

PDB-Eval: An Evaluation of Large Multimodal Models for Description and Explanation of Personalized Driving Behavior

Junda Wu
Computer Science and Engineering
University of California San Diego
La Jolla, USA
juw069@ucsd.edu

Jessica Echterhoff
Computer Science and Engineering
University of California San Diego
La Jolla, USA
jechterh@ucsd.edu

Kyungtae Han
InfoTech Labs
Toyota Motor North America
Mountain View, USA
kt.han@toyota.com

Amr Abdelraouf
InfoTech Labs
Toyota Motor North America
Mountain View, USA
amr.abdelraouf@toyota.com

Rohit Gupta
InfoTech Labs
Toyota Motor North America
Mountain View, USA
rohit.gupta@toyota.com

Julian McAuley
Computer Science and Engineering
University of California San Diego
La Jolla, USA
jmcauley@eng.ucsd.edu

Abstract—Understanding a driver’s behavior and intentions is important for potential risk assessment and early accident prevention. Safety and driver assistance systems can be tailored to individual drivers’ behavior, significantly enhancing their effectiveness. However, existing datasets are limited in describing and explaining general vehicle movements based on external visual evidence. This paper introduces a benchmark, PDB-Eval, for a detailed understanding of Personalized Driver Behavior, and aligning Large Multimodal Models (MLLMs) with driving comprehension and reasoning. Our benchmark consists of two main components, PDB-X and PDB-QA. PDB-X can evaluate MLLMs’ understanding of temporal driving scenes. Our dataset is designed to find valid visual evidence from the external view to explain the driver’s behavior from the internal view. To align MLLMs’ reasoning abilities with driving tasks, we propose PDB-QA as a visual explanation question-answering task for MLLM instruction fine-tuning. As a generic learning task for generative models like MLLMs, PDB-QA can bridge the domain gap without harming MLLMs’ generalizability. Our evaluation indicates that fine-tuning MLLMs on fine-grained descriptions and explanations can effectively bridge the gap between MLLMs and the driving domain, which improves zero-shot performance on question-answering tasks by up to 73.2%. We further evaluate the MLLMs fine-tuned on PDB-X in Brain4Cars’ intention prediction and AIDE’s recognition tasks. We observe up to 12.5% performance improvements on the turn intention prediction task in Brain4Cars, and consistent performance improvements up to 11.0% on all tasks in AIDE.

Index Terms—MLLMs, Personalized Driving Behavior, Question-answering

I. INTRODUCTION

Driver understanding is critical for predicting vehicle movement [1] and assessing potential risks on the road. Recent research has shown advances in recognition tasks of traffic accidents, uncertainty, and vehicle motion prediction [2]. In such tasks, textual explanations for vehicle movement [3] and driver behavior [2] have become increasingly

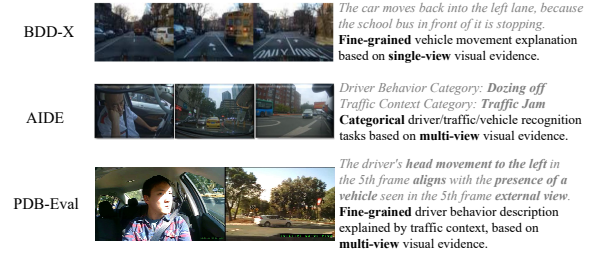


Fig. 1. Existing vehicle movement understanding (e.g., BDD-X [3]) can be fine-grained but limited to single-view explanations. Multi-view and multi-task driving tasks (e.g., AIDE [2]) can provide multi-view visual evidence but lack description granularity. PDB-Eval can provide driver behavior descriptions explained by traffic context based on multi-view visual evidence.

important for more explainable understanding. Since large multimodal models (MLLMs) can generate descriptions and explanations over visual evidence [4]–[6], MLLMs have been recognized as a multimodal reasoner for driving tasks [6], [7]. However, existing MLLMs are limited to fine-tuning and evaluation on general visual understanding and explanation tasks [7]–[9], where a domain gap exists while adapting to driving tasks.

To align MLLMs with driver understanding and reasoning tasks, we propose an evaluation dataset, **PDB-Eval**, for personalized driver behavior understanding. While many existing works focus on analyzing driver behavior from the in-cabin view [1], [2] or the external view of traffic scenes [3], [10], few have tried to understand these viewpoints as parallel time processes, where a video moment (i.e., an event happening between a time boundary) in one viewpoint can be explained based on the other viewpoint’s visual evidence. For example, in Figure 1, existing fine-grained explanations of vehicle movements (e.g., BDD-X [3]) are limited to single-view understanding. Multi-view driving understanding datasets (e.g., AIDE [2]) enable the understanding of

multiple categories of driver behavior and vehicle movement, but they fall short in offering fine-grained descriptions. To unify driving behavior description and explanation from both internal and external visual evidence, we introduce our first evaluation task **PDB-X**. For example, in Figure 1, we show that PDB-X provides fine-grained driving behavior explanations based on multi-view visual evidence. We further propose the visual explanation question-answering task **PDB-QA** to enhance the interpretative and reasoning capabilities of MLLMs in the context of driver behavior.

Extracting personalized driver behavior descriptions can be challenging since human annotators must describe based on knowledge and observations of different types of drivers. We suggest using MLLMs to create personalized driver behavior descriptions by comparative prompting [11], which highlights behavior discrepancies between drivers with the same intention, ensuring specificity and relevance in the descriptions. Despite the potential of this approach, it is crucial to address the inherent challenge of hallucination in MLLM-generated content, where models may produce fabricated or irrelevant information [8], [12]. This issue is particularly significant in tasks requiring detailed and accurate descriptions based on visual evidence [8], [12]. To reduce this risk, we categorize prompted answers for each driver type, and then use these categories to guide more focused MLLM prompts, minimizing irrelevant information.

We evaluate two tasks in PDB-Eval by fine-tuning open-source MLLMs, BLIP-2 and VTimeLLM, which are specialized in image understanding and video understanding, respectively. We also include GPT-4V as a strong zero-shot baseline to demonstrate the performance of the existing MLLMs. In addition, we further evaluate the driver intention prediction tasks in the Brain4Cars dataset [1] (in-domain) and driving recognition tasks in the AIDE dataset [2] (across-domain). The evaluation results showcase the effectiveness and generalizability of MLLMs’ acquired descriptive and explaining abilities via training on PDB-X.

II. RELATED WORK

Recent driving understanding research spans tasks such as time-to-accident, intention, accident prediction, driving anticipation, distraction, and uncertainty estimation [10], [13]. Despite progress, many systems lack fine-grained interpretability and reasoning, limiting trust in applications like self-driving cars. To address this, recent approaches incorporate textual explanations, question-answering based on dash-cam evidence [3], [14], attention maps, and even concept bottleneck frameworks [15] that verbalize drivers’ behavior. In contrast to static, categorical descriptions, our benchmark challenges models to understand personalized driver behavior and correlate it with dynamic vehicle and traffic conditions.

Large multimodal models (MLLMs) have been increasingly applied in driving-related tasks such as planning, navigation, simulation, and command understanding [16], [17]. However, while existing video-based MLLM evaluations primarily focus on single-view analysis, a key challenge remains: reasoning over two causally linked temporal processes—the internal dynamics of driver behavior and

external traffic changes. Our integrated evaluation task thus requires MLLMs to synchronously interpret both streams and provide explanations grounded in visual evidence. This dual-view approach not only enhances the interpretability of driver understanding systems but also paves the way for safer and more transparent human-centered driving technologies.

III. METHODS

A. Pipeline Overview

Our dataset creation pipeline is illustrated in Figure 2, where we identify five data processing steps as follows:

- **Step 1 Comparative Prompting:** To understand personalized driver behavior, we propose comparative prompting in MLLMs for descriptions of driver behavior discrepancies while performing the same action.
- **Step 2 Driver’s Identity & Intention Consistency Filtering:** We extract the intention and driver identities from the generated descriptions as validation measurements for automatic sample filtering.
- **Step 3 Guideline Construction:** Instead of directly using the extracted descriptions from comparative prompting, we extract driving behavior types and categorize descriptions as prompting guidelines, which are post-processed by text-only LLMs.
- **Step 4 Guideline-instruction Prompting:** By prompting MLLMs with the generation guidelines [18], we can focus more effectively on specific driver behavior characteristics. This approach leads to more fine-grained and personalized driver behavior, which reduces the chance of hallucinations.
- **Step 5 Human Annotator Filtering:** However, MLLMs may still hallucinate non-existing visual contexts in the limited visual field of internal and external dash-cams. Therefore, the human annotation involves aspect-level and sample-level filtering for invisible and irrelevant behavior description detection.

With the proposed pipeline, we augment dual-cam videos from the Brain4Cars [1] dataset, with personalized driver behavior descriptions and explanations. The original Brain4Cars dataset was collected from 10 different drivers and segmented into 700 event-based clips with 274 lane changes, 131 turns, and 295 driving straight instances [1].

B. Comparative Prompting

To obtain drivers’ personalized behavior, we propose comparative prompting in MLLMs to extract behavior discrepancies between two drivers with the same intention. For each pair of drivers u, v with the same intention, we first extract N frames from internal dash-cam video frames $\mathbf{I}_u^{in}, \mathbf{I}_v^{in} \in V_i$ as the visual evidence. To construct a compatible visual context for MLLM prompting, temporal frames of drivers’ videos are concatenated (*e.g.*, Step 1 of Figure 2),

$$\mathbf{I}_{u,v}^{in} = \text{vconcat} [\text{hconcat}(\mathbf{I}_u^{in}), \text{hconcat}(\mathbf{I}_v^{in})],$$

in which vconcat and hconcat denote vertical and horizontal concatenations of image frames, respectively.

Since MLLMs are developed from LLMs with multi-modality alignment, MLLMs can have strong textual priors

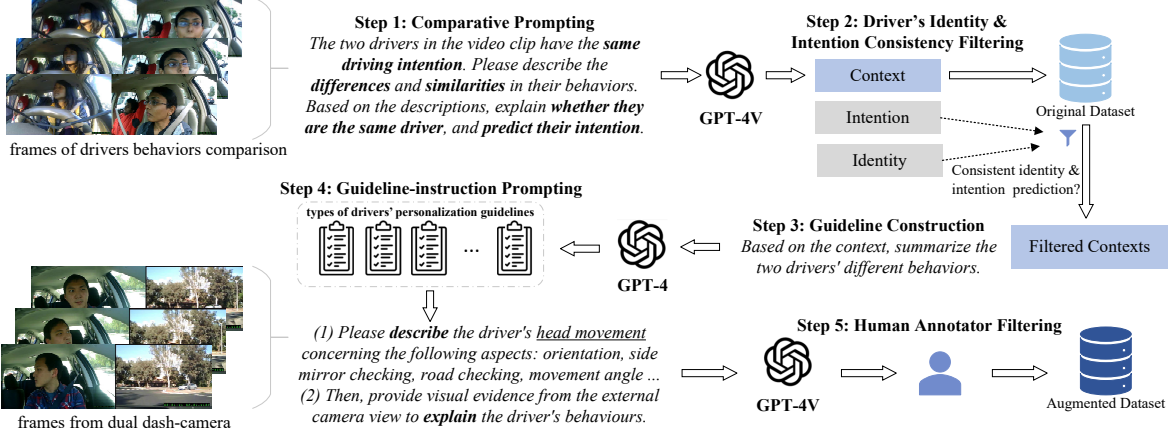


Fig. 2. Illustration of the pipeline for creating a dataset that captures drivers' personalized behavior characteristics, including five key steps: comparative prompting, consistency-based sample filtering, summarizing into generation guidelines, fine-grained behavior generation, and human annotator filtering for quality assurance

(i.e., linguistic bias) [8], which may lead to hallucinated responses neglecting visual evidence [8], [12]. Specifically, previous works point out that failing to instruct the models to focus on fine-grained aspects of visual semantics can result in hallucinations [19]. Thus, in this step, we propose to prompt the MLLMs to generate only comparative descriptions between drivers, which make MLLMs perceive more fine-grained visual details [11] and potentially prevent hallucination problems [19], [20].

The comparative prompt design in our pipeline consists of 4 main instructions. The **explanation** instruction \mathbf{T}_x is to explicitly reason on the sequence of concatenated video frames as videos. Then, the MLLM is prompted with the **comparison** instruction \mathbf{T}_c to answer the discrepancies between the two drivers' behavior. To derive the driver's identity and intention consistency indexes, we prompt the MLLM with both **identity** \mathbf{T}_{id} and **intention** \mathbf{T}_{it} extraction instructions based on the model's previous comparative descriptions. By the driver's identity and intention consistency evaluation, we can improve the accuracy of the MLLM's responses. Together with the visual evidence and textual instructions, we prompt the MLLM to generate initial responses,

$$\mathbf{R}_{u,v} = \text{MLLM}(\mathbf{I}_{u,v}^{in}, [\mathbf{T}_x; \mathbf{T}_c; \mathbf{T}_{id}; \mathbf{T}_{it}]),$$

in which we make sure the drivers' intentions are the same.

C. Identity and Intention Consistency Filtering

We introduce an intermediate sample verification and filtering step based on the driver's identity and intention consistency. The **context** instruction \mathbf{T}'_c is to reformat the comparative descriptions into drivers' personalized behavior types and aspects into dictionaries, which leverages the LLM's structural text generation capacities. We also prompt the textual LLM to extract the **identity** and **intention** predictions following the instructions \mathbf{T}'_{id} and \mathbf{T}'_{it} respectively. By prompting the textual LLM with instructions and the previous response $\mathbf{R}_{u,v}$, we can extract the personalization context $\mathbf{C}_{u,v}$, drivers' identity indicators $\mathbf{1}_{u,v}^{ID}$, and intentions $\mathbf{IT}_u, \mathbf{IT}_v$ respectively. Then, the personalization

contexts are filtered and aggregated into P according to the identity and intention consistency,

$$P = \bigcup_{\substack{\mathbf{I}_u, \mathbf{I}_v \in V_i, \\ \mathbf{1}_{u,v}^{ID} = \text{ID}(u,v), \mathbf{IT}_u = \mathbf{IT}_v}} \mathbf{C}_{u,v},$$

in which the identity is to check if the generated identity-matching result $\mathbf{1}_{u,v}^{ID}$ is the same as the ground truth identities $\text{ID}(u,v)$. Since comparative prompting compares two drivers with the same intention, we validate if the generated intentions $\mathbf{IT}_u, \mathbf{IT}_v$ are also the same.

D. Guideline Instruction Construction and Prompting

The MLLM can generate detailed contrastive characteristics in drivers' behavior by comparative prompting. However, in practice, we observe hallucinations in descriptions of visual details. Due to lacking granularity in visual understanding, the MLLM can respond with visual details that do not exist in the actual visual context [11], [19], [20]. To make the prompt specific in describing a certain type of driver's behavior and prominent visual characteristics of the behavior, we propose aggregating the collected personalization contexts D as the prompting guideline instructions. Following the guideline instructions for each type of personalized behavior, the generated responses contain less irrelevant information and hallucinations [18].

The constructed guideline instructions are types of personalized driver behavior paired with detailed aspects of these types of behavior (e.g., in Step 4 of Figure 2). During guideline instruction prompting, the MLLM is individually prompted with each type of behavior with its descriptions, to prevent irrelevant information and hallucinations. In addition, we also prompt the MLLM to explain the driver's behavior description based on visual evidence from the external dash-cam frames (e.g., in Step 4 of Figure 2). To enable the MLLM to perceive both the internal and external visual evidence based on video frames, we use a similar way to construct the visual evidence. Given a driver's internal $\mathbf{I}_u^{in} \in V_i$ and external $\mathbf{I}_u^{ex} \in V_e$ dash-cam video frames,

$$\mathbf{I}_u^{du} = \text{vconcat}[\text{hconcat}(\mathbf{I}_u^{in}), \text{hconcat}(\mathbf{I}_u^{ex})].$$



Comparative Prompting:

Both individuals are seated in the driver's seat of a vehicle and focused on an activity with their hands. The individual on the top row appears to be leaning forward slightly more than the individual on the bottom row. The individual on the top row has a more visible and expressive face, whereas the individual on the bottom row has a more neutral expression ...



Guideline Instruction Prompting (Head Movement):

Description D: Initially, the driver appears to be looking straight ahead. As the frames progress, his gaze shifts to the left, indicating the checking of the side mirror or a blind spot ...

Explanation E: The traffic light is visible in the sequence, and the driver's initial straight-ahead gaze may correspond to watching the traffic light for changes. The sequence of images seems to capture the driver as he approaches and then waits at an intersection, possibly preparing for a left turn, as indicated by his head turning to check for oncoming traffic or pedestrians.

Fig. 3. A comparative example of comparative prompted description and guideline instruction prompted description. The context (underlined and in red) indicates possible hallucinations by **comparative prompting**. With fine-grained guideline instructions, more detailed visual evidence and temporal information (underlined and in blue) can be prompted from LLMs in **Description D**. Then, the LLMs can capture visual evidence from an external view and explain the driver's behavior (underlined and in green) in **Explanation E**.

For each behavior type $k \in K$ and its description $S(k)$, the prompt $T_g(k, S(k))$ is the guideline instruction for prompting. With the guideline instructions, we can prompt the MLLM to obtain the descriptions \mathbf{D}_u and explanations \mathbf{E}_u ,

$$\mathbf{D}_u, \mathbf{E}_u = \text{MLLM}(\mathbf{I}_u^{du}, T_g(k, S(k))).$$

In Figure 3, we show a comparative example of the difference between the description prompted from comparative prompting and guideline instruction prompting. In comparative prompting, we can observe only coarse-grained aspects (e.g., seating position, facial expression, etc.) are mentioned without further justification. In such descriptions, we can observe possible hallucinations generated from MLLMs (e.g., the first driver is leaning forward more than the other driver). In the guideline instruction prompted responses, we can observe more fine-grained descriptions about the driver's behavior changing with time, which is also specific in a certain aspect without irrelevant information. The MLLM can further try to explain the generated description based on corresponding visual evidence from the external dashcam frames.

E. Human Annotator Filtering

One of the sample filtering tasks is to detect irrelevant or inaccurate types of behavior, namely aspect-level sample filtering. We observe several extracted behaviors (e.g., feet position, leg movement) that cannot be accurately described due to limited visual field from the internal dashcam. In such cases, the MLLM will hallucinate inaccurate visual descriptions and provide false explanations for the descriptions. In addition, some other extracted behaviors (e.g., driver's appearance, attire, etc.) are irrelevant to the

TABLE I
BEHAVIOR TYPES AFTER THE HUMAN ANNOTATOR'S FILTERING.

ACT	action	BOL	body language
DRS	driving style	FAE	facial expression
GEM	gaze/eye movement	HAM	hand movement
HEM	head movement	INT	intention
IWP	interaction w/ passenger		

performance of driving tasks. In such cases, the explanations are not based on visual evidence but on MLLMs' strong textual priors (i.e., linguistic bias) [8].

After aspect-level filtering, we can obtain nine types of behavior in Table I. The sample-level filtering task for human annotators is mainly to filter out failure cases from MLLMs. Since these failure cases of generation follow a similar pattern of expression (e.g., "I'm sorry, but I cannot provide ..."), including such samples in model fine-tuning may result in overfitting problems. Thus, we further conduct the sample-level inspection of these specific failure cases.

IV. PERSONALIZED DRIVING BEHAVIOR EVALUATION (PDB-EVAL)

Based on our data creation pipeline in Section III-A, we develop our evaluation dataset from the existing Brain4Cars dataset [1]. In Figure 4, we illustrate an example in our constructed dataset for these three tasks. In comparative prompting, we sampled 20 pairs of drivers from each annotated driving intention subset: right turn, left turn, right change, left change, and driving straight. Through comparative prompting, the collected comparative descriptions are summarized into **nine types** in Table I with an average of **19.11 guideline instructions** for each type. We summa-

TABLE II
THE NUMBER OF CLIPS, DESCRIPTION-EXPLANATION PAIRS (D/E), QUESTION-ANSWERING PAIRS (QA), AVERAGE WORDS IN DESCRIPTIONS (DESC.), EXPLANATIONS (EXPL.), QUESTIONS (Q), AND ANSWERS (A).

	Sample Number			Average Length			
	Clips	D/E	QA	Desc.	Expl.	Q	A
Train	478	5,084	35,972	125.06	100.45	13.14	37.69
Test	116	1,224	8,881	126.20	100.27	13.13	37.59

size the statistics of PDB-Eval in Table II, where we report the number of clips, description-explanation pairs (D/E), and question-answering pairs (QA). We also report the average number of words in descriptions (Desc.), explanations (Expl.), questions (Q), and answers (A).

In the PDB-X dataset, we observe that the average length of the descriptions is longer than that of the explanations, which is because of the usage of the guideline instructions. Based on each requirement in the guideline, the MLLM will generate an average of **6.54 words per instruction** for the description. As for the PDB-QA dataset, the average number of questions derived from the PDB-X dataset is **7.05 QA pairs per D/E pair**, and the average length of the ground-truth answers is around **3 times the length** of questions. Based on the descriptions and explanations from PDB-X, we can enable automatic question-answer pair generation using text-only LLMs [21], [22]

V. EXPERIMENTS

We evaluate various MLLMs: BLIP-2 [23], VTimeLLM [4], and GPT-4V, on PDB-X and PDB-QA. By fine-tuning (BLIP-2 and VTimeLLM) on PDB-X, we further evaluate the MLLMs on downstream driving benchmarks, Brain4Cars [1] (in-domain) and AIDE [2] (cross-domain), to understand the effectiveness of the generated drivers' personalized behavior descriptions and explanations.

A. Data Preprocessing

We adopt different video preprocessing methods for three baselines to extract compatible visual representations of the input streams. For BLIP-2, we extract 10 image frames with an equal time interval from dual-cam and concatenate them as in Eq. (III-D) before encoding the image using the CLIP encoder from BLIP-2 [23]. For VTimeLLM, we follow the preprocessing method in [4] to extract 10 image frames from each second of the video and encode the video with a ViT-based CLIP model. For GPT-4V, we extract 10 consecutive image frames from the videos and concatenate each frame of the internal and external videos together.

B. Description and Explanation (PDB-X)

We evaluate the MLLM baselines for description and explanation tasks in PDB-X and report the BLEU-4 metric of comparative results in Table III. For open-source MLLMs (BLIP-2 and VTimeLLM) we first fine-tune the models on the training set of PDB-X before evaluation. During inference, we adopt the same instruction design and visual encoding as in the fine-tuning stage. Comparing the performance of fine-tuned MLLMs and GPT-4V, we can observe

a consistent advantage of fine-tuning on PDB-X. Such an observation suggests that prompting MLLMs with simple instructions cannot directly enable fine-grained descriptions and explanations, which showcases the effectiveness of the proposed comparative prompting method and the challenges posed by PDB-X tasks.

TABLE III
PERFORMANCE OF BLEU-4 METRICS ON ALL TYPES OF DRIVER BEHAVIOR (IN TABLE I) AND THE TWO TASKS, DESCRIPTION (DESC.) AND EXPLANATION (EXPL.). WE INDICATE THE BEST PERFORMANCE FOR THE DESCRIPTION AND EXPLANATION TASKS IN BOLD.

Type	BLIP-2		VTimeLLM		GPT-4V	
	Desc.	Expl.	Desc.	Expl.	Desc.	Expl.
ACT	51.10	48.09	47.18	50.41	29.49	28.57
BOL	34.51	52.74	49.13	62.03	33.68	27.22
DRS	72.14	58.50	86.03	54.63	32.33	30.12
FAE	54.52	59.90	72.17	50.40	34.44	28.54
GEM	55.06	57.83	62.72	60.70	39.22	35.73
HAM	46.42	70.50	42.93	44.98	40.96	36.54
HEM	72.20	51.94	72.90	46.71	39.87	36.77
INT	56.52	60.97	53.80	49.22	25.69	26.67
IWP	76.39	66.53	73.80	65.95	37.23	32.04
AVE	57.65	58.55	62.30	53.89	34.77	31.35

Comparing the performance of BLIP-2 and VTimeLLM, we can observe that these two MLLMs are specialized in the explanation and description tasks respectively. We argue that BLIP-2 is pre-trained with enhanced visual knowledge reasoning capacity [23] that benefits visual explanation, while the static modality of the model input limits its description capacity. On the other hand, VTimeLLM is specially pre-trained with the temporally-aware descriptive capacity for fine-grained video moment understanding [4], which explains its better performance on the description task. Nevertheless, creating an MLLM that achieves robust performance in both tasks remains challenging.

C. Visual Explanation QA (PDB-QA)

We further evaluate MLLMs on PDB-QA to showcase their capacities in answering complex questions related to personalized driving behavior. The relatively lower performance on PDB-QA suggests a deficiency of the MLLMs in understanding drivers' personalized information. However, with model fine-tuning, we can still observe consistent improvements in both open-source MLLMs.

Comparing the performance of BLIP-2 and VTimeLLM, we can observe that VTimeLLM has a larger fine-tuning

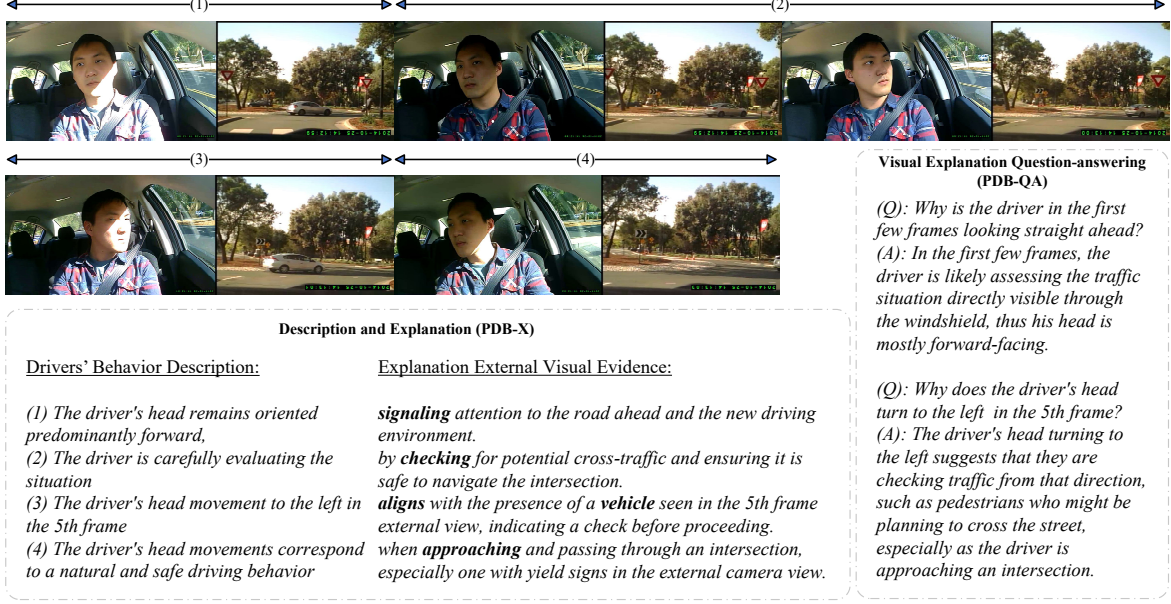


Fig. 4. Illustration of three tasks for drivers' personalized behaviors understanding.

TABLE IV
PERFORMANCE OF BLEU-4 METRICS ON THE PDB-QA DATASET, INCLUDING PRE-TRAINED (PT), FINE-TUNED (FT), AND ZERO-SHOT (ZS).

Metric	BLIP-2		VTimeLLM		GPT-4V
	PT ↑	FT ↑	PT ↑	FT ↑	ZS ↑
BLEU-4	30.29	33.64	28.07	48.61	16.08

improvement compared to its zero-shot performance, which suggests that VTimeLLM has a higher performance upper-bound than BLIP-2. The temporal reasoning capacity in VTimeLLM can enhance the video moment understanding between a certain time boundary [4]. However, how to enable cross-validation of the reasoning in video clips within different time boundaries for better personalization understanding can still be challenging, which requires both temporal understanding and causal reasoning abilities.

D. Driving Task Evaluation

In this section, we evaluate the effectiveness of visual descriptions and explanations in several driving tasks, with the fine-tuned MLLMs' outputs as textual evidence and the corresponding video features as visual evidence. For evaluation consistency, we use the same video features described in Section V-A for each MLLM and textual evidence from MLLMs' fine-tuned on PDB-X.

1) *Evaluation Results on Brain4Cars [1].*: We first evaluate the intention prediction task proposed in the Brain4Cars [1] dataset as the **in-domain** evaluation. To demonstrate the generalizability of MLLMs, we do not incorporate additional low-level signals (*e.g.*, vehicle speed, and GPS information) used by Brain4Cars.

In Table V, we observe that drivers' personalized behavior can be good complementary textual evidence for prediction of drivers' turn maneuvers, given that these maneuvers typically exhibit more discernible behavior from the driver. On the other hand, the original method in Brain4Cars, which fuses various sensory information sources, still outperforms MLLMs with only vision-language information in lane change prediction. We argue that in tasks requiring additional sensory information, our method may still incorporate the information into the general reasoning process.

TABLE V
PERFORMANCE OF PRECISION AND RECALL METRICS ON THE INTENTION PREDICTION TASK [1]. THE BLIP-2 AND VTIMELLM COLUMNS ARE OBTAINED WITH TEXTUAL EVIDENCE OUTPUT FROM LLMs FINE-TUNED ON PDB-X. THE S-RNN COLUMN IS OBTAINED FROM THE METHOD FROM THE ORIGINAL PAPER [1].

	BLIP-2		VTimeLLM		S-RNN	
	Prec.	Recall	Prec.	Recall	Prec.	Recall
Turn	84.57	76.43	83.85	77.18	75.20	75.30
Change	76.47	72.73	76.04	72.11	85.40	86.00

2) *Evaluation Results on AIDE [2].*: We further evaluate the MLLMs that are initially fine-tuned on PDB-X without additional fine-tuning on AIDE as the **cross-domain** evaluation. The tasks proposed by AIDE are driver behavior recognition (DBR), driver emotion recognition (DER), traffic context recognition (TCR), and vehicle condition recognition (VCR). In addition to the dual-cam videos used in our method, AIDE [2] additionally incorporates videos from the left and right cameras, and multimodal annotations of a driver's face, body, gesture, and posture.

TABLE VI

PERFORMANCE OF ACC AND WEIGHTED F1 METRICS ON THE AIDE DATASET [2]. THE BLIP-2 AND VTimeLLM COLUMNS ARE OBTAINED BY THE LMMs THAT ARE INITIALLY FINE-TUNED ON PDB-X AND THEN USED FOR AIDE WITHOUT ADDITIONAL FINE-TUNING.

	BLIP-2		VTimeLLM		AIDE	
	Acc ↑	F1 ↑	Acc ↑	F1 ↑	Acc ↑	F1 ↑
DER	72.74	72.73	74.06	73.41	71.26	68.71
DBR	64.86	64.41	65.52	65.34	65.35	63.29
TCR	90.15	89.86	90.15	90.19	83.74	81.28
VCR	76.35	76.00	78.82	77.67	77.12	75.23

In Table VI, we can observe robust generalizability of the fine-tuned MLLMs in cross-domain driving tasks, where the fine-tuned VTimeLLM achieves better accuracy and weighted F1 than AIDE in all tasks. In the four driving tasks in AIDE, we observe better improvements in driver emotion recognition (DER) and traffic condition recognition (TCR), which benefit from the drivers’ personalized behavior descriptions and external visual explanations respectively.

VI. CONCLUSION

We introduce the Driver’s Personalized Behavior Evaluation dataset (PDB-Eval), which leverages in-cabin and external views for personalized driver behavior analysis. Our benchmark comprises two components—PDB-X and PDB-QA—derived via visual comparative prompting within MLLMs to capture behavioral discrepancies and provide descriptive explanations. Fine-tuning various MLLMs improved performance on these tasks by up to 73.2% over GPT-4V’s zero-shot results. However, the relatively low BLEU-4 scores indicate persistent challenges in fine-grained, temporally aware multimodal reasoning. Future work should focus on enabling MLLMs to identify localized visual evidence with temporal awareness.

REFERENCES

- [1] A. Jain, H. S. Koppula, S. Soh, B. Raghavan, A. Singh, and A. Saxena, “Brain4cars: Car that knows before you do via sensory-fusion deep learning architecture,” *arXiv preprint arXiv:1601.00740*, 2016.
- [2] D. Yang, S. Huang, Z. Xu, Z. Li, S. Wang, M. Li, Y. Wang, Y. Liu, K. Yang, Z. Chen *et al.*, “Aide: A vision-driven multi-view, multi-modal, multi-tasking dataset for assistive driving perception,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 20 459–20 470.
- [3] J. Kim, A. Rohrbach, T. Darrell, J. Canny, and Z. Akata, “Textual explanations for self-driving vehicles,” *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [4] B. Huang, X. Wang, H. Chen, Z. Song, and W. Zhu, “Vtimellm: Empower llm to grasp video moments,” *arXiv preprint arXiv:2311.18445*, 2023.
- [5] J. Wu, H. Lyu, Y. Xia, Z. Zhang, J. Barrow, I. Kumar, M. Mirtaheri, H. Chen, R. A. Rossi, F. Deroncourt *et al.*, “Personalized multimodal large language models: A survey,” *arXiv preprint arXiv:2412.02142*, 2024.
- [6] J. Wu, Z. Zhang, Y. Xia, X. Li, Z. Xia, A. Chang, T. Yu, S. Kim, R. A. Rossi, R. Zhang *et al.*, “Visual prompting in multimodal large language models: A survey,” *arXiv preprint arXiv:2409.15310*, 2024.
- [7] K. Li, Y. Wang, Y. He, Y. Li, Y. Wang, Y. Liu, Z. Wang, J. Xu, G. Chen, P. Luo *et al.*, “Mvbench: A comprehensive multi-modal video understanding benchmark,” *arXiv preprint arXiv:2311.17005*, 2023.
- [8] D. Ko, J. S. Lee, W. Kang, B. Roh, and H. J. Kim, “Large language models are temporal and causal reasoners for video question answering,” *arXiv preprint arXiv:2310.15747*, 2023.
- [9] A. Yan, Z. Yang, J. Wu, W. Zhu, J. Yang, L. Li, K. Lin, J. Wang, J. McAuley, J. Gao *et al.*, “List items one by one: A new data source and learning paradigm for multimodal llms,” *arXiv preprint arXiv:2404.16375*, 2024.
- [10] W. Bao, Q. Yu, and Y. Kong, “Uncertainty-based traffic accident anticipation with spatio-temporal relational learning,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, pp. 2682–2690.
- [11] D. Zhang, J. Yang, H. Lyu, Z. Jin, Y. Yao, M. Chen, and J. Luo, “Cocot: Contrastive chain-of-thought prompting for large multimodal models with multiple image inputs,” *arXiv preprint arXiv:2401.02582*, 2024.
- [12] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds *et al.*, “Flamingo: a visual language model for few-shot learning,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 23 716–23 736, 2022.
- [13] M. Bonyani, M. Rahmadian, S. Jahangard, and M. Rezaei, “Dipnet: Driver intention prediction for a safe takeover transition in autonomous vehicles,” *IET Intelligent Transport Systems*, 2023.
- [14] J. Zhang, F. Ilievski, K. Ma, A. Kollaa, J. Francis, and A. Oltramari, “A study of situational reasoning for traffic understanding,” *arXiv preprint arXiv:2306.02520*, 2023.
- [15] J. Echterhoff, A. Yan, K. Han, A. Abdelraouf, R. Gupta, and J. McAuley, “Driving through the concept gridlock: Unraveling explainability bottlenecks in automated driving,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 7346–7355.
- [16] D. Fu, X. Li, L. Wen, M. Dou, P. Cai, B. Shi, and Y. Qiao, “Drive like a human: Rethinking autonomous driving with large language models,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 910–919.
- [17] Z. Xu, Y. Zhang, E. Xie, Z. Zhao, Y. Guo, K. K. Wong, Z. Li, and H. Zhao, “Drivegpt4: Interpretable end-to-end autonomous driving via large language model,” *arXiv preprint arXiv:2310.01412*, 2023.
- [18] O. Sainz, I. García-Ferrero, R. Agerri, O. L. de Lacalle, G. Rigau, and E. Agirre, “Gollie: Annotation guidelines improve zero-shot information-extraction,” *arXiv preprint arXiv:2310.03668*, 2023.
- [19] H. Liu, W. Xue, Y. Chen, D. Chen, X. Zhao, K. Wang, L. Hou, R. Li, and W. Peng, “A survey on hallucination in large vision-language models,” *arXiv preprint arXiv:2402.00253*, 2024.
- [20] Y. Zhou, C. Cui, J. Yoon, L. Zhang, Z. Deng, C. Finn, M. Bansal, and H. Yao, “Analyzing and mitigating object hallucination in large vision-language models,” *arXiv preprint arXiv:2310.00754*, 2023.
- [21] J. Guo, J. Li, D. Li, A. M. H. Tiong, B. Li, D. Tao, and S. Hoi, “From images to textual prompts: Zero-shot visual question answering with frozen large language models,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 10 867–10 877.
- [22] T. Kim, Y. Cho, H. Shin, Y. Jo, and D. Shin, “Generalizing visual question answering from synthetic to human-written questions via a chain of qa with a large language model,” *arXiv preprint arXiv:2401.06400*, 2024.
- [23] J. Li, D. Li, S. Savarese, and S. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” *arXiv preprint arXiv:2301.12597*, 2023.