Bernhard Jaeger, Kashyap Chitta, and Andreas Geiger. Hidden biases of end-to-end driving models. *arXiv preprint arXiv:2306.07957*, 2023. 5

[19] Xiaosong Jia, Yulu Gao, Li Chen, Junchi Yan, Patrick Langechuan Liu, and Hongyang Li. Driveadapter: Breaking the coupling barrier of perception and planning in end-to-end autonomous driving. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7953–7963, 2023. 3

[20] Xiaosong Jia, Penghao Wu, Li Chen, Jiangwei Xie, Conghui He, Junchi Yan, and Hongyang Li. Think twice before driving: Towards scalable decoders for end-to-end autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21983–21994, 2023. 2

[21] Alex H Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12697–12705, 2019. 5, 7

[22] 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. 6, 1

[23] 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. 1

[24] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 7, 8

[25] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023. 3, 4, 6

[26] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016. 1

[27] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2018. 1

[28] Jiageng Mao, Yuxi Qian, Hang Zhao, and Yue Wang. Gpt-driver: Learning to drive with gpt. *arXiv preprint arXiv:2310.01415*, 2023. 2, 3

[29] Aditya Prakash, Kashyap Chitta, and Andreas Geiger. Multimodal fusion transformer for end-to-end autonomous driving. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 7

[30] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 918–927, 2018. 5, 1

[31] 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. 4

[32] Hao Sha, Yao Mu, Yuxuan Jiang, Li Chen, Chenfeng Xu, Ping Luo, Shengbo Eben Li, Masayoshi Tomizuka, Wei Zhan, and Mingyu Ding. Languagempc: Large language models as decision makers for autonomous driving. *arXiv preprint arXiv:2310.03026*, 2023. 2, 3

[33] Hao Shao, Shengju Qian, and Yu Liu. Temporal interlacing network. *AAAI*, 2020. 6

[34] Hao Shao, Letian Wang, Ruobing Chen, Hongsheng Li, and Yu Liu. Safety-enhanced autonomous driving using interpretable sensor fusion transformer. In *Conference on Robot Learning*, pages 726–737. PMLR, 2023. 2, 4

[35] Hao Shao, Letian Wang, Ruobing Chen, Steven L Waslander, Hongsheng Li, and Yu Liu. Reasonnet: End-to-end driving with temporal and global reasoning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13723–13733, 2023. 1, 2, 6

[36] CARLA team. Carla autonomous driving leaderboard. https://leaderboard.carla.org/, 2020. Accessed: 2021-02-11. 7

[37] Trisha Thadani. Cruise recalls all its driverless cars after pedestrian hit and dragged. *The Washintong Post*. 2

[38] Marin Toromanoff, Emilie Wirbel, and Fabien Moutarde. End-to-end model-free reinforcement learning for urban driving using implicit affordances. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7153–7162, 2020. 3

[39] 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, 6, 7

[40] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023. 3, 7

[41] Jordan Valinsky. 'complete meltdown': Driverless cars in san francisco stall causing a traffic jam. *CNN Business*. 2

[42] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 5

[43] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *European conference on computer vision*, pages 20–36. Springer, 2016. 6

[44] Letian Wang, Jie Liu, Hao Shao, Wenshuo Wang, Ruobing Chen, Yu Liu, and Steven L Waslander. Efficient reinforcement learning for autonomous driving with parameterized skills and priors. *arXiv preprint arXiv:2305.04412*, 2023. 3

[45] Penghao Wu, Xiaosong Jia, Li Chen, Junchi Yan, Hongyang Li, and Yu Qiao. Trajectory-guided control prediction for end-to-end autonomous driving: A simple yet strong baseline. *Advances in Neural Information Processing Systems*, 35:6119–6132, 2022. 3

[46] 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

[47] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11784–11793, 2021. 6

[48] Zhejun Zhang, Alexander Liniger, Dengxin Dai, Fisher Yu, and Luc Van Gool. End-to-end urban driving by imitating a reinforcement learning coach. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 15222–15232, 2021. 3

[49] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric. P Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. 6, 7, 8

[50] Tong Zhou, Letian Wang, Ruobing Chen, Wenshuo Wang, and Yu Liu. Accelerating reinforcement learning for autonomous driving using task-agnostic and ego-centric motion skills. *arXiv preprint arXiv:2209.12072*, 2022. 3

[51] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 3, 4, 6

# LMDrive: Closed-Loop End-to-End Driving with Large Language Models
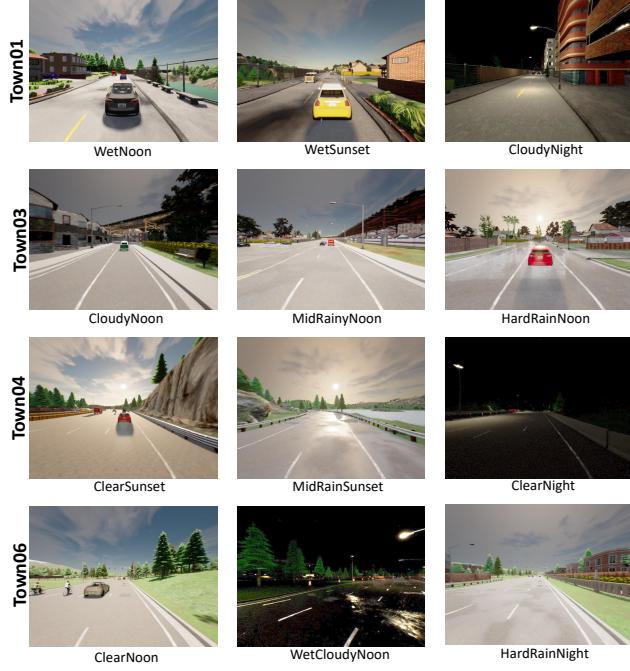
## Supplementary Material



Figure 7. The visualization of some weather and daylight conditions used in the LangAuto Benchmark.



Figure 8. The visualization of eight town maps used in the LangAuto Benchmark.

## A. Implementation Details

**Model details** For the vision encoder, we use $K_{enc} = 1$ encoder layers and $K_{dec} = 3$ decoder layers. The feature of the 5-th stage in the ResNet is employed as the extracted feature map, then we apply an MLP layer to convert its dimension to 768, which is the feature dimension of the following Q-Former. Following LAV [3], we build a simplified version of PointNet [30] with several MLP layers and Batch Normalization layers to encode LiDAR point cloud data. For the Q-Former, we utilize the model architecture and pre-trained weights from the BLIP-2 [22]. In our work, the visual tokens fed into the Q-Former include 400 BEV tokens, 5 future waypoint tokens and 1 traffic light token.

**Sensor configuration** We use one front-facing camera, two side-facing cameras, and one back-facing camera to collect RGB images. Each camera has a resolution of $800 \times 600$ and a $100°$ horizontal field of view (FOV). The two side cameras are angled at $60°$. For the font image, we scale the shorter side of the front camera input to 256 and crop its center patch of $224 \times 224$. For the focusing view image, we directly crop the center of the front camera input to get a $128 \times 128$ patch which can capture distance traffic light status. For the other images, the shorter side of the camera
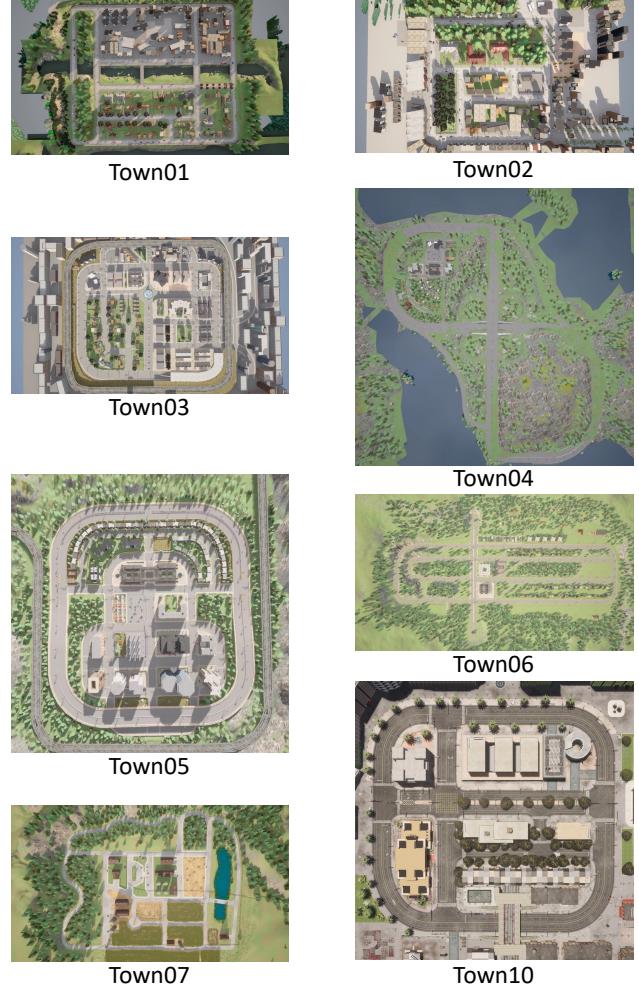
input is scaled to 160 and then takes a $128 \times 129$ center patch. For the LiDAR sensor, the rotation frequency is 10Hz and the upper/lower field-of-view is 10/-30. The number of channels is 64.

**Training details** We first introduce the details of the vision encoder pretraining. We adopt the AdamW optimizer [27] and a cosine learning rate scheduler [26] for the training. The initial learning rate set for the transformer encoder and 3D backbone is $\frac{BatchSize}{512} \times 5e^{-4}$. And $\frac{BatchSize}{512} \times 2e^{-4}$ is the learning rate for the 2D backbone because the backbone is initialized with the ImageNet pre-train weights. We train the models for 35 epochs with the first 5 epochs for warm-up [15]. We used random scaling from $0.9$ to $1.1$ and color

| Sample rate | DS ↑ | RC ↑ | IS ↑ |
|---|---|---|---|
| 1 | 49.5±1.5 | 58.5±2.7 | 0.83±0.03 |
| 2 | **50.6±1.7** | **60.0±3.4** | **0.84±0.04** |
| 4 | 46.0±2.1 | 59.5±2.7 | 0.79±0.03 |

Table 6. Ablation study on the different choices of the sample rate. The experiment is conducted on the LangAuto-Short benchmark with the backbone of LLaVA-v1.5.

jittering to augment the collected RGB images.

Then we introduce the training details in the instruction-finetuning stage. We also adopt the cosine learning rate scheduler and the initial learning rate is $1e^{-4}$ for a batchsize 32. We train the models for 15 epochs with the first 2000 iteration steps for warm-up. The weight decay is set to 0.07. The maximum historic horizon $T_{max}$ is set to 40, and we will truncate the data clip to keep the recent 40 frames if its number of frames exceeds 40.

## B. Additional Experiments

In this section, we delve further into our method by conducting additional ablation studies. First, we assess our methodology using a variety of sample rates in Table 6. The term "sample rate" represents the fixed interval at which training frames are sampled. When the sample rate is set at 1, the horizon becomes narrow, which prevents the application of temporal augmentation, and consequently leads to a decrease in performance. Conversely, setting the sample rate at 4 can create an excessively large gap between two consecutive frames, which might obtain a poor driving score. A sample rate of 2 achieves a good trade-off.

Second, we conduct ablation studies on the usage rate of notice instructions as shown in Table 7. In our method, we randomly removed 75% notice instructions in the data clips to avoid overfitting. It's worth noting that we use the LangAuto-Short here, rather than LangAuto-Notice, where the AV can not receive any notice instruction. We first remove all notice data, and get a worse driving score and infraction score. This suggests that incorporating notice instructions during training may improve the AV's ability to both attend to and understand adverse events, thus reducing collisions. However, when we tried to include all notice data, we found it did not enhance performance. Utilizing 25% of the instructions achieves a good trade-off.

## C. Benchmark Details

In Figure 7, we show the 12 different environmental conditions (16 conditions are used in our benchmark in total). In Figure 8, we show all 8 town maps we used in this work. In Table 8, we list the basic statistical information of the LangAuto benchmarks across various tracks (LangAuto, LangAuto-Short, LangAuto-Tiny).

| Notice Data % | DS ↑ | RC ↑ | IS ↑ |
|---|---|---|---|
| 0 | 45.2±2.8 | **67.1±2.5** | 0.68±0.03 |
| 25 | **50.6±1.7** | 60.0±3.4 | **0.84±0.04** |
| 100 | 49.1±1.9 | 58.2±2.4 | 0.83±0.04 |

Table 7. Ablation study on the usage rate of notice instructions. The experiment is conducted on the LangAuto-Short benchmark with the backbone of LLaVA-v1.5.

| Benchmark Type | LangAuto | LangAuto-Short | LangAuto-Tiny |
|---|---|---|---|
| Avg. Driving Distance (m) | 635.8 | 305.9 | 122.4 |
| Avg. Navigation Instructions | 20.3 | 10.8 | 5.1 |
| Avg. Notice Instructions | 5.8 | 3.3 | 1.7 |

Table 8. Comparative Analysis of Different Benchmark Tracks

## D. Instruction Details

Our work considers 56 different types of navigation and notice instructions, and we use ChatGPT to generate eight different phrases for the same navigation instruction. Table 9 shows an example where we generate eight phrases for the navigation instruction 'Turn Right'. We present the full list of 56 different types of navigation and notice instructions in Table 10. Table 11 demonstrates how we generate the misleading instructions. The first column is the driving scenario in which the agent is located, and the second column represents the possible misleading instructions generated by our framework. For example, when the AV is on a single-lane road, the misleading instruction "*Change your route to the left-hand lane*" violates the traffic rules and raises safety concerns. In Table 12, we show some examples of the connected instructions.

| Instruction | Eight instructions of one kind of instruction |
|---|---|
| **Turn right** | After [x] meters, execute a right turn. After [x] meters, take a right. Right in [x] meters. After [x] meters, hang a right. In [x] meters, prepare to turn right. Continue for [x] meters, then turn right. Go [x] meters, then just take a right. After [x] meters, a right turn is required. |

Table 9. Eight different phrazes of one "turn" instruction generated by the ChatGPT API.

| Type | One randomly chosen instruction of each instruction type |
|---|---|
| **Follow** | Transition to the left lane for travel. |
| | You might want to switch to the right lane. |
| | In about [x] meters, you'll want to switch to the left lane. |
| | In [x] meters, reposition to the right lane. |
| | Keep going on this road, you're doing great! |
| | Continue on the highway. |
| | Maintain your course along this route for precisely [x] meters. |
| | Cruise down the highway for about [x] meters. |
| | Continue in a straight line along your current path. |
| | Keep going straight until you reach the next junction, you're on the right track! |
| | Preserve your current trajectory for exactly [x] meters. |
| | Stay straight for [x] meters until the next intersection. |
| | Veering to the left, prepare to enter the highway. |
| | Execute a right maneuver, prepare for highway exit. |
| | In [x] meters, slide left and plan to hop off the highway. |
| | In [x] meters, proceed to the right, prepare for immediate highway departure. |
| **Turn** | Prepare to turn left up ahead. |
| | Proceed ahead and make a right turn. |
| | Continue for [x] meters, then turn left. |
| | After [x] meters, execute a right turn. |
| | You're going to be turning left at the next junction, alright? |
| | It is mandatory to take a right turn at the imminent intersection. |
| | Straight through the next crossroads. |
| | After navigating [x] meters, a left turn at the intersection is obligatory. |
| | After moving forward [x] meters, prepare to make a right turn at the intersection. |
| | [x] meters more, then straight on at the intersection, piece of cake. |
| | When you reach the next traffic signal, you will need to turn left. |
| | Please execute a right turn upon reaching the upcoming traffic signal. |
| | Please maintain your course straight at the next traffic signal. |
| | Just another [x] meters, then you'll be turning left at the light, okay? |
| | After [x] meters, take a right at the light. |
| | After traversing [x] meters, it's crucial to continue straight at the traffic signal. |
| | Next T-junction, turn left. |
| | At the forthcoming T-intersection, execute a right turn. |
| | Just keep on going straight through the next T-junction, sound good? |
| | In [x] meters, hang a left at the T. |
| | After [x] meters, take a right at the T, no biggie. |
| | After a distance of [x] meters, maintaining a straight course at the T-intersection is imperative. |
| | Find your way out at the first exit on the roundabout, please. |
| | Depart at the second exit on the roundabout. |
| | Shoot out on the third exit. |
| **Others** | Feel free to start driving. |
| | Please implement an immediate reduction in your speed. |
| | Hit the brakes, stop now. |
| | Drive freely. |
| | Please navigate towards the designated point, which is [x] meters in front of you and [y] meters to your left/right. |
| **Notice** | Please watch out for the pedestrians up ahead. |
| | Attention is required for the bicycle ahead. |
| | Watch for the car that's just stopped up front. |
| | Be mindful of the vehicle crossing on a red light to your left. |
| | Car ran red light ahead. |
| | Please be alert of the uneven road surface in the vicinity ahead. |
| | Watch for the tunnel coming up. |
| | Just a heads up, there's a red light ahead. |
| | Green light ahead. |
| | Be mindful of the yellow light ahead. |

Table 10. Full list of the 56 different types of navigation instructions and notice instructions considered in our framework.

| Driving Scenarios | One randomly chosen instruction of each misleading type |
|---|---|
| **Single-lane Road** | Change your route to the left-hand lane. <br> Transition to the right lane for travel. <br> Proceed ahead and make a left turn. <br> A right turn is required up ahead. <br> Hang a left at the next crossroads. <br> Depart at the first exit on the roundabout. <br> Roll out on the second exit. <br> It is necessary for you to take the third exit on the roundabout. |
| **Non-highway Road** | Maintain your course along the highway. <br> Slide left and plan to hop on the highway. <br> Ease on to the right and get set to quit the highway. |
| **Non-roundabout Road** | First exit. <br> No sweat, just hit the second exit on the roundabout. <br> It is necessary for you to take the third exit on the roundabout. |
| **Non-intersection Road** | Just up ahead, take a left. <br> Your next action is a right turn, just ahead. <br> Execute a left maneuver, prepare for highway entry. <br> Right, ready to exit. <br> Carefully navigate to the first exit as you approach the roundabout. <br> You are required to take the second exit on the roundabout. <br> Third exit. |
| **T-Intersection (Left Turn Prohibited)** | You'll be turning left at the next T-junction, alright? <br> Transition to the left lane for travel. <br> You might want to switch to the right lane. |
| **T-Intersection (Right Turn Prohibited)** | A right turn is mandatory at the upcoming T-intersection. <br> Change your route to the left-hand lane. <br> Just head for the right lane. |
| **When turning** | Please adjust your course to the left-most lane. <br> Reposition to the right lane. |

Table 11. Full list of the misleading instructions and their corresponding driving scenarios.

| Driving Scenarios | Examples of connected instructions |
|---|---|
| **Two consecutive instructions** | (1) Prepare to turn right up ahead. (2) Proceed along this route. |
| | (1) Maintain your course along this route. (2) Left, ready to enter. |
| **Three consecutive instructions** | (1) It's critical to keep straight at the forthcoming T-intersection. <br> (2) Keep cruising down this road. (3) At the next traffic signal, you should make a right turn. |
| | (1) At the forthcoming T-intersection, execute a right turn. <br> (2) Just head for the left lane.(3) Maintain your course along this route. |
| | (1) Keep going on this road, you're doing great! (2) Right ahead. (3) Right, ready to exit. |

Table 12. Examples of the connected navigation instructions considered in the LangAuto benchmark.