

Sage Deer: A Super-Aligned Driving Generalist Is Your Copilot

Hao LU^{*1} Jiaqi Tang^{*1} Jiayao Wang¹ Yunfan LU¹ Xu Cao¹ Qingyong Hu¹ Yin Wang² Yuting Zhang¹
 Tianxin Xie¹ Yunpeng Zhang³ Yong Chen⁴ Jiayu Gao⁴ Bin Huang¹ Dengbo He¹ Shiguang Deng²
 Hao Chen¹ Ying-Cong Chen¹

Abstract

The intelligent driving cockpit, an important part of intelligent driving, needs to match different users' comfort, interaction, and safety needs. This paper aims to build a super-aligned and generalist driving agent, *sage deer*. Sage Deer achieves three highlights: (1) Super alignment: It achieves different reactions according to different people's preferences and biases. (2) Generalist: It can understand the multi-view and multi-mode inputs to reason the user's physiological indicators, facial emotions, hand movements, body movements, driving scenarios, and behavioral decisions. (3) Self-Eliciting: It can elicit implicit thought chains in the language space to further increase generalist and super-aligned abilities. Besides, we collected multiple data sets and built a large-scale benchmark. This benchmark measures the deer's perceptual decision-making ability and the super alignment's accuracy.

1. Introduction

The concept of a driving copilot refers to the collaboration between human drivers and vehicles to ensure seamless and efficient operations throughout the driving process (Mao et al., 2023; Cui et al., 2023a; Wang et al., 2020; Hecht et al., 2018). As intelligent vehicles continue to evolve, this concept has become a cornerstone of their development, aiming to enhance safety, improve comfort, increase traffic efficiency, and deliver a superior driving experience for both drivers and passengers (Cui et al., 2023c; Yang et al., 2024). For instance, driver monitoring is essential to proactively assess a driver's health, fatigue, mood, and behavior. Human-machine interaction plays a pivotal role in understanding driver actions, enabling the system to activate

¹HKUST, HKUST-GZ ²ZJU ³Phigent ⁴GEELY. * means equal contribution. Correspondence to: Ying-Cong Chen <yingcongchen@ust.hk>.

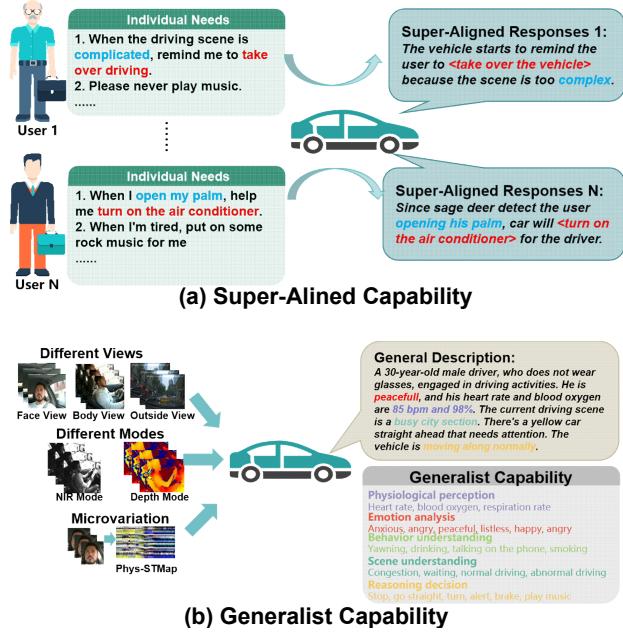


Figure 1. The Capability of Sage Deer. (a) Sage Deer can perform corresponding actions according to the preferences of each user. Users can save their preferences in a document and update them anytime. (b) Sage Deer can integrate multi-view and multi-model inputs for driver and scene understanding. Notably, Sage Deer actually has both super-aligned and generalist capabilities at the same time.

various cockpit service functions (Li et al., 2023c; Yang et al., 2022). Furthermore, driving assistance features, such as lane-keeping and adaptive cruise control, help automate driving tasks, making the driving experience safer and more convenient.

Although research on driving copilots has made progress, much of the prior work has focused on data collection, multi-task learning, and improving generalization performance. However, developing a super-aligned copilot model that addresses individual-specific requirements remains an underexplored area. **Super-aligned**, as defined in this paper, refers to the ability of a model to adapt to individual user requirements and perform corresponding actions, as illustrated in Fig. 1(a). In addition, driving copilots must

serve as generalists, capable of perception, understanding, and decision-making across diverse scenarios, as shown in Fig. 1(b). In essence, Sage Deer seeks to understand multi-view and multi-modal inputs while adapting its behavior to meet user-specific needs.

Multimodal large language models (MLLMs) give a great opportunity to show strong capabilities in integrating multi-modal inputs, handling multitask outputs, and demonstrating generalization abilities. However, most of these approaches are not well-suited for driving copilot applications, with limitations in data, benchmarking, and model design. Specifically, (1) MLLMs must feed the user’s needs and habits into the large model to achieve the super-aligned ability through long text at every inference. (2) Most of the MLLMs can not be satisfied and can accept multi-view and multi-mode input at the same time. (3) Super-aligned and generalist abilities tend to overfit rather than reason through visual instruction tuning. (4) Satisfying both of these requires large-scale tailored training data of the driving copilot, which is also lacking.

To address these limitations, we propose a novel super-aligned and generalist driving copilot framework, called **Sage Deer**. For achieving generalist capabilities, we introduce multi-view tokenizers, multimodal tokenizers, and microvariation tokenizers. For super-aligned capabilities, we design a learnable retrieval-augmented generation (RAG) framework capable of delivering personalized responses based on user-specific needs. Additionally, we propose a Continuous Latent Chain Elicitation (CLCE) mechanism to further enhance both generalist and super-aligned capabilities. The CLCE activates the inherent reasoning capabilities of LLMs by stimulating implicit chains-of-thought (COT) reasoning without requiring explicit COT labels. Lastly, we compile and combine multiple datasets to construct a comprehensive benchmark for evaluating driving copilots. In summary, our key contributions are as follows:

1. We introduce a novel super-aligned generalist driving copilot framework capable of understanding multi-view and multi-modal inputs while providing personalized responses tailored to individual user needs.
2. We design a Continuous Latent Chain Elicitation mechanism to strengthen the copilot’s super-aligned and generalist capabilities by leveraging the LLM’s innate reasoning abilities without relying on additional COT annotations.
3. We establish a multi-view and multi-modal evaluation protocol for driving copilots that unifies the assessment of driver physiology, emotion, behavior, scene understanding, and decision-making.

2. Data Curating and Benchmarking

2.1. Data collection

Our goal is to create a super-aligned driving generalist in the intelligent copilot. We selected two of the most recent multi-view multi-task large-scale driving datasets AIDE (Yang et al., 2023a) and DMD (Ortega et al., 2020). To contactless monitor the user’s health condition, five remote physiological measurement datasets (VIPL-HR (Niu et al., 2020), V4V (Revanur et al., 2021a), PURE (Stricker et al., 2014a), BUAA-rPPG (Xi et al., 2020a) and UBFC (Bobbia et al., 2019a)) is used. To quantitatively evaluate the fatigue state of drivers, YawDD (Abtahi et al., 2014) dataset is used in this paper. Combined with these datasets, sage deer can satisfy multi-view and multi-modal input (RGB, NIR and Depth), generalist capabilities (physiological estimation, emotional estimation, gesture estimation, behavior estimation, driving behavior detection, and driving decision-making), and super-aligned capabilities.

2.2. General Instruction Construction

To construct natural language descriptions for sage deer, we built a data curating pipeline. Firstly, we use existing GPT-4o tools to generate captions for the frame sampled at equal intervals automatically. Likes (Li et al., 2023b; Zhang et al., 2023a), we set the reasonable prompt to merge the information of different frames. Existing image caption methods are often inaccurate or insufficiently annotated for an intelligent copilot system. So, we took advantage of the existing tags (including physiological indicators, emotional indicators, action indicators, behavioral indicators, scene understanding, and reasoning decision-making) in the dataset to correct and supplement the captain. Next, we use GPT-4o as an assistant to build question-answering pairs for different tasks.

2.3. Super-Aligned Reaction

Users have different requirements for the driving copilot, especially interactivity and trustworthiness. Based on interactivity and trustworthiness, we simulate and design a variety of different requirements levels to build requirements documents for each user. In addition, we have built different replies in interactivity and trustworthiness. Interactivity means creating gestures, emotions, and body movements and tailoring copilot feedback to user needs. Trustworthiness refers to personalized warning feedback on user fatigue, bad mood, and bad behavior. Based on different requirements documents, the same scene GPT-4o will be used to generate different interactivity and trustworthiness responses.

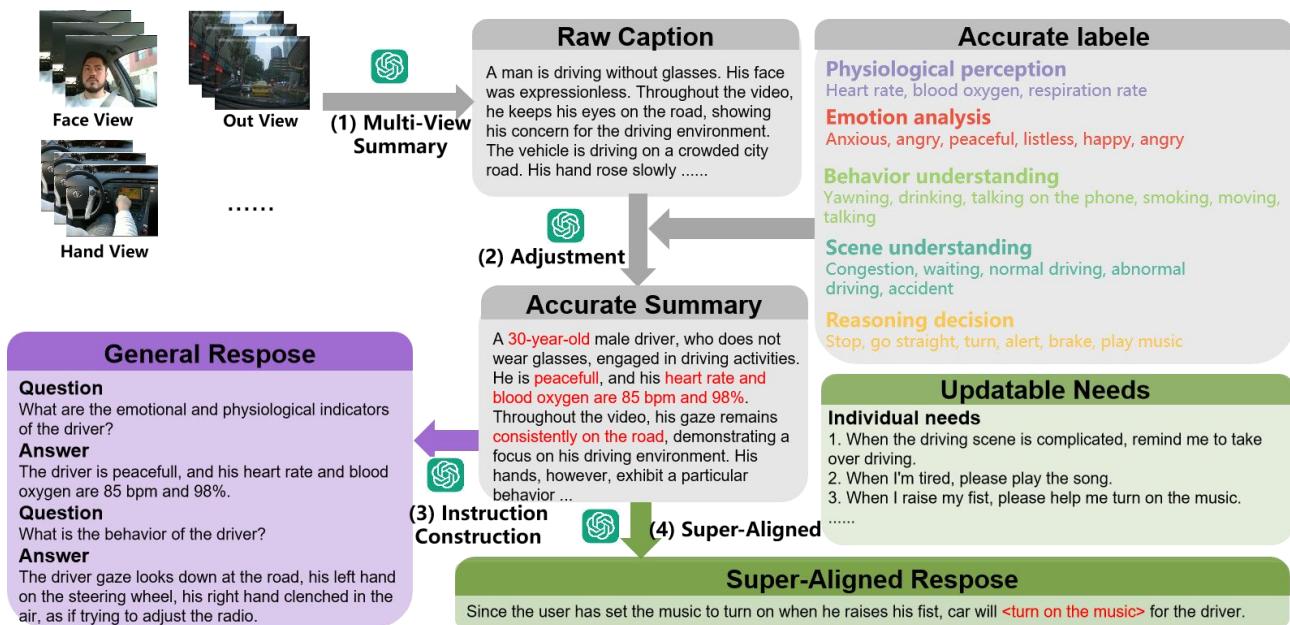


Figure 2. Data construction process of Sage Deer. (1) We use existent GPT-4o tools to generate captions for the videos like (Li et al., 2023b; Zhang et al., 2023a). Then, we set a reasonable prompt to merge the information from different videos. (2) We took advantage of the existing labels (including physiological indicators, emotional indicators, action indicators, behavioral indicators, scene understanding, and reasoning decision-making) to correct and supplement the captain. (3) Next, we use GPT4 (Achiam et al., 2023) as an assistant to build question-answering pairs for different tasks (including physiological indicators, emotion, behavior, and so on.). (4) We design multiple user preferences, and GPT4 responds to the current scenario based on user preferences. All of the above processes are double-checked manually.

3. Framework

We aim to build a super-aligned generalist in the driving copilot as shown in Fig. 3. Specifically, generalist ability is the ability to integrate multi-view and multi-mode information and detect the driver’s physiological indicators, emotions, behaviors, and scene understanding. Super-aligned gives responses based on the user’s preferences. We highlight that Sage Deer’s framework is simple and necessary without a complex design.

3.1. The Generalist Capability

Leveraging multi-view and multi-mode input for driving copilot can significantly enhance the model’s ability to understand complex scenes, especially under challenging conditions such as low-light environments.

Tokenizing Multi-View. Multi-view input is crucial for understanding both the driving scene and the driver against occlusion. Like most MLLMs (Gao et al., 2023a; Lin et al., 2023; Sima et al., 2023; Li et al., 2023b; Wang et al., 2022), we utilized the CLIP as the feature extractor from multi-view images. Using a two-layer linear network, the extracted features were then mapped to the language space features $em \in \mathbb{R}^{C \times L}$, where the L is the number

of language tokens occupied by visual features. Additionally, we employed markers to distinguish between different views. For example, a front-facing RGB video embedding is denoted $em_{front} \in \mathbb{R}^{C \times L}$, formatted as $E_{front} = <\text{Front RGB bos}> em_{front} <\text{Front RGB cos}>$, here $<\text{Front RGB bos}>$ and $<\text{Front RGB cos}>$ are actually tokenized vectors. This format ensures that the LLM can understand exactly which view the feature is coming from.

Tokenizing Multi-Model. For the driving copilot, understanding and reasoning in adverse conditions require a combination of multiple sensing modes, including near-infrared (NIR) and depth images. We uniformly employed the CLIP as the feature extractor for both NIR and depth. Since CLIP is specifically designed for RGB images, we set up the first layer of CLIP as a fine-tuned retraining process. Subsequently, a simple two-layer linear network mapped these mode features into the language space features $em \in \mathbb{R}^{C \times L}$. Here C represents the number of channels characteristic of the language model, and L is a hyperparameter indicating the number of tokens used to represent the video features. To understand different modes, we also incorporate makers. For instance, we appended specific start and end symbols to the NIR video embeddings, formatted as $E + NIR = <\text{NIR bos}> em_{NIR} <\text{NIR cos}>$. Here

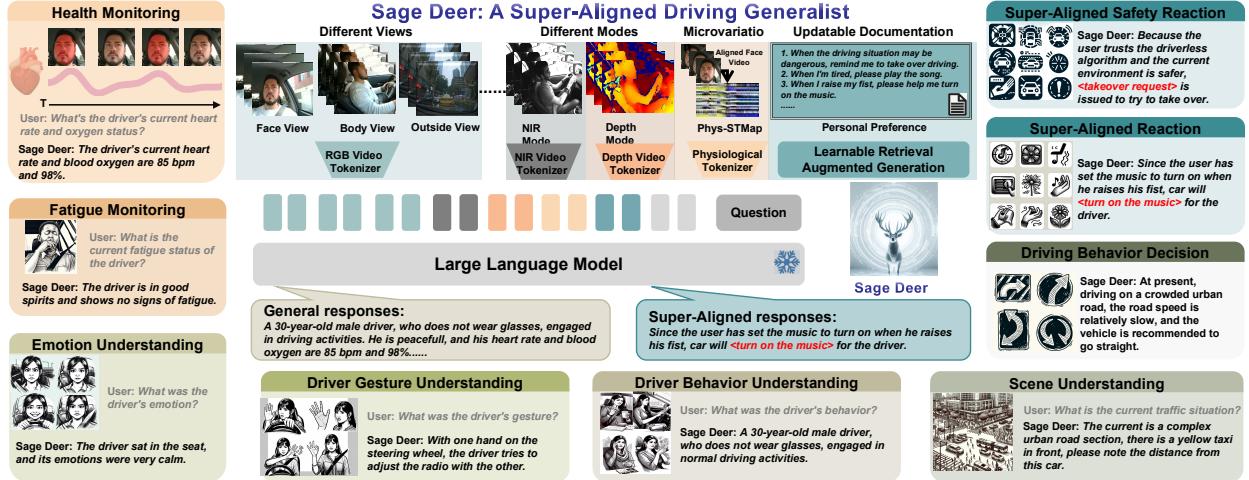


Figure 3. Sage Deer uses pre-trained video encoders (visual tokenizers) to tokenize the different modes and views of the video, especially the physiological encoder used to extract the physiological signals of the face video. Using tokenized visual embedding, large language models can provide general responses for physiological indicators, emotional states, gestures, human behavior, and scene understanding. It is worth emphasizing that the individual needs of the user can be edited in a single document. We then query this information with learnable retrieval augmentation generation, generating results super aligned with user preferences.

< NIR bos > and < NIR cos > are actually tokenized vectors. Similar processing was applied to other video modalities before feeding them into the LLM.

Tokenizing Microvariation. In addition to the ability to observe drivers’ expressions, emotions, and behaviors, Sage Deer also needs to measure further the driver’s health status, such as heart rate (HR), breathing, and heart rate variability (HRV). With remote photoplethysmography (rPPG) technology, we can obtain these indicators through the non-contact RGB video (Wang et al., 2015; Xiao et al., 2024; Lu et al., 2023; Niu et al., 2020; Liu et al., 2024b). Here, we select NEST-rPPG as our physiological tokenizer, which is trained on multiple datasets (Lu et al., 2023). Fellow NEST-rPPG, we pre-trained a neural network that could convert video into HR, breathing, and HRV, and then we kept only the encoder as our physiological tokenizer. Subsequently, a simple two-layer linear network mapped these mode features into the language space features $em_{phys} \in \mathbb{R}^{C \times L}$. Similarly, we set up a special makers $E_{phys} = <\text{Physiological bos}> em_{phys} <\text{Physiological cos}>$.

Along with questions and responses, multi-view, multi-mode, and microvariation tokenization can be jointly fed into LLM, i.e. $\{E_{front}, E_{out}, E_{face}, E_{hand}, E_{NIR}, E_{Depth}, \dots, E_{phys}, E_{rag}, <\text{bos}>, Q, E_{cot}, R, <\text{cos}>\}$. Q and R are the tokenized questions and responses. E_{rag} is the super-aligned feature based on the current scene and the user’s personal preferences. E_{cot} helps the LLM understand and fuse these multiple inputs through implicit COT, which is described in Sec. 4. There’s only one phase of training, and the main

loss is visual instruction tuning \mathcal{L}_{vi} (Li et al., 2023b; Zhang et al., 2023a; Muhammad Maaz & Khan, 2023).

3.2. Super-Aligned Capability

Super-aligned refers to the model’s ability to align with individual requirements to carry out corresponding actions. However, we cannot fine-tune every MLLM model for each user’s different requirements. To this end, we have adopted a simple retrieval-augmented generation (RAG) framework as shown in Fig. 3.

Specifically, we built an up-to-date document that includes personal requirements. Then, we chunked the document in a very concise way to divide it into N sentences. In other words, a chunk is a sentence. Each sentence is then tokenized and then filled with empty to the same length M , i.e., $em_{sa} \in \mathbb{R}^{N \times M \times C}$, where C is the length of the vocabulary table. A four-layer convolutional encoder is then used to compress the length of sentence tokens. Then, we calculate the similarity between visual features and sentence features and the weighted combination of these sentence tokens $E_{rag} \in \mathbb{R}^{\times K \times C}$. em_{rag} reflects the individual requirements, which is also fed into the LLM with other multi-modal information as explained at the end of Sec. 3.1.

4. Continuous Latent Chain Eliciting

To achieve super-aligned generalists, Sage Deer must simultaneously fuse multi-view and multi-modal information, align individual preferences, and model the task relationship. A reasonable chain of thought (COT) can efficiently make

these processes efficient in reasoning in the language space. However, the accuracy label of the COT process for cockpit tasks is missing. Inspired by (Hao et al., 2024), we propose the Continuous Latent Chain Eliciting (CLCE) strategy, as shown in Fig. 4.

Continuous Latent Chain-of-Thought (CL-COT). The CLCE forces the LLM to output the fixed L length CL-COT before the final responses. The fixed CL-COT is expected to learn implicit reasoning without any label supervision. The principle behind this method is that it forces the LLM to think of more COT steps in the hidden space to get the final answer. However, CL-COT input is required because the LLM’s equal length of input and output can speed up training. So, we feed multi-modal, multi-view, and rag features $\{E_{front}, E_{out}, \dots, E_{rag}\}$ into a simple 2-layer convolution network to output a learnable CL-COT embedding E_{cot} . E_{cot} is a self-learning COT embedded in the LLM hidden space, which is expected to activate the LLM’s self-reasoning ability.

Latent Chain Eliciting. The learnable COT embedding E_{cot} exhibits trivial solutions in our experiment as shown in Fig. 5. Specifically, we visualize the Pearson similarity and range of random two token features of Continuous Latent Chain (CLC) embedding. Compared with normal embedding, the similarity between CLC tokens is very high, and the difference between the maximum and minimum values of the same token is very small. This shows that the CLC token E_{cot} tends to predict an invalid token rather than the valid COT process. Here we use a very simple method to solve this problem. We design the latent chain eliciting (LCE) loss \mathcal{L}_{LCE} to make the constraint token valid, i.e., $\mathcal{L}_{LCE} = -\|E_{cot}\|_1$. LCE loss forces the LLM to output an active embedding, and the final response labels ensure that this active embedding is a valid COT process.

5. Experiments

5.1. Datasets.

We found two recent data sets of very comprehensive annotations for assisted driving. The DMD is a driver monitoring dataset, an extensive dataset that includes real and simulated driving scenarios: distraction, eye distribution, drowsiness, handwheel interaction, and contextual data, in 41 hours of RGB, depth, and infrared video from 3 cameras, capturing the faces, bodies, and hands of 37 drivers (Ortega et al., 2020). AIDE proposes an assisted driving awareness dataset that takes into account contextual information both inside and outside the vehicle in a natural scenario (Yang et al., 2023a). AIDE enables overall driver monitoring through multi-perspective Settings of the driver and the scene, multi-modal annotations of the face, body, posture, and gestures, and four practical task designs for driver understanding. we

collect five rPPG face video datasets (VIPL-HR (Niu et al., 2019), PURE (Stricker et al., 2014b), UBFC-rPPG (Bobbia et al., 2019b), V4V (Revanur et al., 2021b), and BUAA-MIHR (Xi et al., 2020b)), mostly with subjects remaining still, and some with head movements.

5.2. Training Details.

We train our model on A6000 for 2 epochs. The learning rate and weight decay are set to 0.001 and 0.02, respectively. The maximum sentence length is set to 64. That is, if the sentence is too long, excess parts will be discarded, and if the sentence is too short, 0 will be filled to make the length uniform. For the feature extraction of physiological signals, we use the pre-trained model of NEST-rPPG (Lu et al., 2023) and do not fine-tune it.

5.3. Baselines and Evaluation Metrics

We compare our proposed framework’s understanding performance with SOTA video understanding baselines. We select five baselines: Video-ChatGPT (Muhammad Maaz & Khan, 2023), VideoChat (Li et al., 2023b), Video-LLaMA (Zhang et al., 2023a), LLaMA-Adapter (Zhang et al., 2023b), and Video-LLaVA (Lin et al., 2023). These methods do not have a multi-perspective, multi-modal understanding, so we migrated our generalist tokenizer method to them and fine-tuned it. The main purpose of comparing these methods is to emphasize the capabilities of our CLCE.

To evaluate our model’s performance accurately, we adopt BLEU Bilingual Evaluation Understudy(BLEU) and SPICE (Papineni et al., 2002) to measure word overlap between the model-generated text and the ground truth. The consistently higher scores across both metrics validate the effectiveness of our approach in generating coherent, accurate, and contextually rich descriptions.

5.4. Generalist Performance

Our model can estimate the driver’s emotion, physiological indicators, gaze, physical behavior, hand behavior, driving scene, and vehicle state. In order to more clearly evaluate the ability of the model in different subtasks, we conducted a systematic evaluation on two multi-task datasets, fatigue and physiological indicators. The AIDE dataset evaluates generalist capabilities across four categories—emotion, behavior, scene, and condition. As shown in Table 1, our model consistently outperforms all baselines, achieving higher BLEU and SPICE scores across all categories. This improvement indicates the model’s enhanced ability to capture nuanced details such as emotional states and contextual anomalies. The integration of expert knowledge through RAG (Retrieval-Augmented Generation) serves as a critical factor in bridging connections between related tasks,

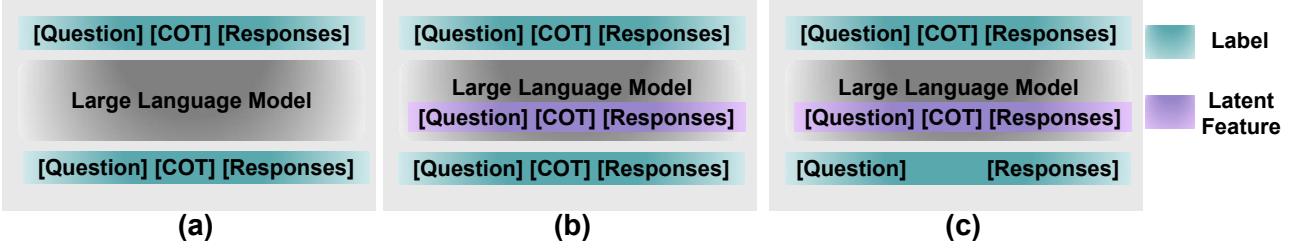


Figure 4. A comparison of our continuous latent chain eliciting with the existent chain-of-thought method. (a) The traditional CoT model generates the reasoning process with the supervision of COT tokens. (b) Coconut uses the LLM to reason in an unrestricted latent space instead of a language space (Hao et al., 2024). But he still distills knowledge from the COT label. (c) Our method tries to activate the implicit language space so that the model learns the implicit CoT from itself.

Table 1. Generalist capabilities on AIDE data sets.

Method	Emotion		Behavior		Scene		Condition	
	BLEU	SPICE	BLEU	SPICE	BLEU	SPICE	BLEU	SPICE
Video-ChatGPT	0.200	0.280	0.170	0.205	0.190	0.330	0.195	0.320
VideoChat	0.205	0.285	0.175	0.210	0.195	0.335	0.200	0.325
Video-Llama	0.217	0.312	0.183	0.223	0.207	0.352	0.211	0.337
Llama-Adapter	0.215	0.310	0.190	0.225	0.210	0.340	0.189	0.321
Video-LLaVA	0.202	0.290	0.172	0.205	0.190	0.330	0.195	0.320
Ours	0.232	0.331	0.194	0.242	0.225	0.369	0.223	0.360

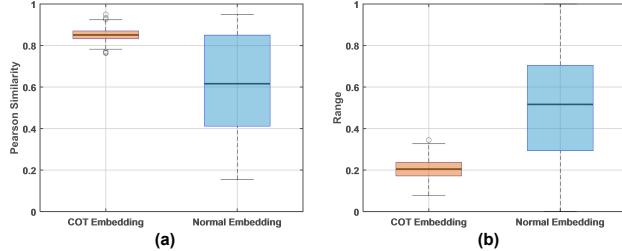


Figure 5. The trivial solutions of the CL-COT embedding E_{cot} . (a) Compared with normal embedding, COT embedding tends to be self-similar. (b) Compared with normal embedding, COT embedding tends to have a smaller range of values. These two phenomena indicate that the COT embedding tends to be invalid.

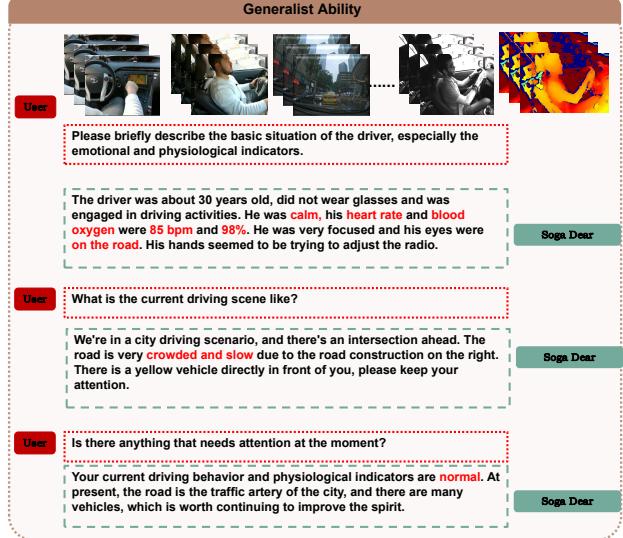


Figure 6. Generalist ability. Soga deer can reply to different tasks and can reply to an open domain.

leading to superior results.

Further, we verify our algorithm on the DMD data set. The DMD dataset focuses on three categories: action, gaze, and hand behavior. As shown in Table 2, our model demonstrates remarkable performance improvements over the baselines, further showcasing its generalist capabilities. Particularly, the model performs exceptionally well in understanding and describing complex interactions involving gaze and hand movement. This suggests that our framework excels at parsing fine-grained details critical for tasks involving driver monitoring and safety analysis.

5.5. Super-Aligned Performance

For driving cockpits, users have diverse and individualized preferences, necessitating personalized super-alignment capabilities. However, the majority of existing algorithms lack support for this level of alignment. A straightforward baseline approach involves encoding all user preferences as text inputs to a large language model (LLM) and leveraging the LLM's long-range reasoning capabilities. In contrast,

Table 2. Generalist capabilities on DMD data sets.

Method	Action		Gaze		Hand	
	BLEU	SPICE	BLEU	SPICE	BLEU	SPICE
Video-ChatGPT	0.175	0.250	0.150	0.185	0.160	0.310
VideoChat	0.190	0.260	0.160	0.185	0.160	0.310
Video-Llama	0.205	0.275	0.170	0.205	0.180	0.325
Llama-Adapter	0.202	0.265	0.170	0.201	0.175	0.320
Video-LLaVA	0.195	0.295	0.155	0.220	0.170	0.315
Ours	0.235	0.315	0.195	0.230	0.210	0.340



Figure 7. Super alignment ability. soga deer can perform different operations according to the individual needs of each person.

Table 3. Baseline compared the performance of our approach on the super alignment protocol. The baseline is to user requirements directly into visual features and feed them into the LLM, taking advantage of the long-range modeling capabilities of the LLM itself.

Method	AIDE		DMD	
	BLEU	SPICE	BLEU	SPICE
Baseline	0.184	0.301	0.191	0.296
Ours	0.231	0.327	0.215	0.314

our method achieves super-alignment through a learned Retrieval-Augmented Generation (RAG) framework. The comparative performance of these two approaches is presented in Table 3.

As shown in Table 3, our approach achieves significant improvements over the baseline. These results highlight the limitations of LLMs in handling long text inputs and making nuanced responses solely based on visual features and prior knowledge. This underscores the necessity of our RAG framework for achieving robust super-alignment.

5.6. Ablation Study

Sage Deer’s core innovation is the CLCE approach. Therefore, we mainly conduct ablation experiments on the hyperparameters of this method. Here we compare the weight of

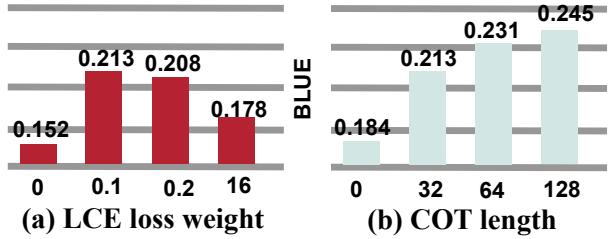


Figure 8. Ablation Study. (a) The impact of LCE loss weight. (b) The impact of COT length.

LCE loss and the compression of implicit COT length for ablation as shown in Fig. 8. The experimental results show that a reasonable weight of LCE loss can enable LLM to stimulate COT ability in the form of hidden language space. The longer the length of the implicit COT, the better the model performance.

6. Related Work

Multimodal Large Language Model Recent research highlights the robust capabilities of large models in video understanding, showcasing significant progress in this area. Beyond the advancements in vision-language large models (Li et al., 2023a; Zhu et al., 2023; Liu et al., 2024a; Lu et al., 2024), there is a growing emphasis on integrating additional modalities, such as audio and sensor data, to further enhance model performance (Lv & Sun, 2024; Li et al., 2023b; Maaz et al., 2023; Ye et al., 2023; Luo et al., 2023). For instance, Bain et al. (Bain et al., 2021) introduced a large-scale dataset offering general video content descriptions, pushing the boundaries of how models can interpret dynamic visual scenes. Several LLM-based approaches (Li et al., 2023b; Maaz et al., 2023; Ye et al., 2023; Luo et al., 2023) aim to effectively interpret the visual elements in videos, unlocking deeper contextual understanding. Furthermore, Video-LLaMa (Zhang et al., 2023a) expands the scope of video comprehension by incorporating both visual and auditory modalities, enhancing the richness of information these models can process. Additionally, Su et

Table 4. Comparison of public driving cockpit datasets. Reaction is the response given by the intelligent cockpit according to different user needs. Question Answering refers to the verbal inquiry interaction between the user and the cockpit.

Dataset	Views	Multimodal	Health	Behavior	Emotion	Traffic Context	Vehicle Condition	Reaction	Question Answering
SEU (Zhao et al., 2012)	1	–	–	✓	–	–	–	–	–
Tran (Tran et al., 2018)	1	–	–	✓	–	–	–	–	–
Zhang (Zhang et al., 2020)	2	✓	–	✓	–	–	–	–	–
StateFarm (sta, 2016)	1	–	–	✓	–	–	–	–	–
AUC-DD (Eraqi et al., 2019)	1	–	–	✓	–	–	–	–	–
LoLi (Saad et al., 2020)	1	✓	–	✓	–	–	–	–	–
Brain4Cars (Jain et al., 2016)	2	–	–	✓	–	–	–	–	–
Drive&Act (Martin et al., 2019)	6	✓	–	✓	–	–	–	–	–
DMD (Ortega et al., 2020)	3	✓	–	✓	–	–	–	–	–
DAD (Kopuklu et al., 2021)	2	✓	–	✓	–	–	–	–	–
DriPE (Guesdon et al., 2021)	1	–	–	–	–	–	–	–	–
LBW (Kasahara et al., 2022)	2	–	–	–	–	–	–	–	–
MDAD (Jegham et al., 2019)	2	✓	–	✓	–	–	–	–	–
3MDAD (Jegham et al., 2020)	2	✓	–	✓	–	–	–	–	–
DEFE (Li et al., 2021a)	1	–	–	–	✓	–	–	–	–
DEFE+ (Li et al., 2021b)	1	✓	–	–	✓	–	–	–	–
Du (Du et al., 2020)	1	✓	–	–	✓	–	–	–	–
KMU-FED (Jeong & Ko, 2018)	1	–	–	–	✓	–	–	–	–
MDCS (Oh et al., 2022)	2	✓	–	–	✓	–	–	–	–
AIDE (Yang et al., 2023a)	4	–	–	✓	✓	✓	✓	–	–
Sage Deer	4	✓	✓	✓	✓	✓	✓	✓	✓

al. (Su et al., 2023) have leveraged multimodal encoders that work across six distinct modalities, pushing the envelope on the versatility of multimodal models. However, despite these exciting advancements, there remains a notable gap in research focused on multimodal large models tailored specifically for driving cockpits.

Multi-modality Large Model in Driving Multi-modality large models have been widely applied in autonomous driving systems, supporting tasks such as automatic planning and control (Mao et al., 2023; Cui et al., 2023a), perception (Wang et al., 2020), and driver health monitoring (Hecht et al., 2018), among others. By integrating multi-modal data (e.g., vision, speech, point clouds, etc.) (Yang et al., 2023b), these models enhance the ability of autonomous driving systems to perceive both the internal and external environments. In user-facing vehicle interfaces, analyzing driver-specific factors such as emotions, health, posture, and actions allows autonomous systems to interact dynamically and adapt to diverse driving styles. Cui et al. (Cui et al., 2023b;a) were the first to propose leveraging multi-modal data from both the vehicle’s interior and exterior to enhance decision-making precision in autonomous driving. Recent advancements (Cui et al., 2023c; Yang et al., 2024) have sought to improve fine-grained comprehension of user-specific language inputs and optimize vehicle control performance. Nonetheless, these approaches overlook the influence of latent factors such as emotional state, health conditions, and physical actions on vehicle safety and operational efficiency.

Retrieval-Augmented Generation Retrieval-Augmented Generation (RAG) enhances the accuracy of knowledge-intensive tasks by retrieving relevant information from

an external knowledge base that can be continuously updated (Gao et al., 2023b). Specifically, RAG begins by collecting external information, archiving previous user queries (Ma et al., 2023), and organizing the knowledge into a structured repository provided by the user. An efficient indexing mechanism is then applied to retrieve image and textual data. Once relevant information is retrieved based on user-provided prompts—through an optimized combination of priority ranking and retrieval algorithms—RAG integrates the query with the selected knowledge, processes it through a large model, and generates a more accurate response. This capability enables RAG to address user-specific queries while simultaneously maintaining a coherent dialogue by aggregating historical interaction contexts. For this paper, we deploy a learnable RAG framework to store and utilize customized user data, including individual driving habits and behaviors. This facilitates the development of a more human-centered vehicle driving agent that can adapt to the unique requirements of specific users.

7. Conclusion

This paper presents the development of a novel driving copilot framework, Sage Deer, designed to provide a super-aligned and generalist solution for intelligent vehicles. The framework integrates multi-view and multi-modal inputs, adapting to individual user preferences and needs while maintaining strong perception, understanding, and decision-making capabilities across diverse driving scenarios. Additionally, a Continuous Latent Chain Elicitation (CLCE) mechanism is proposed to enhance both super-aligned and generalist abilities by tapping into the inherent reasoning capabilities of large language models (LLMs).

References

- State farm distracted driver detection, 2016. <https://www.kaggle.com/c/state-farm-distracted-driver-detection>.
- Abtahi, S., Omidyeganeh, M., Shirmohammadi, S., and Hariri, B. Yawdd: A yawning detection dataset. In *Proceedings of the 5th ACM multimedia systems conference*, pp. 24–28, 2014.
- Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Aleman, F. L., Almeida, D., Altenschmidt, J., Altman, S., Anadkat, S., et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Bain, M., Nagrani, A., Varol, G., and Zisserman, A. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021.
- Bobbia, S., Macwan, R., Benezeth, Y., Mansouri, A., and Dubois, J. Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters*, 124:82–90, 2019a.
- Bobbia, S., Macwan, R., Benezeth, Y., Mansouri, A., and Dubois, J. Unsupervised skin tissue segmentation for remote photoplethysmography. *Pattern Recognition Letters*, 124:82–90, 2019b.
- Cui, C., Ma, Y., Cao, X., Ye, W., and Wang, Z. Drive as you speak: Enabling human-like interaction with large language models in autonomous vehicles. *2024 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*, pp. 902–909, 2023a. URL <https://api.semanticscholar.org/CorpusID:262054629>.
- Cui, C., Ma, Y., Cao, X., Ye, W., and Wang, Z. Receive, reason, and react: Drive as you say with large language models in autonomous vehicles. *ArXiv*, abs/2310.08034, 2023b. URL <https://api.semanticscholar.org/CorpusID:263908840>.
- Cui, C., Yang, Z., Zhou, Y., Ma, Y., Lu, J., and Wang, Z. Personalized autonomous driving with large language models: Field experiments. 2023c. URL <https://api.semanticscholar.org/CorpusID:266335383>.
- Du, G., Wang, Z., Gao, B., Mumtaz, S., Abualnaja, K. M., and Du, C. A convolution bidirectional long short-term memory neural network for driver emotion recognition. *IEEE Transactions on Intelligent Transportation Systems*, 22(7):4570–4578, 2020.
- Eraqi, H. M., Abouelnaga, Y., Saad, M. H., and Moustafa, M. N. Driver distraction identification with an ensemble of convolutional neural networks. *Journal of Advanced Transportation*, 2019, 2019.
- Gao, P., Han, J., Zhang, R., Lin, Z., Geng, S., Zhou, A., Zhang, W., Lu, P., He, C., Yue, X., et al. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv preprint arXiv:2304.15010*, 2023a.
- Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., Dai, Y., Sun, J., Guo, Q., Wang, M., and Wang, H. Retrieval-augmented generation for large language models: A survey. *ArXiv*, abs/2312.10997, 2023b. URL <https://api.semanticscholar.org/CorpusID:266359151>.
- Guesdon, R., Crispim-Junior, C., and Tougne, L. Dripe: A dataset for human pose estimation in real-world driving settings. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2865–2874, 2021.
- Hao, S., Sukhbaatar, S., Su, D., Li, X., Hu, Z., Weston, J., and Tian, Y. Training large language models to reason in a continuous latent space. *arXiv preprint arXiv:2412.06769*, 2024.
- Hecht, T., Feldhütter, A., Radlmayr, J., Nakano, Y., Miki, Y., Henle, C., and Bengler, K. A review of driver state monitoring systems in the context of automated driving. *Advances in Intelligent Systems and Computing*, 2018. URL <https://api.semanticscholar.org/CorpusID:56790391>.
- Jain, A., Koppula, H. S., Soh, S., Raghavan, B., Singh, A., and Saxena, A. Brain4cars: Car that knows before you do via sensory-fusion deep learning architecture. *arXiv preprint arXiv:1601.00740*, 2016.
- Jegham, I., Ben Khalifa, A., Alouani, I., and Mahjoub, M. A. Mdad: A multimodal and multiview in-vehicle driver action dataset. In *International Conference on Computer Analysis of Images and Patterns*, pp. 518–529, 2019.
- Jegham, I., Khalifa, A. B., Alouani, I., and Mahjoub, M. A. A novel public dataset for multimodal multiview and multispectral driver distraction analysis: 3mdad. *Signal Processing: Image Communication*, 88:115960, 2020.
- Jeong, M. and Ko, B. C. Driver’s facial expression recognition in real-time for safe driving. *Sensors*, 18(12):4270, 2018.
- Kasahara, I., Stent, S., and Park, H. S. Look both ways: Self-supervising driver gaze estimation and road scene saliency. In *European Conference on Computer Vision (ECCV)*, pp. 126–142, 2022.

-
- Kopuklu, O., Zheng, J., Xu, H., and Rigoll, G. Driver anomaly detection: A dataset and contrastive learning approach. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pp. 91–100, 2021.
- Li, J., Li, D., Savarese, S., and Hoi, S. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023a.
- Li, K., He, Y., Wang, Y., Li, Y., Wang, W., Luo, P., Wang, Y., Wang, L., and Qiao, Y. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023b.
- Li, W., Cui, Y., Ma, Y., Chen, X., Li, G., Zeng, G., Guo, G., and Cao, D. A spontaneous driver emotion facial expression (defe) dataset for intelligent vehicles: Emotions triggered by video-audio clips in driving scenarios. *IEEE Transactions on Affective Computing*, 2021a.
- Li, W., Zeng, G., Zhang, J., Xu, Y., Xing, Y., Zhou, R., Guo, G., Shen, Y., Cao, D., and Wang, F.-Y. Cogemonet: A cognitive-feature-augmented driver emotion recognition model for smart cockpit. *IEEE Transactions on Computational Social Systems*, 9(3):667–678, 2021b.
- Li, W., Cao, D., Tan, R., Shi, T., Gao, Z., Ma, J., Guo, G., Hu, H., Feng, J., and Wang, L. Intelligent cockpit for intelligent connected vehicles: Definition, taxonomy, technology and evaluation. *IEEE Transactions on Intelligent Vehicles*, 2023c.
- Lin, B., Zhu, B., Ye, Y., Ning, M., Jin, P., and Yuan, L. Video-llava: Learning united visual representation by alignment before projection. *arXiv preprint arXiv:2311.10122*, 2023.
- Liu, H., Li, C., Wu, Q., and Lee, Y. J. Visual instruction tuning. *Advances in neural information processing systems*, 36, 2024a.
- Liu, X., Zhang, Y., Yu, Z., Lu, H., Yue, H., and Yang, J. rppg-mae: Self-supervised pretraining with masked autoencoders for remote physiological measurements. *IEEE Transactions on Multimedia*, 2024b.
- Lu, H., Yu, Z., Niu, X., and Chen, Y.-C. Neuron structure modeling for generalizable remote physiological measurement. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 18589–18599, 2023.
- Lu, H., Niu, X., Wang, J., Wang, Y., Hu, Q., Tang, J., Zhang, Y., Yuan, K., Huang, B., Yu, Z., et al. Gpt as psychologist? preliminary evaluations for gpt-4v on visual affective computing. *2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) workshop*, 2024.
- Luo, R., Zhao, Z., Yang, M., Dong, J., Qiu, M., Lu, P., Wang, T., and Wei, Z. Valley: Video assistant with large language model enhanced ability. *arXiv preprint arXiv:2306.07207*, 2023.
- Lv, H. and Sun, Q. Video anomaly detection and explanation via large language models. *arXiv preprint arXiv:2401.05702*, 2024.
- Ma, X., Gong, Y., He, P., Zhao, H., and Duan, N. Query rewriting for retrieval-augmented large language models. *ArXiv*, abs/2305.14283, 2023. URL <https://api.semanticscholar.org/CorpusID:258841283>.
- Maaz, M., Rasheed, H., Khan, S., and Khan, F. S. Video-chatgpt: Towards detailed video understanding via large vision and language models. *arXiv preprint arXiv:2306.05424*, 2023.
- Mao, J., Qian, Y., Zhao, H., and Wang, Y. Gpt-driver: Learning to drive with gpt. 2023.
- Martin, M., Roitberg, A., Haurilet, M., Horne, M., Reiß, S., Voit, M., and Stiefelhagen, R. Drive&act: A multi-modal dataset for fine-grained driver behavior recognition in autonomous vehicles. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 2801–2810, 2019.
- Muhammad Maaz, Hanoona Rasheed, S. K. and Khan, F. Video-chatgpt: Towards detailed video understanding via large vision and language models. *ArXiv 2306.05424*, 2023.
- Niu, X., Shan, S., Han, H., and Chen, X. Rhythmnet: End-to-end heart rate estimation from face via spatial-temporal representation. *IEEE Transactions on Image Processing*, 29:2409–2423, 2019.
- Niu, X., Shan, S., Han, H., and Chen, X. Rhythmnet: End-to-end heart rate estimation from face via spatial-temporal representation. *IEEE Trans. on Image Process.*, 29:2409–2423, 2020.
- Oh, G., Jeong, E., Kim, R. C., Yang, J. H., Hwang, S., Lee, S., and Lim, S. Multimodal data collection system for driver emotion recognition based on self-reporting in real-world driving. *Sensors*, 22(12):4402, 2022.
- Ortega, J. D., Kose, N., Cañas, P., Chao, M.-A., Unnervik, A., Nieto, M., Otaegui, O., and Salgado, L. Dmd: A large-scale multi-modal driver monitoring dataset for attention and alertness analysis. In *European Conference on Computer Vision (ECCV)*, pp. 387–405, 2020.

-
- Papineni, K., Roukos, S., Ward, T., and Zhu, W.-J. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Revanur, A., Li, Z., Ciftci, U. A., Yin, L., and Jeni, L. A. The first vision for vitals (v4v) challenge for non-contact video-based physiological estimation. In *Proc. CVPR workshop*, pp. 2760–2767, 2021a.
- Revanur, A., Li, Z., Ciftci, U. A., Yin, L., and Jeni, L. A. The first vision for vitals (v4v) challenge for non-contact video-based physiological estimation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 2760–2767, 2021b.
- Saad, M. H., Khalil, M. I., and Abbas, H. M. End-to-end driver distraction recognition using novel low lighting support dataset. In *IEEE International Conference on Computer Engineering and Systems (ICCES)*, pp. 1–6, 2020.
- Sima, C., Renz, K., Chitta, K., Chen, L., Zhang, H., Xie, C., Luo, P., Geiger, A., and Li, H. Drivelm: Driving with graph visual question answering. *arXiv preprint arXiv:2312.14150*, 2023.
- Stricker, R., Müller, S., and Gross, H.-M. Non-contact video-based pulse rate measurement on a mobile service robot. In *Proc. IEEE ISRHIC*, pp. 1056–1062, 2014a.
- Stricker, R., Müller, S., and Gross, H.-M. Non-contact video-based pulse rate measurement on a mobile service robot. In *The 23rd IEEE International Symposium on Robot and Human Interactive Communication*, pp. 1056–1062. IEEE, 2014b.
- Su, Y., Lan, T., Li, H., Xu, J., Wang, Y., and Cai, D. Pandagpt: One model to instruction-follow them all. *arXiv preprint arXiv:2305.16355*, 2023.
- Tran, D., Manh Do, H., Sheng, W., Bai, H., and Chowdhary, G. Real-time detection of distracted driving based on deep learning. *IET Intelligent Transport Systems*, 12(10): 1210–1219, 2018.
- Wang, W., Stuijk, S., and De Haan, G. A novel algorithm for remote photoplethysmography: Spatial subspace rotation. *IEEE transactions on biomedical engineering*, 63(9):1974–1984, 2015.
- Wang, Y., Li, K., Li, Y., He, Y., Huang, B., Zhao, Z., Zhang, H., Xu, J., Liu, Y., Wang, Z., et al. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022.
- Wang, Z., Bian, Y., Shladover, S. E., Wu, G., Li, S. E., and Barth, M. J. A survey on cooperative longitudinal motion control of multiple connected and automated vehicles. *IEEE Intelligent Transportation Systems Magazine*, 12:4–24, 2020. URL <https://api.semanticscholar.org/CorpusID:210931452>.
- Xi, L., Chen, W., Zhao, C., Wu, X., and Wang, J. Image enhancement for remote photoplethysmography in a low-light environment. In *FG*, pp. 1–7. IEEE, 2020a.
- Xi, L., Chen, W., Zhao, C., Wu, X., and Wang, J. Image enhancement for remote photoplethysmography in a low-light environment. In *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*, pp. 1–7. IEEE, 2020b.
- Xiao, H., Liu, T., Sun, Y., Li, Y., Zhao, S., and Avolio, A. Remote photoplethysmography for heart rate measurement: A review. *Biomedical Signal Processing and Control*, 88:105608, 2024.
- Yang, D., Huang, S., Xu, Z., Li, Z., Wang, S., Li, M., Wang, Y., Liu, Y., Yang, K., Chen, Z., 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*, pp. 20459–20470, 2023a.
- Yang, J., Xing, S., Chen, Y., Qiu, R., Hua, C., and Dong, D. A comprehensive evaluation model for the intelligent automobile cockpit comfort. *Scientific Reports*, 12(1): 15014, 2022.
- Yang, Y., Zhang, Q., Li, C., Marta, D. S. o., Batool, N., and Folkesson, J. Human-centric autonomous systems with llms for user command reasoning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV) Workshops*, pp. 988–994, January 2024.
- Yang, Z., Jia, X., Li, H., and Yan, J. Llm4drive: A survey of large language models for autonomous driving. *ArXiv*, abs/2311.01043, 2023b. URL <https://api.semanticscholar.org/CorpusID:264935408>.
- Ye, Q., Xu, H., Xu, G., Ye, J., Yan, M., Zhou, Y., Wang, J., Hu, A., Shi, P., Shi, Y., et al. mplug-owl: Modularization empowers large language models with multimodality. *arXiv preprint arXiv:2304.14178*, 2023.
- Zhang, C., Li, R., Kim, W., Yoon, D., and Patras, P. Driver behavior recognition via interwoven deep convolutional neural nets with multi-stream inputs. *IEEE Access*, 8: 191138–191151, 2020.

Zhang, H., Li, X., and Bing, L. Video-llama: An instruction-tuned audio-visual language model for video understanding. *arXiv preprint arXiv:2306.02858*, 2023a. URL <https://arxiv.org/abs/2306.02858>.

Zhang, R., Han, J., Liu, C., Gao, P., Zhou, A., Hu, X., Yan, S., Lu, P., Li, H., and Qiao, Y. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023b.

Zhao, C., Zhang, B., He, J., and Lian, J. Recognition of driving postures by contourlet transform and random forests. *IET Intelligent Transport Systems*, 6(2):161–168, 2012.

Zhu, D., Chen, J., Shen, X., Li, X., and Elhoseiny, M. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023.