# Vision Language Models in Autonomous Driving and Intelligent Transportation Systems

Xingcheng Zhou ©, Mingyu Liu ©, Bare Luka Zagar ©, Ekim Yurtsever *Member, IEEE* ©,
Alois C. Knoll, *Fellow, IEEE* ©

*Abstract*—The applications of Vision-Language Models (VLMs) in the fields of Autonomous Driving (AD) and Intelligent Transportation Systems (ITS) have attracted widespread attention due to their outstanding performance and the ability to leverage Large Language Models (LLMs). By integrating language data, the vehicles, and transportation systems are able to deeply understand real-world environments, improving driving safety and efficiency. In this work, we present a comprehensive survey of the advances in language models in this domain, encompassing current models and datasets. Additionally, we explore the potential applications and emerging research directions. Finally, we thoroughly discuss the challenges and research gap. The paper aims to provide researchers with the current work and future trends of VLMs in AD and ITS.

*Index Terms*—Vision-Language Model, Large Language Model, Autonomous Driving, Intelligent Transportation Systems.
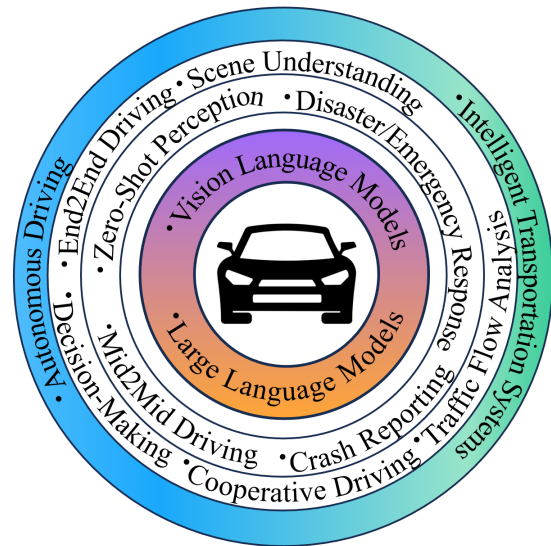
Fig. 1. Vision-Language Models and Large Language Models offer advancements in traditional tasks and pave the way for innovative applications in AD and ITS.

## I. INTRODUCTION

Intelligent mobility is important in modern civilization, driving economic growth, supporting urban development, and strengthening social connections. In recent years, the rapid evolution of deep learning and computing power has profoundly influenced transportation, enhancing its efficiency and intelligence. Two emerging domains central to intelligent mobility are autonomous driving (AD) and Intelligent Transportation Systems (ITS).

Autonomous driving strives to enable vehicles to perceive environments and drive intelligently. The current autonomous driving technologies, especially those related to perception and prediction, have tremendously benefited from advances in computer vision. For instance, perception modules, typically using Convolutional Neural Networks (CNNs) or Transformers [1], process data from sensors like cameras or LiDAR to accurately identify and localize entities in their surroundings. However, despite these technological strides, the current computer vision solutions still struggle in complex and rapidly dynamic environments. They often fail to capture intricate details or understand context, leading to potential safety concerns, and limiting the move toward more advanced autonomous driving. On the other hand, intelligent traffic systems aim to enhance transportation safety and mobility, but challenges persist even

X. Zhou, M. Liu, BL. Zagar, and AC. Knoll are with the Chair of Robotics, Artificial Intelligence and Real-Time Systems, Technical University of Munich, 85748 München, Germany (E-mail: xingcheng.zhou@tum.de, mingyu.liu@tum.de, bare.luka.zagar@tum.de, knoll@in.tum.de)

E. Yurtsever is with the College of Engineering, Center for Automotive Research, The Ohio State University, Columbus, OH 43212, USA (E-mail: yurtsever.2@osu.edu)

as their efficiency and reliability have improved over the years. Real-time traffic flow prediction, for instance, is vulnerable to various environmental factors such as weather, crash events, or road construction.

The emergence of LLMs [2]–[6] and VLMs [7]–[12] provides potential solutions for the inherent limitations of current autonomous driving and intelligent transportation systems. These novel technologies synthesize linguistic and visual data, promising a future where vehicles and systems deeply understand their surroundings. This heralds a new era of intelligent, efficient, and explainable transportation. Besides enhancing traditional tasks in AD or ITS, such as object detection or traffic flow prediction, emerging domains include zero-shot perception, and accident analysis, as shown in Fig. 1. Given the surge in research applying language models to autonomous driving and intelligent systems, a systematic and comprehensive survey is of great importance to the research community. However, existing surveys [12]–[16] focus either on LLMs, VLMs, AD, or ITS individually. To the best of our knowledge, there is still no review that systematically discusses VLMs applications in AD and ITS.

To this end, we present a review of existing algorithms of vision-language models in autonomous driving and intelligent traffic systems, highlighting the recent technological trends in the research community. We illustrate the taxonomy of this
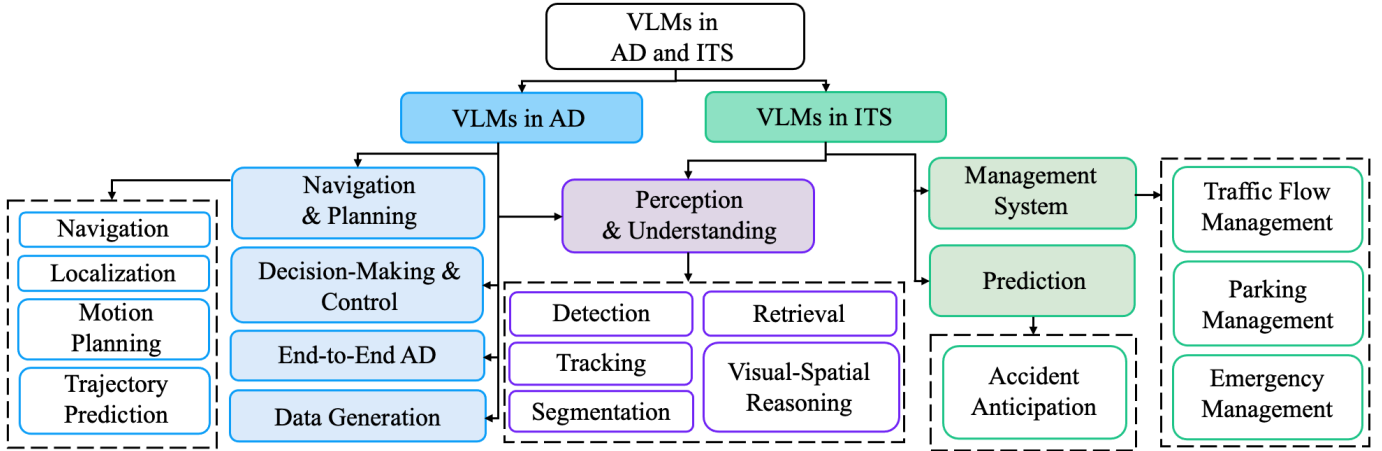
Fig. 2. Overview of taxonomy of VLMs in autonomous driving and intelligent traffic systems.

paper in Fig. 2. The main contributions of this work can be summarized as follows:

- We present the first comprehensive survey on the applications of vision-language models in autonomous driving and intelligent traffic systems.
- We systematically summarize and analyze the existing work and datasets.
- We explore potential applications and technological advances of VLMs in AD and ITS.
- We provide an in-depth discussion on the challenges and research gap in this domain.

## II. BACKGROUND

This section gives a thorough introduction of the related backgrounds, delving into the foundational concepts underlying these technologies: autonomous driving (II-A), intelligent transportation systems (II-B), Large Language Model (II-C), and Vision-Language Model (II-D).

### A. Autonomous Driving

Autonomous driving targets the development of vehicles that can navigate and control themselves without human intervention, reducing incidents and improving traffic efficiency. The driving automation level defined by the Society of Automotive Engineers [17] from Level 0 (No Automation) to Level 5 (Full Automation). As autonomy increases, human intervention reduces, while the requirements for the vehicle to understand its surroundings increase. Currently, most commercial vehicles are at Level 2 or 3, providing partial automation but still requiring driver supervision.

Exiting autonomous driving solutions can be broadly categorized into the classic modular paradigm and the end-to-end approach. However, as mentioned in [13], these schemes all face serious challenges such as interpretability, generalization, causal confusion, robustness, etc. Researchers have attempted to address these issues using various methods, but constructing a safe, stable, and interpretable AD system remains an open topic.

### B. Intelligent Transportation Systems

Intelligent Transportation Systems (ITS) employ advanced technologies to enhance traffic efficiency and safety by optimizing the broader traffic environment. By integrating live data from various sources including road sensors and road users, ITS encompasses a broad range of services and applications, from adaptive traffic signal control to real-time traffic surveillance, accident detection and anticipation, traffic flow prediction, and cooperative vehicle infrastructure systems. Despite the expanding applications of ITS with advancements in sensing, communication, and machine learning technologies, some significant challenges still need to be addressed. As highlighted in [16], driving is a social activity that often requires frequent interaction with other traffic participants, but the intelligence and common sense that humans rely on are still lacking in current systems.

### C. Large Language Models

Large Language Models (LLMs) typically refer to the language models with a massive number of parameters, often in the order of billion or more. The most notable characteristic of LLMs is the exhibition of emergent abilities, such as the capacity for few-shot or zero-shot transfer learning across numerous downstream tasks, strong multi-step reasoning capabilities, and the ability to follow instructions, which are normally not present in smaller models.

ChatGPT, specifically GPT-3.5 [6], serves as a milestone in the development of LLMs. Since its release, GPT-3.5 has consistently drawn attention due to its exceptional performance. An increasing number of researchers are beginning to explore and harness the powerful linguistic understanding, interpretation, analysis, and reasoning capabilities of LLMs to solve problems that were previously hard or even impossible to tackle.

### D. Vision-Language Models

Vision-Language Models (VLMs) bridge the capabilities of Natural Language Processing (NLP) and Computer Vision (CV), breaking down the boundaries between text and visual
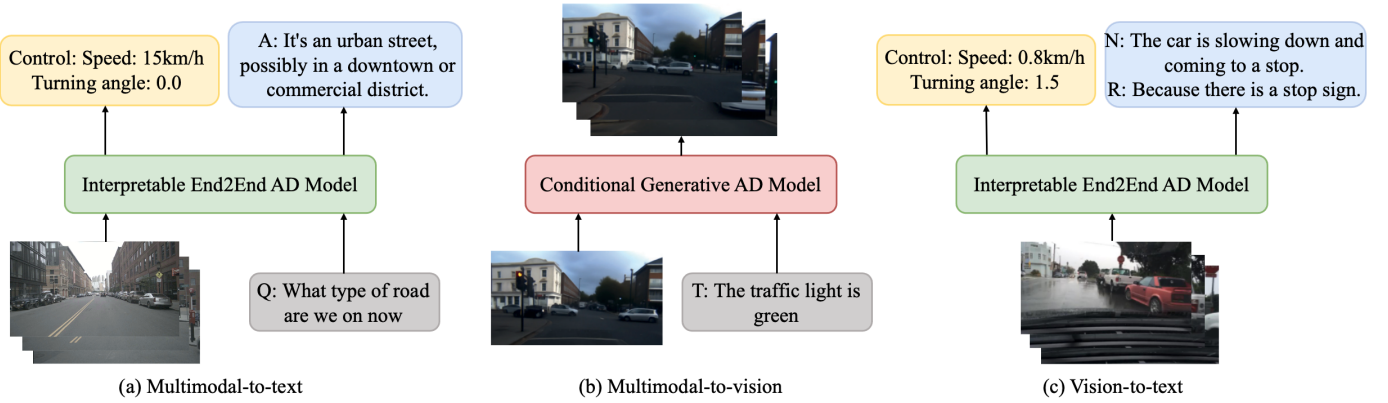
Fig. 3. Overview of mainstream Vision-Language Models in Autonomous Driving. **(a)** Multimodal-to-text models take text and image or video as input and generate text, as in [18]. **(b)** Multimodal-to-vision models take image and text as input and output image or video, depicted with [19]. **(c)** Vision-to-text models accept video or image as input and produce text as output, e.g. GAIA-1 [20].

information to connect multimodal data. With the rise of LLMs, there is also an increasing focus on exploring how to effectively incorporate visual modules into LLMs to perform multimodal tasks.

Mainstream Vision-Language Models in AD can be broadly categorized into Multimodal-to-Text (M2T) [18] [21]–[23], Multimodal-to-Vision (M2V) [20] [24], and Vision-to-Text (V2T) [19] [25] based on input-output modality types, as shown in Fig. 3. M2T typically takes image-text or video-text as input and produces text as output; Correspondingly, M2V accepts image-text as input and generates image or video as output, while V2T takes image or video as input and generates texts as output. As illustrated in Fig. 4, according to the inter-modality information connection approaches, VLMs employed in AD can be divided into Vision-Text-Fusion (VTF) [18] [19] [21] [26]–[28] and Vision-Text-Matching (VTM) [25] [29]–[38]. VTF employs various fusion methods to effectively integrate vision embedding and language embedding, and jointly optimize the feature representation that performs better for the target task. In contrast, VTM, including image-text matching [39], [40] and video-text matching [41], learns a joint representation space by forcing vision-text pairs semantically close to each other and unpaired instances distant to each other, achieving cross-modal semantic alignment, enabling cross-modal semantic propagation. CLIP [39], a milestone image-text matching work in VLMs, captures image feature representations associated with language, and achieves zero-shot transfer ability by training on a vast number of image-text pairs through contrastive learning.

## III. VLMs IN AUTONOMOUS DRIVING

An increasing number of initiatives are endeavoring to implement VLMs across various aspects of AD. In this section, we present the existing work of VLMs within AD consisting of perception and understanding (III-A), navigation and planning (III-B), decision-making and control (III-C), end-2-end AD (III-D), and data generation (III-E) The summarized current methods are shown in Tab. I.
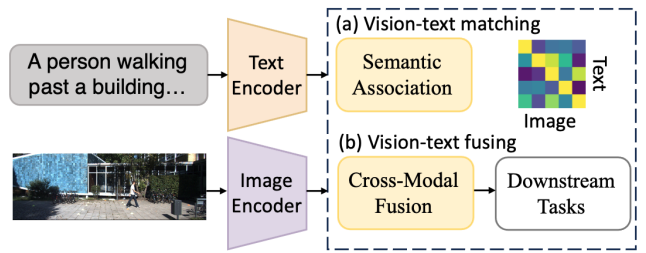


Fig. 4. Two inter-modality connection approaches of Vision-Language Models in Autonomous Driving: **(a)** Vision-text matching. We demonstrate the semantic similarity matching in the top-right of this figure. **(b)** Vision-text fusion. The fused features can be used for downstream tasks. The figure is from KITTI [42].

### A. Perception and Understanding

In autonomous driving perception, VLMs, especially those models pre-trained on large-scale datasets with image-text-matching like [39], has spurred numerous new studies. These studies [30], [33] leverage the substantial prior knowledge of pre-trained large VLMs to boost the performance of perception and understanding, and further introduce many novel tasks in the field.

**Pedestrian Detection.** Human-like object confusion and insufficient border case samples are the inherent challenges in pedestrian detection. To this end, VLPD [29] first proposes a vision-language extra-annotation-free method to enhance model's capability of distinguishing confusing human-like objects. It employs CLIP to acquire the pixel-wise explicit semantic contexts, and distance pedestrian features from the features of other categories through contrastive learning, improving the detection capabilities for broader cases. UMPD [30] also utilizes the zero-shot semantic classes from the CLIP, and proposes a fully-unsupervised multi-view Pedestrian Detection approach without the need of manual annotations.

**Object Referring.** Compared to traditional perception tasks in AD, such as object detection, tracking, and semantic segmentation, the introduction of language enables the model to attain a more finely-grained and comprehensively unconstrained capability to understand the surrounding environment. Object Referring (OR) is a representative task in this area that aims

| Method | Year | Tasks | Type | Contribution |
|---|---|---|---|---|
| **Perception and Understanding** | | | | |
| NIC [25] | 2020 | IC | V2T, VTM | Introduce image captioning task to autonomous driving for traffic scene understanding. |
| VLPD [29] | 2023 | PD | VTM | Propose a vision-language extra-annotation-free method for pedestrian detection. |
| UMPD [30] | 2023 | PD | VTM | Propose a fully-unsupervised multi-view pedestrian detector without manual annotations. |
| MSSG [26] | 2023 | SOR | VTF | Incorporate natural language with 3D LiDAR point cloud and camera image for LiDAR grounding task. |
| TransRMOT [31] | 2023 | MOR-T | VTM | Propose Referring Multi-Object Tracking task, dataset, and benchmark based on KITTI. |
| PromptTrack [32] | 2023 | MOR-T | VTM | Propose Multi-view 3D Referring Multi-Object Tracking dataset based on nuScenes. |
| OpenScene [33] | 2023 | OV-3DSS | VTM | Propose a zero-shot method for extracting 3D dense features from open vocabulary embedding space. |
| CLIP2Scene [34] | 2023 | OV-3DSS | VTM | Distil CLIP knowledge to a 3D network for 3D scene understanding. |
| UP-VL [35] | 2023 | OV-3DOD-T | VTM | Introduce semantic-aware unsupervised detection for objects in any motion state. |
| Zelda [36] | 2023 | LVR | VTM | Introduce a video analytic system using VLMs that delivers semantically diverse, high-quality results. |
| NuScenes-QA [21] | 2023 | VQA | M2T, VTF | Introduce visual question-answering task and baseline model in autonomous driving based on nuScenes. |
| Talk2BEV [23] | 2023 | VSR, OL-DM | M2T | Augment BEV maps with language to enable general-purpose visuolinguistic reasoning for AD scenarios. |
| LLM-AD [43] | 2023 | SAD | LLM | Propose a framework to detect semantic anomalies leveraging LLMs' reasoning abilities. |
| **Navigation and Planning** | | | | |
| TttV [44] | 2019 | LGN | M2T, VTF | Propose a modular framework following Language Instructions to provide accurate waypoints for LGN. |
| GtN [37] | 2022 | LGN | M2T, VTM | Present a vision language navigation tool to controls vehicle movements based on language-driven commands. |
| ALT-Pilot [38] | 2023 | VLL, LGN | M2T, VTM | Propose an autonomous navigation system based on language-augmented topometric maps. |
| GPT-Driver [45] | 2023 | MP | LLM | Transform motion planning task into a language modeling problem. |
| CoverNet-T [28] | 2023 | TP | M2T, VTF | Introduce text and image representation for trajectory prediction. |
| **Decision Making and Control** | | | | |
| Advisable-DM [46] | 2020 | OL-DM | M2T, VTF | Propose a advisable and explainable model for self-driving systems. |
| LanguageMPC [47] | 2023 | OL-DM | LLM | Devised a chain-of-thought framework for LLMs in driving scenarios that divides decision-making process into numerous sub-problems. |
| DaYS [48] | 2023 | OL-DM, MP | LLM | Provide a framework to integrate LLMs into autonomous vehicles. |
| DwLLMs [49] | 2023 | OL-C, VSR | LLM | Propose a model with object-level fusion for Explainable AD. |
| DiLU [50] | 2023 | CL-DM | LLM | Instill knowledge-driven capability into autonomous driving systems from the perspective of how humans drive. |
| SurrealDriver [51] | 2023 | CL-DM | LLM | Develop an LLM-based driver agent for complex urban environments. |
| DLaH [22] | 2023 | CL-DM | M2T, VTF | Show the feasibility and decision-making ability of LLM in driving scenarios in the simulated environment. |
| **End-to-End Autonomous Driving** | | | | |
| DriveGPT4 [18] | 2023 | VQA, OL-C | M2T, VTF | Propose a multimodal model for interpretable AD with multimodal as input and control signal as output. |
| ADAPT [19] | 2023 | VSR, OL-DM | V2T, VTF | Propose a end-to-end action narration and reasoning framework for self-driving vehicles. |
| **Data Generation** | | | | |
| DriveGAN [52] | 2021 | CVG | M2V, VTF | Propose an end-to-end controllable differentiable neural driving simulator for scenario re-creation. |
| GAIA-1 [20] | 2023 | CVG | M2V, VTF | Propose a multi-modal generative model to produce realistic driving scenarios with precise control over ego-vehicle actions and scene attributes. |
| DriveDreamer [24] | 2023 | CVG | M2V, VTF | Introduce the first world model derived from real-world driving scenarios, capable of generating high-quality driving videos and sound driving policies. |
| BEVControl [53] | 2023 | CIG | M2V, VTF | Present a sketch-based street-view images generative model based on nuScenses. |

TABLE I
OVERVIEW OF LLMs AND VLMs IN AUTONOMOUS DRIVING

TASKS: **IC**: IMAGE CAPTIONING **PD**: PEDESTRIAN DETECTION **SOR**: SINGLE-OBJECT REFERRING **MOR-T**: MULTI-OBJECT REFERRING AND TRACKING **OV-3DSS**: OPEN-VOCABULARY 3D SEMANTIC SEGMENTATION **OV-3DOD-T**: OPEN-VOCABULARY 3D OBJECT DETECTION AND TRACKING **LVR**: LANGUAGE-GUIDED VIDEO RETRIEVAL **VQA**: VISUAL QUESTION ANSWERING **VSR**: VISUAL-SPATIAL REASONING **OL-DM**: OPEN-LOOP DECISION MAKING **SAD**: SEMANTIC ANOMALY DETECTION **LGN**: LANGUAGE-GUIDED NAVIGATION **VLL**: VISION-LANGUAGE LOCALIZATION **MP**: MOTION PLANNING **TP**: TRAJECTORY PREDICTION **CL-DM**: CLOSED-LOOP DECISION MAKING **OL-C**: OPEN-LOOP CONTROL **CVG**: CONDITIONAL VIDEO GENERATION **CIG**: CONDITIONAL IMAGE GENERATION
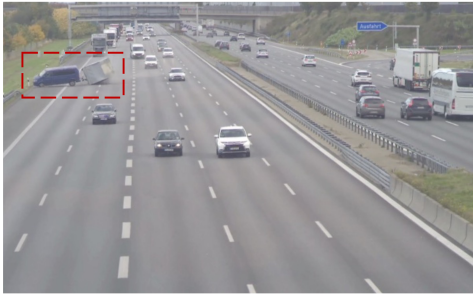TYPE: **V2T**: VISION-TO-TEXT **VTM**: VISION-TEXT MATCHING **VTF**: VISION-TEXT FUSION **M2T**: MULTIMODAL TO TEXT **LLM**: LARGE LANGUAGE MODEL

to localize the described objects using boxes or masks based on the language queries. MSSG [26] proposes a multi-modal 3D single object referring (SOR) task in autonomous driving scenarios. It trains a multi-modal single-shot grounding model by fusing the image, LiDAR, and language features under Bird Eye View (BEV) and predicts the targeted region directly from the detector without any post-processing. TransRMOT [31] extents from the SOR task to Multiple-Object-Referring and Tracking (MOR-T), and construct Refer-KITTI benchmark based on the KITTI dataset. Given a language query, TransRMOT can detect and track an arbitrary number of referent objects in videos. Similarly, PromptTrack [32] proposes a language prompt set for nuScenes dataset, and constructs the NuPrompt benchmark. Compared to ReferKITTI, NuPrompt inherits nuScenes' multi-view attributes, making it applicable for multi-view MOR-T task.

**Open-Vocabulary 3D Object Detection and Semantic Segmentation.** Due to the potent zero-shot transfer and cross-modal mapping capabilities of large VLMs, object detection

and semantic segmentation are endowed with the ability to perform open-vocabulary detection and segmentation on unseen samples. UP-VL framework [35] first presents an unsupervised multi-modal auto-labeling pipeline to generate point-level features and object-level bounding boxes and tracklets for open-vocabulary class-agnostic 3D detector supervision, which is further utilized for proposing 3D bounding box at inference time. Together with the assigned semantic labels through similarity matching, UP-VL framework realizes unsupervised open-vocabulary 3D detection and tracking (OV-3DOD-T) of both static and moving traffic participants in AD scenarios. Opensene [33] utilizes 2D-3D projective correspondence to enforce consistency between 3D point cloud embeddings and the corresponding fused 2D image features. Essentially, it aligns the 3D point cloud representation with CLIP's image-text representation to acquire zero-shot understanding capabilities for dense 3D point features. OpenScene is primarily evaluated in indoor scenarios but also demonstrates satisfactory open vocabulary 3D semantic segmentation (OV-3DSS) capabilities on

Traffic scenario A



**Q**: Describe the traffic scene in this image.
**A**: …there appears to be a significant accident involving a large truck that has tipped over, lying partially on its side and occupying two of the highway lanes. … a few cars have pulled over to the shoulder, likely as a result of the accident or to avoid it. …
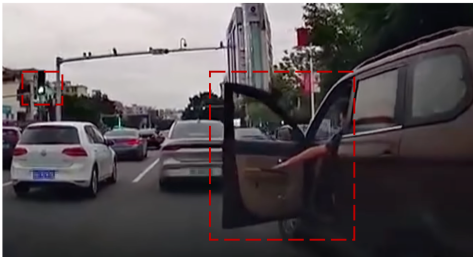**Q**: Based on the description, list the three most possible reasons for this accident.
**A**: 1. Loss of Control: The truck driver may have lost control of the vehicle due to a mechanical failure, a sudden maneuver to avoid an obstacle or another vehicle, or even due to a tire blowout.
    2. High Speed or Improper Lane Change: The truck could have been traveling at a speed not suitable for its size and weight, leading to tipping when trying to change lanes or make a sudden movement.
    3. External Factors: There might have been external factors, such as an object on the road, another vehicle's unexpected behavior, or even factors like a brief distraction of the truck driver, which led to the accident.

Traffic scenario B



**Q**: Describe the traffic scene in this image.
**A**: A multi-lane city road with several cars is shown. …traffic light displays a green light, suggesting that the traffic is allowed to move forward. … on the right side of the image, an open car door …appears to be swung wide open while the vehicle is in the lane…

**Q**: Based on the description, can the ego-car drive forward?
**A**: …the ego-car (the car from which the image is taken) appears to have a clear path directly ahead, with the traffic light displaying a green signal. Thus, it seems the ego-car can drive forward. However, it's essential to note the unusual situation with the open car door on the right. The driver should proceed with caution and be aware of the surrounding environment to ensure safety.

Fig. 5. Evaluation of GPT-4V [54] on traffic anomaly scene interpretability and reasoning. The traffic scenario A (first row) is sourced from Providentia++. The traffic accident is highlighted in the red rectangle dotted box at the left-top of this figure. In the second row, we demonstrate traffic scenario B from a video, where a car's door is open while the car is in the middle of the lane. The red dotted boxes show the opened door and the traffic light.

nuScenes. Similarly, CLIP2Scene [34] explores how to utilize CLIP to assist with 3D scene understanding in autonomous driving. By seeking connections in modality between pixel-text mapping and pixel-point mapping, CLIP2Scene constructs point-text pairs and pixel-point-text pairs for contrastive learning training respectively. The objectiveness is also to ensure that the 3D point feature and its corresponding language achieve semantic consistency, thereby facilitating OV-3DSS. Experiments indicate that using CLIP2Scene as pre-training greatly outperforms other self-supervised methods.

**Traffic Scene Understanding.** A correct and high-level understanding of the traffic scene is critical for driving safety. In Fig. 5, we illustrate examples of the understanding ability of GPT-4V [54] in a traffic accident scenario and an urban road scene with potential risk. There are several exploratory works utilizing VLMs to understand traffic scenes through specific downstream tasks. [25] attempts to understand traffic scenes by describing images of the scenes through Image Captioning (IC). [43] transforms visual information into language descriptions and then leverages the strong reasoning capabilities of LLMs to address Semantic Anomaly Detection (SAD). Zelda [36] employs the semantical similarity capabilities of VLMs for Video Retrieval (VR), achieving performance that surpasses other state-of-the-art video analytics systems. NuScenes-QA [21], establishes a benchmark for Visual Question Answering (VQA) tasks in the AD scene based on the nuScenes dataset, providing a foundation for subsequent research. Talk2BEV [23] employs Bird's Eye View (BEV) detection model, dense captioning model, and text recognition model to construct Ground-truth language-enhanced BEV map, and evaluate the model's performance in visual and spatial understanding based on VQA tasks.

We note that most aforementioned works, including OV-

3DOD, OV-3DSS, SOR, MOR-T in AD are still at the early stage, but we believe this is a promising direction and anticipate an increasing number of interesting works emerge in the future.

### B. Navigation and Planning

In the field of navigation, with the advancement of VLMs, especially the proposal of CLIP [39], Language-Guided Navigation (LGN) task begins to extend from specific pre-defined location descriptions to free and arbitrary instructions, which also prompts the development of Language-Augmented Maps [38].

**Language-Guided Navigation.** Talk to the Vehicle [27] proposes a waypoint generator network (WGN) that maps semantic occupancy and pre-defined natural language encodings (NLE) to local waypoints. The planning module then takes local waypoints to predict the trajectory for execution. Ground then Navigatie [37] solves the language-guided navigation tasks with the help of CLIP. It proposes a pipeline that takes video frames, historical trajectory context, and language command as input, and output predicted navigation mask as well as trajectory at each timestamp. ALT-Pilot [38] enhances OpenStreetMap (OSM) road networks by incorporating linguistic landmarks, including street signs, traffic signals, and other prominent environmental features that aid in localization to substitute traditional memory and compute expensive HD LiDAR maps. ALT-Pilot also leverages CLIP to precompute the feature descriptors for each of these landmarks and match them with the pixel-wise vision descriptors using cosine similarity at inference time, which facilitates the correspondence from language navigation instructions to map locations, thereby assisting multimodal localization and navigation.

**Prediction and Planning.** Some works have also begun to explore how to leverage LLMs to enhance the performance of motion planning and trajectory predictions. GPT-driver [45] reformulates motion planning as a language modeling problem and transforms the GPT-3.5 model into a motion planner for autonomous driving, capitalizing on its strong reasoning and generalization capabilities. CoverNet-T [28] proposes to train a joint encoder with text-based scene descriptions and rasterized scene images for trajectory predictions. It shows that text-based scene representation has complementary strengths to image encodings, and a joint encoder outperforms its individual counterparts.

### C. Decision-Making and Control

In the decision-making and control area of autonomous driving, several works seek to leverage the powerful common sense comprehending and reasoning capabilities of LLMs to assist drivers [46] [48], to emulate and entirely substitute drivers [22] [47] [49]–[51]. By utilizing LLMs for closed-loop control in AD, most works [22] [48] [50] [51] introduce a memory module to record driving scenarios, experiences, and other essential driving information.

LanguageMPC [47] employs LLMs as a decision-making component to solve complex AD scenarios that require human commonsense understanding. Drive as You Speak [48] proposes a framework that integrates LLMs into AD and orchestrates other modules accordingly. Drivers can directly communicate with the vehicle through LLMs. The framework includes a memory module to save past driving scenario experiences in a vector database, which includes decision cues, reasoning processes, and other valuable information. LLMs then make decisions based on acquired experience and common sense. DiLU [50] studies the driving methods of human drivers and proposes a paradigm that uses reasoning, memory, and reflection modules to facilitate interaction between LLMs and the environment. This approach embeds such knowledge-driven capabilities of human drivers into AD system. DwLLMs [49] encodes traffic participants and environment into object-level vectors. It introduces a new paradigm with a two-stage pre-training and fine-tuning approach, allowing the model to understand driving scenarios and generate driving actions. SurrealDriver [51] proposes a human-like AD framework based on LLMs within the CARLA simulator. Through memory and safety mechanisms, LLMs can accomplish situation understanding, decision-making, and action generation. It also learns the driving habits of human drivers and continuously optimizes the driving skills in the closed loop. DLaH [22] introduces reasoning, interpretation, and memory modules to construct an AD system based on GPT-3.5 [6] and LLaMA-Adapter v2 [55]. It demonstrates strong capabilities in scenario understanding and addressing long-tail problems within simulations.

Although the existing control and decision-making work in AD rely solely on LLMs, they can easily get connected with the perception module utilizing visual-LLMs connectors [55]–[58], achieving either mid-to-mid or end-to-end AD. Furthermore, designing a vision-LLMs connector specifically suitable for AD systems is a promising direction. We encourage exploration in this area and believe that a substantial amount of work in this segment will emerge in the near future.

### D. End-to-End Autonomous Driving

As defined in [13], end-to-end AD system is a fully differentiable program that takes raw sensor data as input and produces a plan and/or low-level control actions as output, which aligns well with the structure of M2T model in VLMs. Due to such a natural synergy, some studies start to explore the feasibility of applying M2T VLMs models to end-to-end AD. Compared to traditional end-to-end AD systems, large VLMs-based end-to-end AD systems possess potent interpretability, trustworthiness, and complex scene comprehending ability, paving the way for the practical application and implementation of end-to-end AD.

DriveGPT4 [18] is the pioneering work leveraging large VLMs for end-to-end AD tasks, which takes raw sensor data and human questions as input, and outputs predicted control signals and corresponding answers. It retains the powerful zero-shot generation capability of LLMs and is able to handle unseen scenarios. ADAPT [19] proposes an end-to-end AD pipeline based on the transformer model. With video input, ADAPT continuously outputs control signals as well as narration and reasoning descriptions for actions. Unlike the DriveGPT4, ADAPT does not incorporate VQA module but instead transforms interpretable end2end AD into the vision captioning task.

### E. Data Generation

Benefiting from the advancement and success of generative networks [59]–[65], the application of conditional generative models in AD allows the generation of large-scale high-quality data, thereby promoting the development of data-driven AD.

DriveGAN [52] learns the sequences of driving videos and their corresponding control signals. By disentangling scene components into action-dependent and action-independent features, it can control the vehicle behaviors in the generated video. This capability enables high-fidelity, controllable neural simulations and AD data generation. BEVControl [53] takes a sketch-style BEV layout and text prompts as inputs to generate street-view multi-view images. It introduces controller and coordinator elements to ensure the geometrical consistency between the sketch and output and appearance consistency across the multi-view images. This approach facilitates the possibility for controllable AD scenario sample generation based on BEV sketch.

Some works incorporate world models [66] into AD data generation for a more reasonable, predictable, and structured environment simulation. DriveDreamer [24] is a pioneering world model in AD entirely learned from real-world driving scenarios. It undergoes two stages of training: it initially understands and models driving scenarios from real-world driving videos, thereby acquiring structured traffic information. In the second stage, it constructs a driving world model through a video prediction task, gaining the ability to forecast future events and interact with the environment. DriveDreamer

generates realistic and controllable driving scenarios, which can be used for AD model training. GAIA-1 [20] takes video, action, and text description as inputs, utilizing the powerful capability of world models to learn structured representation and understand the environment, encoding the inputs into a sequence of tokens. It then employs denoising video diffusion models as the video decoder to achieve highly realistic video.

## IV. VLMs IN INTELLIGENT TRANSPORTATION SYSTEMS

Intelligent transportation systems are becoming more and more prevalent and work as an important step to achieve fully self-driving capabilities [75]. ITS usually consists of a multisensor setup in which cameras are almost always present due to its rich and dense semantic information. Similar to AD, VLMs can also facilitate the development of ITS in many aspects including ITS perception and understanding (IV-A), and ITS management system (IV-B). We further analyze the potential applications of VLMs in ITS (IV-C). Current works are demonstrated in Tab. II.

### A. ITS Perception and Understanding

The rapid increase in the number of multi-sensor systems mounted on traffic infrastructure significantly empowers the capabilities of ITS. It facilitates a more comprehensive perception and understanding of the traffic environment, empowering ITS to identify and interpret complex traffic scenarios accurately.

**Language-Guided Vehicle Retrieval.** Vehicle Retrieval serves as a crucial component of ITS perception and understanding. To draw the attention of researchers in this direction, AI City Challenge [76] [77] has been hosting the Tracked-Vehicle Retrieval by Natural Language Descriptions as a challenge track.

As part of this collective effort, [67] propose a multi-granularity retrieval approach for natural language-based vehicle retrieval. The key concept is the introduced multi-query retrieval module based on language augmentation. The idea behind this module is to leverage multiple imperfect language descriptions to achieve higher robustness and accuracy. An interesting approach for solving out-of-distribution input data for vehicle retrieval has been proposed by [68]. The key contribution of this work is the introduced domain adaptive training method which transfers knowledge from labeled data to unseen data by generating pseudo labels. MLVR [69] proposes a multi-modal language vehicle retrieval framework that employs text and image extractors for feature encoding, subsequently generating video vector sequences through a video recognition module. By integrating modules that combine various vehicle characteristics, MLVR creates more informative vehicle vectors for matching control and accomplishes language-guided retrieval.

**Traffic Visual-Scene Reasoning.** Another emerging field of using VLMs in ITS is visual scene event understanding, which commonly forms as Visual Question Answering (VQA) task. [70] proposes a weakly supervised Traffic-domain Video Question Answering with Automatic Captioning method. The core contribution is the usage of automatically generated synthetic captions for online available urban traffic videos. The automatically generated video-caption pairs are then used for fine-tuning, and thus injecting additional traffic domain knowledge into the trained model. [71] propose a Cross-Modal Question Reasoning framework to identify temporal causal context for event-level question reasoning. An attention-based module enables the learning of temporal causal scenes and question pairs. [72] introduces Tem-Adapter to minimize the gap between image and video domains from the temporal aspect by learning temporal dependencies. It shows great performance in traffic video question-answering tasks. AnomalyCLIP [73] employs the CLIP model for video anomaly detection. By specifying anomaly categories and using context optimization [78], it distinguishes between normal and abnormal samples, enabling the model to identify anomalous instances. AnomalyCLIP achieves good results in various datasets, including roadside anomaly detection. The VLMs enhanced semantic anomaly recognition algorithms can be further extended for disaster or emergency response in ITS.

### B. ITS Management System

ITS Management System enhances the operational safety and efficiency of the transportation system based on real-time perception and scene understanding. It facilitates smooth traffic flow through the management of traffic signals, offers timely updates on road conditions, and provides prompt accident alerts.

**Traffic Flow Management.** Pioneering work explores the use of LLMs in the realm of traffic flow management, specifically in traffic signal control. PromptGAT [74] introduces a prompt-based grounded action transformation method, marking the first application of LLMs to mitigate the sim-to-real transfer problem real-world in traffic signal control. Leveraging LLMs' to understand of the impacts of weather conditions on traffic states and road types, it enhances the applicability of policies in actual scenarios, effectively narrowing the sim-to-real gap.

### C. Potential Applications

Although many approaches have explored leveraging the capabilities of VLMs in various transportation scenarios, the full potential of VLMs in ITS remains untapped.

**Accident Anticipation and Detection.** The capability to promptly detect and anticipate accidents is critical in ensuring road safety. As [79] points out, current vision-based Traffic Accident Detection (TAD) and Traffic Accident Anticipation (TAA) still face various challenges, including its long-tailed and safe-critical properties, complex scene evolution, harsh environments, and determination uncertainty. Large VLMs, thanks to their exceptional abilities in zero-shot generalization, profound scene comprehension, border case recognition, and multi-step reasoning, serve as a promising solution to address current challenges.

**Crash Reporting.** Rapid crash analysis enhances traffic efficiency and prevents further congestion. M2V conditional generative VLMs represent a potential solution. By utilizing video footage recorded by drivers, along with their descriptions, the model can instantly generate narratives of the accident scene

| Method | Year | Tasks | Type | Contribution |
|---|---|---|---|---|
| **ITS Perception and Understanding** | | | | |
| MGNLVR [67] | 2022 | LGVR | VTM | Introduce a multi-granularity system for natural language-based vehicle retrieval. |
| SSDA-CLIP [68] | 2023 | LGVR | VTM | Propose a domain-adaptive CLIP model with semi-supervised training for vehicle retrieval. |
| MLVR [69] | 2023 | LGVR | VTM | Propose Multi-modal Language Vehicle Retrieval system for retrieving the trajectory of tracked vehicles. |
| TRIVIA [70] | 2023 | IC,VQA | VTF | Present a novel approach termed Traffic-domain Video Question Answering with Automatic Captioning. |
| CMQR [71] | 2023 | VQA | VTM, VTF | Propose a causality-aware event-level visual question reasoning framework to achieve robust VQA. |
| Tem-adapter [72] | 2023 | VQA | VTM | Propose an adapter to enable the learning of temporal dynamics and complex semantics by a visual Temporal Aligner and a textual Semantic Aligner. |
| AnomalyCLIP [73] | 2023 | VAR | VTM | Propose the first method for VAR based on LLV models to detect and classify anomalous events. |
| **ITS Management System** | | | | |
| PromptGAT [74] | 2023 | TSC | LLM | Propose method using LLMs to mitigate the sim-to-real transfer in traffic signal control. |

TABLE II

OVERVIEW OF VISION-LANGUAGE MODEL IN INTELLIGENT TRAFFIC SYSTEM.

TASKS: **LGVR**: LANGUAGE-GUIDED VEHICLE RETRIEVAL **IC**: IMAGE CAPTIONING **VQA**: VISUAL QUESTION ANSWERING **VAR**: VIDEO ANOMALY RECOGNITION **TSC**: TRAFFIC SIGNAL CONTROL
TYPE: **VTM**: VISION-TEXT MATCHING **VTF**: VISION-TEXT FUSION **LLM**: LARGE LANGUAGE MODEL

and the driver's account of the incident, significantly enhancing the response time in handling accidental circumstances.

**Parking Management System.** Smart parking solutions can take advantage of the planning capabilities of LLMs exhibited within the realms of language-guided navigation and motion planning [38], [45] to significantly reduce the time of finding parking space in urban areas. By integrating with parking space management systems, it is feasible to provide vehicles with linguistic guidance for viable routes, assisting in parking management. This approach can be further enhanced by interfacing with vehicle-side language-guided navigation systems, potentially actualizing autonomous parking solutions.

## V. DATASETS

Datasets play a fundamental role in ensuring the robustness and generalizability of the intelligent transportation. Beyond traditional vision-based datasets, integrating language modality into data offers advantages for driving and transportation systems. This section exhibits and analyzes the datasets foundational to autonomous driving (V-A) and those integrating language in autonomous and intelligent transportation contexts (V-B and V-C). The overview of datasets is presented in Tab. III.

### A. Autonomous Driving Dataset

In the autonomous driving domain, datasets serve as one of the key points for developing safe and efficient perception, prediction, and planning systems.

Several datasets like KITTI [42], nuScenes [85], BDD100K [86], and Waymo [87] span multiple tasks, such as object detection, tracking, and segmentation, with various data modalities. Cityscapes [81] provides precise annotated image data for object detection and semantic segmentation. In contrast to the versatile datasets, Caltech Pedestrian Detection [80] offers annotated images for pedestrian detection within urban traffic scenarios. Meanwhile, as a subset of [81], CityPersons focuses on image-based pedestrian detection from varied city environments. Other task-specific datasets, such as SemanticKITTI [83] provides labeled LiDAR point clouds for semantic segmentation. Data given by CityFlow [84] can be utilized to solve object tracking and re-identification.

### B. Language-Enhanced AD Dataset

As autonomous driving advances, combining linguistic information with visual data enriches semantic and contextual comprehension. By facilitating better recognition of road agencies and a deeper understanding of driving scenarios, natural language assistance enhances the safety and interaction capabilities of autonomous vehicles.

The prior work [89] provides a potential opportunity to enhance the capability of the perception system in autonomous vehicles by introducing language understanding into the detector. For object tracking task, CityFlow-NL [93], Refer-KITTI [31], and NuPrompt [32] extend [84] [42] [85] with language prompts, respectively. TOUCHDOWN [90], LCSD [44], and CARLA-NAV [37] generate language-guided navigation datasets. Talk2Car [92] is proposed for single traffic object referring task. Safe autonomous driving requires reliable scene understanding, [21] [94] [95] evaluate the understanding and reasoning capabilities of autonomous vehicles by providing question-answer pairs. Talk2BEV [23] focuses on visual-spatial reasoning (VSR). Beyond image and video data, Rank2Tell [96] takes LiDAR point clouds into account for multi-modal importance ranking and reasoning. BDD-X [88] offers textual explanations for improving the explainability of AD algorithms. HAD [91] proposes a human-to-vehicle advice dataset for developing advisable autonomous driving models.

### C. Language-Enhanced ITS Dataset

Purely vision-based intelligent transportation systems can exhibit vulnerabilities in challenging scenarios such as traffic incidents or heavy traffic flow. Hence, integrating natural language with visual data can elevate the robustness and analytical process of the transportation system.

Anomaly recognition is a critical task in intelligent transportation systems. ShanghaiTech [97] provides a dataset with 437 videos for video-level anomaly detection. UCF-Crime [98] presents 128 hours of real-world surveillance videos, including 13 realistic anomalies. SUTD-TrafficQA [99] comprises 10,080 in-the-wild traffic accident videos, annotated in a video Question-Answer (QA) format. AerialVLN [100] introduces a synthetic dataset that aims to fill the gap of the UAV-based vision-language navigation in complex urban environments.

| Dataset | Year | Tasks | Source Datasets | Data Modalities | | | |
|---|---|---|---|---|---|---|---|
| | | | | Image | Video | Point Cloud | Text |
| Autonomous Driving | | | | | | | |
| Caltech Ped Det [80] | 2009 | 2D OD | - | ✓ | - | - | - |
| KITTI [42] | 2012 | 2D/3D OD, SS, OT | - | ✓ | - | ✓ | - |
| Cityscapes [81] | 2016 | 2D/3D OD, SS | - | ✓ | - | - | - |
| CityPersons [82] | 2017 | 2D OD | - | ✓ | - | - | - |
| SemnaticKITTI [83] | 2019 | 3D SS | - | - | ✓ | - | - |
| CityFlow [84] | 2019 | OT, ReID | - | ✓ | ✓ | - | - |
| nuScenes [85] | 2020 | 2D/3D OD, 2D/3DSS, OT, MP | - | ✓ | ✓ | ✓ | - |
| BDD100K [86] | 2020 | 2D OD, 2D SS, OT | - | ✓ | ✓ | - | - |
| Waymo [87] | 2020 | 2D/3D OD, 2D SS, OT | - | ✓ | ✓ | ✓ | - |
| Language-Enhanced Autonomous Driving Dataset | | | | | | | |
| BDD-X [88] | 2018 | TE | - | ✓ | ✓ | - | ✓ |
| Cityscapes-Ref [89] | 2018 | OD | Cityscapes | ✓ | - | - | ✓ |
| TOUCHDOWN [90] | 2019 | VSR, VLN | - | ✓ | - | - | ✓ |
| LCSD [44] | 2019 | VLN | - | ✓ | - | - | ✓ |
| HAD [91] | 2019 | H2V-Advice | - | ✓ | ✓ | - | ✓ |
| Talk2Car [92] | 2020 | SOR | nuScenes | ✓ | ✓ | - | ✓ |
| CityFlow-NL [93] | 2021 | VR, OT | CityFlow | ✓ | ✓ | - | ✓ |
| CARLA-NAV [37] | 2022 | VLN | - | ✓ | - | - | ✓ |
| NuPrompt [32] | 2023 | OT | nuScenes | ✓ | ✓ | - | ✓ |
| NuScenes-QA [21] | 2023 | VQA | nuScenes | ✓ | ✓ | ✓ | ✓ |
| Refer-KITTI [31] | 2023 | OT | KITTI | ✓ | ✓ | - | ✓ |
| Talk2BEV [23] | 2023 | VSR, DM | nuScenes | ✓ | ✓ | - | ✓ |
| Driving LLMs [94] | 2023 | VQA | - | ✓ | - | - | ✓ |
| DRAMA [95] | 2023 | IC, VQA | - | ✓ | ✓ | - | ✓ |
| Rank2Tell [96] | 2023 | IR,VSR | - | ✓ | ✓ | ✓ | ✓ |
| Language-Enhanced Intelligent Transportation Systems Dataset | | | | | | | |
| ShanghaiTech [97] | 2018 | AD | - | ✓ | ✓ | - | - |
| UCF-Crime [98] | 2018 | AD | - | ✓ | ✓ | - | - |
| SUTD-TrafficQA [99] | 2021 | VQA | - | ✓ | ✓ | - | ✓ |
| AerialVLN [100] | 2023 | VLN | - | ✓ | ✓ | - | ✓ |

TABLE III

OVERVIEW OF DATASET. **OD**: OBJECT DETECTION **SS**: SEMANTIC SEGMENTATION **OT**: OBJECT TRACKING **ReID**: RE-IDENTIFICATION **MP**: MOTION PLANNING **TE**: TEXTUAL EXPLANATION **VSR**: VISUAL-SPATIAL REASONING **VLN**: VISION-LANGUAGE NAVIGATION **H2V-ADVICE**: HUMAN-TO-VEHICLE ADVICE **SOR**: SINGLE OBJECT REFERRING **VR**: VEHICLE RETRIEVAL **VQA**: VISUAL QUESTION ANSWERING **DM**: DECISION-MAKING **IC**: IMAGE CAPTIONING **IR**: IMPORTANCE RANKING **AD**: ANOMALY DETECTION

## VI. DISCUSSION

Given the aforementioned summary of existing work, we deeply discuss challenges and research gap related to language models in autonomous driving and intelligent transportation systems in this section, and outline potential directions for future research.

**Autonomous Driving Foundation Model.** Existing foundation models—including vision foundation models [101]–[103], language foundation models [2]–[4], and multi-modal foundation models [7]–[9]—have set the stage for the feasibility of Autonomous Driving Foundation Models (ADFMs). We formulate ADFMs as models pre-trained on vast and diverse datasets, excelling in interpretability, reasoning, forecasting, and introspection, and effective in various autonomous driving tasks, such as perception, understanding, planning, control, and decision-making. Some studies have made preliminary attempts [18] [20] [45] [47], while how to adapt existing foundation models to ADFMs, to align the objectiveness of autonomous driving remains a relatively uncharted domain.

**Data Availability and Formatting.** Although there are already many on-the-shelf large-scale autonomous driving datasets [85] [87] available, they are not suitable and optimal for direct adaptation of LLMs in AD and ITS. For example, how to generate instruction tuning datasets and design instruction formats based on AD datasets for ADFMs adaptations remains barely investigated. Besides, a large-scale image-text traffic-specific pairs dataset can also be of great help for the development of AD and ITS, especially for the approaches relying on VTM pre-training models in object detection, semantic segmentation, language-guided navigation, and language-guided retrieval.

**Safe Driving Alignment.** LLMs can generate toxic, biased, harmful content that may conflict with human values, necessitating alignment tuning. Similarly, when training

autonomous driving foundation models, it's also imperative to align their controlling strategy, decision-making, and response mechanisms with safety standards to ensure adherence to stable, safe, and sound driving values. Existing techniques in LLMs alignment tuning, such as reinforcement learning from human feedback (RLHF) [104] and supervised alignment tuning, are all worth trying in this domain.

**Multi-Modality Adaptation.** As mentioned in III-C, current approaches utilizing LLMs for motion planning, control, and decision-making often transform sensor data into text formulation directly, either through existing perception algorithms or by direct extraction from the simulator. While this modular approach simplifies experiments, it can lead to the loss of context and environment information and is heavily dependent on the performance of the perception algorithms. In light of this, exploring the establishment of vision-language connections through VTM or VTF or a hybrid of both, specifically for autonomous driving scenarios, as alternatives to simplistic manual reformulations, is a direction worth pursuing.

**Temporal Scene Understanding.** Scene understanding in autonomous driving and ITS typically requires temporal information from video to continuously perceive and comprehend the dynamics and causality of traffic environment and traffic participants. Merely using image-level VLMs is insufficient for the demands. For instance, it's impossible to determine the specific cause of a car accident from a single image Fig. 5. Therefore, how to process temporal sensor data for traffic scenarios is an issue that still needs exploration. One possible approach is to train a video-language model, e.g. [18] as an initial attempt, with all existing video-language adapters [105] [106] being potentially applicable in this regard. Another possible route involves converting video data into a paradigm that can be processed by image-language models [72], integrating temporal adapter layers and fine-tuning accordingly by necessity, thereby enhancing the model's understanding of spatial-temporal information in the traffic environment.

**Computation Resource and Processing Speed.** Real-time processing and limited computational resources pose significant challenges for model deployment in autonomous driving and intelligent traffic systems. Current LLMs typically contain billion-scale parameters, making both fine-tuning and inference highly resource-intensive, and failing to meet real-time requirements. Several existing techniques can alleviate these problems. For instance, Parameter-Efficient Fine-Tuning (PEFT) [107]–[109] reduces the number of trainable parameters while maintaining satisfactory model performance, thereby minimizing resource consumption for fine-tuning. Besides, unlike general LLMs, the knowledge required for autonomous driving is often specialized and domain-specific, and much of the knowledge contained within LLMs is actually redundant for AD. Hence, employing knowledge distillation [110] [111] to train a smaller, more tailored model suitable for autonomous driving presents

a feasible approach. Other common model compression techniques in deep learning, e.g. quantization [112] [113] and pruning, are also applicable in this context.

## VII. Conclusion

This survey provided an overview of the background, current advancements, potential applications and future trajectories for vision-language models in autonomous driving and intelligent traffic systems. It comprehensively summarized and analyzed the notable tasks, methods and datasets available in this field up to the present date. Drawing upon current studies, this work expounded on the prevailing challenges, potential solutions, and possible directions for future exploration. We hope that this paper can draw interest and attention within the research community to this area and facilitate more meaningful investigations.

## References

[1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," vol. 30, 2017.

[2] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale, *et al.*, "Llama 2: Open foundation and fine-tuned chat models," *arXiv preprint arXiv:2307.09288*, 2023.

[3] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, *et al.*, "Llama: Open and efficient foundation language models," 2023.

[4] T. L. Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić, D. Hesslow, R. Castagné, A. S. Luccioni, F. Yvon, M. Gallé, *et al.*, "Bloom: A 176b-parameter open-access multilingual language model," *arXiv preprint arXiv:2211.05100*, 2022.

[5] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann, *et al.*, "Palm: Scaling language modeling with pathways," *arXiv preprint arXiv:2204.02311*, 2022.

[6] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

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

[8] Z. Peng, W. Wang, L. Dong, Y. Hao, S. Ma, and F. Wei, "Kosmos-2: Grounding multimodal large language models to the world," *arXiv preprint arXiv:2306.14824*, 2023.

[9] T. Gong, C. Lyu, S. Zhang, Y. Wang, M. Zheng, Q. Zhao, K. Liu, W. Zhang, P. Luo, and K. Chen, "Multimodal-gpt: A vision and language model for dialogue with humans," *arXiv preprint arXiv:2305.04790*, 2023.

[10] S. Liu, L. Fan, E. Johns, Z. Yu, C. Xiao, and A. Anandkumar, "Prismer: A vision-language model with an ensemble of experts," *arXiv preprint arXiv:2303.02506*, 2023.

[11] C. Wu, S. Yin, W. Qi, X. Wang, Z. Tang, and N. Duan, "Visual chatgpt: Talking, drawing and editing with visual foundation models," *arXiv preprint arXiv:2303.04671*, 2023.

[12] J. Zhang, J. Huang, S. Jin, and S. Lu, "Vision-language models for vision tasks: A survey," *arXiv preprint arXiv:2304.00685*, 2023.

[13] L. Chen, P. Wu, K. Chitta, B. Jaeger, A. Geiger, and H. Li, "End-to-end autonomous driving: Challenges and frontiers," *arXiv preprint arXiv:2306.16927*, 2023.

[14] J. Gu, Z. Han, S. Chen, A. Beirami, B. He, G. Zhang, R. Liao, Y. Qin, V. Tresp, and P. Torr, "A systematic survey of prompt engineering on vision-language foundation models," *arXiv preprint arXiv:2307.12980*, 2023.

[15] W. X. Zhao, K. Zhou, J. Li, T. Tang, X. Wang, Y. Hou, Y. Min, B. Zhang, J. Zhang, Z. Dong, *et al.*, "A survey of large language models," *arXiv preprint arXiv:2303.18223*, 2023.

[16] T. Yuan, W. D. R. Neto, C. E. Rothenberg, K. Obraczka, C. Barakat, and T. Turletti, "Machine learning for next-generation intelligent transportation systems: A survey," *Transactions on Emerging Telecommunications Technologies*, vol. 33, 4 2022.

[17] S. International, "Taxonomy and definitions for terms related to driving automation systems for on-road motor vehicles," 2021.

[18] 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.

[19] B. Jin, X. Liu, Y. Zheng, P. Li, H. Zhao, T. Zhang, Y. Zheng, G. Zhou, and J. Liu, "Adapt: Action-aware driving caption transformer," *arXiv preprint arXiv:2302.00673*, 2023.

[20] A. Hu, L. Russell, H. Yeo, Z. Murez, G. Fedoseev, A. Kendall, J. Shotton, and G. Corrado, "Gaia-1: A generative world model for autonomous driving," *arXiv preprint arXiv:2309.17080*, 2023.

[21] T. Qian, J. Chen, L. Zhuo, Y. Jiao, and Y.-G. Jiang, "Nuscenes-qa: A multi-modal visual question answering benchmark for autonomous driving scenario," *arXiv preprint arXiv:2305.14836*, 2023.

[22] 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," *arXiv preprint arXiv:2307.07162*, 2023.

[23] V. Dewangan, T. Choudhary, S. Chandhok, S. Priyadarshan, A. Jain, A. K. Singh, S. Srivastava, K. M. Jatavallabhula, and K. M. Krishna, "Talk2bev: Language-enhanced bird's-eye view maps for autonomous driving," *arXiv preprint arXiv:2310.02251*, 2023.

[24] X. Wang, Z. Zhu, G. Huang, X. Chen, and J. Lu, "Drivedreamer: Towards real-world-driven world models for autonomous driving," *arXiv preprint arXiv:2309.09777*, 2023.

[25] W. Li, Z. Qu, H. Song, P. Wang, and B. Xue, "The traffic scene understanding and prediction based on image captioning," *IEEE Access*, vol. 9, pp. 1420–1427, 2020.

[26] W. Cheng, J. Yin, W. Li, R. Yang, and J. Shen, "Language-guided 3d object detection in point cloud for autonomous driving," *arXiv preprint arXiv:2305.15765*, 2023.

[27] N. Sriram, T. Maniar, J. Kalyanasundaram, V. Gandhi, B. Bhowmick, and K. M. Krishna, "Talk to the vehicle: Language conditioned autonomous navigation of self driving cars," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5284–5290, IEEE, 2019.

[28] A. Keysan, A. Look, E. Kosman, G. Gürsun, J. Wagner, Y. Yu, and B. Rakitsch, "Can you text what is happening? integrating pre-trained language encoders into trajectory prediction models for autonomous driving," *arXiv preprint arXiv:2309.05282*, 2023.

[29] M. Liu, J. Jiang, C. Zhu, and X.-C. Yin, "Vlpd: Context-aware pedestrian detection via vision-language semantic self-supervision," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6662–6671, 2023.

[30] M. Liu, C. Zhu, S. Ren, and X.-C. Yin, "Unsupervised multi-view pedestrian detection," *arXiv preprint arXiv:2305.12457*, 2023.

[31] D. Wu, W. Han, T. Wang, X. Dong, X. Zhang, and J. Shen, "Referring multi-object tracking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14633–14642, 2023.

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

[33] S. Peng, K. Genova, C. Jiang, A. Tagliasacchi, M. Pollefeys, T. Funkhouser, *et al.*, "Openscene: 3d scene understanding with open vocabularies," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 815–824, 2023.

[34] R. Chen, Y. Liu, L. Kong, X. Zhu, Y. Ma, Y. Li, Y. Hou, Y. Qiao, and W. Wang, "Clip2scene: Towards label-efficient 3d scene understanding by clip," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7020–7030, 2023.

[35] M. Najibi, J. Ji, Y. Zhou, C. R. Qi, X. Yan, S. Ettinger, and D. Anguelov, "Unsupervised 3d perception with 2d vision-language distillation for autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 8602–8612, 2023.

[36] F. Romero, C. Winston, J. Hauswald, M. Zaharia, and C. Kozyrakis, "Zelda: Video analytics using vision-language models," *arXiv preprint arXiv:2305.03785*, 2023.

[37] K. Jain, V. Chhangani, A. Tiwari, K. M. Krishna, and V. Gandhi, "Ground then navigate: Language-guided navigation in dynamic scenes," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4113–4120, IEEE, 2023.

[38] M. Omama, P. Inani, P. Paul, S. C. Yellapragada, K. M. Jatavallabhula, S. Chinchali, and M. Krishna, "Alt-pilot: Autonomous navigation with language augmented topometric maps," *arXiv preprint arXiv:2310.02324*, 2023.

[39] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," 2 2021.

[40] M. Cherti, R. Beaumont, R. Wightman, M. Wortsman, G. Ilharco, C. Gordon, C. Schuhmann, L. Schmidt, and J. Jitsev, "Reproducible scaling laws for contrastive language-image learning," 2022.

[41] H. Xu, G. Ghosh, P.-Y. Huang, D. Okhonko, A. Aghajanyan, F. Metze, L. Zettlemoyer, and C. Feichtenhofer, "Videoclip: Contrastive pre-training for zero-shot video-text understanding," 2021.

[42] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.

[43] A. Elhafsi, R. Sinha, C. Agia, E. Schmerling, I. Nesnas, and M. Pavone, "Semantic anomaly detection with large language models," *arXiv preprint arXiv:2305.11307*, 2023.

[44] N. Sriram, T. Maniar, J. Kalyanasundaram, V. Gandhi, B. Bhowmick, and K. M. Krishna, "Talk to the vehicle: Language conditioned autonomous navigation of self driving cars," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5284–5290, IEEE, 2019.

[45] J. Mao, Y. Qian, H. Zhao, and Y. Wang, "Gpt-driver: Learning to drive with gpt," *arXiv preprint arXiv:2310.01415*, 2023.

[46] J. Kim, S. Moon, A. Rohrbach, T. Darrell, and J. Canny, "Advisable learning for self-driving vehicles by internalizing observation-to-action rules," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9661–9670, 2020.

[47] H. Sha, Y. Mu, Y. Jiang, L. Chen, C. Xu, P. Luo, S. E. Li, M. Tomizuka, W. Zhan, and M. Ding, "Languagempc: Large language models as decision makers for autonomous driving," *arXiv preprint arXiv:2310.03026*, 2023.

[48] C. Cui, Y. Ma, X. Cao, W. Ye, and Z. Wang, "Drive as you speak: Enabling human-like interaction with large language models in autonomous vehicles," *arXiv preprint arXiv:2309.10228*, 2023.

[49] 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," *arXiv preprint arXiv:2310.01957*, 2023.

[50] 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.

[51] 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.

[52] S. W. Kim, J. Philion, A. Torralba, and S. Fidler, "Drivegan: Towards a controllable high-quality neural simulation," pp. 5820–5829, 2021.

[53] K. Yang, E. Ma, J. Peng, Q. Guo, D. Lin, and K. Yu, "Bevcontrol: Accurately controlling street-view elements with multi-perspective consistency via bev sketch layout," *arXiv preprint arXiv:2308.01661*, 2023.

[54] OpenAI, "Gpt-4v(ision) system card," 2023.

[55] P. Gao, J. Han, R. Zhang, Z. Lin, S. Geng, A. Zhou, W. Zhang, P. Lu, C. He, X. Yue, H. Li, and Y. Qiao, "Llama-adapter v2: Parameter-efficient visual instruction model," 2023.

[56] P. Gao, J. Han, R. Zhang, Z. Lin, S. Geng, A. Zhou, W. Zhang, P. Lu, C. He, X. Yue, H. Li, and Y. Qiao, "Llama-adapter v2: Parameter-efficient visual instruction model," 2023.

[57] H. Liu, C. Li, Q. Wu, and Y. J. Lee, "Visual instruction tuning," 2023.

[58] H. Liu, C. Li, Y. Li, and Y. J. Lee, "Improved baselines with visual instruction tuning," 2023.

[59] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. L. Denton, K. Ghasemipour, R. Gontijo Lopes, B. Karagol Ayan, T. Salimans, *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *Advances in Neural Information Processing Systems*, vol. 35, pp. 36479–36494, 2022.

[60] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," pp. 10684–10695, 2022.

[61] A. Nichol, P. Dhariwal, A. Ramesh, P. Shyam, P. Mishkin, B. McGrew, I. Sutskever, and M. Chen, "Glide: Towards photorealistic image generation and editing with text-guided diffusion models," *arXiv preprint arXiv:2112.10741*, 2021.

[62] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents, 2022," *URL https://arxiv. org/abs/2204.06125*, vol. 7, 2022.

[63] J. Ho, W. Chan, C. Saharia, J. Whang, R. Gao, A. Gritsenko, D. P. Kingma, B. Poole, M. Norouzi, D. J. Fleet, *et al.*, "Imagen video: High definition video generation with diffusion models," *arXiv preprint arXiv:2210.02303*, 2022.

[64] L. Khachatryan, A. Movsisyan, V. Tadevosyan, R. Henschel, Z. Wang, S. Navasardyan, and H. Shi, "Text2video-zero: Text-to-image diffusion models are zero-shot video generators," *arXiv preprint arXiv:2303.13439*, 2023.

[65] X. Wang, H. Yuan, S. Zhang, D. Chen, J. Wang, Y. Zhang, Y. Shen, D. Zhao, and J. Zhou, "Videocomposer: Compositional video synthesis with motion controllability," *arXiv preprint arXiv:2306.02018*, 2023.

[66] Y. LeCun, "A path towards autonomous machine intelligence," 2022.

[67] J. Zhang, X. Lin, M. Jiang, Y. Yu, C. Gong, W. Zhang, X. Tan, Y. Li, E. Ding, and G. Li, "A multi-granularity retrieval system for natural language-based vehicle retrieval," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3216–3225, 2022.

[68] H. D.-A. Le, Q. Q.-V. Nguyen, D. T. Luu, T. T.-T. Chau, N. M. Chung, and S. V.-U. Ha, "Tracked-vehicle retrieval by natural language descriptions with multi-contextual adaptive knowledge," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5510–5518, 2023.

[69] D. Xie, L. Liu, S. Zhang, and J. Tian, "A unified multi-modal structure for retrieving tracked vehicles through natural language descriptions," in *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 5419–5427, 2023.

[70] E. Qasemi, J. M. Francis, and A. Oltramari, "Traffic-domain video question answering with automatic captioning," *arXiv preprint arXiv:2307.09636*, 2023.

[71] Y. Liu, G. Li, and L. Lin, "Causality-aware visual scene discovery for cross-modal question reasoning," *arXiv preprint arXiv:2304.08083*, 2023.

[72] G. Chen, X. Liu, G. Wang, K. Zhang, P. H. Torr, X.-P. Zhang, and Y. Tang, "Tem-adapter: Adapting image-text pretraining for video question answer," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 13945–13955, 2023.

[73] L. Zanella, B. Liberatori, W. Menapace, F. Poiesi, Y. Wang, and E. Ricci, "Delving into clip latent space for video anomaly recognition," *arXiv preprint arXiv:2310.02835*, 2023.

[74] L. Da, M. Gao, H. Mei, and H. Wei, "Llm powered sim-to-real transfer for traffic signal control," *arXiv preprint arXiv:2308.14284*, 2023.

[75] Z. Mahrez, E. Sabir, E. Badidi, W. Saad, and M. Sadik, "Smart urban mobility: When mobility systems meet smart data," 2021.

[76] S. Wang, D. C. Anastasiu, S. Sclaroff, and A. Li, "The 6th AI City Challenge," pp. 3347–3356, 2022.

[77] S. Wang, D. C. Anastasiu, M. S. Arya, and S. Sclaroff, "The 7th AI City Challenge," 2023.

[78] K. Zhou, J. Yang, C. C. Loy, and Z. Liu, "Learning to prompt for vision-language models," *International Journal of Computer Vision*, vol. 130, no. 9, pp. 2337–2348, 2022.

[79] J. Fang, iahuan Qiao, J. Xue, and Z. Li, "Vision-based traffic accident detection and anticipation: A survey," pp. 1–16, 2023.

[80] P. Dollár, C. Wojek, B. Schiele, and P. Perona, "Pedestrian detection: A benchmark," in *2009 IEEE conference on computer vision and pattern recognition*, pp. 304–311, IEEE, 2009.

[81] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3223, 2016.

[82] S. Zhang, R. Benenson, and B. Schiele, "Citypersons: A diverse dataset for pedestrian detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 3213–3221, 2017.

[83] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, "Semantickitti: A dataset for semantic scene understanding of lidar sequences," in *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9297–9307, 2019.

[84] Z. Tang, M. Naphade, M.-Y. Liu, X. Yang, S. Birchfield, S. Wang, R. Kumar, D. Anastasiu, and J.-N. Hwang, "Cityflow: A city-scale benchmark for multi-target multi-camera vehicle tracking and re-identification," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8797–8806, 2019.

[85] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.

[86] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving dataset for heterogeneous multitask learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2636–2645, 2020.

[87] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, *et al.*, "Scalability in perception for autonomous driving: Waymo open dataset," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2446–2454, 2020.

[88] J. Kim, A. Rohrbach, T. Darrell, J. Canny, and Z. Akata, "Textual explanations for self-driving vehicles," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 563–578, 2018.

[89] A. B. Vasudevan, D. Dai, and L. Van Gool, "Object referring in videos with language and human gaze," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4129–4138, 2018.

[90] H. Chen, A. Suhr, D. Misra, N. Snavely, and Y. Artzi, "Touchdown: Natural language navigation and spatial reasoning in visual street environments," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12538–12547, 2019.

[91] J. Kim, T. Misu, Y. T. Chen, A. Tawari, and J. Canny, "Grounding human-to-vehicle advice for self-driving vehicles," in *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, vol. 2019-June, pp. 10583–10591, 2019.

[92] T. Deruyttere, S. Vandenhende, D. Grujicic, L. Van Gool, and M.-F. Moens, "Talk2car: Taking control of your self-driving car," *arXiv preprint arXiv:1909.10838*, 2019.

[93] Q. Feng, V. Ablavsky, and S. Sclaroff, "Cityflow-nl: Tracking and retrieval of vehicles at city scale by natural language descriptions," *arXiv preprint arXiv:2101.04741*, 2021.

[94] 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," *arXiv preprint arXiv:2310.01957*, 2023.

[95] S. Malla, C. Choi, I. Dwivedi, J. H. Choi, and J. Li, "Drama: Joint risk localization and captioning in driving," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 1043–1052, 2023.

[96] E. Sachdeva, N. Agarwal, S. Chundi, S. Roelofs, J. Li, B. Dariush, C. Choi, and M. Kochenderfer, "Rank2tell: A multimodal driving dataset for joint importance ranking and reasoning," *arXiv preprint arXiv:2309.06597*, 2023.

[97] W. Liu, W. Luo, D. Lian, and S. Gao, "Future frame prediction for anomaly detection–a new baseline," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6536–6545, 2018.

[98] W. Sultani, C. Chen, and M. Shah, "Real-world anomaly detection in surveillance videos," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 6479–6488, 2018.

[99] L. Xu, H. Huang, and J. Liu, "Sutd-trafficqa: A question answering benchmark and an efficient network for video reasoning over traffic events," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9878–9888, 2021.

[100] S. Liu, H. Zhang, Y. Qi, P. Wang, Y. Zhang, and Q. Wu, "Aerialvln: Vision-and-language navigation for uavs," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15384–15394, 2023.

[101] X. Wang, X. Zhang, Y. Cao, W. Wang, C. Shen, and T. Huang, "Seggpt: Segmenting everything in context," 2023.

[102] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, "Segment anything," 2023.

[103] X. Wang, W. Wang, Y. Cao, C. Shen, and T. Huang, "Images speak in images: A generalist painter for in-context visual learning," 2023.

[104] P. Christiano, J. Leike, T. B. Brown, M. Martic, S. Legg, and D. Amodei, "Deep reinforcement learning from human preferences," 2023.

[105] H. Zhang, X. Li, and L. Bing, "Video-llama: An instruction-tuned audio-visual language model for video understanding," 2023.

[106] M. Maaz, H. Rasheed, S. Khan, and F. S. Khan, "Video-chatgpt: Towards detailed video understanding via large vision and language models," 2023.

[107] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," 2021.

[108] B. Lester, R. Al-Rfou, and N. Constant, "The power of scale for parameter-efficient prompt tuning," 2021.

[109] Q. Zhang, M. Chen, A. Bukharin, P. He, Y. Cheng, W. Chen, and T. Zhao, "Adaptive budget allocation for parameter-efficient fine-tuning," 2023.

[110] Y. Gu, L. Dong, F. Wei, and M. Huang, "Knowledge distillation of large language models," 2023.

[111] C. Liu, C. Tao, J. Feng, and D. Zhao, "Multi-granularity structural knowledge distillation for language model compression," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, (Dublin, Ireland), pp. 1001–1011, Association for Computational Linguistics, May 2022.

[112] G. Xiao, J. Lin, M. Seznec, H. Wu, J. Demouth, and S. Han, "Smoothquant: Accurate and efficient post-training quantization for large language models," in *International Conference on Machine Learning*, pp. 38087–38099, PMLR, 2023.

[113] H. Bai, L. Hou, L. Shang, X. Jiang, I. King, and M. R. Lyu, "Towards efficient post-training quantization of pre-trained language models," *Advances in Neural Information Processing Systems*, vol. 35, pp. 1405–1418, 2022.