# A Novel MLLM-based Approach for Autonomous Driving in Different Weather Conditions

Sonda Fourati, Wael Jaafar, *Senior Member, IEEE,* and Noura Baccar

*Abstract*—Autonomous driving (AD) technology promises to revolutionize daily transportation by making it safer, more efficient, and more comfortable. Their role in reducing traffic accidents and improving mobility will be vital to the future of intelligent transportation systems. Autonomous driving in harsh environmental conditions presents significant challenges that demand robust and adaptive solutions and require more investigation. In this context, we present in this paper a comprehensive performance analysis of an autonomous driving agent leveraging the capabilities of a Multi-modal Large Language Model (MLLM) using GPT-4o within the LimSim++ framework that offers close loop interaction with the CARLA driving simulator. We call it MLLM-AD-4o. Our study evaluates the agent's decision-making, perception, and control under adverse conditions, including bad weather, poor visibility, and complex traffic scenarios. Our results demonstrate the AD agent's ability to maintain high levels of safety and efficiency, even in challenging environments, underscoring the potential of GPT-4o to enhance autonomous driving systems (ADS) in any environment condition. Moreover, we evaluate the performance of MLLM-AD-4o when different perception entities are used including either front cameras only, front and rear cameras, and when combined with LiDAR. The results of this work provide valuable insights into integrating MLLMs with AD frameworks, paving the way for future advancements in this field.

*Index Terms*—Autonomous driving, multimodal large language models, harsh environment, CARLA, LimSim++, GPT-4o.

## I. INTRODUCTION

Autonomous driving systems (ADS) are set to play a pivotal role in the future, necessitating advanced technology to ensure their safe and efficient integration into transportation networks. Artificial Intelligence (AI), in particular large language models (LLMs) and multimodal LLMs (MLLMs), can be used to enhance the decision-making capabilities of ADS [1] [2].

Typically, MLLMs combine multi-modal data including images, video, and audio data with the advanced reasoning capabilities of LLMs [3]. Hence, they can act as decision-making agents for various applications including medicine, education, and intelligent transport systems (ITS). More accurate and timely responses to dynamic driving conditions can be achieved when integrated into ADS. Specifically, they can process vast amounts of driving data to understand and predict complex traffic patterns, thus enabling autonomous vehicles (AVs) to learn from their environment and enhance

their driving performances. In addition, MLLMs contribute significantly to the development of ADS by providing the computational power required to analyze the driving scene in real time and make accurate decisions [4].

In this context, various approaches have been proposed to use MLLMs for AD. Indeed, several strategies have been proposed based on prompt engineering [5]–[7], fine-tuning [8]–[12], and reinforcement learning with human feedback (RLHF) [13]–[15]. Furthermore, MLLMs toward AD provisioning could be utilized for perception and scene understanding [16]–[18], question answering [19], planning and control [9], [20], and multitasking [21]. Despite their diverse methods and applications, MLLMs experience limitations in ADS. For instance, fine-tuning-based approaches are computationally expensive and may result in models overfitting for specific tasks or environments. Also, RLHF can align models with human preferences, however, it still faces scaling and generalization issues beyond particular training data instances. Finally, most works focused on using MLLM for ADS in perfect and clear weather conditions, which does not reflect the reality of many areas around the globe.

Consequently, in this work, we introduce MLLM-AD-4o, a novel prompt engineering MLLM-based approach for AD, and capable of handling autonomous driving in any weather condition. Our approach integrates different sensor modalities, including LiDAR and cameras, allowing the system to adapt according to available sensor data. Based on prompt engineering, MLLM-AD-4o provides enhanced environment perception, scene understanding, reasoning, and decision-making in ADS. Hence, our contributions can be summarized as follows:

- Unlike previous works, we integrate for the first time the harsh environment driving conditions into the CARLA simulator by leveraging functions and modules from the LimSim++ framework. This work has been documented and disseminated in Github for open public access and reproducibility [22].
- We propose MLLM-AD-4o, a novel MLLM-based AD agent that leverages GPT-4o [23] for driving perception and control decision-making.
- Through extensive experiments, we conduct a performance analysis of the proposed MLLM-AD-4o in terms of safety, comfort, efficiency, and speed score metrics, and under different weather conditions and types of involved sensor data (e.g., front and/or rear cameras and/or LiDAR).

S. Fourati and N. Baccar are with the Computer Systems Engineering Department, Mediterranean Institute of Technology (MedTech), Tunis, Tunisia. W. Jaafar is with the Software and IT Engineering Department, Ecole de Technologie Supérieure (ÉTS), Montreal, QC H3C 1K3, Canada. E-mails: {sonda.fourati, noura.baccar}@medtech.tn; wael.jaafar@etsmtl.ca.

The remainder of the paper is organized as follows. In Section II, we present the related works that integrate LLMs/M-LLMs into driving agents in a closed-loop using the CARLA simulator. In Section III, we explain the basic concepts of ADS and LLM/MLLMs with a focus on the GPT-4o architecture and functions. Section IV presents the architecture, main modules, and technical details for integrating an MLLM agent driver in a closed-loop environment using Limsim++. Section V evaluates the performance of MLLM-AD-4o in diverse scenarios and conditions. Finally, Section VI closes the paper.

## II. RELATED WORK

Recently, there has been a significant effort to enhance the efficiency and reliability of ADS using LLMs and MLLMs. Among proposed driving agents, we focus here on the contributions based on prompt engineering. Authors of [24] proposed "LLM-driver", a framework integrating numeric vector modalities, into pre-trained LLMs for question-answering (QA). LLM-driver uses object-level 2D scene representations to fuse vector data into a pre-trained LLM with adapters. A language generator (lanGen) is used to ground vector representations into LLMs. The model's performance was evaluated on perception and action prediction using mean absolute error (MAE) for predictions, accuracy of traffic light detection, and normalized errors for acceleration, brake pressure, and steering. In [25], the SurrealDriver framework has been proposed. It consists of an LLM-based generative driving agent simulation framework with multitasking capabilities that integrate perception, decision-making, and control processes to manage complex driving tasks. In [26], the authors developed LLM-based driver agents with reasoning and decision-making abilities, aligned with human driving styles. Their framework utilizes demonstrations and feedback to align the agents' behaviors with those of humans, utilizing data from human driving experiments and post-driving interviews. Also, the authors of [27] proposed the Co-driver framework for planning & control and trajectory prediction tasks. Co-driver utilizes prompt engineering to understand visual inputs and generate driving instructions, while it uses deep reinforcement learning (DRL) for planning and control. In [28], the authors introduced the PromptTrack framework, an approach to 3D object detection and tracking, by integrating cross-modal features within prompt reasoning. PromptTrack uses language prompts that act as semantic guides to enhance the contextual understanding of a scene. Finally, authors in [6] proposed HiLM-D as an efficient technique to incorporate high-resolution information in ADS for hazardous object localization.

Despite the variety of proposed driving agents, to our knowledge, no study has fully assessed the performance of MLLM-based driving agents within a closed-loop framework under harsh environmental conditions. Table I summarizes the aforementioned contributions with a comparison to this work.

## III. BACKGROUND

The key modules associated with an AV are *perception*, *planning*, *localization*, *decision-making*, and *action or control* as highlighted in Fig. 1. The *perception* module ensures sensing of the surroundings environment and identifies the driving scene. The main goal of the *localization* module is to estimate with high precision and accuracy the vehicle position on the map. *Path planning and decision* modules establish the waypoints that the vehicle should follow when moving through the surroundings. The set of waypoints corresponds to the vehicle trajectory. The *action and control* module deals with the system's actuation such as braking, steering, and acceleration.

Several MLLMs have been developed [2], with a subset of them suitable for ADS as discussed in Section II. In this work, we opt for GPT-4o [23], an optimized version of GPT-4 built around the latter's foundational transformer architecture with enhancements in understanding, generating, and maintaining context across interactions. GPT-4o includes various specialized modules to improve its operation, namely (i) language understanding module, (ii) language generation module, (iii) context management module, (iv) task-specific modules, and (v) optimization and fine-tuning modules. Programming GPT-4o involves using APIs or SDKs, custom prompts, or fine-tuning and customization. In our work, interaction with GPT-4o is realized via API calls, where we utilize a customized prompt to guide the model response. In Table II, we provide a comparison of the most recent GPT models.

## IV. PROPOSED MLLM-BASED DRIVING AGENT IN A CLOSED-LOOP ENVIRONMENT

In this section, we present the tools used to integrate the proposed MLLM-based driving agent, a.k.a., MLLM-AD-4o.

### A. CARLA Simulator

Car Learning to Act (CARLA) is an open-source simulator designed for AD research [29]. It facilitates the development, training, prototyping, and validation of autonomous urban driving systems, including perception and control. CARLA offers open-source code and protocols, and a repository of digital assets (e.g., urban layouts, multitude of buildings, vehicles, pedestrians, street signs, etc.). The platform allows for customized sensor configurations and versatile environmental settings. CARLA can be used to evaluate the effectiveness of three distinct methods to ADS: (i) A traditional modular pipeline, including a vision-based perception module, a rule-based planner, and a maneuver controller; (ii) A deep network that maps sensory inputs to driving commands, trained end-to-end via imitation learning; and (iii) an end-to-end model trained using RL.

CARLA has been built as an open-source layer over Unreal Engine 4 (UE4). It operates as a server-client system. Specifically, the server handles the simulation and scene rendering running on port 2000, while the client, implemented in Python, manages the interaction with the autonomous agent through socket communication. The client sends commands and meta-commands to the server and receives sensor data in return. In addition, the Python API serves as the primary client interface for CARLA enabling interaction with the simulation (including vehicle control), sensor data retrieval, and environment manipulation. Also, CARLA provides diverse

TABLE I: Summary of related works

| Ref. | Contribution | Used Models | Used Data | AD Tasks | Performance Metrics | Harsh Env. |
|---|---|---|---|---|---|---|
| [24] | Design and validation of an LLM Driver that interprets and reasons about driving situations to generate adequate actions. | GPT-3.5; RL agent trained with PPO, lanGen, and trainable LoRA modules. | 160k QA driving pairs dataset; Control Commands Dataset. | Perception and action prediction. | Traffic light detection accuracy; Acceleration, brake, and steering errors. | No |
| [25] | LLM-based agent in a simulation framework to manage complex driving tasks. | GPT-3 and GPT-4. | Driving Behavior Data; Simulation data from CARLA simulator; NLP libraries. | Perception, control, and decision-making. | Collision rate. | No |
| [26] | LLM-based driver agents with reasoning and decision-making, aligned with human driving styles. | GPT-4. | Private dataset from a human driver; NLP library; RL library. | Behavioral alignment; Human-in-the-Loop system. | Collision rate, average speed, throttle percentage, and brake percentage. | No |
| [27] | Co-driver agent, based on a vision language model, for planning, control, and trajectory prediction. | Qwen-VL (9.6 billion parameters), including a visual encoder, a vision-language adapter, and the Qwen LLM. | CARLA simulation data; ROS2; Customized dataset. | Adjustable driving behaviors; Trajectory and lane prediction; Planning and control. | Fluctuations frequency, and running time. | Foggy/gloomy; Rainy/gloomy |
| [28] | PromptTrack framework for 3D object detection and tracking by integrating cross-modal features within prompt reasoning. | VoVNetV2 to extract visual features; RoBERTa to embed language prompts. | NuPrompt for 3D perception in AD. | 3D object detection and tracking; Scene understanding. | Average multiple object tracking precision (AMOTA), and identity switches (IDS). | No |
| [6] | High-resolution scene understanding using proposed HiLM-D. | BLiP-2; Q-Former; MiniGPT-4. | DRAMA dataset; Pytorch; Visual Encoder; Query detection module; ST-adapter module. | Risky object detection; Vehicle intentions' and motions' prediction. | Mean intersection over union (mIoU) for detection; BLEU-4, METEOR, CIDER, and SPICE for captioning. | No |
| This work | Implementation of harsh environment scenarios and performance evaluation of an MLLM-based driving agent. | GPT-4o. | Dedicated Prompt. | Perception and decision-making. | Safety, comfort, efficiency, and speed scores. | Heavy rain; Storm; Foggy; Wetness; Good |

TABLE II: Comparison of recent GPT models

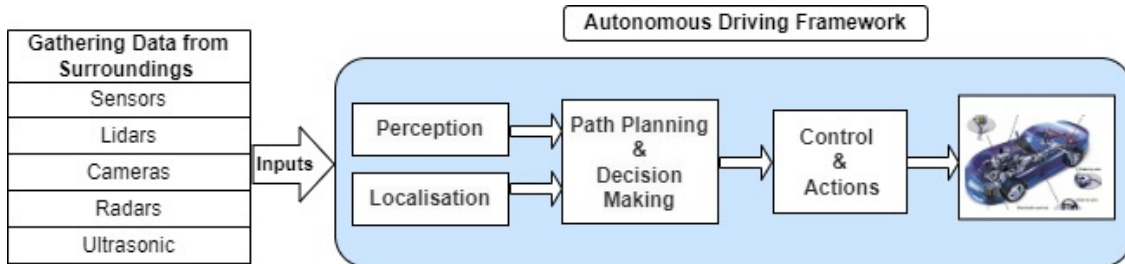| Attribute | GPT-4o | GPT-4 | GPT-4 Turbo | GPT-3.5 |
|---|---|---|---|---|
| Initial Release Date | May 13, 2024 | March 14, 2023 | July 31, 2023 | March 15, 2022 |
| Modality | Audio, video, text & images | Text & images | Text & images | Text |
| Context window | 128000 | 8192 | 25600 | 4096 |
| Parameters | Yet-to-Disclose (YTD) | 1.76 trillion | YTD | 175 billion |
| Cost per 1M token | Input: $5, Output: $15 | Input: $30, Output: $60 | Input: $10, Output: $20 | Input: $1.5, Output: $2 |
| Decoder layers | YTD | 120 | YTD | 96 |
| Prompting Method | YTD | Chain-of-Thought, n-shot | Chain-of-Thought, n-shot | n-shot |
| Training Data | YTD | 13T tokens | YTD | Over 570 GB of text data |
| Response Time | 65.8 ms per token | 94 ms per token | 50 ms per token | 35 ms per token |
| Type of Optimizer | YTD | AdamW | AdamW | Adam |
| Fine-Tuning Capability | Yes | Yes | Yes | Yes |
| Special Features | Multi-modal integration | Complex task handling | Optimized for performance and efficiency | Basic text understanding |



Fig. 1: Typical architecture of ADS.

testing environments, e.g., 14 towns, each with unique layouts, road networks, and features. In our work, we focused on using Town-06, which is a low-density town that has long highways/roads with multiple lanes (4 to 6) per direction, several highway entrances and exits, and special junctions.

*B. SUMO Simulator*

Simulation of Urban MObility (SUMO) is an open-source modular traffic simulation tool designed to handle large-scale traffic scenarios [30]. It enables the simulation of vehicular traffic, pedestrian movements, and multimodal traffic flows across complex road networks. It offers detailed modeling of various traffic scenarios including vehicle behaviors and environments. It can simulate large-scale traffic systems involving autonomous and traditional vehicles in both urban and highway scenarios. SUMO is integrated with *LimSim* as the engine

managing traffic flows and vehicle behaviors, thus enabling *LimSim* to simulate realistic and dynamic environments that reflect the complexity of urban traffic systems.

*C. LimSim++ Framework*

LimSim++ is a closed-loop framework for deploying MLLMs for autonomous driving, where multimodal data, such as text, audio, and video, can be analyzed, and then accurately react to driving scenarios [31]. LimSim++ integrates algorithms for real-time decision-making, obstacle detection, and navigation, making it an interesting tool for designing, developing, testing, and refining ADS.

To do so, LimSim++ simulates a closed-loop system with specifications of the environment including road topology, dynamic traffic flow, navigation, and traffic control. The MLLM's agent is based on prompts that provide real-time scenario

information in visuals or text. The MLLM-supported driving agent system can process information, use tools, develop strategies, and assess itself. In Fig. 2, we depict the architecture of Limsim++ and explain its main modules as follows:

- *Simulation system*: It gathers the scenario information from SUMO and CARLA, and packages it as input for the MLLM, such as visual content, scenario cognition, and task description.
- *Driver Agent*: It communicates with the "simulation system" to make driving judgments and employs MLLMs to interpret them. Data is represented as prompts for MLLMs to make suitable driving decisions. The outputs of MLLMs influence vehicle behavior and are kept in the case log system for knowledge accumulation. Following the simulation, the decisions are evaluated as follows: Those that performed well are immediately added to the memory, while poor decisions are added after reflection. The memory helps MLLMs make better driving decisions and improve the agent's performance.

When designing our MLLM-AD-4o, we customized the aforementioned modules as follows:

1) We modified the original "MPGUI.py" file within Lim-Sim++ and its related dependencies to include the visualization of six cameras (rather than the default three front cameras) in the autonomous driving simulation. Hence, a comprehensive view is provided, which is crucial to evaluating the perception capability of the driving system, analyzing/debugging the AV's sensor inputs, and decision-making.
2) We updated the default "VLM-Driver-Agent" to allow changes in the scenario and integrate new sensors and environmental conditions.
3) We created a novel function in the CARLA simulator to automate the setup of different weather conditions.

## V. EXPERIMENTAL RESULTS

### A. Definition of Metrics and Parameters

To accurately assess the AD performance of the proposed MLLM-AD-4o, we introduce the following performance metrics (called scores):

- *Safety score:* The safety level is commonly assessed through Time-to-Conflict (TTC). When TTC falls below a specific threshold, a potential risk warrants penalties. The safety score can be given by

$$\text{Safety\_score} = \begin{cases} 1, & \text{if } \tau_e \geq \tau_{\text{th}} \\ \tau_e/\tau_{\text{th}}, & \text{otherwise,} \end{cases} \quad (1)$$

where $\tau_e$ is the TTC of the AV and $\tau_{\text{th}}$ is the threshold value of TTC. Higher values reflect safer traveling.
- *Comfort score:* It evaluates the smoothness of the AV ride. Using empirical data, reference values can be determined for different driving styles, such as cautious, normal, and aggressive. It is computed as the average of the summation of acceleration score (acc_score), jerk score,

lateral acceleration score (lat_acc_score), and lateral jerk score (lat_jerk_score), as shown below

$$\begin{aligned} \text{Comfort\_score} \quad = \quad & \frac{1}{4}\left(\text{acc\_score} + \text{jerk\_score} \right. \quad (2) \\ & + \quad \text{lat\_acc\_score} + \text{lat\_jerk\_score}\left.\right). \end{aligned}$$

A higher comfort score means that traveling is smoother.
- *Efficiency score:* Driving efficiency is computed according to the vehicle's speed. In regular traffic conditions, the AV should maintain a speed at least equivalent to the average speed of surrounding vehicles. In sparse traffic conditions, the AV should approach the road's speed limit for efficiency. This score is computed as 3.

$$\text{Efficiency\_score} = \begin{cases} 1.0, & \text{if } v_e \geq v^* \\ \frac{v_e}{v^*}, & \text{otherwise,} \end{cases} \quad (3)$$

where $v_e$ represents the speed of the AV and $v^* \in \{v_{avg}, v_{limit}\}$ is the targeted speed for efficiency, being $v_{avg}$ as the average speed of surrounding vehicles, and $v_{limit}$ as the road's speed limit. A high value means efficient traveling.
- *Speed score:* The speed limit score (or speed score) penalizes the vehicle for exceeding the speed limit. This score is set to 0.9 when the AV exceeds the speed limit and 1 otherwise. This score is expressed as

$$\text{Speed\_score} = 0.9^{\left(\frac{\text{Nbr\_frames\_with\_speeding}}{\text{Total\_nbr\_frames}} \times 10\right)}, \quad (4)$$

where "Nbr_frames_with_speeding" is the number of captured frames while speeding, and "Total_nbr_frames" is the total number of captured frames.

In the following simulations, we set up $\alpha_1 = \alpha_2 = 0.25$, $\alpha_3 = 0.5$, and $v_{limit} = 50$ km/h (i.e., 13.89 m/s).

### B. Interaction with GPT-4o API

To use the GPT-4o API with prompts to process the driving environment, with data mainly collected from CARLA and SUMO, we conduct the following steps [22]:

1) Install OpenAI package and get OpenAI API key.
2) Capture images from the front and back cameras (in CARLA), save them in a suitable format, and then convert them for API submission.
3) Prepare the API request. It corresponds to sending the prompt with the image data to the GPT-4o API. The prompt describes what we want the model to do with the images and environmental data. In our work, we ask the model to analyze the surroundings, then, provide a driving decision. Fig. 3 presents an example of used prompts for captured images with the AV's cameras.
4) Process the response received from the API.

### C. Setup of Weather Conditions

In this work, we considered, in addition to "Good weather", four harsh weather conditions, namely "Heavy rain", "Storm", "Foggy", and "Wetness". We developed a new function in the CARLA simulator, named "set_weather", to configure different weather conditions. This function has eight parameters:
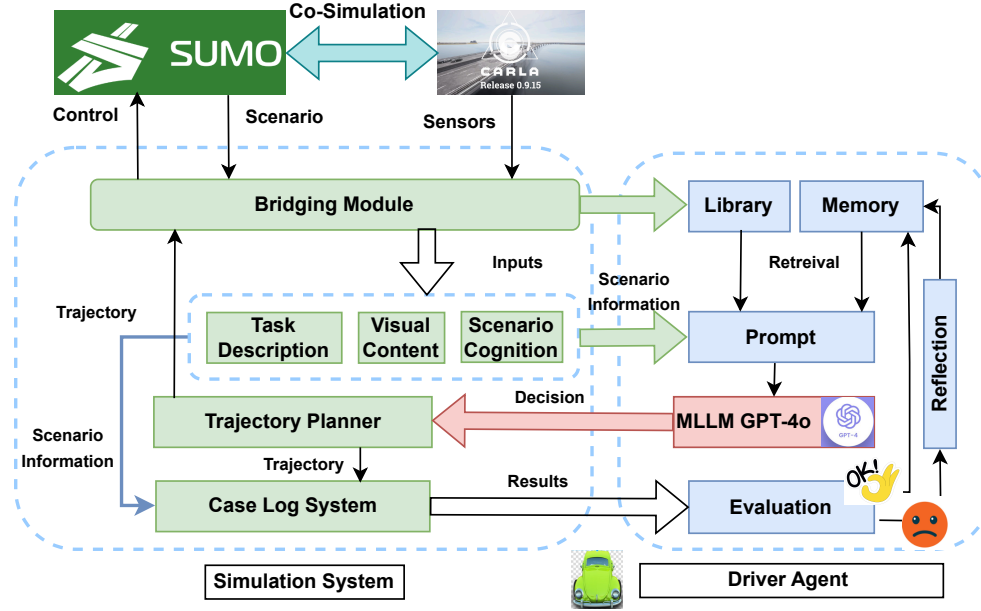
Fig. 2: Integration of MLLM-AD-4o in LimSim++.



Fig. 3: Example of used prompt for GPT-4o.

- *Cloudiness:* The amount of cloud cover in the sky. A high value means more clouds.
- *Precipitation:* The intensity of rain. A high value means heavier rain.
- *Precipitation_deposits:* The amount of water accumulating on the ground due to precipitation. A high value means important accumulation.
- *Wind_intensity:* The strength of wind. A high value means stronger wind.
- *Sun_altitude_angle:* The angle of the sun above the horizon. Lower values can simulate dawn, dusk, or low-light conditions.
- *Fog_density:* The density of fog in the environment. A high value implies thicker fog.
- *Fog_distance:* The distance at which the fog starts to become noticeable.
- *Wetness:* The wetness of the road surfaces. A high value makes the roads wetter.

In Table III, we summarize the parameters' values to set up each weather condition using the function "set_weather" [22].

### D. Integration of Semantic LiDAR in LimSim++

CARLA supports various sensing tools, including traditional LiDAR and semantic LiDAR. Traditional LiDAR provides raw distance measurements in the form of a point cloud, with no inherent understanding of the objects in the scene. In contrast, semantic LiDAR provides both the point cloud and an understanding of what each point represents, thus simplifying tasks that involve scene understanding, object recognition, and data processing reduction. We opted here for semantic

TABLE III: Parameters to set up the weather conditions

| Case | Set_weather parameters values |
|---|---|
| Heavy Rain | self.set_weather(cloudiness=80.0, precipitation=70.0, precipitation_deposits=60.0, wind_intensity=30.0, sun_altitude_angle=45.0, fog_density=10.0, fog_distance=10.0, wetness=80.0) |
| Storm | self.set_weather(cloudiness=80.0, precipitation=100.0, precipitation_deposits=100.0, wind_intensity=100.0, sun_altitude_angle=20.0, fog_density=20.0, fog_distance=10.0, wetness=80.0) |
| Fog | self.set_weather(cloudiness=40.0, precipitation=5.0, precipitation_deposits=5.0, wind_intensity=10.0, sun_altitude_angle=60.0, fog_density=70.0, fog_distance=3.0, wetness=10.0) |
| Wetness | self.set_weather(cloudiness=30.0, precipitation=0.0, precipitation_deposits=0.0, wind_intensity=0.0, sun_altitude_angle=70.0, fog_density=0.0, fog_distance=0.0, wetness=100.0) |
| Good Weather | self.set_weather(cloudiness=00.0, precipitation=00.0, precipitation_deposits=00.0, wind_intensity=00.0, sun_altitude_angle=60.0, fog_density=00.0, fog_distance=20.0, wetness=00.0) |

LiDAR. To add the latter in the perception environment within Limsim++, we follow several steps [22]:

1) Attach "Semantic_LIDAR" function to the AV and gather data in CARLA.
2) Save and process LiDAR data.
3) Add "Semantic_LIDAR" data to the prompt request.

*E. Simulation Results*

In this section, we assess the performances of the proposed MLLM-AD-4o autonomous driving agent in different conditions. Each experiment is built with the same scenario of our AV driving between points A and B, where it travels on multilane main roads and has to maneuver in curves. The autonomous driving agent has to take one decision among five possible decisions in each time frame, including idle (no change in current behavior), acceleration, deceleration, turn left, or turn right.

Figs. 4a-4d present the cumulative distribution functions (CDFs) for the safety, comfort, efficiency, and speed scores/metrics, respectively, under various weather conditions, when 3 cameras are used for autonomous driving by the proposed MLLM-AD-4o method[1]. In Fig. 4a, we see that when the weather is good (blue line), the agent takes a riskier behavior that may cause collisions than in a harsh condition. This also impacts the passengers' comfort during travel as shown in Fig. 4b. Nevertheless, in good weather, the driving behavior is the most efficient one as shown in Fig. 4c and it aligns with the speed performance depicted in Fig. 4d, This trade-off highlights that while the AV may be driven more aggressively and with less comfort, its efficiency improves as it aligns more closely with the surrounding traffic's speed. In contrast, in harsh weather conditions, e.g., heavy rain (green line) or wetness (dark yellow line), safety is enforced through slower driving, which translates into a better comfort score, offering a smoother and more cautious ride, at the expense of degraded efficiency and speed performances.
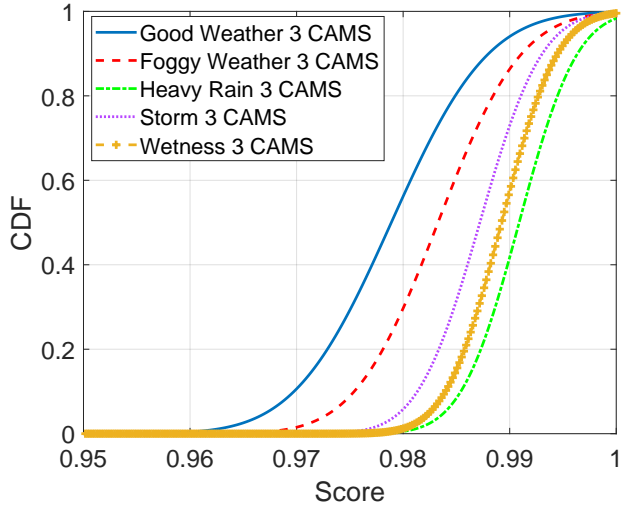
Among the weather conditions, we notice in Fig. 4 that the best safety CDF is achieved in heavy rain (green line), the best comfort CDF in stormy weather (pink line), the best efficiency CDF in good weather, and the best speed CDF in all of the foggy, heavy rain, and stormy weather conditions (red

[1]Fig. 4d illustrates also speed score results for the 6 cameras case, which will be discussed later.
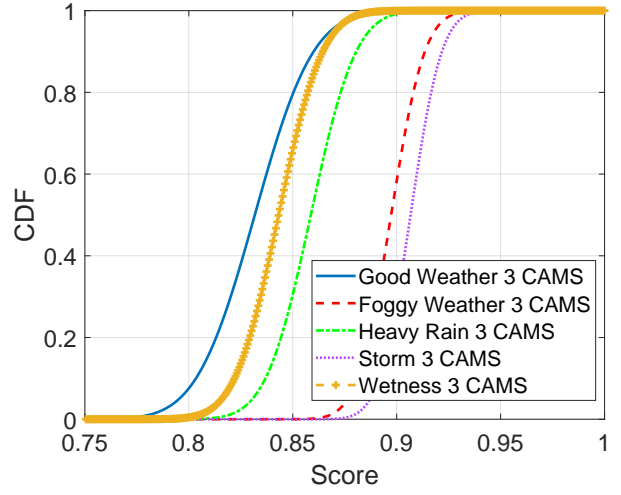
triangle). This suggests that the MLLM-AD-4o driving agent considers heavy rain the most unsafe environment requiring low-speed traveling. Also, traveling can be smooth in stormy weather, while driving is most efficient in good weather. Finally, in foggy/heavy rain/storm conditions, the agent is more likely to respect the speed limits as it favors slower driving to guarantee security, which is not the case in good and wet weather conditions. These results reflect the trade-offs the MLLM-AD-4o agent makes between various performance scores, under different weather patterns.

In Figs. 5a-5c, we depict the CDFs of the safety, comfort, and efficiency scores when 6 cameras (3 front + 3 rear) are used by the AD agent. In good weather, MLLM-AD-4o achieves now the best safety CDF performance (blue triangle), at the expense of comfort and efficiency. Indeed, the additional information brought by the 3 rear cameras makes the agent more precautious regarding the behavior of the surrounding vehicles, which was not possible when it only relied on 3 cameras. This is confirmed with the speed CDF result (purple line) in Fig. 4d which is higher than that of the 3 cameras' CDF. According to Fig. 5b, the best CDFs are realized in stormy/foggy weather (with a preference for stormy) at the expense of efficiency and safety. Indeed, in these situations, the AD agent is more likely to drive smoothly with minor behavior changes in response to poor visibility, with a higher safety risk and reduced efficiency, due to the misalignment with the surrounding road traffic behavior. Moreover, the best efficiency CDF is obtained in wet conditions, as shown in Fig. 5c. Indeed, with the knowledge of its 6 cameras, the AV maintains a speed closer to the average of other vehicles by continuously adapting to the dynamics of the road and surrounding traffic. This is validated through the results in Fig. 4d where the 6 cameras' speed CDF is better than the 3 cameras' one in wetness. However, as illustrated in Figs. 5a-5b, this behavior negatively impacts the safety and comfort scores. In summary, the addition of the rear cameras enhances the vehicle's ability to perceive and respond to potential collision threats, particularly in challenging weather conditions.
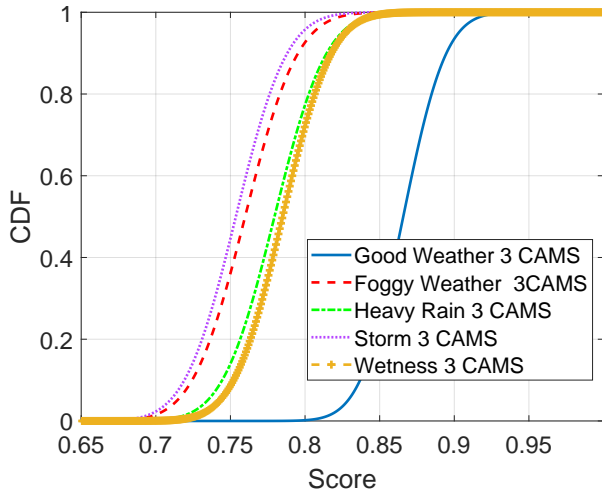
Figs. 6a-c illustrate the CDFs of the safety, comfort, and efficiency scores, when the AD agent is equipped with 3 (front) cameras only, 6 (3 front + 3 rear) cameras only, 3 (front) cameras + LiDAR, and 6 (3 front + 3 rear) cameras + LiDAR, respectively, operating in a heavy rain condition. The best safety CDF is obtained when the LiDAR is used with 3 cameras (green triangle in Fig. 6a). Indeed, activating the LiDAR in adverse weather conditions enhances the system's awareness of its environment leading to more responsible driving actions, and thus to an improved safety. Nevertheless, safety and efficiency CDFs degrade when LiDAR is activated with 6 cameras. This suggests that rear cameras' information might distort the agent's comprehension of its surroundings when combined with LiDAR. The best comfort score is also achieved when 3 cameras are activated with LiDAR (green line in Fig. 6b), while the best efficiency score is realized with a 3-camera system (blue line in Fig. 6c). Nevertheless, the activation of LiDAR with the 3 cameras (green line in Fig. 6c) realizes similar efficiency CDF performance to
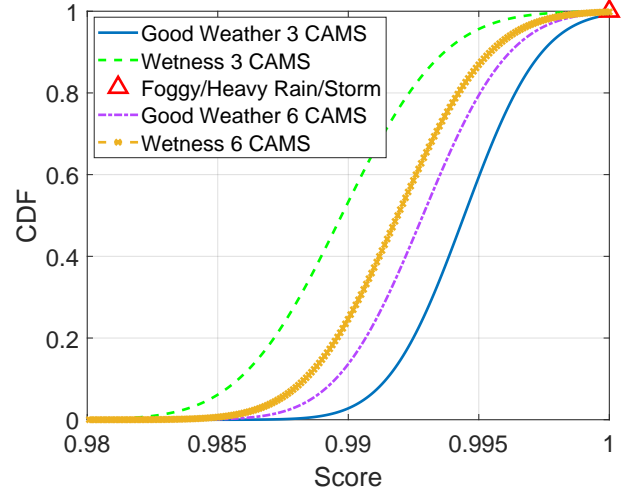
(a) CDF of safety score

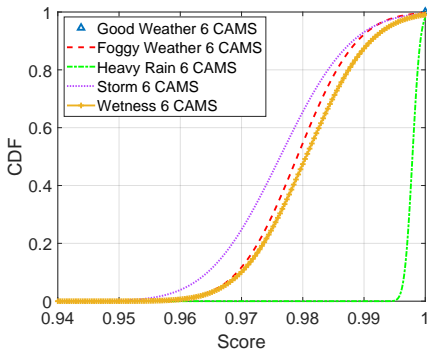(b) CDF of comfort score

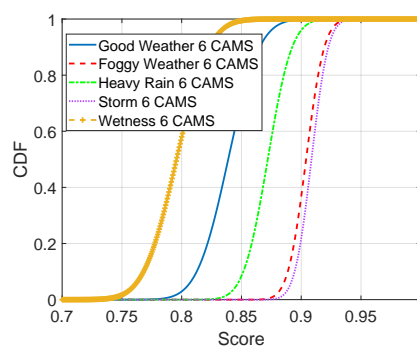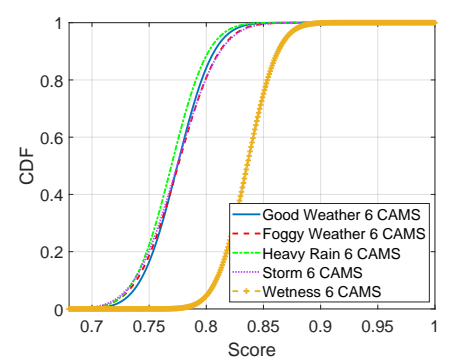(c) CDF of efficiency score

(d) CDF of speed score

Fig. 4: CDFs of MLLM-AD-4o scores (Front cameras).



(a) CDF of safety score

(b) CDF of comfort score

(c) CDF of efficiency score

Fig. 5: CDFs of MLLM-AD-4o scores (Front+rear cameras).

(a) CDF of safety score     (b) CDF of comfort score     (c) CDF of efficiency score
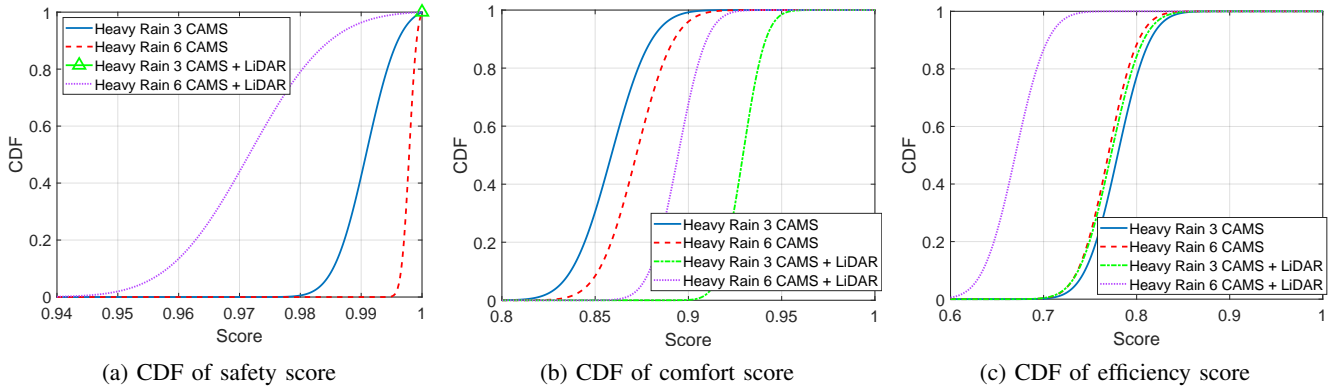
Fig. 6: CDFs of MLLM-AD-4o scores (Combinations of cameras with LiDAR; Heavy rain).

the latter. Based on these results, the activation of LiDAR with cameras should be triggered in a timely manner, in particular conditions, and with specific configurations, as the full combination of cameras with LiDAR does not necessarily lead to better AD performances.

## VI. CONCLUSION

Following integrating the harsh environment conditions into the CARLA/Limsim++ simulation framework, we proposed a novel AD agent based on the GPT-4o MLLM model for driving perception and control decision-making, a.k.a., MLLM-AD-4o. Through extensive experiments, we assessed the performances of MLLM-AD-4o in terms of safety, comfort, efficiency, and speed score CDFs, for different weather conditions, and given diverse combinations of cameras and LiDAR activated at the AV. Our findings indicate that activating both front and rear cameras significantly enhances the system's performance, compared to the case where only the front cameras are used, especially in good weather conditions. Moreover, activating LiDAR with only 3 cameras provides the best performance compared to the cases where only cameras are used. However, when LiDAR is used with all 6 cameras, the AV's performance degrades, suggesting that the right combination of cameras and LiDAR should be activated in specific configurations and weather conditions. The obtained results draw important insights into the efficient design of AD systems, where careful data fusion should be accurately realized to optimize autonomous driving behavior.

## ACKNOWLEDGMENT

## REFERENCES

[1] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang *et al.*, "A survey on evaluation of large language models," *ACM Trans. Intelli. Syst. Technol.*, vol. 15, no. 3, pp. 1–45, 2024.

[2] S. Fourati, W. Jaafar, N. Baccar, and S. Alfattani, "XLM for autonomous driving systems: A comprehensive review," *arXiv preprint arXiv:2409.10484*, 2024.

[3] Y. Wang, W. Chen, X. Han, X. Lin, H. Zhao, Y. Liu, B. Zhai, J. Yuan, Q. You, and H. Yang, "Exploring the reasoning abilities of multimodal large language models (MLLMs): A comprehensive survey on emerging trends in multimodal reasoning," *arXiv preprint arXiv:2401.06805*, 2024.

[4] C. Cui, Y. Ma, X. Cao, W. Ye, Y. Zhou, K. Liang, J. Chen, J. Lu, Z. Yang, K.-D. Liao *et al.*, "A survey on multimodal large language models for autonomous driving," in *Proceed. of the IEEE/CVF Winter Conf. Appl. Comput. Vis.*, 2024, pp. 958–979.

[5] L. Wen, D. Fu, X. Li, X. Cai, T. Ma, P. Cai, M. Dou, B. Shi, L. He, and Y. Qiao, "Dilu: A knowledge-driven approach to autonomous driving with large language models," *arXiv preprint arXiv:2309.16292*, 2023.

[6] X. Ding, J. Han, H. Xu, W. Zhang, and X. Li, "HiLM-D: Towards high-resolution understanding in multimodal large language models for autonomous driving," *arXiv preprint arXiv:2309.05186*, 2023.

[7] J. Yuan, S. Sun, D. Omeiza, B. Zhao, P. Newman, L. Kunze, and M. Gadd, "RAG-Driver: Generalisable driving explanations with retrieval-augmented in-context learning in multi-modal large language model," in *Proc. Robotics: Science and Syst. (RSS) Conf.*, 2024, pp. 1–14.

[8] T. Wang, E. Xie, R. Chu, Z. Li, and P. Luo, "Drivecot: Integrating chain-of-thought reasoning with end-to-end driving," *arXiv preprint arXiv:2403.16996*, 2024.

[9] W. Wang, J. Xie, C. Hu, H. Zou, J. Fan, W. Tong, Y. Wen, S. Wu, H. Deng, Z. Li *et al.*, "DriveMLM: Aligning multi-modal large language models with behavioral planning states for autonomous driving," *arXiv preprint arXiv:2312.09245*, 2023.

[10] Y. Huang, J. Sansom, Z. Ma, F. Gervits, and J. Chai, "DriVLMe: Exploring foundation models as autonomous driving agents that perceive, communicate, and navigate," in *Proc. Vis. and Langua. for Autonom. Driv. and Roboti. Wrkshp.*, 2024.

[11] G. Liao, J. Li, and X. Ye, "VLM2Scene: Self-supervised image-text-LiDAR learning with foundation models for autonomous driving scene understanding," in *Proc. AAAI Conf. on Artifi. Intelli.*, vol. 38, no. 4, 2024, pp. 3351–3359.

[12] S. Wang, Z. Yu, X. Jiang, S. Lan, M. Shi, N. Chang, J. Kautz, Y. Li, and J. M. Alvarez, "OmniDrive: A holistic LLM-agent framework for autonomous driving with 3D perception, reasoning and planning," *arXiv preprint arXiv:2405.01533*, 2024.

[13] H. Tian, K. Reddy, Y. Feng, M. Quddus, Y. Demiris, and P. Angeloudis, "Enhancing autonomous vehicle training with language model integration and critical scenario generation," *arXiv preprint arXiv:2404.08570*, 2024.

[14] J. Mao, J. Ye, Y. Qian, M. Pavone, and Y. Wang, "A language agent for autonomous driving," in *Proc. Conf. Langua. Model. (COMS)*, 2024.

[15] J. Z. Z. H. A. Ray and E. Ohn-Bar, "Feedback-guided autonomous driving," 2024.

[16] H. I. Ashqar, T. I. Alhadidi, M. Elhenawy, and N. O. Khanfar, "Leveraging multimodal large language models (MLLMs) for enhanced object detection and scene understanding in thermal images for autonomous driving systems," *Automat.*, vol. 5, no. 4, pp. 508–526, 2024.

[17] S. Luo, W. Chen, W. Tian, R. Liu, L. Hou, X. Zhang, H. Shen, R. Wu, S. Geng, Y. Zhou *et al.*, "Delving into multi-modal multi-task foundation models for road scene understanding: From learning paradigm perspectives," *IEEE Trans. Intelli. Veh. (Accepted)*, 2024.

[18] R. Li, Z. Zhang, C. He, Z. Ma, V. M. Patel, and L. Zhang, "Dense multimodal alignment for open-vocabulary 3D scene understanding," *arXiv preprint arXiv:2407.09781*, 2024.

[19] Z. Xu, Y. Zhang, E. Xie, Z. Zhao, Y. Guo, K.-Y. K. Wong, Z. Li, and H. Zhao, "DriveGPT4: Interpretable end-to-end autonomous driving via large language model," *IEEE Robot. Automat. Lett.*, vol. 9, no. 10, pp. 8186–8193, Oct. 2024.

[20] C. Cui, Y. Ma, X. Cao, W. Ye, and Z. Wang, "Receive, reason, and react: Drive as you say, with large language models in autonomous vehicles," *IEEE Intelli. Transport. Syst. Mag.*, 2024.

[21] Z. Li, W. Wang, Y. Cai, X. Qi, P. Wang, D. Zhang, H. Song, B. Jiang, Z. Huang, and T. Wang, "UnifiedMLLM: Enabling unified representation for multi-modal multi-tasks with large language model," *arXiv preprint arXiv:2408.02503*, 2024.

[22] "MLLM applied to autonomous driving across various weather conditions," https://github.com/SondaFourati/MLLM-applied-to-autonomous-driving-across-various-weather-conditions/tree/main, accessed: 2024-11-10.

[23] R. Islam and O. M. Moushi, "GPT-4o: The cutting-edge advancement in multimodal LLM," *Authorea Preprints*, 2024.

[24] L. Chen, O. Sinavski, J. Hünermann, A. Karnsund, A. J. Willmott, D. Birch, D. Maund, and J. Shotton, "Driving with LLMs: Fusing object-level vector modality for explainable autonomous driving," in *Proc. IEEE Int. Conf. Robot. Automat. (ICRA)*, May 2024, pp. 14 093–14 100.

[25] Y. Jin, X. Shen, H. Peng, X. Liu, J. Qin, J. Li, J. Xie, P. Gao, G. Zhou, and J. Gong, "Surrealdriver: Designing generative driver agent simulation framework in urban contexts based on large language model," *arXiv preprint arXiv:2309.13193*, 2023.

[26] R. Yang, X. Zhang, A. Fernandez-Laaksonen, X. Ding, and J. Gong, "Driving style alignment for llm-powered driver agent," *arXiv preprint arXiv:2403.11368*, 2024.

[27] Z. Guo, A. Lykov, Z. Yagudin, M. Konenkov, and D. Tsetserukou, "Co-driver: VLM-based autonomous driving assistant with human-like behavior and understanding for complex road scenes," *arXiv preprint arXiv:2405.05885*, 2024.

[28] D. Wu, W. Han, T. Wang, Y. Liu, X. Zhang, and J. Shen, "Language prompt for autonomous driving," *arXiv preprint arXiv:2309.04379*, 2023.

[29] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proc. Conf. Robot Learn.* PMLR, 2017, pp. 1–16.

[30] A. Kusari, P. Li, H. Yang, N. Punshi, M. Rasulis, S. Bogard, and D. J. LeBlanc, "Enhancing SUMO simulator for simulation based testing and validation of autonomous vehicles," in *Proc. IEEE Intelli. Veh. Symp. (IV)*, 2022, pp. 829–835.

[31] D. Fu, W. Lei, L. Wen, P. Cai, S. Mao, M. Dou, B. Shi, and Y. Qiao, "LimSim++: A closed-loop platform for deploying multimodal LLMs in autonomous driving," in *Proc. IEEE Intelli. Veh. Symp. (IV)*, 2024, pp. 1084–1090.